

# EE5600: Intro to AI and ML

## 1. Linear Regression:

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

~~(d+1) \times N~~  
 $N \times (d+1)$

$N \rightarrow$  number of training examples

$d \rightarrow$  length of input vector.

$$\bar{y} = \begin{bmatrix} y_1^{(1)} \\ y_2^{(1)} \\ \vdots \\ y_N^{(1)} \end{bmatrix} \rightarrow \text{outputs (labels)} \quad \bar{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \rightarrow \text{weights}$$

$N \times 1$        $(d+1) \times 1$

Cost function .  $E(\bar{w}) = \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2$   
(sum of squared error)

Here  $\hat{y}^{(i)}$  is predicted output.

$$\hat{\bar{y}} = X \cdot \bar{w} \rightarrow (N \times 1)$$

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

$$= (\bar{y} - X \bar{w})^T (\bar{y} - X \bar{w})$$

optimal weight,  $\bar{\omega}^* = \underset{\bar{\omega}}{\operatorname{argmin}} (E(\bar{\omega}))$   
 (weights that minimize the sum of squared error)

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

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

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

$$\Rightarrow -2 X^T (\bar{y} - X \bar{\omega}) = 0$$

$$\Rightarrow \bar{\omega}^* = (X^T X)^{-1} X^T \bar{y}$$

To show that this is minimum,

check second derivative ~~of~~.  $\Rightarrow 2 X^T X$

This is positive definite (if  $X$  has full column rank)

So, 
$$\bar{\omega}^* = (X^T X)^{-1} X^T \bar{y}$$

2. Basis functions:  $\phi_j(\bar{x}^{(i)}) \rightarrow 0 \leq j \leq M$

$$\bar{\omega} = \begin{bmatrix} \omega_0 \\ \omega_1 \\ \vdots \\ \omega_M \end{bmatrix}_{(M+1) \times 1}$$

$$\phi_0(\bar{x}) = 1$$

assuming  $\phi_j(\cdot)$  are  $\mathbb{R}^{(d+1)} \rightarrow \mathbb{R}$

Prediction, 
$$\hat{y}^{(i)} = \sum_{j=0}^M \phi_j(\bar{x}^{(i)}) \omega_j$$

Cost function,  $E(\bar{\omega}) = \sum_{i=1}^N \left[ y^{(i)} - \sum_{j=0}^M (\phi_j(\bar{x}^{(i)}) \omega_j) \right]^2$

Notation:-  $\Phi_j(x) = \begin{bmatrix} \phi_j(\bar{x}^{(1)}) \\ \vdots \\ \phi_j(\bar{x}^{(N)}) \end{bmatrix}$

Notation

$\hookrightarrow \Phi = \Phi(x) = \begin{bmatrix} \phi_0(\bar{x}^{(1)}) & \phi_1(\bar{x}^{(1)}) & \dots & \phi_M(\bar{x}^{(1)}) \\ \phi_0(\bar{x}^{(2)}) & \phi_1(\bar{x}^{(2)}) & \dots & \phi_M(\bar{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\bar{x}^{(N)}) & \phi_1(\bar{x}^{(N)}) & \dots & \phi_M(\bar{x}^{(N)}) \end{bmatrix}$

$$E(\bar{\omega}) = (\bar{y} - \Phi \bar{\omega})^T (\bar{y} - \Phi \bar{\omega})$$

minimizing the cost function,

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

$$\Rightarrow -2 \Phi^T \bar{y} + 2 \Phi^T \Phi \bar{\omega} = 0$$

$$\Rightarrow \boxed{\bar{\omega}^* = (\Phi^T \Phi)^{-1} \Phi^T \bar{y}}$$

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

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

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

$$\hat{y}(x, \underline{\omega}) = \omega_0 + \sum_{i=1}^M \omega_i \cdot \sigma\left(\frac{x - u_i}{s}\right)$$

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

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

Take  $2x = z$ ,  $2\underline{u} = \underline{v} \Rightarrow u_j = \frac{k_j}{2}$

$\Rightarrow \hat{y}\left(\frac{z}{2}, \frac{\underline{v}}{2}\right) = \frac{v_0}{2} + \sum_{j=1}^M u_j \sigma\left(\frac{z - k_j}{s}\right) - \frac{v_j}{2}$

$\hat{y}(x, \underline{v})$

$k_j = \text{mean of } u_j$   
 $[z \text{ replaced by } x]$

$$w_0 = \frac{1}{2} \left( v_0 - \sum_{j=1}^M v_j \right) = u_0 - \sum_{j=1}^M u_j$$

$$(w_1, w_2, \dots, w_M) = (2u_1, 2u_2, \dots, 2u_M)$$

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

$\rightarrow$  Inputs

outputs,  $Y = \begin{bmatrix} y_1^{(1)} & y_2^{(1)} & \dots & y_k^{(1)} \\ y_1^{(2)} & y_2^{(2)} & \dots & y_k^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ y_1^{(N)} & y_2^{(N)} & \dots & y_k^{(N)} \end{bmatrix}_{N \times k}$

~~$\hat{y}$~~  =

$\bar{x}^{(i)}$   $\rightarrow$  input vector,  $\bar{y}^{(i)}$   $\rightarrow$  output vector,  $1 \leq i \leq N$

Weights:  $W = \begin{bmatrix} w_0^{(1)} & w_1^{(1)} & \dots & w_d^{(1)} \\ w_0^{(2)} & w_1^{(2)} & \dots & w_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ w_0^{(K)} & w_1^{(K)} & \dots & w_d^{(K)} \end{bmatrix}_{(d+1) \times K}$

Cost function  $E(W) = \sum_{i=1}^N \| \bar{y}^{(i)} - \hat{y}^{(i)} \|^2$

$\hat{y}^{(i)}$  → predicted output vector.

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

~~$\hat{y}^{(i)} = x^{(i)} W = \sum_{p=0}^d x_p^{(i)} w_p^{(i)}$~~

$$E(W) = \sum_{i=1}^N \sum_{j=1}^K (y_j^{(i)} - \hat{y}_j^{(i)})^2$$

$$= \sum_{i=1}^N \sum_{j=1}^K \left( y_j^{(i)} - \sum_{p=0}^d (x_p^{(i)} w_p^{(j)}) \right)^2$$

$$= \sum_{j=1}^K \sum_{i=1}^N \left[ y_j^{(i)} - \sum_{p=0}^d (x_p^{(i)} w_p^{(j)}) \right]^2$$

let's write this as  $E(W) = \sum_{j=1}^K E(\bar{w}^{(j)})$

$$E(W) = \sum_{j=1}^K E(\bar{w}^{(j)})$$

where  $\bar{w}^{(j)} = (w_0^{(j)}, w_1^{(j)}, \dots, w_d^{(j)})^T$

i.e.  $j$ th ~~row~~ column in  $W$ .

Minimizing  $E(W)$  is minimizing each term

in  $\sum_{j=1}^k E(\bar{w}^{(j)})$

[each of  $E(\bar{w}^{(j)})$  is

a positive quantity  
non-negative]

i.e.

$$\bar{w}^{(j)*} = (X^T X)^{-1} X^T \bar{y}_j$$

optimizing,

$$\bar{w}^{(j)*} = (X^T X)^{-1} X^T \bar{y}_j$$

where  $\bar{y}_j$  is  $j$ th column in  $Y$ .

So optimizing weights matrix:

$$W = \begin{bmatrix} ((X^T X)^{-1} X^T \bar{y}_1)^T & (X^T X)^{-1} X^T \bar{y}_2 & \dots & (X^T X)^{-1} X^T \bar{y}_k \end{bmatrix}$$

$$= (X^T X)^{-1} X^T \begin{bmatrix} \bar{y}_1 & \bar{y}_2 & \dots & \bar{y}_k \end{bmatrix}$$

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

5. Simple linear regression with weighted sum of squared error:

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

$$r_i > 0$$

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

$$= (\bar{y} - \hat{\bar{y}})^T R (\bar{y} - \hat{\bar{y}})$$

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

diagonal matrix

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

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

minimizing this  $\Rightarrow \nabla E(\bar{w}) = 0$

$$\cancel{\nabla E(\bar{w}) = -2(RX)^T \bar{y} +}$$

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

$$\Rightarrow \bar{w}^* = (X^T R X)^{-1} (RX)^T \bar{y}$$

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

6. Regularization:-

put a constraint on the  $L_2$ -norm of the weight vector  $\bar{w}$ .

\* Bias is excluded.

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_p^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(N)} & x_2^{(N)} & \dots & x_p^{(N)} \end{bmatrix}_{N \times p}$$

$$\bar{\omega} = [\omega_1, \omega_2, \dots, \omega_d]^T$$

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

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

$$\Rightarrow \nabla E(\bar{\omega}) = -2X^T(\bar{y} - X\bar{\omega}) + 2\lambda I \bar{\omega} = 0$$

$$\Rightarrow -2X^T\bar{y} + 2X^TX\bar{\omega} + 2\lambda I \bar{\omega} = 0$$

$$\Rightarrow \bar{\omega} \neq X$$

$$(XX^T + \lambda I) \bar{\omega} = X^T \bar{y}$$

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

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

Regularization helps in avoiding overfitting.

8. When training data is noisy, the model won't be able to generalize well to future data. Regularization helps in shrinking these learned estimates to zero.

7. Cost function :

$$E(\bar{\omega}) = \sum_{i=1}^N (y_i - \omega_1 x_i - \omega_0 - \omega_2 z_i)^2$$

model:

$$\begin{aligned} \hat{y} &= \omega_1 x' + \omega_0 \\ &= \omega_1 (x + z) + \omega_0 \end{aligned}$$



$$\Rightarrow E(\bar{w}) = \sum_{i=1}^N \left\{ (y_i - w, x_i - w_0)^2 + (w, z_i)^2 - 2(y_i - w, x_i - w_0)(w, z_i) \right\}$$

$$= \cancel{N E[y_i - w]}$$

$$= \sum_{i=1}^N (y_i - w, x_i - w_0)^2 + N E[(w, \bar{z})^2] - 2 N E[(y_i - w, x_i - w_0)(w, z_i)]$$

$$\Rightarrow E(\bar{w}) = \sum_{i=1}^N (y_i - w, x_i - w_0)^2$$

$$+ N w_1^2 E[(\bar{z})^2]$$

$$- 2 N E[(y_i - w, x_i - w_0)] E[\bar{z}] \cdot w_1$$

$E[\cdot] \rightarrow \text{expectation}$

$$E[x] = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\bar{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}$$

$$= \sum_{i=1}^N (y_i - w, x_i - w_0)^2 + N^2 w_1^2 \sigma^2$$

$$| E[\bar{z}] = 0$$

this is noise-free cost function

$$\lambda w_1^2$$

$\lambda^2$  norm regularization term.

8. We have  $p(\bar{y} | x, \bar{w}, \sigma^2) \sim \mathcal{N}\left(\sum_{j=0}^d \bar{x}_j \bar{w}_j, \sigma^2 I\right)$

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

~~$$p(\bar{w} | x, \bar{y}, \alpha, \sigma) \sim \mathcal{N}(0, \alpha^2 I)$$~~

By

we

need to find  $\bar{w}$  that maximizes

posterior distribution of  $\bar{w}$  given  $X, \bar{y}, \alpha, \sigma$

i.e.  $p(\bar{w} | X, \bar{y}, \alpha, \sigma)$

By Bayes' theorem,

$$p(\bar{w} | x, \bar{y}, \alpha, \sigma) = \frac{p(\bar{y} | x, \bar{w}, \sigma) \cdot p(\bar{w} | \alpha)}{p(x, \bar{w}, \sigma, \alpha)}$$

$$\Rightarrow p(\bar{w} | x, \bar{y}, \alpha, \sigma) \propto p(\bar{y} | x, \bar{w}, \sigma) p(\bar{w} | \alpha)$$

$$p(\bar{y} | x, \bar{w}, \sigma) p(\bar{w} | \alpha) = \frac{1}{(\sqrt{2\pi\sigma^2})^d} \exp\left(-\frac{(\bar{y} - \hat{\bar{y}})^T (\bar{y} - \hat{\bar{y}})}{2\sigma^2}\right).$$

$$\hat{\bar{y}} = \sum_{j=0}^d x_j \bar{w}_j$$

$$\frac{1}{(\sqrt{2\pi\alpha^2})^d} \exp\left(-\frac{\bar{w}^T \bar{w}}{2\alpha^2}\right)$$

$$= \frac{1}{(\sqrt{2\pi\sigma^2})^d} \frac{1}{(\sqrt{2\pi\alpha^2})^d} \exp\left[-\left\{\frac{(\bar{y} - \hat{\bar{y}})^T (\bar{y} - \hat{\bar{y}})}{2\sigma^2} + \frac{\bar{w}^T \bar{w}}{2\alpha^2}\right\}\right]$$

$$= \frac{1}{(\sqrt{2\pi}\sigma^2)^N} \frac{1}{(\sqrt{2\pi}\alpha^2)^d} \exp \left[ - \frac{[(\bar{y} - \hat{\bar{y}})^T (\bar{y} - \hat{\bar{y}}) + \frac{\sigma^2}{\alpha^2} \bar{w}^T \bar{w}]}{2\sigma^2} \right]$$

maximizing this probability is minimizing

the expression in exponent (without -ve sign)

i.e. minimizing  $(\bar{y} - \hat{\bar{y}})^T (\bar{y} - \hat{\bar{y}}) + \frac{\sigma^2}{\alpha^2} \bar{w}^T \bar{w}$

This expression is similar to  $\lambda_2$  regularization

so with  ~~$\lambda = \frac{\sigma^2}{\alpha^2}$~~   $\lambda = \frac{\sigma^2}{\alpha^2}$