

Assignment 4

$$1. E_D(\underline{\omega}) = \frac{1}{2} \sum_{n=1}^N g_n (t_n - \underline{\omega}^T \phi(\underline{x}_n))^2$$

Consider $G_i = \begin{bmatrix} \sqrt{g_1} \\ \sqrt{g_2} \\ \vdots \\ \sqrt{g_n} \end{bmatrix}$ (possible $\because g_i > 0, t_i$)

$$Y = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}, X = \begin{bmatrix} (\phi(\underline{x}_1))^T \\ (\phi(\underline{x}_2))^T \\ \vdots \\ (\phi(\underline{x}_n))^T \end{bmatrix}$$

a) We can rewrite $E_D(\underline{\omega})$ as

$$E_D(\underline{\omega}) = \frac{1}{2} \| G^T (Y - X \underline{\omega}) \|^2$$

It's convex, so we can differentiate w.r.t $\underline{\omega}$, to get minima

$$\frac{\partial E_D(\underline{\omega})}{\partial \underline{\omega}} = \frac{\partial}{\partial \underline{\omega}} \left[\frac{1}{2} \| G^T (Y - X \underline{\omega}) \|^2 \right]$$

$$\|G^T(Y - X\omega)\|^2 = [G^T(Y - X\omega)]^T [G^T(Y - X\omega)]$$

$$= (Y - X\omega)^T G G^T (Y - X\omega)$$

$$= Y^T G G^T Y - 2\omega^T X^T G G^T Y + \omega^T X^T G G^T X \omega$$

$$\therefore \frac{\partial E_D(\omega)}{\partial \omega} = \frac{\partial}{\partial \omega} \left[\frac{1}{2} (Y^T G G^T Y - 2\omega^T X^T G G^T Y + \omega^T X^T G G^T X \omega) \right]$$

~~$$\frac{\partial}{\partial \omega} \left[\frac{1}{2} (Y^T G G^T Y - 2\omega^T X^T G G^T Y + \omega^T X^T G G^T X \omega) \right]$$~~

$$\frac{\partial E_D(\omega)}{\partial \omega} = 0 - X^T G G^T Y + X^T G G^T X \omega$$

$$\text{So, @ minima } \therefore \frac{\partial E_D(\omega)}{\partial \omega} \geq 0$$

$$\therefore X^T G G^T X \omega^* = X^T G G^T Y$$

$$\therefore \underline{\omega^*} = (X^T G G^T X)^+ X^T G G^T Y$$

b) (i) data-dependent noise variance.

$$\tilde{t}_i \sim \omega \cdot x_i + N(0, \sigma_{\tilde{t}_i}^2)$$

$$\sim N(\omega \cdot x_i, \sigma_{\tilde{t}_i}^2)$$

Maximising log-likelihood $\left(\begin{array}{l} = \text{Minimising} \\ \text{negative} \\ \text{log-likelihood} \end{array} \right)$

$$= \underset{\theta}{\operatorname{argmax}} \log(L(\theta | x))$$

$$= \underset{\mu, \sigma^2}{\operatorname{argmax}} \log \left(\prod_i P_{\mu, \sigma^2}(x_i) \right)$$

$$= \underset{\mu, \sigma^2}{\operatorname{argmax}} \sum_{i=1}^n \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$g_i(\mu, \sigma^2)$$

$$\text{Maximising } \mu \Rightarrow \mu = \frac{1}{n} \sum x_i$$

$$\text{Maximising } \sigma^2 \Rightarrow \sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

$$\begin{aligned}\therefore \text{Log } L(w/s) &= \sum_i \log P(t_i^*/w) \\ &= \sum_{i=1}^n (w x_i - t_i^*)^2 \\ &\quad \text{(if } \sigma_i = \text{fixed)} \end{aligned}$$

By setting $\frac{1}{\sigma_i^2} = g_i$, we get the given expression.

Date _____
Page _____

(ii) repeated data points in list.

If our data is repeated, we can simply add a term g_i , to the cost, if x_i 's will be unique now.

g_i corresponds to frequency of occurrence of every x_i .

2) Bayes Optimal Classification:

$$P(F) = \sum P(F/h_i) P(h_i/D)$$
$$= 1(0.4) + 0(0.2) + 0(0.1) + 0(0.1) + 0(0.2)$$

$$P(L) = \sum P(L/h_i) P(h_i/D)$$

$$= 0(0.4) + 1(0.2) + 0(0.1) + 1(0.1) + 1(0.2)$$
$$= 0.2 + 0.1 + 0.2 = 0.5$$

$$P(R) = \sum P(R/h_i) P(h_i/D)$$

$$= 0(0.4) + 0(0.2) + 1(0.1) + 0(0.1) + 0(0.2)$$

Bayes optimal classifier recommends

to turn left. (as $P(L)$ is maximum).

MAP estimate.

$$\underset{h \in H}{\text{MAP}} = \underset{h \in H}{\operatorname{argmax}} P(h/D) \quad (\text{definition})$$

Clearly, from definition, & data given

MAP recommends to go forward.

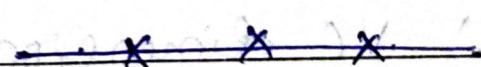
($\because P(h_i/D)$ is max for $h_1 \in H$, $P(F/h_1) = 1$)

Hence, they are NOT SAME.

3) VC-dimension: Model class \mathcal{F} has VC dimension d iff d is the largest set of points $x_1, \dots, x_d \in X$ such that for all labellings of x_1, \dots, x_d, \exists some $f \in \mathcal{F}$ that achieves that labelling.

Given, P -dimensional data.

Model class \mathcal{H} , characterised by $\{P, q\}$

Consider 
Case of 3 points.

label 1: $\{p \in \mathcal{C} : q \in \mathcal{C}\}$

label 0: else.

Clearly, for the following configuration

$$\begin{array}{ccc} 1 & 0 & 1 \\ \hline x_1 & x_2 & x_3 \end{array}$$

this class of models H_1 , can't achieve the labellings.

$\because x_1$ labelled as 1 $\Rightarrow p \subset x_1 \subset q$

$\therefore x_3$ " " $\Rightarrow p \subset x_3 \subset q$

Also $x_1 \subset x_2 \subset x_3$

$\Rightarrow p \subset x_2 \subset q$

Contradiction, as x_2 is labelled

0. i.e., $x_2 \not\subset q$ (Or) $x_2 \not\subset p$.

\therefore VC dimension is ~~not~~ ~~2~~ 3 or more.

for 2 points.

$$\begin{array}{cccc} + & 1 & 0 & + \\ \hline x_1 \cap x_2 & x_1 \cup x_2 \cap & p \subset x_1 \cap x_2 & x_1 \cap x_2 \end{array}$$

for all combinations, H can achieve the labellings,

\therefore VC dimension = 2

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

$$E(\underline{w}) = \frac{1}{2} \sum_{i=1}^n (y(\underline{x}_i, \underline{w}) - t_i)^2$$

If Gaussian noise $\epsilon_k \sim N(0, \sigma^2)$ is added to each i/p variable x_k , then

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

$$\text{i.e., } y'(\underline{x}_i, \underline{w}) = w_0 + \sum_{k=1}^D w_k (x_{ik} + \epsilon_{ik}),$$

$$= y(\underline{x}_i, \underline{w}) + \sum_{k=1}^D w_k \epsilon_{ik}$$

Also,

$$E[y'(\underline{w})] = \frac{1}{2} \sum_{i=1}^n \left\{ (y(\underline{x}_i, \underline{w}) - t_i)^2 + 2(y(\underline{x}_i, \underline{w}) - t_i) \left(\sum_{k=1}^D w_k \epsilon_{ik} \right) + \left(\sum_{k=1}^D w_k \epsilon_{ik} \right)^2 \right\}$$

Applying expectation and then

$$E[E[y'(\underline{w})]] = \frac{1}{2} \sum_{i=1}^n \left\{ (y(\underline{x}_i, \underline{w}) - t_i)^2 + 2(y(\underline{x}_i, \underline{w}) - t_i) \left(\sum_{k=1}^D w_k E[\epsilon_{ik}] \right) + E\left[\left(\sum_{k=1}^D w_k \epsilon_{ik} \right)^2\right] \right\}$$

→

Assuming

$$E[\varepsilon_i] = 0 \quad E[\varepsilon_i \varepsilon_{i'}] = \begin{cases} \sigma^2, & i \neq i' \\ 0, & i = i' \end{cases}$$

we get,

$$E[(\sum w_k \varepsilon_{ik})^2] = E\left[\sum_{k=1}^D \sum_{k'=1}^D w_k w_{k'} \varepsilon_{ik} \varepsilon_{ik'}\right]$$

$$= \sum_{k=1}^D \sum_{k'=1}^D w_k w_{k'}' E [\varepsilon_{ik} \varepsilon_{ik'}]$$

$$= \sum_{k=1}^D w_k^2 \sigma^2$$

$$= \sigma^2 \sum_{k=1}^D w_k^2$$

Hence,

$$E[\hat{E}(\underline{w})] = E(\underline{w}) + \frac{N}{2} \sum_{k=1}^D w_k^2$$

∴ we could get L₂ regularization term, without bias parameter w_0 .