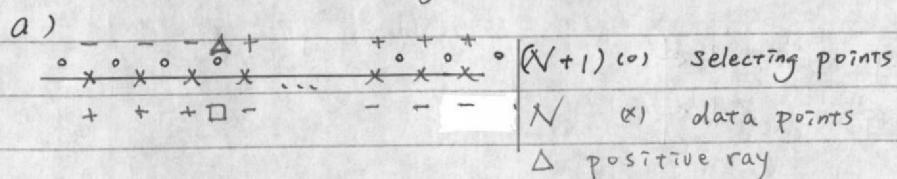


1. Problem 2.3 Compute the maximum number of dichotomies, $m_H(N)$, for these learning models, and consequently compute dvc , the VC dimension.

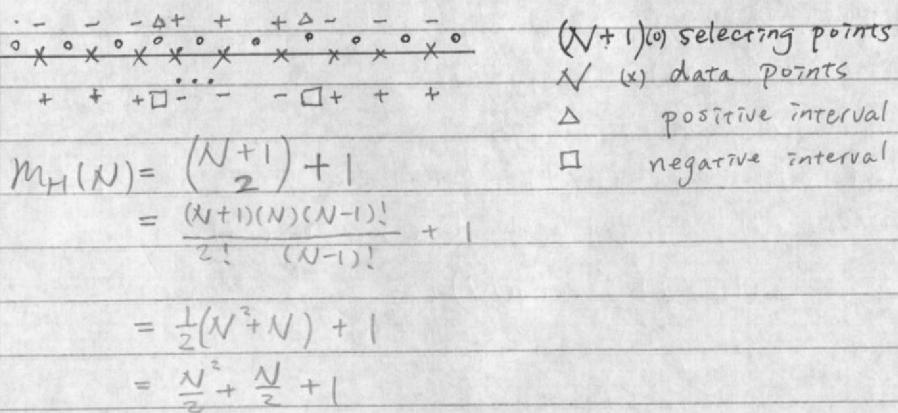
- (a) Positive or negative ray : H contains the functions which are $+1$ on $[a, \infty]$ (for some a) together with those that are $+1$ on $(-\infty, a]$ (for some a)



$$m_H(N) = N + 1 \quad \square \text{ negative ray}$$

break point, $k = 2 = dvc + 1$
 $dvc = 1$

- (b) Positive or negative interval : H contains the functions which are $+1$ on an interval $[a, b]$ and -1 elsewhere or -1 on an interval $[a, b]$ and $+1$ elsewhere.



break point $k = 3 = dvc + 1$

$dvc = 2$

- (c) Two concentric spheres in \mathbb{R}^d : H contains the functions which are $+1$ for $a \leq (x_1^2 + \dots + x_d^2)^{1/2} \leq b$



in this case, no data sets can be shattered by the spheres in \mathbb{R}^d

\Rightarrow break point $k = 3 = dvc + 1$

$dvc = 2$

$$m_H(N) = \binom{N+1}{2} + 1$$

$$= \frac{N^2}{2} + \frac{N}{2} + 1$$

$$\begin{array}{c} a+b \\ -a+b \\ -a-b \\ -b+a \end{array} \rightarrow \|X\|$$

$$= (x_1^2 + x_2^2 + \dots + x_d^2)^{1/2}$$

$$x^a x x x^b x x x \rightarrow$$

$$1 + N + \frac{N^{\frac{3}{2}} - 3N^2 + 2N}{6}$$

2. Problem 2.8 Which of the following are possible growth function $m_H(N)$ for $N \in \mathbb{N}$.

Some hypothesis set: $1 + N$; $1 + N + \frac{N(N-1)}{2}$; 2^N ; 2^{LN^2} ; $1 + N + \frac{N(N-1)(N-2)}{6}$

$1 + N$ is possible, e.g. $m_H(N)$ of positive ray.

$1 + N + \frac{N(N-1)}{2}$ is possible, e.g. $m_H(N)$ of positive interval

2^N is possible, e.g. $m_H(N)$ of convex set.

2^{LN^2} and 2^{LN^2} are not possible because they're not monotonically increasing function.

3. Problem 2.10 Show that $m_H(2N) \leq m_H(N)^2$, and hence obtain a generalization bound which only involves $m_H(N)$.

If we separate $2N$ to N and N ,

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

and if these two sets are independent

then $m_H(2N) = m_H(N)^2$

if in other dependent cases

$m_H(2N) < m_H(N)^2$ because there're some of possible dichotomies in $m_H(2N)$ intersect with those in $m_H(N)^2$.

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

(a) Let $H = \{h_1, h_2, \dots, h_M\}$ for a finite M . Prove that $dvc(H) \leq \log_2 M$

4. Problem 2.13
 (b) For hypothesis sets H_1, H_2, \dots, H_K with finite VC dimension $dvc(H_k)$, derive and prove the tightest upper and lower bound that you can get on $dvc(\bigcap_{k=1}^K H_k)$

(c) For hypothesis sets H_1, H_2, \dots, H_K with finite VC dimension $dvc(H_k)$, derive and prove the tightest upper and lower bounds that you can get on $dvc(\bigcup_{k=1}^K H_k)$

$$(a) M \leq m_H(N) \leq 2^K \quad dvc \rightarrow \infty$$

k : break point
 $m_H(N) \leq \sum_{i=0}^{dvc} \binom{N}{i}$

$$m_{H+1}(N) \leq (N) dvc + 1$$

$k-1$ points shattered

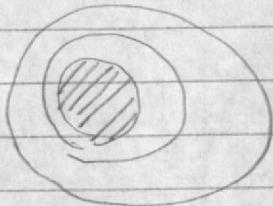
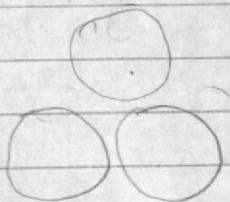
$$B(N_0+1, k) \leq B(N_0, k) + B(N_0, k-1)$$

$$2^{k-1} \leq m_H(N) \leq 2^k$$

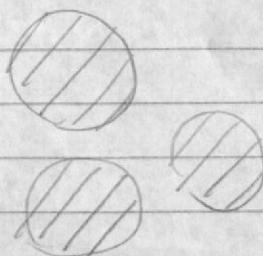
$$2^{dvc} \leq m_H(N) \leq M$$

$$\Rightarrow dvc(H) \leq \log_2 M.$$

$$(b) \emptyset \leq dvc(\bigcap_{k=1}^K H_k) \leq \min \{ dvc(H_k) \}_{k=1}^K$$



$$(c) \max \{ dvc(H_k) \}_{k=1}^K \leq dvc(\bigcup_{k=1}^K H_k) \leq \sum_{k=1}^K dvc(H_k)$$



5. Problem 2.22 When there is noise in data, $E_{out}(g^{(D)}) = \mathbb{E}_{x,y}[(g^{(D)}(x) - y(x))^2]$, where $y(x) = f(x) + \epsilon$, if ϵ is a zero-mean noise random variable with variance σ^2 , show that the bias-variance decomposition becomes $\mathbb{E}_D[E_{out}(g^{(D)})] = \sigma^2 + \text{bias} + \text{var}$

$$\begin{aligned}
\mathbb{E}_D[E_{out}(g^{(D)})] &= \mathbb{E}_D[\mathbb{E}_{x,y}[(g^{(D)}(x) - y(x))^2]] \\
&= \mathbb{E}_D[\mathbb{E}_{x,y}[(g^{(D)}(x) - (f(x) + \epsilon))^2]] \\
&= \mathbb{E}_{x,y}[\mathbb{E}_D[g^{(D)}(x)^2] - 2g^{(D)}(x)f(x) - 2g^{(D)}(x)\epsilon + f(x)^2 + 2f(x)\epsilon + \epsilon^2]] \\
&= \mathbb{E}_{x,y}[\mathbb{E}_D[g^{(D)}(x)^2] - \bar{g}(x)^2 \\
&\quad - 2\mathbb{E}_D[g^{(D)}(x)]\mathbb{E}[\epsilon] + 2f(x)\mathbb{E}_D[\epsilon] + \mathbb{E}_D[\epsilon^2]] \\
&= \mathbb{E}_{x,y}[\mathbb{E}_D[g^{(D)}(x)^2] - \bar{g}(x)^2 \\
&\quad + \bar{g}(x)^2 - 2\bar{g}(x)f(x) + f(x)^2 \\
&\quad + \sigma^2 + \text{bias} + \mathbb{E}_D[\epsilon^2]] \\
&= \text{Var} + \text{bias} + \sigma^2
\end{aligned}$$

6.

$X \in \mathbb{R}^1 : [-1, 1]$, data set: $\{x_1, x_2\}$, target function $f(x) = x^2$

Problem 2.24 $D = \{(x_1, x_1^2), (x_2, x_2^2)\}$. The learning algorithm returns the line fitting these two points as $g(\cdot)$ consists of functions of the form $h(x) = ax + b$. We are interested in the best performance (E_{out}) of our learning system with respect to the squared error measure, the bias and the var

(a) $\bar{g}(x)$ analytic expression

(b) Experiment that you could run to determine (numerically) $\bar{g}(x)$, E_{out} , bias, and var.

(c) Run your experiment and report the results. Compare E_{out} (expectation is with respect to data set) with bias+var.
Provide a plot of your $\bar{g}(x)$ and $f(x)$ (on the same plot)

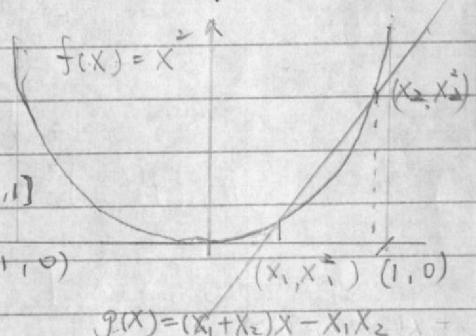
(d) Compute analytically what E_{out} , bias and var should be

$$(a) \bar{g}(x) = \bar{a}x + \bar{b}, \text{ where } \bar{a} = E[x_1 + x_2]$$

$$= E[x_1] + E[x_2] = 0 + 0 = 0, \therefore x_1, x_2 \in [-1, 1]$$

$$\bar{b} = E[x_1 x_2] = E[x_1] E[x_2] = 0 \cdot 0 = 0$$

$\because X_1 \& X_2 \text{ are independent} \Rightarrow \bar{g}(x) = 0$



(b) The experiment would be:

1 step: Generating K (input variable, e.g. 100) numbers of data set D_1, \dots, D_k , where $D = \{(x_1, x_1^2), (x_2, x_2^2)\}$ and x_1, x_2 are generated randomly $[-1, 1]$, in implementation D would be a $(K \times 2)$ matrix.

2 step: Compute $\bar{g}(x) = \bar{a}x + \bar{b}$, where $\bar{a} = \frac{1}{K} \sum_{i=1}^K (x_i + x_2)$ and function $g(x) = ax + b$

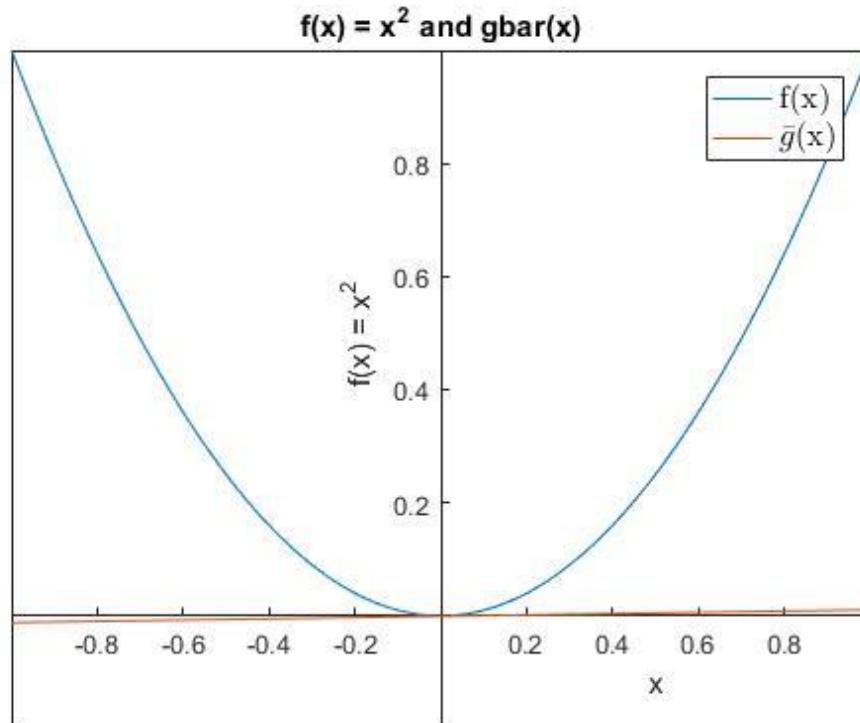
3 step: Compute bias = $\frac{1}{2} \int_{-1}^1 ((\bar{a}x + \bar{b}) - x^2)^2 dx$

$$\text{Compute Var} = E_x \left[\frac{1}{K} \sum_{i=1}^K (g(x_i^{D_i}) - \bar{g}(x_i^{D_i}))^2 \right] + \frac{1}{K} \sum_{i=1}^K (g(x_i^{D_i}) - \bar{g}(x_i^{D_i}))^2$$

$$E_x \left[\text{mean} \left[(x_i^{D_i} + x_2^{D_i}) x_i^{D_i} - (x_i^{D_i} x_2^{D_i}) - (\bar{a} x_i^{D_i} + \bar{b}) \right]^2 \right]; \\ \text{mean} \left[(x_i^{D_i} + x_2^{D_i}) x_2^{D_i} - (x_i^{D_i} x_2^{D_i}) - (\bar{a} x_2^{D_i} + \bar{b}) \right]^2 \right]$$

$$E_D[E_{out}] = \frac{1}{5} + \frac{1}{3} E_D[(x_1 + x_2)^2] - \frac{2}{3} E_D[x_1 x_2] + E_D[(x_1 x_2)^2]$$

Problem 2.24 (c)



	Analytical computation (from part(d))	Numerical computation (experiment with $k=10000$)
$E[E_{out}]$	$\frac{8}{15} \cong 0.5333$	0.5274
Bias + Var	$\frac{8}{15} \cong 0.5333$	0.53012
Bias	$\frac{1}{5} = 0.2$	0.2004
Var	$\frac{1}{3} \cong 0.3333$	0.3297
$\bar{g}(x) = \bar{a}x + \bar{b}$	0	$0.0113x - 0.0006$

Another way to compute Var is on the next page.

$$(d) \text{Var}(X_1) = E_D[X_1^2] - (E_D[X_1])^2 = \frac{1}{3} - 0 = \frac{1}{3}$$

$$\text{Var}(X_2) = E_D[X_2^2] - (E_D[X_2])^2 = \frac{1}{3} - 0 = \frac{1}{3}$$

$$E_D[X] = \int_{-1}^1 \frac{1}{2} X^2 dX = \frac{1}{3}$$

$$E_D[X] = 0$$

$$\text{Var} = \frac{1}{2} [\text{Var}(X_1) + \text{Var}(X_2)] = \frac{1}{3}$$

$$E_{\text{out}} = E_D[\int_{-1}^1 \frac{1}{2} [g(x) - f(x)]^2 dx]$$

$$= \int_{-1}^1 \frac{1}{2} [(ax+b) - x^2]^2 dx$$

$$= \int_{-1}^1 \frac{1}{2} [x^4 - 2ax^3 + (a^2 - 2b)x^2 + 2abx + b^2] dx$$

$$= \frac{1}{5} + \frac{a^2 - 2b}{3} + b^2$$

$$E_D[E_{\text{out}}] = \frac{1}{5} + \frac{1}{3} E_D[X_1 + X_2] - \frac{2}{3} E_D[X_1 X_2] + E_D[(X_1 X_2)^2]$$

$$E_D[X_1 X_2] = E_D[X_1] E_D[X_2] = 0 \cdot 0 = 0$$

$$E_D[X_1 + X_2] = E_D[X_1] + E_D[X_2] = 0 + 0 = 0$$

$$E_D[(X_1 + X_2)^2] = E_D[X_1^2] + 2E_D[X_1] E_D[X_2] + E_D[X_2^2] \\ = \frac{1}{3} + 2 \cdot 0 \cdot 0 + \frac{1}{3}$$

$$= \frac{2}{3}$$

$$\Rightarrow E_D[E_{\text{out}}] = \frac{1}{5} + \frac{1}{3} \cdot \frac{2}{3} - \frac{2}{3}(0) + \frac{1}{3} \cdot \frac{1}{3}$$

$$= \frac{8}{15}$$

$$\text{Bias} = \int_{-1}^1 \frac{1}{2} [\bar{g}(x) - f(x)]^2 dx$$

$$= \frac{1}{5} + \frac{\bar{a}^2 - 2\bar{b}}{3} + \bar{b}^2 \rightarrow \text{from result of (a), } \bar{a} = \bar{b} = 0$$

$$= \frac{1}{5}$$

~~Y~~

6. Problem 2.24(d)

$$\text{Var} = \mathbb{E}_x[\mathbb{E}_b[(g^b(x) - \bar{g}(x))^2]]$$

$$\begin{aligned}& \mathbb{E}_x[\mathbb{E}_b[(ax+b)^2]] \\& \mathbb{E}_x[\mathbb{E}_b[a^2x^2 + 2abx + b^2]] \\& \mathbb{E}_x[X^2\mathbb{E}_b[(x_1+x_2)^2] + x \cdot 2\mathbb{E}_b[x_1+x_2]\mathbb{E}_d[x_1x_2] + \mathbb{E}_d[x_1x_2^2]] \\& \mathbb{E}_x\left[\frac{2}{3}x^2 + 2x \cdot 0 \cdot 0 + \frac{1}{9}\right] \\& = \int_{-1/2}^{1/2} \frac{1}{2} \left[\frac{2}{3}x^2 + \frac{1}{9} \right] dx \\& = \frac{1}{2} \left[\frac{2}{9}x^3 + \frac{1}{9}x \right] \Big|_{-1/2}^{1/2} \\& = \frac{2}{9} + \frac{1}{9} \\& = \frac{1}{3}\end{aligned}$$

7. $H = X(X^T X)^{-1} X^T$, $X \in \mathbb{R}^{N \times d+1}$ is an N by $d+1$ matrix, $X^T X$ invertible

Exercise 3.3(a) Show that H is symmetric

$$\text{first}, X^T X \text{ is invertible} \Rightarrow [(X^T X)^{-1}]^T = [(X^T X)^T]^{-1} \\ = (X^T X)^{-1} \quad \text{using result (1)}$$

$$\begin{aligned} \text{Second, } H^T &= [X(X^T X)^{-1} X^T]^T \\ &= ((X^T X)^{-1} X^T)^T X^T \\ &= X[(X^T X)^{-1}]^T X^T \quad \text{using result (1)} \\ &= X(X^T X)^{-1} X^T \\ &= H \end{aligned}$$

$\Rightarrow H$ is symmetric

(b) Show that $H^k = H$ for any positive integer k

$$k=1, H^1 = H$$

$$\begin{aligned} k=2, H^2 &= H \cdot H = [X(X^T X)^{-1} X^T][X(X^T X)^{-1} X^T] \\ &= X(X^T X)^{-1}(X^T X)(X^T X)^{-1} X^T \\ &= X(X^T X)^{-1} X^T \\ &= H \end{aligned}$$

$$\begin{aligned} n \in \mathbb{Z}^+, k=n, H^n &= \underbrace{H \cdot H \cdots H}_n = \underbrace{(X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T)}_n \cdots \underbrace{(X(X^T X)^{-1} X^T)}_n \\ &= X \cdot \underbrace{I \cdot I \cdot \cdots \cdot I}_{n-1} \cdot (X^T X)^{-1} X^T, \text{ where } I \text{ is } (d+1) \text{ by } (d+1) \text{ identity matrix} \\ &= X(X^T X)^{-1} X^T \end{aligned}$$

Therefore, $H^k = H$, $k \in \mathbb{Z}^+$

(c) If I is the identity matrix of size N , show that $(I - H)^k = I - H$ for any positive integer k .

$$k=1, (I - H)^1 = (I - H)$$

$$\begin{aligned} k=2, (I - H)^2 &= (I - H)(I - H) = I^2 - IH - HI + H^2 \quad (\text{using result from (b)}) \\ &= I - H - H + H = (I - H) \end{aligned}$$

$$\begin{aligned} k=n, (I - H)^n &= \underbrace{(I - H)(I - H) \cdots (I - H)}_n \\ &= \underbrace{(I - H)(I - H) \cdots (I - H)}_{n-1} \\ &\vdots \\ &= (I - H) \end{aligned}$$