

Given $E_D(w) = \frac{1}{2} \sum_{n=1}^N g_n (t_n - w^T \phi(x_n))^2$
 Make $\nabla E_D(w) = 0$ for getting w^*
 $\nabla E_D(w) = \frac{1}{2} \sum_{n=1}^N g_n (t_n - w^T \phi(x_n)) \phi(x_n)^T$

lets denote $\sqrt{g_n} \phi(x_n) = \phi'(x_n)$ & $\sqrt{g_n} t_n = t'_n$.

$$\Rightarrow \sum_{n=1}^N g_n t_n \phi(x_n)^T - \sum_{n=1}^N g_n w^T \phi(x_n) \phi(x_n)^T = 0$$

$$\Rightarrow \sum_{n=1}^N \phi'(x_n)^T t'_n - \sum_{n=1}^N w^T \phi'(x_n) \phi'(x_n)^T = 0$$

$$\Rightarrow \sum_{n=1}^N \phi'(x_n)^T t'_n = w^T \sum_{n=1}^N \phi'(x_n) \phi'(x_n)^T$$

From this we can simply derive a result for w^* which minimizes the above error function.

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

But t is defined as:

$$t = [\sqrt{g_1} t_1, \sqrt{g_2} t_2, \dots, \sqrt{g_N} t_N]^T$$

lets define ϕ as $N \times M$ matrix, with element $\phi_{ij} = \sqrt{g_i} \phi_j(x_i)$

1b. Interpretation:-

(i) data dependent noise variance:-

Considering a gaussian noise model, let us assume target variable t is given by deterministic function $y(x, w)$ with additive gaussian noise so that $t = y(x, w) + \epsilon$

where ϵ is a zero mean Gaussian RV with precision

(converse variance) β .

$$p(t/x, w, \beta) = \mathcal{N}(t/y(x, w), \beta^{-1})$$

$$E(t/x) = \int t p(t/x) dt = y(x, w)$$

$$p(t/x, w, \beta) = \prod_{i=1}^n \mathcal{N}(t_n/w^T \phi(x_n), \beta^{-1})$$

Now considering logarithm of likelihood.

$$\ln P(t/x, w, \beta) = \sum_{n=1}^N \ln \left(\frac{1}{\sqrt{\beta}} \exp \left(-\frac{1}{2\beta} (t_n - w^T \phi(x_n))^2 \right) \right)$$

$$= \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \frac{1}{2\beta} E_0(w)$$

Here sum of squares Error is given by

$$E_0(w) = \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

Consider gradient of log of likelihood

$$\nabla \ln P(t/w, \beta) = \sum_{n=1}^N (t_n - w^T \phi(x_n)) \phi(x_n)^T$$

Setting gradient to zero,

$$0 = \sum_{n=1}^N t_n (\phi(x_n))^T - w^T \sum_{n=1}^N \phi(x_n) \phi(x_n)^T$$

$$\Rightarrow \boxed{w = (\Phi^T \Phi)^{-1} \Phi^T t} \longrightarrow (2)$$

(ii) From given question, q_n can also be viewed as. effective no. of observation; (x_n, t_n) . i.e. (x_n, t_n) can be treated as repeatedly occurring q_n times (replication of datapoints)

②

Bayes estimate.

$$\sum_{h_i \in H} P(F/h_i) P(h_i/D) = 0.4.$$

$$\sum_{h_i \in H} P(L/h_i) P(h_i/D) = 0.2 + 0.1 + 0.2 = 0.5$$

$$\sum_{h_i \in H} P(R/h_i) P(h_i/D) = 0.1.$$

[Thus Bayes optimal recommend the robot turn left]

MAP Hypothesis is defined as follows.

$$h_{MAP} = \operatorname{argmax}_{h_i \in H} P(h_i/D).$$

$$\Rightarrow \begin{aligned} & \text{---argmax} \\ P(h_1/D) &= 0.4 ; P(h_2/D) = 0.2 ; P(h_3/D) = 0.1 ; P(h_4/D) = 0.1 \\ P(h_5/D) &= 0.2. \end{aligned}$$

\Rightarrow Max values occur at hypothesis 1 which is h_1

$$\Rightarrow \text{in } P(h_1/D) \Rightarrow P(F/h_1)=1, P(L/h_1)=0, P(R/h_1)=0$$

\Rightarrow [The robot should go forward]

3. Given is one dimensional data $\in \mathbb{R}^2$. The parameters are $\{p, q\}$ where x is classified as 1 if $p < x < q$ then the VC dimensions of H is given by
- Iff (a) let us suppose that the training points are in sphere of radius R then let $g(x) = \text{sign}[f(x)] = \text{sign}[a + bx]$
- (b) So, the class of functions $\{g(x) | \|B\| \leq \Delta y\}$, it has VC dimension h satisfying $|h| \leq R^2 \Delta^2$
- Hence VC dimension of $H \leq R^2 \Delta^2$

4. Given D dimensional data $= (x_1, x_2, \dots, x_D)$ of linear model

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

5. N such data samples with labels $(x_i, t_i) = i=1, 2, \dots, N$

the sum of square function is given by

$$E_0(w) = \frac{1}{2} \sum_{i=1}^N (y(x_i, w) - t_i)^2$$

3. so rearranging:

$$E_0(w) = \frac{1}{2} \sum_{i=1}^N \left\{ \left[w_0 + \sum_{k=1}^D w_k (x_{ik} + \epsilon_{ik}) - t_i \right]^2 \right\}$$

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

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

4. Where we have used $y(x_i, w)$ to denote the output of the linear model when the input variable is x_i without noise added for the 2nd term equation we have.

$$E_E \left[\left(\sum_{i=1}^D w_i \epsilon_i \right)^2 \right] = E_E \left[\sum_{i=1}^D \sum_{j=1}^D w_i w_j \epsilon_i \epsilon_j \right]$$

$$\Rightarrow \sum_{i=1}^D \sum_{j=1}^D w_i w_j E_E [\epsilon_i \epsilon_j] = \sigma^2 \sum_{i=1}^D \sum_{j=1}^D w_i w_j \delta_{ij}$$

which gives

$$E_E \left[\left(\sum_{i=1}^D w_i \epsilon_i \right)^2 \right] = \sigma^2 \sum_{i=1}^D w_i^2$$

5. for the third term we can obtain:-

$$E_{\epsilon} \left[2 \left(\sum_{i=1}^D w_i \epsilon_i \right) (y(x_i, w) - t_i) \right] = 2(y(x_i, w) - t_i) E_{\epsilon} \left(\sum_{i=1}^D w_i \epsilon_i \right) \\ = 2(y(x_i, w) - t_i) \sum_{i=1}^D E_{\epsilon} [w_i \epsilon_i] = 0$$

Therefore if we calculate the expectation of $E_D(w)$ with respect to ϵ we can obtain;

$$E_{\epsilon} [E_D(w)] = \frac{1}{2} \sum_{i=1}^N (y(x_i, w) - t_i)^2 + \frac{\sigma^2}{2} \sum_{i=1}^D w_i^2$$