

FOML

Assignment - 5

AIZ BTECH 11011

Gunjit Mittal

5.1 Q1 Loss function  $L(h)$  is given by

$$L(h) = |\{i \in \{1, \dots, m\} : h(x_i) \neq y_i\}|$$

ERM problem for domain  $\mathcal{Z}$  and input sample  $S \in \mathcal{Z}^m$  is given by

$$L_S(h) = \frac{1}{|S|} \sum_{p \in S} L(h, p)$$

The sample complexity of learning a finite class is upper bounded by  $m_H(\epsilon, \delta) = c \log(C|H|/\delta)/\epsilon$

where  $c=1$  in realizable case and  $c=2$  in the non realizable one

performing exhaustive search the run time is  $K|H|c \log(C|H|/\delta)/\epsilon$

Let us consider the example in which  $H_1, H_2, \dots, H_K$  and  $|H_i| = 2^i \forall i \in K$ . Learning  $H_K$  in Agnostic-loc model provides the following bound for an ERM hypothesis  $h$ :

$$L_D(h) = \min_{h \in H_K} L_D(h) + \sqrt{\frac{2(KH + \log(1/\delta))}{n}}$$

Using model selection, assuming  $j$  is the minimal index which contains  $h^* \in \arg\min_{h \in H} L_D(h)$  and fixing  $\epsilon \in (0, K)$

By Hoeffding's inequality we have

$$|L_D(\hat{h}_n) - L_V(\hat{h}_n)| \leq \sqrt{\frac{1}{2\alpha m} \log\left(\frac{4}{\delta}\right)}$$

is true with probability  $\geq 1 - \frac{\delta}{2k}$

Applying the union bound we have  
with probability  $\geq 1 - \frac{\delta}{2}$

$$\begin{aligned} L_D(\hat{h}) &\leq L_V(\hat{h}) + \sqrt{\frac{1}{2\alpha m} \log\left(\frac{4k}{\delta}\right)} \\ &\leq L_V(\hat{h}_n) + \sqrt{\frac{1}{2\alpha m} \log\left(\frac{4k}{\delta}\right)} \\ &\leq L_D(\hat{h}_n) + 2\sqrt{\frac{1}{2\alpha m} \log\left(\frac{4k}{\delta}\right)} \\ &\leq L_D(\hat{h}_n) + \sqrt{\frac{2}{\alpha m} \log\left(\frac{4k}{\delta}\right)} \end{aligned}$$

for  $n > j$  we have

$$L_D(\hat{h}) \leq L_D(\hat{h}_j) + \sqrt{\frac{2}{\alpha m} \log\left(\frac{4k}{\delta}\right)}$$

Similarly we have with probability  $> 1 - \frac{\delta}{2}$

$$L_D(\hat{h}_j) \leq L_D(h^*) + \sqrt{\frac{2}{(1-\alpha)m} \log\left(\frac{4|H_j|}{\delta}\right)}$$

Combining last two we have with  
probability  $> 1 - \delta$

$$L_D(\hat{h}) \leq L_D(h^*) + \sqrt{\frac{2}{\alpha m} \log\left(\frac{4k}{\delta}\right)} + \sqrt{\frac{2}{(1-\alpha)m} \log\left(\frac{4|H_j|}{\delta}\right)}$$



Substituting  $|H_j| = 2^j$  we get

$$L_D(\hat{h}) \leq L_D(h^*) + \sqrt{\frac{2}{\alpha m} \log\left(\frac{4K}{\delta}\right)} + \sqrt{\frac{3}{(1-\alpha)m} \left(j + \log\left(\frac{4}{\delta}\right)\right)}$$

From this we can observe that when the optimal index  $j$  is much smaller than  $K$  using model selection is the better approach.

5.3 Q1 Q6.3 According to the representer theorem the minimiser of training error lies in ~~span~~  $\text{span}\{\psi(x_1), \psi(x_2), \dots, \psi(x_m)\}$ .

$\therefore$  The ERM problem can be rewritten as

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i \psi(x_i) \right\|^2 + \frac{1}{2m} \sum_{i=1}^m \left( \left\langle \sum_{j=1}^m \alpha_j \psi(x_j), \psi(x_i) \right\rangle - y_i \right)^2$$

Using gram matrix we can write it as

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2} \alpha^T G \alpha + \frac{1}{2m} \sum_{i=1}^m \left( \langle \alpha, G_i \rangle - y_i \right)^2$$

Since it is convex we differentiate it and get

$$(\lambda' G + G G^T) \alpha - G y = 0 \quad (\lambda' = mI)$$

As  $G$  is symmetric

$$G (\lambda' I + G^T) \alpha = G y$$

If  $G$  is invertible

$$(\lambda' I + G^T) \alpha = y$$

Since  $G$  is positive semi-definite and  $\lambda > 0$   
 $\lambda' I + G$  is positive definite

$\therefore (A'I + G^T)$  is invertible

$\therefore$  the minimiser will be

$$\alpha = (A'I + G^T)^{-1} y$$

Q16.4 Let  $\Psi = \{1, 2, \dots, N\} \rightarrow \mathbb{R}^N$   
 $\Psi_j = (1^j, 0^{N-j})$

where  $1^j$  is the one vectors in  $\mathbb{R}^j$   
 and  $0^{N-j}$  is the zeroes vectors in  $\mathbb{R}^{N-j}$   
 then

$$\begin{aligned} \langle \Psi_i, \Psi_j \rangle &= \langle (1^i, 0^{N-i}), (1^j, 0^{N-j}) \rangle \\ &= \min\{i, j\} \\ &= K(i, j) \end{aligned}$$

Q2 Q6.1 The cauchy schwarz inequality for PDS kernels state that for a PDS kernel  $K$ , for any  $x, x' \in X$

$$K(x, x')^2 \leq K(x, x) K(x', x')$$

For the kernel  $K(x, y) = \frac{K(x, y)}{\alpha(x) \alpha(y)}$

~~Let  $G$  be kernel of~~

Let  $G$  be the gram matrix of the above kernel

considering any column matrix  $G$   
 with column vectors  $G_1, G_2, \dots, G_m \in \mathbb{R}^{m \times 1}$



$$C^T G C = \sum_{i,j=1}^m \frac{C_i C_j K(x_i, y_j)}{\alpha(x_i) \alpha(y_j)}$$

since both the numerators and denominators ~~evaluate~~ evaluate to  $> 0$

$$G^T G C > 0$$

which implies  $K'(x, y)$  is PDS

Q6.2 a)  $K(x, y) = \cos(x-y)$  over  $\mathbb{R} \times \mathbb{R}$   
 $\cos(x-y) = \cos x \cos y + \sin x \sin y$

$K(x, y)$  can be written as inner product of  
 $\phi(x) = \begin{bmatrix} \cos x \\ \sin x \end{bmatrix}$  and  $\phi(y) = \begin{bmatrix} \cos y \\ \sin y \end{bmatrix}$

b)  $K(x, y) = \cos(x^2 - y^2)$  over  $\mathbb{R} \times \mathbb{R}$

$$C^T G C = \sum_{i,j=1}^m C_i C_j \cos(x_i^2 - y_j^2)$$

taking  $x_i^2 = x_i'$  &  $y_j^2 = y_j'$

$$= \sum_{i,j=1}^m C_i C_j \cos(x_i' - y_j')$$

which is the same as proving  $\cos(x-y)$  is PDS

c)  $K(x, y) = \sum_{i=1}^N \cos^n(x_i^2 - y_i^2)$

Since each term in  $K(x, y)$  is PDS  
 $K(x, y)$  will also be PDS

e)  $K(x, x') = \cos L(x, x')$

$$C^T G C = \sum_{i,j=1}^N C_i C_j \cos(L(x_i, x_j))$$

which can be written as  $\cos(x_i - x_j)$  for some  $x_i \geq x_j$

so it becomes the same as a)

hence  $K(x, x') = \cos L(x, x')$  is PDS