

Ans 1. (a)

Given: $E_D(w) = \frac{1}{2} \sum_{n=1}^N g_n (t_n - w^T \phi(x_n))^2$

→ So taking gradient of $E_D(w)$ and set it to zero (For minimum)

$$\therefore \frac{\partial}{\partial w} E_D(w) = - \sum_{n=1}^N g_n (t_n - w^T \phi(x_n)) \phi(x_n) = 0 \quad \text{--- (1)}$$

∴ From eqⁿ (1) →

$$\sum_{n=1}^N g_n t_n \phi(x_n) = \left(\sum_{n=1}^N g_n \phi(x_n) \phi(x_n)^T \right) w$$

$$\therefore w = \frac{\left(\sum_{n=1}^N g_n t_n \phi(x_n) \right)}{\left(\sum_{n=1}^N g_n \phi(x_n) \phi(x_n)^T \right)} \quad \underline{\underline{as}}$$

Ans 1. (b) The linear model for the output eqⁿ →

$$\hat{y}_i = w^T x_i + \hat{w}_i$$

where, \hat{w}_i is $\mathcal{N}(0, \sigma_i^2)$

∴ i.e. - Normalize the distribution with 0 and σ_i^2 as parameters then making $\sigma_i^2 = \frac{1}{2g_i}$ we can get the original one.

~~we can~~

It can also be derived from making each i^{th} element g_i times so that the $T(w)$ derived.

Hence the above ~~are~~ are the two alternative interpretations of weighted sum of sq. error function in terms of data-dependent noise variance and replicated data points.

Ans 2. MAP estimate :-

$$i > \arg \max_{h_i} P\left(\frac{h_i}{D}\right) \Rightarrow h_1 \quad (\because P\left(\frac{h_1}{D}\right) = 0.4)$$

also, $P\left(\frac{F}{h_1}\right) = 1 \Rightarrow$ this implies that

MAP estimate is "Forward".

Bayes Optimal Estimate :-

$$P\left(\frac{K=F}{D}\right) = \sum_{h_i} P\left(\frac{F}{h_i}\right) P\left(\frac{h_i}{D}\right) = (1) * (0.4) = 0.4 \quad \text{--- (1)}$$

$$P\left(\frac{K=L}{D}\right) = \sum_{h_i} P\left(\frac{L}{h_i}\right) P\left(\frac{h_i}{D}\right) = (0.2)(1) + (0.1)(1) + (0.2)(1) = 0.5 \quad \text{--- (2)}$$

$$P\left(\frac{K=R}{D}\right) = \sum_{h_i} P\left(\frac{R}{h_i}\right) P\left(\frac{h_i}{D}\right) = 0.1 \quad \text{--- (3)}$$

\therefore From eqⁿ (1), (2), (3) i.e. Bayes Optimal Estimate is Left . as

∴ MAP Estimate is Forward (F)
and Bayes Optimal Estate is Left (L)

So they are not same, because MAP chooses the most probable hypothesis function then evaluate function of predicted data.

$$\text{i.e. } \rightarrow h_{\text{MAP}} = \operatorname{argmax} (P(D|h)) P(h)$$

$$\text{Forward (F)} \Rightarrow \operatorname{argmax} (0.4 \times 1) = 0.4$$

$$\text{Left (L)} \Rightarrow \operatorname{argmax} (0.2 \times 1, 0.1 \times 1, 0.2 \times 1)$$

$$\Rightarrow \underline{0.2}$$

$$\text{Right (R)} \Rightarrow \operatorname{argmax} (0.1 \times 1) = 0.1$$

~~Bayes~~ → Bayes Optimal Classifier makes most probable prediction using entire training data for making predictions.

It also uses hypothetical space for making predictions.

Ans 3. ~~Q3~~ Let two data points y_1 and y_2
and $y_1 < y_2$.

So they can always be shattered by H , no matter how they are labeled.

Explanation \rightarrow

(a) Case-a : if y_1 is positive, and y_2 negative
then we'll choose $a < y_1 < b < y_2$.

Case-b : if y_1 is negative, and y_2 is positive,
then choose $\rightarrow y_1 < a < y_2 < b$.

Case-c : if y_1 is positive and y_2 is also positive
then choose $\rightarrow a < y_1 < y_2 < b$.

Case d if y_1 is negative & y_2 is also negative
then choose $\rightarrow a < b < y_1 < y_2$.

Now, if we have three points $y_1 < y_2 < y_3$ and if they
are labeled as y_1 as positive, y_2 as negative,
and y_3 as positive.

Then in above scenario they cannot be
shattered by H .

Hence $VC(H) = 2$ as.

Ans 4. Given: $y(x, w) = w_0 + \sum_{k=1}^D w_k x_k$

and $E(w) = \frac{1}{2} \sum_{i=1}^N (y(x_i, w) - t_i)^2$

∴ According to question, independent noise is added to each dimension of each input variable x_i

Hence our new model becomes →

$$y'(x_i, w) = w_0 + \sum_{k=1}^D w_k (x_{ik} + E_{ik})$$

$$y'(x_i, w) = w_0 + \underbrace{\sum_{k=1}^D w_k x_{ik} + \sum_{k=1}^D w_k E_{ik}}_{\text{(here } E_{ik} \text{ is independent of } i \text{ and } k)}$$

Hence our new error function is →

$$E'_D(w) = \frac{1}{2} \sum_{i=1}^N (y'(x_i, w) - t_i)^2$$

$$= \frac{1}{2} \sum_{i=1}^N \left(y(x_i, w) + \sum_{k=1}^D w_k E_{ik} - t_i \right)^2$$

$$= \frac{1}{2} \sum_{i=1}^N \left[(y(x_i, w) - t_i)^2 + 2 (y(x_i, w) - t_i) \left(\sum_{k=1}^D w_k E_{ik} \right) + \left(\sum_{k=1}^D w_k E_{ik} \right)^2 \right]$$

Taking the expectation of above eqⁿ and using linearity of expectation $\rightarrow \dots E[\epsilon_{ik}] = 0$

$$E[E'_D(\omega)] = \frac{1}{2} \sum_{i=1}^N \left[(y(x_i, \omega) - t_i)^2 + 2 \left(y(x_i, \omega) - t_i \right) \left(\sum_{k=1}^D \omega_k E[\epsilon_{ik}] \right) + E \left[\left(\sum_{k=1}^D \omega_k \epsilon_{ik} \right)^2 \right] \right]$$

Now,

$$E \left[\left(\sum_{k=1}^D \omega_k \epsilon_{ik} \right)^2 \right] = E \left[\sum_{k=1}^D \sum_{k'=1}^D \omega_k \omega_{k'} \epsilon_{ik} \epsilon_{ik'} \right]$$

$$= \sum_{k=1}^D \sum_{k'=1}^D \omega_k \omega_{k'} E[\epsilon_{ik} \epsilon_{ik'}]$$

$$= \sum_{k=1}^D \sum_{k'=1}^D \omega_k \omega_{k'} \delta_{kk'}$$

$$= \sum_{k=1}^D \omega_k^2 \quad \text{--- (1)}$$

Using result of eqⁿ (1) we get \rightarrow

$$E[E'_D(\omega)] = \frac{1}{2} \sum_{i=1}^N \left[(y(x_i, \omega) - t_i)^2 + \sum_{k=1}^D \omega_k^2 \right]$$

$$= E_D(\omega) + \frac{N}{2} \sum_{k=1}^D \omega_k^2$$

↓

Hence here we get a $L_2^{(norm)}$ regularization term without the bias parameter ω_0 .

Hence Proved.