

Theory :-

1) Say

$$\hat{y} = f(x; \vec{w})$$

where $\vec{x} = [1 \ x_1 \ x_2 \ \dots \ x_d]$

& $\vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$ be the ~~best~~

relation between y & x

We try to minimize the error

$$\sum_{i=1}^N (y^i - \hat{y}^i)^2 \quad \text{with respect to}$$

\vec{w} so that we can find the

best fit for $\hat{y} = f(x; \vec{w})$

Let the optimal \vec{w} be \vec{w}^* .

& the matrix of inputs be

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

where $N =$ no. of training samples
 and $d =$ dimension of input data (no. of parameters)

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

where $\hat{y}^{(i)} = \sum_{j=0}^d x_j^{(i)} w_j$

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

$$\Rightarrow \frac{\partial L(w)}{\partial w_k} = -2 \sum_{i=1}^N \left(y^{(i)} - \sum_{j=0}^d x_j^{(i)} w_j \right) x_k^{(i)}$$

$$\Rightarrow \sum_{i=1}^N x_k^{(i)} \left(y^{(i)} - \sum_{j=0}^d x_j^{(i)} w_j \right) = 0$$

$$\Rightarrow \sum_{i=1}^N (X^T)_{ki} \cdot (y - Xw)_i = 0$$

$$\Rightarrow X^T (y - Xw) = 0$$

$$\Rightarrow \boxed{w = (X^T X)^{-1} \cdot X^T y}$$

$\therefore \vec{w^*}$ minimizes $L(w)$ if $X^T X$ is positive definite

$$2.) \hat{y}^{(i)} = \sum_{j=0}^M \phi_j^{(i)} \cdot w_j$$

where $\phi_j^{(i)}$ represents $\phi(x_j^{(i)})$, ϕ being basis function.

$$\text{i.e., } \Phi = \begin{bmatrix} \phi_0(x^{(1)}) & \phi_1(x^{(1)}) & \phi_2(x^{(1)}) & \dots & \phi_H(x^{(1)}) \\ \phi_0(x^{(2)}) & \phi_1(x^{(2)}) & \phi_2(x^{(2)}) & \dots & \dots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi_0(x^{(N)}) & \dots & \dots & \dots & \phi_H(x^{(N)}) \end{bmatrix}$$

$N \times (H+1)$

$$L(w) = \sum_{i=1}^N \left(y^{(i)} - \sum_{j=0}^M \phi_j^{(i)} \cdot w_j \right)^2$$

$$\frac{\partial L(w)}{\partial w_k} = 0 \Rightarrow \sum_{i=1}^N \phi_k^{(i)} \left(y^{(i)} - \sum_{j=0}^M \phi_j^{(i)} \cdot w_j \right) = 0$$

$$\Rightarrow \sum_{i=1}^N (\Phi^T)_{ki} (y - \Phi w)_i = 0$$

$$\Rightarrow \phi^T y = (\phi^T \phi) w$$

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

3.)

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

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

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

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

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

$$\boxed{\therefore \tanh(x) = 2\sigma(2x) - 1}$$

$$\hat{y}(x, \vec{u}) = u_0 + \sum_{j=1}^M u_j \left(2\sigma\left(\frac{2(x-\mu_i)}{s}\right) - 1 \right)$$

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

$$- \sum_{j=1}^M u_j$$

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

Comparing this with

$$\hat{y}(x, w) \quad , \quad (\text{also } 2(x-\mu_i) = \cancel{x} - \mu'_i)$$

$$\Rightarrow \boxed{w_0 = u_0 - \sum_{j=1}^M u_j}$$

$$\& \boxed{w_i = 2u_i}$$

4.) This is a similar argument as in
 ①, ② except

$$\hat{y}_{ik} = \sum_{j=0}^d x_{ij} w_{jk} \Rightarrow \mathbf{Y}_{N \times K} = \mathbf{X}_{N \times (d+1)} \mathbf{W}_{(d+1) \times K}$$

The weight matrix will be a two dimensional matrix (dimensions :- $(d+1) \times K$)

$$L(w) = \sum_{i=1}^N \sum_{k=1}^K (y_{ik} - f_k(x_i))^2$$

~~$$= \text{tr} [(\mathbf{Y} - \mathbf{XW})^T (\mathbf{Y} - \mathbf{XW})]$$~~

$$= \text{tr} [(\mathbf{Y} - \mathbf{XW})^T (\mathbf{Y} - \mathbf{XW})]$$

Minimizing this is similar to
 minimize individual terms as seen
 in prob ①, prob ② except \mathbf{Y} is
 a $N \times K$ dimensional matrix

$$\text{So } \mathbf{W}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

W^* being optimal W

(for basis functions, $X = \phi$)
 ϕ is the same matrix as in prob 2)

$$\Rightarrow W^* = (\phi^T \phi)^{-1} \phi^T y$$

$$b) L(w) = \sum_{i=1}^N r_i \left(y^{(i)} - \sum_{j=0}^d x_j^{(i)} w_j \right)^2, \quad r_i > 0$$

$$\frac{\partial L(w)}{\partial w_k} = 0$$

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

$$\Rightarrow \sum_{i=1}^N r_i \left(y^{(i)} x_k^{(i)} - x_k^{(i)} \sum_{j=0}^d x_j^{(i)} w_j \right) = 0$$

$$\Rightarrow \sum_{i=1}^N$$

$$\frac{\partial L(w)}{\partial w_k}$$

$$L(w) = \sum_{i=1}^N \left(\sqrt{r_i} y_i - \sum_{j=0}^d \sqrt{r_i} x_j^{(i)} w_j \right)^2$$

$$\text{Let } y_i' = \sqrt{r_i} y_i \text{ \& } (x_j^{(i)})' = \sqrt{r_i} x_j^{(i)}$$

$$\Rightarrow L(w) = \sum_{i=1}^N \left(y_i' - \sum_{j=0}^d (x_j^{(i)})' w_j \right)^2$$

Clearly from prob ①, ②

the \vec{w}^* which minimizes above function is

$$\vec{w}^* = \left((X')^T (X') \right)^{-1} (X')^T \cdot y'$$

where

$$X' = \cancel{R_{d \times d}} \cdot \cancel{X_{d \times d}} \cdot R_{N \times N} \cdot X_{N \times (d+1)}$$

$$\& y' = \cancel{R_{d \times d}} \cdot \cancel{y_{d \times 1}} \cdot R_{N \times N} \cdot y_{N \times 1}$$

$$\cancel{R_{d \times d}} = \begin{bmatrix} \sqrt{r_1} & & & \\ & \sqrt{r_2} & & \\ & & \sqrt{r_3} & \\ & & & \ddots \end{bmatrix}$$

$$R_{N \times N} = \begin{bmatrix} \sqrt{r_1} & 0 & 0 & 0 & \dots \\ 0 & \sqrt{r_2} & 0 & 0 & \dots \\ 0 & 0 & \sqrt{r_3} & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

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

$$= (X^T (R^2) X)^{-1} \cdot X^T (R^2) y$$

$$R^2 = \begin{bmatrix} r_1 & 0 & 0 & 0 & \dots \\ 0 & r_2 & 0 & 0 & \dots \\ 0 & 0 & r_3 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}_{N \times N}$$

6.)

$$E(\vec{w}) = (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w}) + \lambda \underbrace{\vec{w}^T \cdot \vec{w}}_{\text{regularizer!}}$$

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

$$\Rightarrow (X^T X + \lambda I) \vec{w} = X^T \vec{y}$$

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

where $X =$

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

$$\& y = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{pmatrix}_{N \times 1}$$

Use of regularization :-

> Clearly $X^T X + \lambda I$ adds a non-zero term to ~~column~~ diagonals of $X^T X$

this makes the matrix $X^T X + \lambda I$ invertible for any matrix X .

> Also, this prevents weights from taking extreme values and reduces the high variance of model (better generalization) by constraining the norm of weight vector.

7. ~~★~~ Here $X_i = x_i + n_i$ is a random variable

$\Rightarrow \vec{X}$ is a random vector

$\Rightarrow \vec{X} \cdot \vec{w}$ is a random variable

and so ~~$y = \vec{X} \cdot \vec{w}$~~ is

$$y = w_0 + \vec{X} \cdot \vec{w}$$

the dimension of \vec{X}, \vec{w} is $1 \times d$ & ~~consider~~

The loss function

$L(w)$ is given by

$$\begin{aligned}
 L(w) &= E \left((y - w_0 - \vec{x}^T \cdot \vec{w})^2 \right) \\
 &= E \left((y - w_0 - \vec{x} \cdot \vec{w} - \vec{N} \cdot \vec{w})^2 \right) \\
 &= E \left((y - w_0 - \vec{x} \cdot \vec{w})^2 - 2(y - w_0 - \vec{x} \cdot \vec{w}) \vec{N} \cdot \vec{w} + (\vec{N} \cdot \vec{w})^2 \right)
 \end{aligned}$$

where \vec{N} is a random vector where each $n_i \sim \mathcal{N}(0, \sigma^2)$.

& all are independent i.e.,

Same as normal regression

$$E(n_i n_j) = 0.$$

$$\begin{aligned}
 \Rightarrow L(w) &= E \left((y - w_0 - \vec{x} \cdot \vec{w})^2 \right) - 2 \left((y - w_0 - \vec{x} \cdot \vec{w}) \cdot E(\vec{N} \cdot \vec{w}) \right) + E \left((\vec{N} \cdot \vec{w})^2 \right)
 \end{aligned}$$

~~a constant~~
not random.

$$E(\vec{N} \cdot \vec{w}) = E \left(\sum_{i=1}^d n_i w_i \right) = 0$$

[$\because E(n_i) = 0$]

$$E((\vec{N} \cdot \vec{w})^2)$$

$$= E\left(\left(\sum_{i=1}^d n_i w_i\right)^2\right)$$

$$= E\left((n_1 w_1 + n_2 w_2 + \dots + n_d w_d)^2\right)$$

$$= E\left(\sum_{i=1}^d n_i^2 w_i^2 + k'\right)$$

Here k' is $2 \sum_{j=1}^d \sum_{i=1}^{j-1} n_i n_j w_i w_j$

$$\Rightarrow E(k') = 0$$

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

$$\Rightarrow E((\vec{N} \cdot \vec{w})^2) = \sum_{i=1}^d E(n_i^2) \cdot w_i^2$$

$$= (\sigma^2) \|\vec{w}\|^2$$

$$\therefore L(w) = E((y - w_0 - xw)^2)$$

$$+ \sigma^2 \|\vec{w}\|^2$$

$$\|\vec{w}\|^2 = \sum_{i=1}^d w_i^2$$

which is same as ~~loss~~ ridge regression loss
function with $\lambda = \sigma^2$
Hence proved

* (Adding noise to input is same as regularization)

Q. $p(\omega) \sim \mathcal{N}(0, \alpha^2 I)$

& $p(\vec{y} | \omega, X) \sim \mathcal{N}(X\vec{\omega}, \sigma^2 I)$

we find $\vec{\omega}$
⊗ maximize

$p(\vec{\omega} | X, y)$

$$8.) \quad p(\omega) \sim \mathcal{N}(0, \sigma^2 I)$$

we find ω such that

$$\Rightarrow \cancel{p(\omega)} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\omega^2}{2\sigma^2}}$$

$p(\omega | x, y)$ is maximized.

we know that

$$p(\vec{y} | x, \omega) \sim \mathcal{N}(X\omega, \sigma^2 I)$$

& also $y_i \in \vec{y}$ are independent.

$$\Rightarrow p(\vec{y} | x, \vec{\omega}) = \prod_{i=1}^N \frac{1}{(\sqrt{2\pi})\sigma} e^{-\frac{(y_i^{(i)} - \sum_{j=0}^d x_j \omega_j)^2}{2\sigma^2}}$$

where d - dimension of features

& N - No. of train samples.

$$p(\vec{\omega} | x, \vec{y}) \propto p(\vec{y} | x, \vec{\omega}) \cdot p(\omega)$$

(According to Baye's theorem)

~~$p(\omega | x, y) \propto p(y | x, \omega) \cdot p(\omega)$~~

$$p(y) \cdot p(w|y) = p(w|y) \cdot p(y)$$

$$p(y|w) \cdot p(w)$$

$$\Rightarrow p(\vec{w} | x, \vec{y}) \propto p(\vec{y} | x, \vec{w}) \cdot p(\vec{w})$$

~~$$\Rightarrow p(\vec{w} | x, \vec{y}) \propto \left(\frac{1}{\pi} \prod_{i=1}^N \frac{1}{\sigma} e^{-\frac{(y^i - xw)^2}{2\sigma^2}} \right)$$~~

~~$$\frac{1}{\pi} \prod_{j=1}^d \frac{1}{\sqrt{2\pi}\alpha} e^{-\frac{w_j^2}{2\alpha^2}}$$~~

~~$$\Rightarrow p(w | x, y) \propto \frac{1}{2\pi\sigma^2} e^{-\frac{(\sum_{i=1}^N (y^i - xw)^2 + w^2)}{2\sigma^2}}$$~~

$$\Rightarrow p(\vec{w} | x, \vec{y}) \propto \left[\prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^i - xw)^2}{2\sigma^2}} \right] \times \left[\prod_{j=1}^d \frac{1}{\sqrt{2\pi}\alpha} e^{-\frac{w_j^2}{2\alpha^2}} \right]$$

$$\Rightarrow p(\vec{w} | x, \vec{y}) \propto \frac{1}{(2\pi)^{N/2} \cdot \sigma^N \cdot \alpha^d \cdot (2\pi)^{d/2}} \cdot e^{-\left(\frac{\|\vec{y} - x\vec{w}\|^2}{2\sigma^2} + \frac{\|\vec{w}\|^2}{2\alpha^2} \right)}$$

\Rightarrow maximizing $p(\vec{w} | x, \vec{y})$ is the same as minimizing the mean which is

$$L(w) = \frac{\|\vec{y} - X\vec{w}\|^2}{2\sigma^2} + \frac{1}{2\alpha^2} (\|w\|^2)$$

~~Scalar~~

$$= \left(\frac{1}{2} \right) \left(\frac{\|\vec{y} - X\vec{w}\|^2}{\sigma^2} + \frac{\|w\|^2}{\alpha^2} \right)$$

$$= \left(\frac{1}{2\sigma^2} \right) \left(\|\vec{y} - X\vec{w}\|^2 + \left(\frac{\sigma}{\alpha} \right)^2 \|w\|^2 \right)$$

Scalar

This is same as ridge regression

loss: function with

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

It can also be observed that this loss function is actually the mean of the posterior distribution

$$p(\vec{w} | X, \vec{y}) \text{ with } \lambda = \left(\frac{\sigma}{\alpha} \right)^2$$

The same $L(\omega)$ is also the mode of
the distribution (peak value in distribution)