

HW-1

① Formulate and solve simplest case of linear regression

Consider N samples for X each of d dimension

$$X = \begin{bmatrix} 1 & x_1^1 & x_1^2 & x_1^3 & \dots & x_1^d \\ 1 & x_2^1 & x_2^2 & \dots & \dots & x_2^d \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_N^1 & x_N^2 & \dots & \dots & x_N^d \end{bmatrix}$$

$N \times d+1$

and the labels Y be

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$N \times 1$

Consider the weights W

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$d+1 \times 1$

then we have the set of equations as

$$XW = \hat{Y}$$

where \hat{Y} is the estimated labels for the given weights W .

Consider the cost function $E(W)$

$$E(W) = \sum_{i=1}^N \|y^{(i)} - \hat{y}^{(i)}\|^2 = \text{sum of squared errors}$$

$$\Rightarrow E(W) = \|y - \hat{y}\|^2 \quad \| \cdot \| \text{ being the norm}$$

for optimal weights we have

$$\nabla E(W) = 0 \quad \left(\text{Gradient wrt } W \right)$$

$$\Rightarrow \sqrt{\|y - \hat{y}\|^2} = 0$$

$$\Rightarrow \sqrt{\|Y - XW\|^2} = 0$$

$$\Rightarrow \frac{\partial}{\partial w_j} \left(y^i - \sum_{j=0}^d x_j^i w_j \right)^2 = 0$$

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

In matrix consolidated form, we have

$$2X^T(Y - XW^*) = 0$$

W^* = optimal weights

$$X^T Y = X^T X W^*$$

now, the first assumption is $X^T X$ is positive definite

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

② ~~Bias~~ Basis function is a function ϕ which we apply on the matrix X to further continue and estimate weights.

So given data X , after applying bias function, we have

$$\phi(X) = \begin{bmatrix} \phi_0'(x_0) & \phi_1'(x_1) & \dots & \phi_M'(x'_M) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0^N(x_0^N) & \phi_1^N(x_1^N) & \dots & \phi_M^N(x_M^N) \end{bmatrix}$$

$$\phi(X) = \begin{bmatrix} \phi_0' & \phi_1' & \dots & \phi_M' \\ \phi_0^2 & \phi_1^2 & \dots & \phi_M^2 \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0^N & \phi_1^N & \dots & \phi_M^N \end{bmatrix}_{N \times (M+1)}$$

Consider the weight matrix W

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix}_{(M+1) \times 1}$$

then we have

$$\hat{y} = \phi(x) W$$

for simplicity sake, drop x in $\phi(x)$

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

Considering cost function as SSE, we have

$$E(w) = \|y - \hat{y}\|^2$$

$$E(w) = \|y - \phi w\|^2$$

for optimal weights w^*

$$\nabla E(w) = 0$$

Similar to vanilla regression,

$\nabla E(w)$ doesn't depend on

X (or) in this case ϕ

hence

$$\nabla E(w) = 0$$

$$\Rightarrow 2\phi^T(Y - \phi w^*) = 0$$

$$\Rightarrow \phi^T Y = \phi^T \phi w^*$$

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

$$(3) \quad \sigma(x) = \frac{1}{1+e^{-x}}$$

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

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

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

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

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

Consider

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

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

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

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

$$= \underbrace{u_0 - \sum_{j=1}^M u_j}_{k_0} + \sum_{j=1}^M \underbrace{2u_j}_{k_j} \sigma\left(\frac{2x - 2\mu_j}{s}\right)$$

$$\hat{y}(x, k) = k_0 + \sum_{j=1}^M k_j \sigma\left(\frac{2x - 2\mu_j}{s}\right)$$

$$\text{take } 2x = z \quad \& \quad 2\mu_j = d_j$$

$$\hat{y}(z, k) = k_0 + \sum_{j=1}^M k_j \sigma\left(\frac{z - \alpha_j}{s}\right)$$

which is equivalent to

$$\hat{y}(x, u)$$

$$\& k_0 = w_0 = u_0 - \sum_{j=1}^M u_j$$

$$\& k_j = 2u_j = w_j$$

$$\therefore w_i = \begin{cases} u_0 - \sum_{j=1}^M u_j & i=0 \\ 2u_i & \text{else} \end{cases}$$

(4) Consider the SSE cost function

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

for i^{th} row of y and k^{th} term (dimension)

$$CF_{ik} = \left(y_k^{(i)} - \sum_{j=0}^d x_j^{(i)} w_{jk} \right)^2$$

for optimal w^*

$$\frac{\partial CF_{ik}}{\partial w_{jk}} = 0 \Rightarrow -2x_j^{(i)} \left(y_k^{(i)} - \sum_{j=0}^d x_j^{(i)} w_{jk} \right) = 0$$

\Rightarrow In matrix format

$$X^T (Y - XW^*) = 0$$

which yields $\boxed{W^* = (X^T X)^{-1} X^T Y}$

In case of basis functions;

$$\boxed{W^* = (\Phi^T \Phi)^{-1} \Phi^T Y}$$

their orders being $d+1 \times K$ and

$(M+1) \times K$ respectively

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

$$= \sum_{i=1}^N \left(Y^{(i)} - \sum_{j=0}^d X_j^{(i)} w_j \right)^2$$

where $Y^{(i)} = \sqrt{r_i} y^{(i)}$

$X_j^{(i)} = \sqrt{r_i} x_j^{(i)}$

for $w^*, \nabla E(w^*) = 0$

$$\Rightarrow 2 X^T (Y - X w^*) = 0$$

$$\Rightarrow \boxed{w^* = (X^T X)^{-1} X^T Y}$$

where $X = R x$

$Y = R y$

$$(x^T R^2 x)^{-1} x^T R^2 y$$

$$R = \begin{bmatrix} \sqrt{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \sqrt{\sigma_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \sqrt{\sigma_N^2} \end{bmatrix}$$

$$Rx = \begin{bmatrix} \sqrt{\sigma_1^2} x_1 & \sqrt{\sigma_1^2} x_2 & \dots \\ \sqrt{\sigma_2^2} x_1 & \sqrt{\sigma_2^2} x_2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$$Ry = \begin{bmatrix} \sqrt{\sigma_1^2} y_1 \\ \sqrt{\sigma_2^2} y_2 \\ \vdots \\ \sqrt{\sigma_N^2} y_N \end{bmatrix}$$

$$⑥ \quad E(w) = \|y - Xw\|^2 + \lambda w^T w.$$

where $\|y - Xw\|^2$ is the standard SSE

and $\lambda w^T w$ is the l_2 norm &

λ is the Lagrangian Multiplier

$$\nabla E(w) = 0 \Rightarrow -2X^T(Y - Xw) + 2\lambda Iw = 0$$

$$\cancel{2}X^TY = \cancel{2}X^TXw^* + \cancel{2}\lambda Iw^*$$

$$\Rightarrow w^* = (X^TX + \lambda I)^{-1} X^TY$$

$$\boxed{w^* = (X^TX + \lambda I)^{-1} X^TY}$$

Regularization is applied when

→ To prevent over fitting of data

→ When data is noisy

⑦ Let $N \sim (0, \sigma^2)$ be added to x

$$X = x + N$$

If $x = \text{matrix } N \times d+1$

$N = \text{matrix of } N \times d+1$

where $N^{(i)} \sim (0, d)$

Now consider SSE cost function for this altered data

$C_F(X, W) = \text{cost function on } X \text{ and } W.$

$$C_F(X, W) = \sum_{i=1}^N \left(\left(y^i - \sum_{j=0}^d x_j^i w_j \right)^2 \right)$$

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

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

$$\text{Let } y^i - \sum_{j=0}^d x_j^i w_j = \alpha$$

$$\sum_{j=0}^d N_j^i w_j = \beta$$

$$C_F(X, W) = \sum_{i=1}^N (\alpha - \beta)^2$$

$$C_F(X, W) = \sum_{i=1}^N (\alpha^2 + \beta^2 - 2\alpha\beta)$$

Now $\sum_{i=1}^N (R.V) = N E(R.V)$

where R.V is a random variable

$E_i \Rightarrow$ expectancy over iterations of i

$$\Rightarrow C_F(X, W) = N E_i (\alpha^2 + \beta^2 - 2\alpha\beta)$$

$$= N E(\alpha^2) + N E(\beta^2) - 2N E(\alpha\beta)$$

note that $N E(\alpha^2) = C_F(\alpha, W)$

i.e., $N E(\alpha^2) =$ cost function of

data without gaussian noise

$$C_F(X, W) = N E \left(y^i - \sum_{j=0}^d x_j^i w_j \right)^2 \\ + N E \left(\left(\sum_{j=0}^d N_j^i w_j \right)^2 \right) - 2N E \left(\left(y^i - \sum_{j=0}^d x_j^i w_j \right) \left(\sum_{j=0}^d N_j^i w_j \right) \right)$$

α is basically a constant

$$\Rightarrow E(\alpha \beta) = \alpha E(\beta)$$

$$\begin{aligned} \& E(\beta) &= E \left(\sum_{j=0}^d N_j^i w_j \right) \\ &= \sum_{j=0}^d w_j E(N_j^i) = 0 \end{aligned}$$

~~$C_F(X, W)$~~

and

$$\begin{aligned} E \left(\left(\sum_{j=0}^d N_j^i w_j \right)^2 \right) &= \|w\|^2 E((N_j^i)^2) \\ &= \sigma^2 \|w\|^2 \\ &= \sigma^2 w^T w. \end{aligned}$$

(By expanding)

$$\Rightarrow C_F(x, w) = C_F(x, w) + \sigma^2 w^T w$$

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

$\therefore C_F(x, w)$ is equivalent to
 l_2 regularized cost function of
 $C_F(x, w)$

⑧ Consider the maximum a posterior probability expression

$$p(W|x, Y, \alpha, \beta)$$

We need to maximize this to get optimal weights

but,

$$p(W|x, Y, \alpha, \beta) \propto p(Y|x, W, \alpha, \beta) \cdot p(W|\alpha)$$

from Bayes theorem

and

$$p(Y|x, W, \alpha, \beta) \sim \mathcal{N}(\hat{Y}, \sigma^2 I)$$

$$p(W|\alpha) \sim \mathcal{N}(0, \alpha^2 I)$$

$$\text{Max}_{\text{aposterior prob}} \propto \frac{1}{(2\pi)^{p/2} \sigma^p} \exp\left(\frac{-(y-\hat{y})^T (y-\hat{y})}{2\sigma^2}\right) \frac{1}{(2\pi)^{p/2} \alpha^p}$$

$$\exp\left(\frac{-W^T W}{2\alpha^2}\right)$$

these expressions follow from PDF's of a Gaussian Random Vector PDF with independent entries.

$$\text{Max}_{\text{Apost Prob}} \propto \frac{1}{(2\pi)^{p/2} \sigma^p \alpha^p} \exp\left(-\left[\frac{(y-\hat{y})^T (y-\hat{y})}{2\sigma^2} + \frac{W^T W}{2\alpha^2}\right]\right)$$

$$\text{Max}_{\text{aposterior prob}} = K \exp\left(-\left[\frac{(y-\hat{y})^T (y-\hat{y})}{2\sigma^2} + \frac{W^T W}{2\alpha^2}\right]\right)$$

For maximum aposterior probability gradient w.r.t W must be zero.

$$\Rightarrow \nabla \left(K \exp \left(- \left[\frac{(y - \hat{y})^T (y - \hat{y})}{2\sigma^2} + \frac{w^T w}{2\alpha^2} \right] \right) \right) = 0$$

$$\Rightarrow \cancel{K \exp \left(- \left[\right.} \right.$$

$$\text{Let } \frac{(y - \hat{y})^T (y - \hat{y})}{2\sigma^2} + \frac{w^T w}{2\alpha^2} = CF$$

$$\nabla K \exp(-CF) = 0$$

$$\Rightarrow K \exp(-CF) \cdot (-\nabla(CF)) = 0$$

$$\Rightarrow \nabla(CF) = 0$$

$$\Rightarrow \nabla \left(\frac{(y - \hat{y})^T (y - \hat{y})}{2\sigma^2} + \frac{w^T w}{2\alpha^2} \right) = 0$$

$$\Rightarrow \nabla \left((y - \hat{y})^T (y - \hat{y}) + \left(\frac{\sigma^2}{\alpha^2} \right) w^T w \right) = 0$$

$$\Rightarrow \nabla \left((y - \hat{y})^T (y - \hat{y}) + \lambda w^T w \right) = 0$$

this is similar to L_2 normalized
cost function for ridge regression

$$\rightarrow \boxed{\lambda = \frac{\sigma^2}{\alpha^2}}$$