

Problem 1

$$t = h(x) + \epsilon$$

Here x is k -dimensional vector, $h(x)$ is a deterministic function of x , where x includes a bias & ϵ is random noise that has Gaussian probability with 0 mean & variance σ^2 .

Linear function approximation with vector w , we model the data as

$$y = \sum_{i=0}^k w_i x_i$$

Prove, $w^* = \arg \min \sum_{n=1}^N (t^n - y^n)^2$ that minimizes the SSE.

$$\rightarrow P(t^* | x^*) = \prod_{k=1}^N P(t^k | x^k)$$

Let there be N ^{BR-1 categories} examples, so, the combined probability of one data is $P(t|x)$ & if all distributions are independent, then probabilities can be multiplied.

Now, since, $t = h(x) + \epsilon$, where ϵ is Gaussian noise

$$P(\epsilon_k) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{(-\epsilon_k^2 / 2\sigma^2)} \quad - 1)$$

Now, we are modelling function $h_k(x)$ by neural network with outputs $y_n(x; w)$, using $t_k = h_k(x) + \epsilon_k$. - 2)

& 1) we get

$$\begin{aligned} P(t_k | x) &= \frac{1}{(2\pi\sigma^2)^{1/2}} e^{(-\frac{(y_n(x; w) - t_k)^2}{2\sigma^2})} \\ &= \frac{1}{(2\pi\sigma^2)^{1/2}} e^{(-\frac{(y_k - t_k)^2}{2\sigma^2})} \end{aligned}$$

{ If we replace $h_k(x)$ by $y_n(w, x)$ the value our model predicts }

Now, Likelihood of data

$$L = \prod_n P(x^n, t^n)$$

$$\& E = -\ln L = -\sum_n \ln P(t^n | x^n) = -\sum_n \ln P(t^n) - \sum_n \ln P(x^n)$$

Since second term is independent of network parameters we can eliminate it.

$$E = \frac{1}{2\sigma^2} \sum_n \sum_k \{ (f_k(x^n; w) - t_k)^2 + N_c \ln \sigma + \frac{N_c}{2} \ln(2\pi) \}$$

→ we can remove the last 2 terms for purpose of minimization as they are independent of w & also factor σ^2

$$E = \frac{1}{2} \sum_n \sum_k y_k (x^n, w) - t_k)^2$$

$$E = \frac{1}{2} \sum_n (y^n - t^n)^2 \quad - 3)$$

Finding minimum set of w^* which minimizes the error function E in 3). is same as

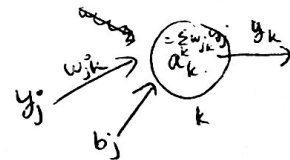
$$w^* = \arg \min_w \sum_{n=1}^N (t^n - y^n)^2.$$

2) & for output layer δ_k , output layer is represent as k , hidden layer as j .

$$\delta_k = -\frac{\partial E}{\partial a_k}, \quad a_k = \sum_j w_{jk} y_j$$

Now, by chain rule

$$\delta_k = -\frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial a_k} \quad - 1)$$



Since, a_k is the input to a unit k after applying the ~~act~~ summation. & y_k is the output of a unit in k .

$$\text{So } \frac{\partial y_k}{\partial a_k} = y_k' \quad - 2)$$

$$\text{Now, } \frac{\partial E}{\partial y_k} = \frac{\partial E}{\partial y_k} \frac{\partial w_{jk}}{\partial w} = \frac{\partial E}{\partial w_{jk}} \frac{\partial w_{jk}}{\partial y_k}$$

For logistic & softmax units

$$\frac{\partial E}{\partial y_k} = \frac{y_k - t_k}{y_k} \quad - 3)$$

$$\text{Using 2) & 3) } -\frac{\partial E}{\partial a_k} = t_k - y_k.$$

Derivation of equation of $\frac{\partial E}{\partial a_k}$ for softmax function

$$\frac{\partial y_l}{\partial a_m} = \frac{\sum_x e^{a_x}}{\sum_x e^{a_x}} \quad \text{if } m = l$$

, then

$$m = l, \frac{\partial y_l}{\partial a_l} = \frac{e^{a_l} \sum_x e^{a_x} - e^{a_l} e^{a_l}}{(\sum_x e^{a_x})^2} \quad \left\{ \text{by } \frac{d(uv)}{dv} = \frac{u'v - uv'}{v^2} \right\}$$

$$= \frac{e^{a_l}}{\sum_x e^{a_x}} \left[1 - \frac{e^{a_l}}{\sum_x e^{a_x}} \right] = (y_l)(1 - y_l)$$

if $m \neq l$

$$\frac{\partial y_l}{\partial a_m} = \frac{0 - e^{a_l} e^{a_m}}{(\sum_x e^{a_x})^2} = -y_l y_m$$

Now, cross entropy function is given by

$$E = -\sum_{l=1}^C t_l \log y_l$$

$$\frac{\partial E}{\partial a_k} = -\frac{\partial}{\partial a_k} \sum_{l=1}^C t_l \log y_l = -\sum_{l=1}^C t_l \frac{\partial y_l}{y_l} \frac{\partial y_l}{\partial a_k}$$

$$= \frac{t_k (y_k)}{y_k} (1 - y_k) \frac{\partial y_k}{\partial a_k} + \sum_{l=1, l \neq k}^C \frac{-t_l y_l}{y_l} (-y_l)(y_l) \frac{\partial y_l}{\partial a_k}$$

$l = k$

$$= -t_k + t_k y_k + \sum_{l=1, l \neq k}^C t_l y_l$$

$$= -t_k + \sum_{l=1}^C t_l y_l$$

$$= -t_k + y_k \quad \left\{ \sum_{l=1}^C t_l = 1 \right\}$$

$$= y_k - t_k$$

Now, $\frac{\partial E}{\partial a_k} = \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial a_k} = y_k - t_k$

for output layer.

Q For hidden layer

$$\delta_j^o = -\frac{\partial E}{\partial a_j^o} = -\frac{\partial E}{\partial y_j^o} \frac{\partial y_j^o}{\partial a_j^o} \quad \{ \text{by chain rule} \}$$

$$= -\sum_k \frac{\partial E}{\partial y_j^o} \frac{\partial a_k}{\partial a_j^o} \frac{\partial y_j^o}{\partial a_j^o}$$

$$= -\sum_k \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial y_j^o} \frac{\partial y_j^o}{\partial a_j^o} \quad - 1)$$

since, y_j^o is the output of a unit in hidden layer
 & a_j^o is its input

$$\frac{\partial y_j^o}{\partial a_j^o} = y_j^o \quad - 2)$$

also, $\frac{\partial a_k}{\partial y_j^o} = \frac{\sum_x w_{jk}^o y_x}{\partial y_j^o} = w_{jk}^o \text{ when } x=j \quad - 3)$

since, $a_k = \sum_j w_{jk}^o y_j^o$ {where y_j^o is the output of previous layer}

using 2) & 3) in 1)

we get

$$\delta_j^o = y_j^o \sum_k w_{jk}^o \delta_k^o \quad \left\{ \because \frac{\partial E}{\partial a_k} = \delta_k^o \right\}$$

so, $\delta_j^n = y_j^n \sum_k w_{jk}^n \delta_k^n$

b) Update Rule :-

$$w_{ij}^o = w_{ij}^o - \eta \frac{\partial E}{\partial w_{ij}^o}$$

$$\frac{\partial E}{\partial w_{ij}^o} = \frac{\partial E}{\partial a_j^o} \frac{\partial a_j^o}{\partial w_{ij}^o}$$

$$= \delta_j^o \frac{\partial a_j^o}{\partial w_{ij}^o} \quad - 1) \{ \text{By part a} \}$$

$$a_j^o = \sum_i w_{ij}^o z_i^o$$

$$\frac{\partial a_j^o}{\partial w_{ij}^o} = z_i^o \quad \text{when } x=i \quad - 2)$$

Using this 1)

$$\frac{\partial E}{\partial w_{ij}^o} = -\delta_j^o z_i^o$$

$$\begin{aligned} \text{So, } w_{ij}^o &= w_{ij}^o + \eta \delta_j^o z_i^o \\ &= w_{ij}^o + \eta \sum_k \delta_k^o z_i^o w_{jk}^o \end{aligned}$$

$$\begin{aligned} &= w_{ij}^o + \eta \sum_k (t_k - y_k) z_i^o w_{jk}^o y_j^o \quad (\text{Part a}) \\ &= w_{ij}^o + \eta \sum_{n=1}^N y_j^n (1 - y_j^n) z_i^o \sum_{k=1}^C (t_k^n - y_k^n) w_{jk}^o \end{aligned}$$

considering the hidden unit as sigmoid

C) Vectorization

* In vectorized form, we can represent our input to be matrix of $N \times I$ where N is the total no. of examples & I is the no. of units in each example. In our case of MNIST data, $I = 784 + 1$.
 $784(28 \times 28)$ & 1 more for bias.

* The weights for hidden ^{layer} ~~unit~~ can be represented as w_{ij} as a matrix of size $I \times H$ where H is the no. of ~~weights~~ hidden units.

→ The weights for output layer can also be represented (as w_{jk}^o) as a matrix of size $H \times C$ where C is the no. of classifications (10) in ^{case of} MNIST (0-9)

* For each layer input is represented as a matrix Z
for input layer Z has dimensions $N \times I$.

for hidden layer Z is $N \times H$

for output layer Z is $N \times C$

* For each layer output is represented as a matrix Y
for input layer Y 's dimensions $N \times I$
hidden Y 's dimension $N \times H$
output layer Y 's dimension $N \times C$

* For each layer, the corresponding matrix is denoted by
 Z^l, Y^l, δ^l, W^l .

Now, for weight update

$$W^l = W^l + \eta (\delta^l (Z^l)^T)$$

For hidden layer, δ

$$\delta^l = (\sigma'(Z^l)) * (\delta^{l+1} (W^{l+1})^T)$$

~~dot product~~ or
elementwise multiplication

{ \bullet : matrix multiplication
 $*$: dot product or elementwise multiplication }

For output layer

$$\delta^k = T - Y^k$$

{ T : Target matrix containing
one-hot representation of
labels }