

problem 1.

(a) $E_{in}(h)$ 이 최소가 되는 곳에서, $\nabla E_{in}(h) = 0$ 이다.
 $\therefore \nabla E_{in}(h) = \sum_{n=1}^N 2(h - y_n) = 0, \sum_{n=1}^N h = \sum_{n=1}^N y_n$, 양변을 N 으로 나누면, $h_{mean} = \frac{1}{N} \sum_{n=1}^N y_n$ 이다.

(b) $E_{in}(h)$ 에 대해, $h > y_n$ 에서 $|h - y_n| = h - y_n$,
 $h < y_n$ 에서 $|h - y_n| = -h + y_n$ 이다. h 에 대해 미분하면,
 $\frac{dE_{in}(h)}{dh} = \begin{cases} 1, & \text{for } h > y_n \\ -1, & \text{for } h < y_n \end{cases}$ 이고, $E_{in}'(h) = 0$ 이기
 위해서는 $h < y_n$ 과 $h > y_n$ 의 개수가 같아야 한다.
 $\therefore \frac{dE_{in}}{dh} = 0$ 인 h 는 h_{med} 이다.

problem 2.

(a) $1 - y_n W^T x_n > 0$ 이면, $e_n(w) = 1 - y_n W^T x_n$.
 $= 1 - \sum_{i=1}^M w_i x_{ni}$ 에서, 상수 y_n, x_{ni} 에 대해
 미분 가능하다. $1 - y_n W^T x_n < 0$ 에서 $e_n(w) = 0$ 이므로
 연속이고 미분 가능하다.

$1 - y_n W^T x_n = 0$ 일 때, $\nabla e_n(w) = -y_n x_n$ 에 대해,
 미분 가능하기 위해 모든 방향에서의 미분 극한값이
 같아야 하므로, $-y_n x_n = 0$ 이어야 한다. 이 때, y_n 은
 $+1, -1$ 이므로, $x_n = 0$ 이고, 이 경우 $y_n W^T x_n \neq 1$ 이다.
 따라서 $y_n = W^T x_n$ 에서는 미분 가능하기 위한 조건이 존재하므로,
 이 점에서만 미분 불가능하다.

(b) $y_n = \text{sign}(W^T x_n)$ 이면, $[\text{sign}(W^T x_n) \neq y_n] = 0$ 이고, $y_n W^T x_n > 0$ 에 대해,
 $0 \leq e_n(w) = \max(0, 1 - y_n W^T x_n)$ 이므로 이 때는
 upper bound이다.
 $y_n \neq \text{sign}(W^T x_n)$ 이면, $[\text{sign}(W^T x_n) \neq y_n] = 1$ 이고,
 $y_n W^T x_n < 0$ 에 대해, $1 - y_n W^T x_n > 1$ 이므로
 이때의 $e_n(w) > 1$ 에 대해 두 경우 모두 upper
 bound임을 알 수 있다.
 $\therefore E_{in}(w) = \frac{1}{N} \sum_{n=1}^N [\text{sign}(W^T x_n) \neq y_n] \leq \frac{1}{N} \sum_{n=1}^N e_n(w)$
 즉 in sample error의 upper bound임을 알 수 있다.
 (모든 n 에 대해 $[\text{sign}(W^T x_n) \neq y_n] \leq e_n(w)$ 이므로).

(c) minimum $\frac{1}{N} \sum_{n=1}^N e_n(w)$ 를 찾기 위해 gradient
 descent method를 사용하라,
 $\nabla e_n(w) = \begin{cases} 0, & \text{for } y_n W^T x_n > 1 \\ -y_n x_n, & \text{for } y_n W^T x_n < 1 \end{cases}$
 $\therefore \frac{1}{N} \sum_{n=1}^N e_n(w)$ 의 gradient는 위 식을 n 에 대해 더하고 N 으로
 나눈 것과 같다. 이를 기반으로 $S(w)$ 를 짜면 다음과 같다.
 Random value $W(0)$ 에 대해,
 for t in $(0, T)$:
 $V_t = \nabla \left(\frac{1}{N} \sum_{n=1}^N e_n(w) \right) = \frac{1}{N} \sum_{n=1}^N \nabla e_n(w(t))$
 $= \frac{1}{N} \sum_{n=1}^N -y_n x_n$.
 (이때 $\langle S \rangle$ 는 $y_n W^T x_n < 1$ 인 집합 n).
 $W(t+1) = W(t) - \eta V_t$. (η 는 선택된 learning rate)
 return W .

problem 3. 뒤쪽에 프린트하여 첨부합니다.

problem 4.

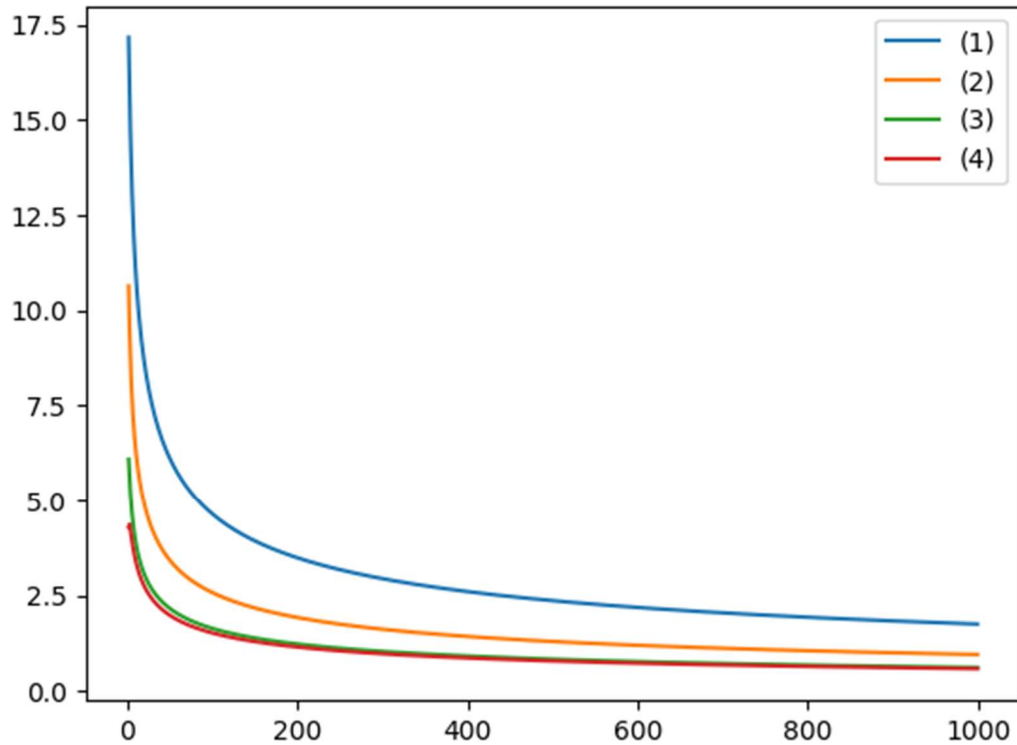
(a) D 에 대해 $g^{(0)}(x)$ 들의 mean을 구하는 것만
 족, 다음과 같이 나타낼 수 있다.

$$\bar{g}(x) = E_D[g^{(0)}(x)] = \frac{1}{K} \sum_{i=1}^K g^{(0)}(x).$$

$$\begin{aligned} (b) E_D[(g^{(0)}(x) - f(x))^2] &= E_D[(g^{(0)} - \bar{g} + \bar{g} - f)^2] \\ &= E_D[(g^{(0)} - \bar{g})^2 + (\bar{g} - f)^2 - 2(g^{(0)} - \bar{g})(\bar{g} - f)] \\ &\text{이때, } (g^{(0)} - \bar{g}) \text{의 경우 } D \text{와 independent하므로} \\ &E_D \text{를 취해도 같다. 또한 } E_D[g^{(0)} - \bar{g}] = E[g^{(0)}] - \bar{g} \\ &= \bar{g} - \bar{g} = 0 \text{이므로, 준식은 다음과 같다.} \\ &= E_D[(g^{(0)} - \bar{g})^2] + (\bar{g} - f)^2 = \text{Var}(x) + \text{bias}(x). \end{aligned}$$

$$\begin{aligned} (c) E_D[E_{out}(g^{(0)})] &= E_x[E_D[(g^{(0)}(x) - f(x))^2]] \\ &= E_x[\text{Var}(x) + \text{bias}(x)] = E_x[\text{Var}(x)] + E_x[\text{bias}(x)] \\ &= \text{bias} + \text{var} \quad (\text{b)에서 증명한 것을 사용}) \end{aligned}$$

Problem 3



위 그림은 python의 matplotlib module을 이용하여 각 error bound를 plot한 것이다. $d_{VC} = 50$, $\delta = 0.05$ 에 대해, perceptron을 가정하여 $m_H(2N) \leq (2N)^{50} + 1$ 을 이용하였다.

그림에서 확인할 수 있듯, Devroye bound가 가장 빠르게 수렴하는 것을 확인할 수 있다. 이때 Parrondo and van den Broeke와 비슷한 수준을 보인다.

problem 5.

(a). $Xw + b = \begin{bmatrix} -0.5 \\ -4.5 \\ 0.7 \end{bmatrix}$ 이므로,

$y_n(w^T x_n + b)$ 가 $n=1$ 에서 0.5 이다.

$\therefore \rho = 0.5$.

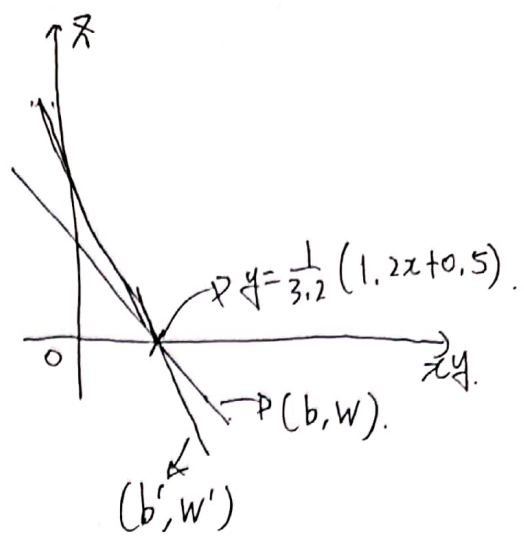
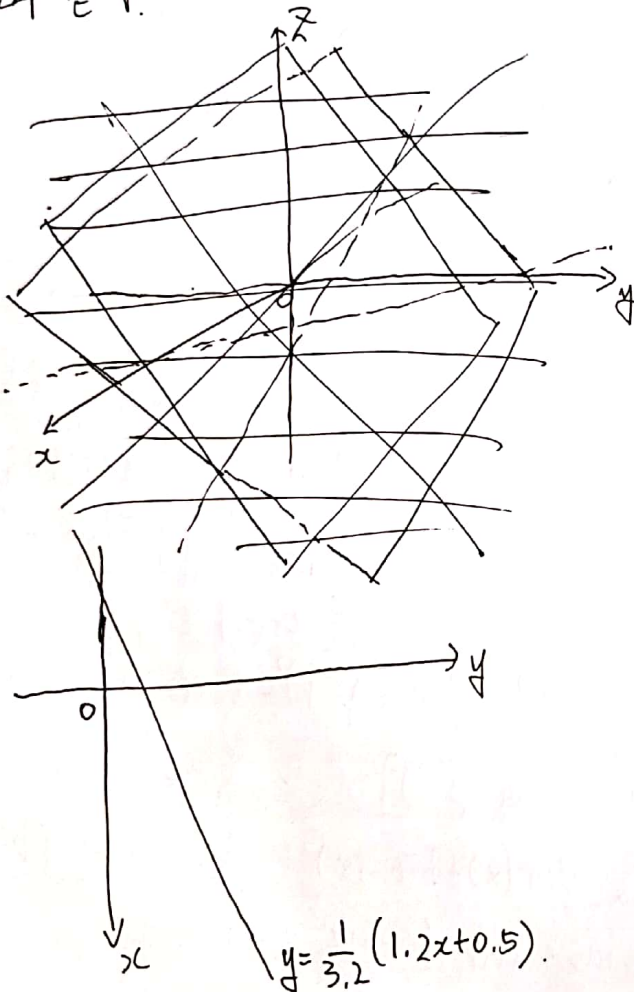
(b). ~~(b, w)~~ $(b', w') = \frac{1}{\rho} (b, w)$ 할 때,

$b' = -1, w' = \begin{bmatrix} 2.4 \\ -6.4 \end{bmatrix}$ 이고,

$Xw' + b' = \begin{bmatrix} -1 \\ -9 \\ 1.4 \end{bmatrix}$ 이다.

\therefore 이때 $\rho = \min y_n(w^T x_n + b) = 1$ 이다.

(c). ~~$y = w^T x + b$~~ $y = w^T x + b$ 에 대해
 (w, b) 와 (w', b') 에 대해 hyperplane은
 다음과 같다.



위 그림은 각각 $x-y$ 평면, plane을 뿔에서
 본 그림이다. 그림과 같이 x_n 에 대해 같은
 separator로 작용하며, $w^T x_n + b$ 값의 상수배가
 되는 것을 확인할 수 있다.