

1) Given,
 data points $(x_n, t_n) \quad n = 1 \dots N$
 weighing factor $g_n > 0$
 error function

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

$\phi(\cdot)$ is any representation of data.

a) we have to find an expression for the solution w^* that minimizes the above error function -

~~we~~ we can take gradient and set to zero.

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

Now,

let solve for w :-

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

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

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

b) let us assume a linear model of the output
 $y_i = w^T x_i + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ is the noise.

i> For data dependent noise variance, the above objective function (E_D) can be derived by minimizing the negative log-likelihood of the output if we set $\sigma^2 = \frac{1}{g_i}$ or,

$$= \frac{1}{2g_i}$$

ii> For replicated data points, the above objective function (E_D) can be derived if we create g_i copies of the i^{th} data point.

2 Bayes Optimal Estimate

$$y = \{F, L, R\}$$

$$= \arg \max_{h_i} \sum_{F, L, R} P(y_j | h_i) P(h_i | D)$$

$$= \arg \max \left(\sum (0.4 * 1), \sum (0.2 * 1 + 0.1 * 1 + 0.2 * 1), \sum (0.1 * 1) \right)$$

$$= \arg \max (0.4, 0.5, 0.1)$$

$$= \arg (0.5)$$

$$= L$$

So according to Bayes Optimal Estimator for a new data instance, the most probable prediction is Left (L) move by the robot.

MAP estimate

$$= \arg \max_{\theta} P(\theta | D)$$

$$= \arg \max (P(h_1 | D), P(h_2 | D), P(h_3 | D), \dots, P(h_5 | D))$$

$$= \arg \max (0.4, 0.2, 0.1, 0.1, 0.2)$$

$$= \arg (0.4)$$

$$= h_1$$

According to MAP estimator, h_1 ~~hypothesis~~ is the most probable hypothesis that describes the training dataset.

Hence, MAP estimate \neq Bayes Optimal Estimate.

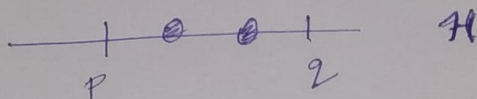
[3] \mathbb{R}^1 - 1D data

we can represent the data in a number line.

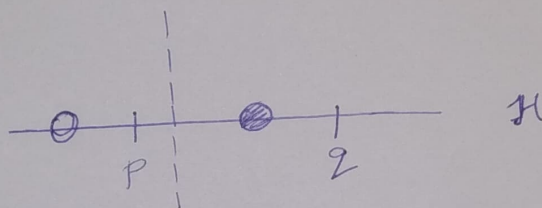
let \bullet denote class 1 when $P < x < Q$
 \circ denote class 0 when $x \leq P$ or, $x \geq Q$.

~~let's~~ let's take two points and place them such that they are correctly classified.

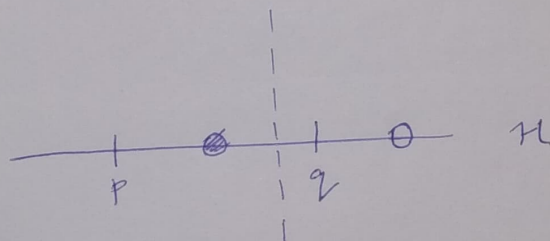
case 1: Both points are between P and Q



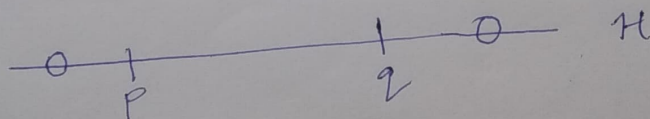
case 2: One point $\leq P$ and the other between P & Q



case 3: one point $\geq Q$ and other between P & Q

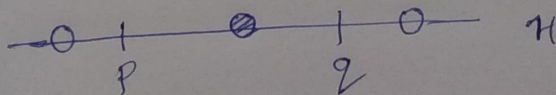


case 4: One point $\leq P$ and the other $\geq Q$



Hence, $\boxed{VC(H) = 2}$

Also, for three points we can prove that the classifier cannot separate all points. For ex -
 (for atleast one labelling order)



$VC(H) \neq 3$
 proved.

Noisy data

$$\hat{x}_i = x_i + \varepsilon_i$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\hat{y}(x_i, w) = y(x_i, w) + \sum_{k=1}^D w_k \varepsilon_k$$

$$\hat{y}(x_i, w) = y(x_i, w) + w^T \varepsilon$$

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

$$\frac{\partial}{\partial w} \hat{E}(w) = \sum_{i=1}^N (\hat{y}(x_i, w) - t_i) (x + \varepsilon)$$

Setting $\frac{\partial}{\partial w} \hat{E}(w)$ to 0 for minimizing,

$$\sum_{i=1}^N (\hat{y}(x_i, w) - t_i) (x + \varepsilon) = 0 \quad \text{--- (1)}$$

Noise-free data

Similar to above, we can directly write it in derivative form —

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

$$\sum_{i=1}^N (y(x_i, w) - t_i) (x) = 0 \quad \text{--- (2)}$$

Let us subtract (2) from (1) to find a relation —

$$\sum_{i=1}^N \left[(\hat{y}(x_i, w) - t_i) (x + \varepsilon) - (y(x_i, w) - t_i) (x) \right] = 0$$

$$\Rightarrow \sum_{i=1}^N \left[(y(x_i, w) + w^T \varepsilon - t_i) (x + \varepsilon) - (y(x_i, w) - t_i) (x) \right] = 0$$

$$\Rightarrow \sum_{i=1}^N \left[x w^T \varepsilon + \underbrace{y(x_i, w) \varepsilon + w^T \varepsilon^2 - \varepsilon t_i}_{l_2 \text{ norm}} \right] = 0$$

$$\Rightarrow \sum_{i=1}^N (y(x_i, w) - t_i) \varepsilon = -x w^T \varepsilon - w^T \varepsilon^2$$

$$\Rightarrow \sum_{i=1}^N (y(x_i, w) - t_i) = -(x w^T + w^T \varepsilon)$$

$$= -w^T (x + \varepsilon)$$

or,

$$= -\sum_{k=1}^D w_k (x_k + \varepsilon_k)$$

— (3)

Putting this result into

(1), we get -

$$\sum_{i=1}^N (\hat{y}(x_i, w) - t_i) (x + \varepsilon) = 0$$

or,

$$\sum_{i=1}^N (y(x_i, w) + w^T \varepsilon - t_i) (x + \varepsilon) = 0$$

$$\text{on } \cancel{\sum_{i=1}^N} (w^T \varepsilon - w^T (x + \varepsilon)) (x + \varepsilon) = 0$$

$$\text{on } \cancel{\sum_{i=1}^N} (x + \varepsilon) w^T (-x) = 0$$

$$\text{on } \cancel{\sum_{i=1}^N} x(x + \varepsilon) w^T = 0$$

Similarly in (2), we get

$$-w^T (x + \varepsilon) (x) = 0$$

$$\text{or, } -x(x + \varepsilon) w^T = 0$$

Hence, both (1) and (2) are equal, on substituting the result we found in (3).