

① Theory(1) Simplest case of Linear regression

Linear regression is where the expression ^(function) for finding the labels from input is linear in parameter and labels are from a continuous set.

First,
in linear regression

$$\hat{\vec{y}} = X \vec{w}$$

where

$\hat{\vec{y}}$ - estimated labels ^{vector} from input ^{matrix} and weights

X - Matrix containing N input ^{vectors} \vec{x}_i

$\{i \in \mathbb{N}[1, N]\}$
 \mathbb{N} - natural numbers

\vec{w} - weights vector

$$\vec{y} = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(N)}]^T$$

$$\vec{w} = [w_0 \quad w_1 \quad \dots \quad w_d]^T$$

$$X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(N)} \\ \vdots & \vdots & \ddots & \vdots \\ x_d^{(1)} & x_d^{(2)} & \dots & x_d^{(N)} \end{bmatrix}$$

where $x_i^{(j)} \in \mathbb{R}$ and $y^{(j)} \in \mathbb{R} \quad \forall i$

To find the optimal weights vector,

we use sum of squared error as cost function,

$$E(\vec{w}) = \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2$$

Here $y^{(i)}$ is ground truth label associated with $\vec{x}^{(i)}$, $y^{(i)} \in \mathbb{R} \quad \forall i$

$E(\vec{w})$ is thus the error fn (SSE).

$$E(\vec{w}) = \sum_{i=1}^N \left[y^{(i)} - \left(\sum_{j=1}^d x_j^{(i)} w_j + w_0 \right) \right]^2$$

In Matrix notation,

$$E(\vec{w}) = (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w}) \quad \text{--- (1)}$$

$$\vec{w}^* = \underset{\vec{w}}{\text{argmin}} (E(\vec{w}))$$

where \vec{w}^* is the optimal weight vector.

$$\text{Set } \nabla E(\vec{w}) = 0$$

from (1)

$$E(\vec{w}) = \vec{y}^T \vec{y} - \vec{w}^T X^T \vec{y} - \vec{y}^T X \vec{w} + \vec{w}^T X^T X \vec{w}$$

$$\nabla E(\vec{w}) = -X^T \vec{y} - X^T \vec{y} + 2(X^T X) \vec{w} = 0$$

$$\nabla E(\vec{w}) = -2X^T \vec{y} + 2(X^T X) \vec{w} = 0$$

$$X^T \vec{y} = X^T X \vec{w}^*$$

$$\boxed{\vec{w}^* = (X^T X)^{-1} X^T \vec{y}}$$

By taking 2nd derivative,

we get $2X^T X$, which should be positive definite for the cost function to be minimum, which is possible if X has full column rank.

A matrix is said to be positive definite, if for any \vec{a}

$$\boxed{\vec{a}^T X \vec{a} > 0}$$

The best estimate for y , would be

$$\boxed{\hat{y}_{\text{best}} = X \vec{w}^{\text{opt}}}$$

where $\vec{w}^{\text{opt}} = (X^T X)^{-1} X^T y$

(2) Using basis function with $(M+1)$ parameters

$$\hat{y}^{(i)} = \sum_{j=0}^M \phi_j(\vec{x}^{(i)}) w_j \quad \text{--- (1)}$$

where

$$\phi_0(\vec{x}^{(i)}) = 1 \quad \text{for all } i \text{ in } [1, N] \quad i \in \mathbb{N}$$

we can represent (1) in Matrix form ^(Natural numbers)

$$\vec{y} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(N)} \end{bmatrix} \quad \Phi = \begin{bmatrix} 1 & \phi_1^{(1)} & \dots & \phi_M^{(1)} \\ \vdots & \phi_1^{(2)} & \dots & \phi_M^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1^{(N)} & \dots & \phi_M^{(N)} \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix}$$

Cost function (SSE):

$$E(\vec{w}) = \sum_{i=1}^N \left(y^{(i)} - \hat{y}^{(i)} \right)^2 \quad \sum_{i=1}^N \left(y^{(i)} - \sum_{j=0}^M \phi_j(\vec{x}^{(i)}) w_j \right)^2$$

In Matrix form,

$$E(\vec{\omega}) = (\vec{y} - \phi \vec{\omega})^T (\vec{y} - \phi \vec{\omega})$$

Set,
 $\nabla E(\vec{\omega}) = 0$

to find $\vec{\omega}^*$,

$$\nabla E(\vec{\omega}) = -2\phi^T \vec{y} + 2\phi^T \phi \vec{\omega}$$

Setting
 $\nabla E(\vec{\omega}) = 0$

$$-2\phi^T \vec{y} + 2\phi^T \phi \vec{\omega}^* = 0$$

$$\vec{\omega}^* = (\phi^T \phi)^{-1} \phi^T \vec{y}$$

The best estimate

$$\hat{y}_{\text{best}} = \phi \vec{\omega}^*$$

(6) Linear ridge regression with L_2 Regularisation

→ Firstly as we do this,
we subtract the mean from the ~~data~~ ^{label}
data then we can ignore 'dc' bias term.

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_M^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_M^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(N)} & x_2^{(N)} & \dots & x_M^{(N)} \end{bmatrix}$$

$$\vec{\hat{y}} = X \vec{\omega}$$

$$\vec{\omega} = [\omega_1 \quad \omega_2 \quad \dots \quad \omega_M]^T$$

$$\vec{\hat{y}} = [\hat{y}^{(1)} \quad \hat{y}^{(2)} \quad \dots \quad \hat{y}^{(N)}]^T$$

where $\vec{\hat{y}}$ is estimated label vector from given inputs

$$E(\vec{\omega}) = (\vec{y} - X\vec{\omega})^T (\vec{y} - X\vec{\omega}) + \lambda \vec{\omega}^T \vec{\omega}$$

$$\nabla E(\vec{\omega}) = -2X^T(\vec{y} - X\vec{\omega}) + 2\lambda\vec{\omega}$$

Setting

$$\nabla E(\vec{\omega}) = 0$$

$$2X^T(\vec{y} - X\vec{\omega}) = 2\lambda\vec{\omega}$$

$$\vec{\omega}^* = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

The best estimate

$$\vec{\hat{y}} = X \vec{\omega}^*$$

→ In most general cases No of observations (N) will be less than the total dimension ($d+1$)

$$N < (d+1)$$

which leads to overfitting.

Thus the model performs poorly on unseen data.

By using L_2 regularization by trying to put a constraint on the squared norm of weights vector, along with $SS\bar{E}$ we generate weights as close to origin as possible.

[Note :- We do this after subtracting the mean^{of labels} from the^{each} labels].

By varying λ , we can regulate the^{amount of} shrinkage of weights.

→ A proper value of λ should be chosen to ensure both less error and lesser squared norm.

$$b) \sum_{i=1}^N r_i (y^{(i)} - \hat{y}^{(i)})^2 = \sum_{i=1}^N r_i (y^{(i)} - \sum_{j=0}^d x_j^{(i)} w_j)^2$$

and $r_i > 0$

now,

$$E(\vec{w}) = \sum_{i=1}^N r_i \left(y^{(i)} - \sum_{j=0}^d x_j^{(i)} w_j \right)^2$$

In Matrix form

$$E(\vec{w}) = (\vec{y} - X\vec{w})^T (R\vec{y} - R X \vec{w})$$

where

$$R = \begin{bmatrix} r_1 & 0 & \dots & 0 \\ 0 & r_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & r_N \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_d^{(N)} \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}$$

$$E(\vec{w}) = \left(\vec{y}^T - \vec{w}^T X^T \right) (R\vec{y} - R X \vec{w})$$

$$= \vec{y}^T R \vec{y} - \vec{y}^T R X \vec{w} - \vec{w}^T X^T R \vec{y} + \vec{w}^T X^T R X \vec{w}$$

$$\nabla(E\vec{w}) = -2 X^T R \vec{y} + 2 (X^T R X) \vec{w} = 0$$

$$\boxed{\vec{w}^* = (X^T R X)^{-1} X^T R \vec{y}}$$

where $R = \begin{bmatrix} r_1 & 0 & 0 & \dots & 0 \\ 0 & r_2 & 0 & \dots & 0 \\ \vdots & & & & \\ 0 & 0 & 0 & \dots & r_N \end{bmatrix}$

(4) Linear regression
In multidimensional labels

$$\boxed{\hat{y}_{ik} = \sum_{j=0}^d x_{ij} w_{jk}}$$

$$\boxed{x_{i0} = 1 \forall i}$$

(or)

$$\boxed{\hat{y}_k^{(i)} = \sum_{j=0}^d x_j^{(i)} w_{jk}}$$

→ the weights matrix will be of dimension

$$(d+1) \times K$$

$$\vec{y}^{(i)} \in \mathbb{R}^K$$

$$W = \begin{bmatrix} w_{01} & w_{02} & \dots & w_{0k} \\ w_{11} & w_{12} & \dots & w_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1} & \vdots & \dots & w_{dk} \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1^{(1)} & \hat{y}_2^{(1)} & \dots & \hat{y}_k^{(1)} \\ \hat{y}_1^{(2)} & \hat{y}_2^{(2)} & \dots & \hat{y}_k^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{y}_1^{(n)} & \hat{y}_2^{(n)} & \dots & \hat{y}_k^{(n)} \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}$$

$$\hat{y} = X \cdot W$$

$\begin{matrix} N \times K & K \times d+1 & d+1 \times K \end{matrix}$

Applying the ~~RSS~~ cost function,

$$E(W) = \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - b_k(x_i))^2$$

$$= \text{tr} [(y - XB)^T (y - XB)]$$

Minimizing this is similar to minimizing prob 1 for each column of y matrix

$$\phi = \sum_{i=1}^N \left[(y_i^{(1)} - \sum_{j=0}^d \pi_j^{(1)} \omega_{j1})^2 + \dots + (y_i^{(d)} - \sum_{j=0}^d \pi_j^{(d)} \omega_{ji})^2 \right]$$

$$E(\omega) = \left(y_1^{(1)} - \sum_{j=0}^d \pi_j^{(1)} \omega_{j1} \right)^2 + \dots + \left(y_1^{(N)} - \sum_{j=0}^d \pi_j^{(N)} \omega_{j1} \right)^2 + \left(y_2^{(1)} - \sum_{j=0}^d \pi_j^{(1)} \omega_{j2} \right)^2 + \dots + \left(y_2^{(N)} - \sum_{j=0}^d \pi_j^{(N)} \omega_{j2} \right)^2 + \dots + \left(y_K^{(1)} - \sum_{j=0}^d \pi_j^{(1)} \omega_{jK} \right)^2 + \dots + \left(y_K^{(N)} - \sum_{j=0}^d \pi_j^{(N)} \omega_{jK} \right)^2$$

$$\nabla E(\omega) = \begin{bmatrix} \frac{\partial E}{\partial \omega_{01}} & \frac{\partial E}{\partial \omega_{02}} & \dots & \frac{\partial E}{\partial \omega_{0K}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial E}{\partial \omega_{d1}} & \frac{\partial E}{\partial \omega_{d2}} & \dots & \frac{\partial E}{\partial \omega_{dK}} \end{bmatrix}$$

$$= \begin{bmatrix} -2 \sum_{k=1}^N \left(y_k^{(1)} - \sum_{j=0}^d \pi_j^{(1)} \omega_{j1} \right) \pi_0^{(1)} & \dots & -2 \sum_{k=1}^N \left(y_k^{(d)} - \sum_{j=0}^d \pi_j^{(d)} \omega_{jd} \right) \pi_d^{(d)} \end{bmatrix}$$

$$= -2 X^T [Y - XW]$$

on taking the $\nabla(E(w)) = 0$ similar to

we get

$$W^* = (X^T X)^{-1} X^T Y$$

For Basis function, $M+1$ parameters.

For multidimensional labels,

$$\hat{y}_{ik} = \sum_{j=0}^d \phi_j(\vec{x}^{(i)}) w_{jk}$$

$$\hat{y}_{ik}^{(i)} = \sum_{j=0}^d \phi_j(\vec{x}^{(i)}) w_{jk}$$

The weights matrix is of dimension $(d+1) \times (k)$

$$\vec{y}^{(i)} \in \mathbb{R}^k$$

w, \hat{y} are having similar matrix notation as above

$$\Phi = \begin{bmatrix} 1 & \phi_1^{(1)} & \dots & \phi_M^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1^{(N)} & \dots & \phi_M^{(N)} \end{bmatrix}$$

Now

$$\hat{y} = \phi w$$

$$E(w) = \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - \phi_{ik}(w))^2$$

$$= \text{tr} \left[(Y - \phi w)^T (Y - \phi w) \right]$$

Minimizing this is similar to minimizing
 prob(2) for each column of Y matrix.

on taking $\nabla(E(w)) = 0$ similar to prob(2)
 we get

$$w^* = (\phi^T \phi)^{-1} \phi^T Y$$

(7) Now

$$X_{ij} = x_j^{(i)} + n_{ij} \rightarrow \textcircled{1}$$

(or)

$$X_i = x^{(i)} + n_i$$

$i = 1 \text{ to } N$

$$L(\vec{\omega}) = \sum_{i=1}^N \left(y^{(i)} - \omega_0 - \sum_{j=1}^d \hat{x}_j^{(i)} \omega_j \right)^2$$

$L(\vec{\omega})$ is cost function.

$$\hat{x}_j^{(i)} = x_j^{(i)} = n_{ij}$$

$$\begin{aligned} L(\vec{\omega}) &= \sum_{i=1}^N \left(y^{(i)} - \omega_0 - \sum_{j=1}^d x_j^{(i)} \omega_j - \sum_{j=1}^d n_{ij} \omega_j \right)^2 \\ &= \sum_{i=1}^N \left[\left(y^{(i)} - \omega_0 - \sum_{j=1}^d x_j^{(i)} \omega_j \right) - \left(\sum_{j=1}^d n_{ij} \omega_j \right) \right]^2 \\ &= \sum_{i=1}^N \left[\left(y^{(i)} - \omega_0 - \sum_{j=1}^d x_j^{(i)} \omega_j \right) \right]^2 + \sum_{i=1}^N \left(\sum_{j=1}^d n_{ij} \omega_j \right)^2 \\ &\quad - 2 \sum_{i=1}^N \left[\left(y^{(i)} - \omega_0 - \sum_{j=1}^d x_j^{(i)} \omega_j \right) \left(\sum_{j=1}^d n_{ij} \omega_j \right) \right] \end{aligned}$$

now

$$\sum_{i=1}^N f_i(x) = NE(f(x))$$

where E is expectation

$$\begin{aligned} L(\vec{\omega}) &= N \left[E \left(y - \omega_0 - \sum_{j=1}^d x_j \omega_j \right)^2 \right] + NE \left[\sum_{j=1}^d n_{ij} \omega_j \right]^2 \\ &\quad - 2NE \left(y - \omega_0 - \sum_{j=1}^d x_j \omega_j \right) E \left[\sum_{j=1}^d n_{ij} \omega_j \right] \end{aligned}$$

Taking term wise

$$\rightarrow 2NE \left(y - \omega_0 - \sum_{j=1}^d \omega_j x_j \right) E \left(\sum_{j=1}^d n_j \omega_j \right)$$

$$= 2NE \left(y - \omega_0 - \sum_{j=1}^d \omega_j x_j \right) E \left(\vec{N} \cdot \vec{\omega} \right)$$

where \vec{N} is a random vector where

each $n_j \sim N(0, \sigma^2)$ & all are

independent: $E(n_i n_j) = E(n_i) \cdot E(n_j) = 0$

$$= 2N \left[E \left(y - \omega_0 - \sum_{j=1}^d \omega_j x_j \right) E \left(\sum_{j=1}^d n_j \omega_j \right) \right]$$

$$\left\{ \begin{aligned} E(n_1 \omega_1 + n_2 \omega_2 + \dots + n_d \omega_d) &= E(n_1 \omega_1) + \dots + E(n_d \omega_d) \\ &= 0 + \dots + 0 = 0 \end{aligned} \right\}$$

0

$$\rightarrow NE(\vec{N} \cdot \vec{\omega})^2 = NE \left(\left(\sum_{i=1}^d n_i \omega_i \right)^2 \right)$$

$$= NE \left((n_1 \omega_1 + n_2 \omega_2 + \dots + n_d \omega_d)^2 \right)$$

$$= NE \left(\sum_{i=1}^d n_i^2 \omega_i^2 + 2 \sum_{j=1}^d \sum_{i=1, i \neq j}^d n_i n_j \omega_i \omega_j \right)$$

$$= NE \left(\sum_{i=1}^d n_i^2 \omega_i^2 \right) + E \left(2 \sum_{j=1}^d \sum_{i=1, i \neq j}^d n_i n_j \omega_i \omega_j \right)$$

$$= 2N \sum_{i=1}^d E(n_i^2) \omega_i^2$$

$$\because E(n_i n_j) = E(n_i) E(n_j) = 0$$

$$= N \sigma^2 \sum_{i=1}^d \omega_i^2 = N \sigma^2 \|\omega\|^2$$

$\rightarrow N \sum \left(y - w_0 - \sum_{j=1}^d x_j w_j \right)^2 = \text{cost fn for input without noise}$

overall

$$\begin{aligned} L(\vec{w}) &= \sum_{i=1}^N \left(y^{(i)} - \sum_{j=1}^d x_j^{(i)} w_j \right)^2 + N\sigma^2 \|\vec{w}\|^2 \\ &= \sum_{i=1}^N \left(y^{(i)} - \sum_{j=1}^d x_j^{(i)} w_j \right)^2 + N\sigma^2 \vec{w}^T \vec{w} \end{aligned}$$

\rightarrow
 (b) Given, $P(\vec{w} | x, y, \frac{1}{\sigma^2}, \frac{1}{\alpha^2})$ is maximised if

$$\boxed{\vec{w}^* = \underset{\vec{w} \in \mathbb{R}^{d+1}}{\text{argmax}} \{ P(\vec{w} | x, y, \frac{1}{\sigma^2}, \frac{1}{\alpha^2}) \}} \quad \text{--- (1)}$$

now given information

$$P(\vec{w} | \frac{1}{\alpha^2}) \sim N(0, \alpha^2 I)$$

$$P(\vec{y} | x, \vec{w}, \frac{1}{\sigma^2}) \sim N(x\vec{w}, \sigma^2 I)$$

$$P(\vec{w} | x, \vec{w}, \frac{1}{\sigma^2}, \frac{1}{\sigma^2}) \propto P(\vec{y} | x, \vec{w}, \frac{1}{\sigma^2}) \cdot P(w | \frac{1}{\sigma^2}) \quad (2)$$

According to Bayes theorem.

Taking ① :-

① is maximised for if $P(\vec{y} | x, \vec{w}, \frac{1}{\sigma^2}) \cdot P(w | \frac{1}{\sigma^2})$ is maximised by ②.

Taking log on both of ②.

$$\log \left(P(\vec{y} | x, \vec{w}, \frac{1}{\sigma^2}) \cdot P(w | \frac{1}{\sigma^2}) \right)$$

$$= \log \left(P(\vec{y} | x, \vec{w}, \frac{1}{\sigma^2}) \right) + \log P(w | \frac{1}{\sigma^2}) \quad (3)$$

Now

$$\text{Since } P(\vec{y} | x, \vec{w}, \frac{1}{\sigma^2}) \sim \mathcal{N}(x\vec{w}, \sigma^2 I)$$

which says all elements of \vec{y} given x, w are independent since the covariance matrix is a diagonal matrix ($\sigma^2 I$)

Thus,

$$P(\vec{y} | x, \omega, \frac{1}{\sigma^2}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - (x\vec{\omega})_i)^2}{2\sigma^2}} \quad (3)$$

Similarly,

$$P(\vec{\omega} | \frac{1}{\alpha^2}) \sim N(0, \alpha^2 I)$$

which says all elements of $\vec{\omega}$ are independent

Since the covariance matrix is a diagonal matrix ($\alpha^2 I$).

$$P(\vec{\omega} | \frac{1}{\alpha^2}) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\alpha^2}} e^{-\frac{(\omega_i)^2}{2\alpha^2}} \quad (4)$$

Sub. (3), (4) in (eq 1).

$$\vec{\omega} = \arg \max_{\vec{\omega} \in \mathbb{R}^{d+1}} \left\{ \log(P(\vec{y} | x, \omega, \frac{1}{\sigma^2})) + \log(P(\vec{\omega} | \frac{1}{\alpha^2})) \right\}$$

Thus we have to maximise

$$\begin{aligned} \log(\cdot) &= 2 \left[\sum_{i=1}^N \left(-\frac{1}{2} \right) \log(2\pi) \right] - \sum_{i=1}^N \left(\frac{1}{2} \right) \log(\sigma^2) \\ &\quad - \sum_{i=1}^N \frac{(y_i - (x\vec{\omega})_i)^2}{2\sigma^2} - \sum_{i=1}^d \frac{1}{2} \log(\alpha^2) \\ &\quad - \sum_{i=1}^d \frac{(\omega_i)^2}{2\alpha^2} \end{aligned}$$

maximising $f(w)$ means minimising

$$J(w) = \frac{\|\vec{y} - X\vec{w}\|^2}{2\sigma^2} + \frac{\|\vec{w}\|^2}{2\lambda}$$

$$\because \sum_{i=1}^n (y_i - (Xw)_i)^2 = \|\vec{y} - X\vec{w}\|^2$$

$$= \frac{1}{2\sigma^2} \left[\|\vec{y} - X\vec{w}\|^2 + \left(\frac{\sigma^2}{\lambda}\right) \|\vec{w}\|^2 \right] \quad \left[\sum_{i=1}^d (w_i)^2 = \|\vec{w}\|^2 \right]$$

(scalar)

This is same as ridge regression loss function with

$$\lambda = \left(\frac{\sigma}{\alpha}\right)^2$$

It can also be seen that maximising the posterior distribution is minimising the

L_2 regression of linear parameters or Ridge regression with

$$\lambda = \left(\frac{\sigma}{\alpha}\right)^2$$

$$(3) \quad \tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^x [1 - e^{-2x}]}{e^x [1 + e^{-2x}]} = \frac{1 - e^{-2x}}{1 + e^{-2x}}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma(2x) = \frac{1}{1 + e^{-2x}}$$

$$2\sigma(2x) = \frac{2}{1 + e^{-2x}}$$

$$2\sigma(2x) - 1 = \frac{1 - e^{-2x}}{1 + e^{-2x}} = \tanh(x)$$

$$\tanh(x) = 2\sigma(2x) - 1$$

$$\hat{y}(x; u) = u_0 + \sum_{j=1}^M u_j \tanh\left(\frac{x - u_j}{s}\right)$$

$$= u_0 + \sum_{j=1}^M u_j \left[2\sigma\left(\frac{2(x - u_j)}{s}\right) - 1 \right] = \sum_{j=1}^M u_j$$

$$= \left[u_0 - \sum_{j=1}^M u_j \right] + \sum_{j=1}^M (2u_j) \sigma\left(\frac{2(x - u_j)}{s}\right)$$

comparing with

$$\hat{y}(x; w) = w_0 + \sum_{j=1}^M w_j \left(\sigma\left(\frac{x - u_j}{s}\right) \right)$$

$$\hat{q}(x, \mu) = \left(\mu_0 - \sum_{j=1}^M \mu_j \right) + \sum_{j=1}^M (2\mu_j) \left(\sigma \left(\frac{x - \mu_0}{t} \right) \right)$$

where $t = s/2$

thus it is equivalent to

$$w_0 = \mu_0 - \sum_{j=1}^M \mu_j$$

$$w_j = 2\mu_j$$