

HW - 6

Tanoglu

Q1. Contd.

$N < k_0$ where k_0 is the smallest break point

$$\therefore m_H(N) = 2^N$$

$$\Rightarrow m_H(N)^2 = (2^N)^2 = 2^{2N}$$

when $N < \frac{k_0}{2} \Rightarrow$ i.e., $2N < k_0$

$$m_H(2N) = 2^{2N} = m_H(N)^2$$

when $N > \frac{k_0}{2} \Rightarrow$ i.e., $2N > k_0$

$$m_H(2N) < 2^{2N} = m_H(N)^2$$

$$\therefore \boxed{m_H(2N) \leq m_H(N)^2}$$

Generalization Bound :

$$\left\{ \begin{array}{l} m_H(N) \leq m_H\left(\frac{N}{2}\right)^2 \\ \text{or } m_H(N) \geq \sqrt{m_H(2N)} \end{array} \right. , \quad N \text{ is even}$$

$$m_H(N) \leq m_H\left(\frac{N+1}{2}\right)^2 , \quad N \text{ is odd}$$

Q1 Let $m_H(N) = k$

Now, if we partition any set of $2N$ points into two sets of N points each, each of these partitions will produce at best k dichotomies. Now, combining the two sets, the maximum number of dichotomies possible will be the cross product of the two sets of dichotomies (with N points each) i.e.,

$$m_H(2N) \leq k^2 = m_H(N)^2$$

$$\therefore m_H(2N) \leq m_H(N)^2$$

Combining this result with the VC generalization bound, we get,

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \frac{4m_H(2N)}{\delta}}$$

$$\therefore E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \frac{4m_H(N)^2}{\delta}}$$

(contd. on last page)

$$m_H(N) = N+1$$

Q2. (a) $N=100$

Then, VC generalization bound,

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{100} \cdot \ln \frac{4(2 \cdot 100 + 1)}{0.1}}$$

$$\therefore E_{\text{out}}(g) \leq E_{\text{in}}(g) + 0.8481596$$

with probability at least 90%

$$N = 10000$$

Then, VC generalization ~~without~~, bound,

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{10000} \cdot \ln \frac{4(2 \cdot 10000 + 1)}{0.1}}$$

$$\therefore E_{\text{out}}(g) \leq E_{\text{in}}(g) + 0.1042782$$

with probability at least 90%.

$$(b) N = 100,$$

$$d_{VC} = 1$$

Then, VC generalization bound,

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4((2N)^{d_{VC}} + 1)}{8} \right)}$$

$$\leq E_{\text{in}}(g) + \sqrt{\frac{8}{100} \ln \left(\frac{4 \cdot (200 + 1)}{0.1} \right)}$$

$$\leq E_{\text{in}}(g) + \sqrt{\frac{8}{100} \ln \frac{4 \cdot (200 + 1)}{0.1}}$$

$$\leq E_{\text{in}}(g) + 0.8487596$$

with probability at least 90%.

This is the same bound as previously seen in (a).

$$N = 10000$$

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{100} \ln \left(\frac{4 \cdot (2 \cdot 10000 + 1)}{0.1} \right)}$$

$$\leq E_{\text{in}}(g) + 0.1042782$$

with probability at least 90%.

This is the same bound as previously seen in (a).

$$(c) d_{VC} = 10$$

$$\delta = 0.05$$

$$\epsilon = 0.05$$

Then, $N \geq \frac{8}{(0.05)^2} \ln \left(\frac{4 \left[(2N)^{10} + 1 \right]}{0.05} \right)$

Let $N = 1000$ in RHS, then,

$$\begin{aligned} N &\geq \frac{8}{(0.05)^2} \ln \left(\frac{4 \left[(2 \cdot 1000)^{10} + 1 \right]}{0.05} \right) \\ &\approx 2.57251 \times 10^5 \end{aligned}$$

Substituting $N = 2.57251 \times 10^5$ in RHS, and continuing this process till convergence, we get,

$$N \approx 4.52957 \times 10^5$$

Q3.

(a) Let $d_{VC}(H) = d$, then, we have, $m_H(d) = 2^d$

so,

$$\begin{aligned} 2^d &= m_H(d) = \max_{x_1, \dots, x_d} |H(x_1, \dots, x_d)| \\ &= \max_{x_1, \dots, x_d} |\{f(h(x_1), \dots, h(x_d)) : h \in H\}| \\ &\leq |H| = M \end{aligned}$$

So, we can say that $d \leq \log_2(M)$
so, basically,

if the VC dimension of a hypothesis set is d , then we need at least 2^d hypothesis to run the d points through in order to get them to shatter. So,

$$d \leq \log_2(M).$$

(b) At worst, $\bigcap_{k=1}^K \mathcal{H}_k = \{h\}$

In this case VC dimension = 0 as $m_H(N) = 1 \forall N$

So,

$$d_{VC}(\bigcap_{k=1}^K \mathcal{H}_k) \geq 0$$

Now, let's assume,

$$d_{VC}(\bigcap_{k=1}^K \mathcal{H}_k) > \min_{1 \leq k \leq K} d_{VC}(\mathcal{H}_k) = d$$

which means that

$\bigcap_{k=1}^K \mathcal{H}_k$ can shatter $d+1$ points.

Let these points be x_1, \dots, x_{d+1}

Then,

$$\begin{aligned} \{ -1, +1 \}^{d+1} &= \bigcap_{k=1}^K \mathcal{H}_k(x_1, \dots, x_{d+1}) \\ &= \{ (h(x_1), \dots, h(x_{d+1})) : h \in \bigcap_{k=1}^K \mathcal{H}_k \} \end{aligned}$$

$$\begin{aligned} &\subset \{ (h(x_1), \dots, h(x_{d+1})) : h \in \mathcal{H}_k \} \\ &= \mathcal{H}_k(x_1, \dots, x_{d+1}) \quad \forall k=1, \dots, K \end{aligned}$$

So,

$$\begin{aligned} 2^{d+1} &\leq |\{ (h(x_1), \dots, h(x_{d+1})) : h \in \mathcal{H}_k \}| \leq 2^{d+1} \\ \Rightarrow |\{ (h(x_1), \dots, h(x_{d+1})) : h \in \mathcal{H}_k \}| &= 2^{d+1} \quad \forall k=1, \dots, K \end{aligned}$$

so, any \mathcal{H}_k can shatter $d+1$ points

If we let $\min_{1 \leq k \leq K} d_{VC}(\mathcal{H}_k) = d_{VC}(\mathcal{H}_{k_0})$, we get,

$$d = d_{VC}(\mathcal{H}_{k_0}) \geq d+1 \text{ which is not possible}$$

So,

$$0 \leq d_{VC}(\bigcap_{k=1}^K \mathcal{H}_k) \leq \min_{1 \leq k \leq K} d_{VC}(\mathcal{H}_k)$$

(c) Let $d_{VC}(\mathcal{H}_K) = d_K \quad \forall k = 1, \dots, K$

so, \mathcal{H}_K shatters d_K points x_1, \dots, x_{d_K}

$$\therefore \mathcal{H}_K(x_1, \dots, x_{d_K}) = \{-1, +1\}^{d_K}$$
$$\quad \forall k = 1, \dots, K$$

so,

$$\{-1, +1\}^{d_K} = \{(h(x_1), \dots, h(x_{d_K})): h \in \mathcal{H}_K\}$$

$$\subset \{(h(x_1), \dots, h(x_{d_K})): h \in \bigcup_{k=1}^K \mathcal{H}_k\}$$

So, $2^{d_K} \leq |\{(h(x_1), \dots, h(x_{d_K})): h \in \bigcup_{R=1}^K \mathcal{H}_k\}| \leq 2^{d_K}$

so, we can say,

$$\Rightarrow |\{(h(x_1), \dots, h(x_{d_K})): h \in \bigcup_{R=1}^K \mathcal{H}_k\}| = 2^{d_K}$$

~~This is similar to proof done~~

$$\therefore m \bigcup_{k=1}^K \mathcal{H}_k(d_k) = 2^{d_K} \quad \forall k = 1, \dots, K$$

$$\Rightarrow d_{VC}\left(\bigcup_{k=1}^K \mathcal{H}_k\right) \geq d_k \quad \forall k$$

$$\Rightarrow d_{VC}\left(\bigcup_{k=1}^K \mathcal{H}_k\right) \geq \max_{1 \leq k \leq K} d_k = \max_{1 \leq k \leq K} d_{VC}(\mathcal{H}_k)$$

• For upper bound, by induction proof:

Let $K = 2$ then, $d_{VC}(\mathcal{H}_1) = d_1$
 $d_{VC}(\mathcal{H}_2) = d_2$

So, the number of dichotomies generated by $H_1 \cup H_2$ is the sum of the dichotomies generated by H_1 and H_2 . So,

$$\begin{aligned}
 m_{H_1 \cup H_2}(N) &\leq m_{H_1}(N) + m_{H_2}(N) \\
 &\leq \sum_{i=0}^{d_1} \binom{N}{i} + \sum_{i=0}^{d_2} \binom{N}{i} \\
 &\leq \sum_{i=0}^{d_1} \binom{N}{i} + \sum_{i=0}^{d_2} \binom{N}{N-i} \\
 &\leq \sum_{i=0}^{d_1} \binom{N}{i} + \sum_{i=N-d_2}^N \binom{N}{i} \\
 &\leq \sum_{i=0}^{d_1} \binom{N}{i} + \sum_{i=d_1+1}^{N-d_2-1} \binom{N}{i} + \sum_{i=N-d_2}^N \binom{N}{i} \\
 &= \sum_{i=0}^N \binom{N}{i} \\
 &= 2^N
 \end{aligned}$$

$\forall N$

$$\begin{aligned}
 \text{s.t. } d_1 + 1 &\leq N - d_2 - 1 \\
 \Rightarrow N &\geq d_1 + d_2 + 2
 \end{aligned}$$

So, we can say,

$$\begin{aligned}
 d_{VC}(H_1 \cup H_2) &\leq d_1 + d_2 + 1 \\
 &\leq 1 + \sum_{k=1}^2 d_{VC}(H_k)
 \end{aligned}$$

So, by induction, if this is correct for $k-1$, Then,

$$d_{VC}\left(\bigcup_{k=1}^K H_k\right) = d_{VC}\left(\left(\bigcup_{k=1}^{K-1} H_k\right) \cup H_K\right)$$

$$\begin{aligned}
 d_{VC} \left(\bigcup_{k=1}^K \mathcal{H}_k \right) &\leq 1 + d_{VC} \left(\bigcup_{k=1}^{K-1} \mathcal{H}_k \right) + d_{VC} (\mathcal{H}_K) \\
 &\leq 1 + (K-2) + \sum_{k=1}^{K-1} d_{VC} (\mathcal{H}_k) + d_{VC} (\mathcal{H}_K)
 \end{aligned}$$

$\therefore d_{VC} \left(\bigcup_{k=1}^K \mathcal{H}_k \right) \leq K-1 + \sum_{k=1}^K d_{VC} (\mathcal{H}_k)$

So,

$\max_{1 \leq k \leq K} d_{VC} (\mathcal{H}_k) = d_{VC} \left(\bigcup_{k=1}^K \mathcal{H}_k \right) \leq K-1 + \sum_{k=1}^K d_{VC} (\mathcal{H}_k)$

Q5. Tikhonov regularization constraint

$$w^\top \Gamma^\top \Gamma w \leq C$$

(a) To obtain: $\sum_{q=0}^Q w_q^2 \leq C$

In order to obtain $w^\top \Gamma^\top \Gamma w = \sum_{q=0}^Q w_q^2$

we should use $\Gamma = I$

where I is identity matrix.

(b) To obtain: $\left(\sum_{q=0}^Q w_q^2 \right)^2 \leq C$

In order to obtain $w^\top \Gamma^\top \Gamma w = \left(\sum_{q=0}^Q w_q^2 \right)^2$

we should use $\Gamma = [1 \ 1 \ 1 \ 1 \dots 1]$

i.e., Γ should be a row of ones.

This implies that

$$w^T \Gamma^T = \sum_{q=0}^Q w_q$$

and thus,

$$w^T \Gamma^T \Gamma w = \left(\sum_{q=0}^Q w_q \right)^2$$

$$\text{Q6. (a) } p(Y) = \prod_{j=1}^D \text{Ber}(Y_j | \pi_0)$$

$$= \pi_0^{\|Y\|_0} (1 - \pi_0)^{D - \|Y\|_0}$$

$$\text{Here, } \pi_0 = 0.2 = \frac{1}{5}$$

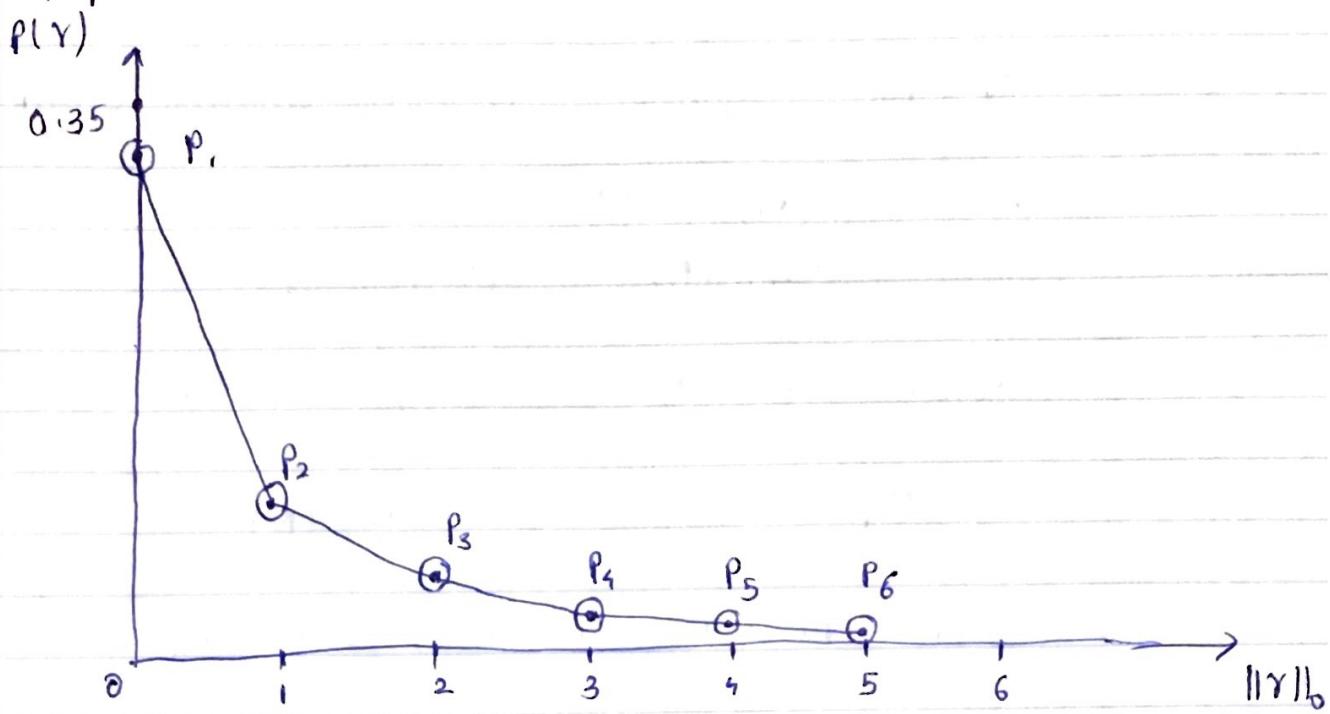
$$D = 5$$

$$\text{so, } p(Y) = \left(\frac{1}{5}\right)^{\|Y\|_0} \left(\frac{4}{5}\right)^{5 - \|Y\|_0}$$

Substituting different values, we get,

$\ Y\ _0$	$p(Y)$
0	0.32768 (P_1)
1	0.08192 (P_2)
2	0.02048 (P_3)
3	0.00512 (P_4)
4	0.00128 (P_5)
5	0.00032 (P_6)

So, plot :



- (b) Adding a constraint that the norm of gamma should be greater than 0 would prevent the feature selection in this case from eliminating all features (i.e., resulting in $\gamma_i = 0 \forall i$).