

Assignment 5 Part 1

$$P(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

derivative of y_k with respect to a_j

$$\frac{\partial y_k}{\partial a_j} = y_k(\delta_{kj} - y_j)$$

Next we need likelihood function using the 1-of-K coding scheme

$$P(T | w_1, \dots, w_K) = \prod_{n=1}^N \prod_{k=1}^K P(C_k | a_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_k(\phi_n)^{t_{nk}}$$

where $y_{nk} = y_k(\phi_n)$ and T is an $N \times K$ matrix of target variables with element t_{nk}

Taking negative logarithm gives the cross entropy error function:

$$E(w_1, \dots, w_K) = -\ln P(T | w_1, \dots, w_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln(y_{nk})$$

Gradient of the error function wrt to one parameter vector w_j

$$\nabla_{w_j} E(w_1, \dots, w_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

$$\nabla_{w_K} \nabla_{w_j} E(w_1, \dots, w_N) = -\sum_{n=1}^N y_{nk} (\delta_{kj} - y_{nj}) \phi_n \phi_n^T$$

Differentiating the softmax function First we differentiate $y_k = h(a_k) = \text{softmax}(a_k)$ with respect to a_K

$$\frac{\partial y_k}{\partial a_k} = \frac{\partial}{\partial a_k} \text{softmax}(a_k) = \frac{\partial}{\partial a_k} \left(\frac{\exp a_k}{\sum_j \exp a_j} \right) = \frac{\exp a_k \sum_j \exp a_j - \exp a_k \exp a_k}{(\sum_j \exp a_j)^2} \quad (\text{product rule})$$

$$= y_k \frac{\sum_{j \neq k} \exp a_j}{\sum_j \exp a_j} = y_k \left(1 - \frac{\exp a_k}{\sum_j \exp a_j} \right) = y_k (1 - y_k)$$

we differentiate $y_k = \text{softmax}(a_k)$ with respect to a_i to obtain

$$\frac{\partial y_k}{\partial a_i} = \frac{\partial}{\partial a_i} \text{softmax}(a_k) = \frac{-\exp(a_k) \exp(a_i)}{(\sum_j \exp a_j)^2} = -y_i y_k$$

These two results may be summarized using the Kronecker delta function δ_{ki} as

$$\frac{\partial y_k}{\partial a_i} = \frac{\partial}{\partial a_i} \text{softmax}(a_k) = y_k (\delta_{ki} - y_i)$$

Given $E(w_1, \dots, w_K) = -\ln P(T | w_1, \dots, w_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$ mathematically show that the derivative E will take the following form:

minimize the negative log likelihood function

$$E(w) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} + (1-t_{nk}) \ln(1-y_{nk})$$

Taking the gradient of the error function:

$$\nabla_w E(w) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n = \phi^T (y - T)$$

minimize error function efficient

$$E(w) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} + (1-t_{nk}) \ln(1-y_{nk})$$

$$\nabla_w E(w) = \phi^T (y - T)$$

$$\nabla_w \nabla_w E(w) = H = \phi^T R \phi$$

Using Newton-Raphson

$$w^{t+1} = (R^T R)^{-1} R^T y$$

define Z as

$$Z = \phi^T R^{-1} (y - T)$$

The likelihood of an example (x, t) is

$$P(t|x) = \prod_{k=1}^K y_k^{t_k}$$

where

$$y_k = y_k(x) = \frac{\exp(a_k(x))}{\sum_j \exp(a_j(x))}$$

The likelihood of a data set $D = \{(x_n, t_n)\}_{n=1}^N$ is

$$P(D|w) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

Maximizing the data likelihood is equivalent to minimizing the negative log data likelihood

$$E(w) = -\log \prod_{n=1}^N P(t_n|x_n)$$

$$= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_{nk}$$

$$= \sum_{n=1}^N E_n(w)$$

note that

$$E_n = -\sum_{k=1}^K t_{nk} \log y_{nk} \text{ is the cross entropy between } t_n \text{ and } y_n$$

From $a_j = w_j^T \phi$ and $y_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$ we have

$$\nabla_{w_j} a_j = \delta_{j1} \phi \text{ and } \frac{\partial y_k}{\partial a_j} = y_k (\delta_{kj} - y_j)$$

It follows that

$$\nabla_{w_j} y_k = \sum_{l=1}^K \frac{\partial y_k}{\partial a_l} (\nabla_{w_l} a_j)$$

$$= \sum_l y_k (\delta_{kl} - y_l) \delta_{jl} \phi$$

(Consider the first-order derivatives of the error $E = -\sum_{k=1}^K t_{nk} \log y_{nk}$ of point (x, t))

$$\nabla_{w_j} \left(-\sum_{k=1}^K t_{nk} \log y_{nk} \right) = -\sum_{k=1}^K t_{nk} \left(\frac{1}{y_{nk}} \nabla_{w_j} y_k \right)$$

$$= -t_j \phi + \sum_{k=1}^K t_{nk} y_k \phi$$

$$= -t_j \phi + y_j \phi$$

$$= (y_j - t_j) \phi$$

Gradient of $E(w)$ of data set D is

$$\nabla_w E(w) = \nabla_w \left(\sum_{n=1}^N E_n(w) \right) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n = \Phi^T (y - T)$$

$$\nabla_w E(w) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$