Ans: Considering linear model

$$y(x,w) = w_0 + \sum w_i x_i \quad \text{for } i = 1 \text{ to } D$$

∴ D → features

Error, $$E_D(w) = \frac{1}{2} \sum_{n=1}^{N} \{y_n(x_n, w) - t_n\}^2$$

Note: This Sol'n is formulated in terms of

matrix / vector calculations.

→ Denotions:

Data → $$\overline{X} = \begin{bmatrix} \overline{x^1} \\ \vdots \\ \overline{x^N} \end{bmatrix}$$
$$N \times (D+1)$$

where $\overline{x_*^n}$ → row vector

$$\overline{x^n} = \begin{bmatrix} x_0^n & x_1^n & \cdots & x_{D+1}^n \end{bmatrix}$$

Here, $\boxed{x_0^n = 1 \quad \forall \, n = 1 \text{ to } N}$

∴ $$\overline{x^n} = \begin{bmatrix} 1, & x_1^n & \cdots & x_{D+1}^n \end{bmatrix}$$

Here, for every data point $x_1, x_2, x_3, ... x_n$ is added → which is Gaussian distributed

$\bar{x}_k^n = x_k^n + \epsilon_k^n$

Such that

$\bar{\epsilon}^n = [\epsilon_1^0, \epsilon_2^0, \epsilon_3^0 ... \epsilon_D^0]$

Here, $\boxed{E(\epsilon_i^0) = 0}$

$\boxed{E(\epsilon_i^0 \epsilon_j^0) = \sigma^2 \delta_{ij}}$

→ cbs NOTE!

Since here $\bar{\epsilon}^n$ is a row vector,
where, it has $(D+1)$ elements, Since
no error is added to $x_0^n$, $x_0^n \to \boxed{\epsilon_0^n = 0}$, $i \in [1, D]$

→ Each element among other $D$ ⇒ $i \in [1, D]$

$\epsilon_i$ are generated from a gaussian distribution of mean '0' and variance: $\sigma^2$

$\boxed{\epsilon_0^n = 0}$

$\therefore \quad E(\varepsilon_i) = 0 \quad \forall \; i$

$\Rightarrow \quad E(\varepsilon_i, \varepsilon_j) = \sigma^2 \delta_{ij} \; ; \quad i,j \in \{1, D\}$

$E(\varepsilon_0 \varepsilon_j) = 0 \quad \Rightarrow \quad \boxed{\varepsilon_0 = 0}$

$\Rightarrow \quad cov(\varepsilon) = \sigma^2 \begin{bmatrix} 0 & 0 & \cdots & \cdots & 0 \\ 0 & 1 & & \cdots & 0 \\ 0 & & 1 & & 0 \\ 0 & & & 1 & 0 \\ 0 & & & & 1 \\ 0 & & & & \end{bmatrix}$

$(D+1) \times (D+1)$

Note: First & row of $cov(\varepsilon) = \bar{0}$

---

(a) Minimising $E_D$ averaged over the noise distribution.

$\Rightarrow J(w) = \underset{\varepsilon_i}{E}\left( \left(\frac{1}{2}\right) \sum_{n=1}^{N} \{y_n \{\bar{x}_n, w\} - t_n\}^2 \right)$

$\Rightarrow \quad y_n\{\bar{x}_n, w\} = \underset{1\times(D+1)}{\bar{x}_n^T} \cdot \underset{(D\times 1, 1)}{\bar{w}}$

$t_n \rightarrow$ ground truth.

Matrix form

$W = \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix} \quad \rightarrow \quad$ weight vector

$(D+1, 1)$

$$J(w) = E_{\varepsilon}\left( \tfrac{1}{2} \| \hat{x}w - t \|^2 \right)$$

Here $t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}_{N \times 1} \longrightarrow$ ground truth.

Here $\hat{x}$ is the noised data

$$\hat{x} = \bar{x} + \bar{\varepsilon}$$
$$(N, D+1)$$

$$\Rightarrow J(w) = E_{\varepsilon_i}\left( \tfrac{1}{2} \| (x+\varepsilon)w - t \|^2 \right)$$

We need to minimise $J(w)$.
Gradient for $k$ th feature

$$\Rightarrow \frac{\partial J(w)}{\partial w_k} = \frac{\partial}{\partial w_k} E_{\varepsilon_i}\left( \tfrac{1}{2} ((x+\varepsilon)w - t)^T ((x+\varepsilon)w - t) \right)$$

$$\frac{\partial J(w)}{\partial w} = E_{\varepsilon}\left( (x+\varepsilon)^T ((x+\varepsilon)w - t) \right)$$

$$\left( \because \frac{\partial(\|A\|^2}{\partial A} \cdot \frac{\partial A^T A}{\partial A} = 2A \right)$$

$$\Rightarrow \frac{\partial J}{\partial W} = E_\varepsilon \left( (x+\varepsilon)^T (x+\varepsilon) W - (x+\varepsilon)^T t \right)$$

$$= E_\varepsilon \left( (x^T x + \varepsilon^T x + x^T \varepsilon + \varepsilon^T \varepsilon) W - x^T t - \varepsilon^T t \right)$$

$$\Rightarrow = \left( x^T x + E_\varepsilon(\varepsilon^T) x + x^T E_\varepsilon(\varepsilon) + E_\varepsilon(\varepsilon^T \varepsilon) \right) W$$
$$- p x^T t - E(\varepsilon^T t)$$

∵ x, t. are independent of ε

$$\Rightarrow \quad E(x\varepsilon) = x E(\varepsilon)$$

$$\frac{\partial J}{\partial W} = \left( x^T x + 0 + 0 + Cov(\varepsilon) \right) W$$
$$- x^T t - 0$$

$$\boxed{\frac{\partial J}{\partial W} = \left( x^T x + Cov(\varepsilon) \right) W - x^T t}$$

To minimize J, $\frac{\partial J}{\partial W} = 0$

$$\Rightarrow \left( x^T x + Cov(\varepsilon) \right) W = x^T t$$

$$\boxed{W = \left( x^T x + Cov(\varepsilon) \right)^{-1} (x^T t)} \rightarrow ①$$

↳ Result -1.

## Part b :

Minimising sum of squares error for noise free variables with weight decay regularisation

i.e, $$J(w) = \frac{1}{2} \sum_{n=1}^{N} \{ y_n(x_n, w) - t_n y \}^2 + \lambda \tilde{w}^T w$$

$$J(w) = \frac{1}{2} \| Xw - t \|^2 + \underbrace{\frac{\lambda}{2} \| \tilde{w} \|^2}_{}$$

$\downarrow$ noise-free data        $\downarrow$ regularisation

Here, $\cdot w = \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix}_{(D+1 \, \times \, N)}$

Note: The regularisation has a omitted $w_0$

i.e, $$\tilde{w} = \begin{bmatrix} 0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}$$

$$J(w) = \frac{1}{2} \| Xw - t \|^2 + \frac{\lambda}{2} \| \tilde{w} \|^2$$

$$\frac{\partial J(\cdot)}{\partial w} = \left(\frac{1}{2}\right)^{(2)} \left( (Xw-t)^T (X) \right) + \frac{\lambda}{2} \left( 2 (w^T) \tilde{I} \right) = 0$$

Nok. $\frac{\partial \| \tilde{w} \|^2}{\partial w}$, $2 w^T \tilde{I}$, $\tilde{I} = \begin{bmatrix} 0 & 0 & 0 & \cdots \\ 0 & 1 & 0 & \cdots \\ & & & \\ & & & 1 \end{bmatrix}$

NEED
$(D+1 \times D+1)$

→ First row of $\tilde{I} = \bar{0}$

→ $\frac{\partial J}{\partial w}$, $(Xw-t)^T X + \lambda (w^T \tilde{I})$

→ To minimise, $\frac{\partial J}{\partial w} = 0$

② $(Xw-t)^T X + \lambda (w^T \tilde{I}) = 0$

→ $X^T (Xw-t) + \lambda (\tilde{I} w) = 0$

→ $(X^T X + \lambda \tilde{I}) w = X^T t$

→ $\boxed{w = (X^T X + \lambda \tilde{I})^{-1} (X^T t)} \longrightarrow ②$

## Conclusion:

From ①:

$$W = (X^TX + Cov(\varepsilon))^{-1}(X^T t)$$

From ②:

$$W = (X^TX + \lambda \tilde{I})^{-1}(X^T t)$$

we know that $Cov(\varepsilon) = \begin{bmatrix} 0 & 0 & & & 0 \\ 0 & \sigma^2 & & & 0 \\ 0 & 0 & \sigma^2 & \cdots & 0 \\ & & & & \sigma^2 \end{bmatrix}$

$$Cov(\varepsilon) = \sigma^2 \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 1 & & 0 \\ & & & 1 \\ & & & 1 \end{bmatrix} \quad (D+1 \times D+1)$$

$$Cov(\varepsilon) = \sigma^2 \tilde{I}$$

∴ Both results are the same

where $\boxed{\lambda = \sigma^2}$

— X —