

005230642 Rui Xu

1.(a) A: if $\exists x_i \in \text{Strain} \neq 0 \quad i=1, 2, \dots, m$

$$A(\text{strain}) = h_z, z = x_i$$

else

$$A(\text{strain}) = h^-$$

(b). ① If $h^* = h^-$, the algorithm gives h^- .

$$\Rightarrow L(D, h^*)(A(\text{strain})) = \mathbb{E}_{x \sim D} (0-1\text{-Loss}(h^*(x), A(\text{strain})(x))) = P_{x \sim D}(h^*(x) \neq A(\text{strain})(x)) = 0$$

$$\Rightarrow P(L(D, h^*)(A(\text{strain})) > \varepsilon) = 0 \Rightarrow m > 0 \text{ is enough}$$

② If $h^* \neq h_z$:

If the algorithm gives h_z , $\Rightarrow h_z = h^*$,

$$\Rightarrow L(D, h^*)(A(\text{strain})) = 0$$

If the algorithm gives h^- , $\Rightarrow h^- \neq h^*$.

$$\begin{aligned} \Rightarrow L(D, h^*)(A(\text{strain})) &= \mathbb{E}_{x \sim D} (0-1\text{-Loss}(h^*(x), A(\text{strain})(x))) = P_{x \sim D}(h^*(x) \neq A(\text{strain})(x)) \\ &= \frac{1}{N}(0+0+\dots+0+1) = \frac{1}{N} > \varepsilon \end{aligned}$$

$\Rightarrow P(L(D, h^*)(A(\text{strain})) > \varepsilon) \Leftrightarrow$ The possibility that the algorithm gives h^-
($h^* = h_z$)

\Leftrightarrow The possibility that $z \notin \text{Strain}$

$$\Rightarrow P(L(D, h^*)(A(\text{strain})) > \varepsilon) = \left(1 - \frac{1}{N}\right)^m \leq \alpha$$

$$\Rightarrow \left(1 - \frac{1}{N}\right)^m \leq e^{-\frac{m}{N}} < e^{-m\varepsilon} \leq \alpha$$

$$\Rightarrow m > \frac{1}{\varepsilon} \log \frac{1}{\alpha}$$

\Rightarrow If training set Strain has a sample size larger than $\frac{\log(1/\alpha)}{\varepsilon}$, then
the probability that the generalization error of $A(\text{strain})$ larger than
 ε is at most α .

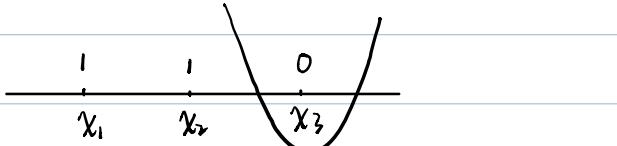
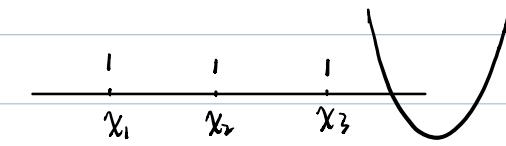
2. VC dimension of H is 3.

① Suppose we have three points x_1, x_2, x_3 and $x_1 < x_2 < x_3$. They can always be shattered by H , no matter how they are labeled.

if x_1, x_2, x_3 all labeled 1,

if x_1, x_2 labeled 1, x_3 labeled 0,

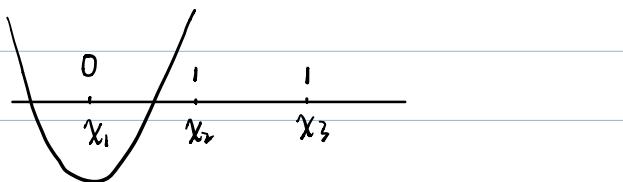
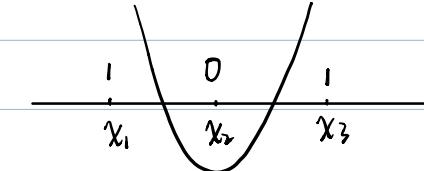
choose a, b, c to form the curve below: choose a, b, c to form the curve below:



if x_1, x_3 labeled 1, x_2 labeled 0

if x_1 labeled 0, x_2, x_3 labeled 1,

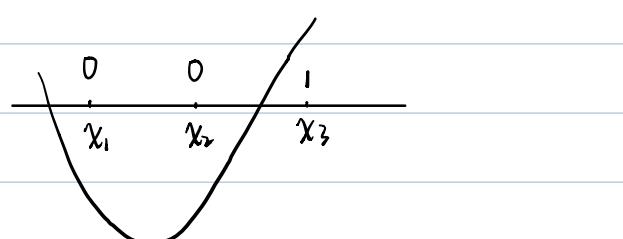
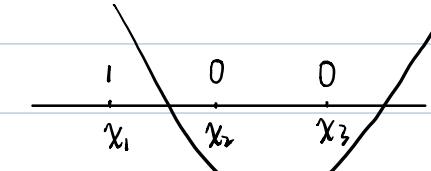
choose a, b, c to form the curve below: choose a, b, c to form the curve below:



if x_1 labeled 1, x_2, x_3 labeled 0,

if x_1, x_2 labeled 0, x_3 labeled 1,

choose a, b, c to form the curve below: choose a, b, c to form the curve below:

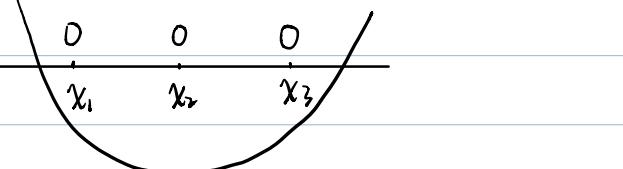
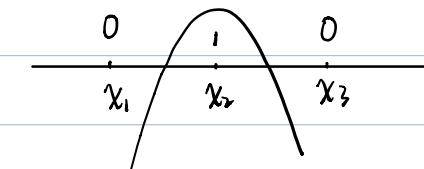


if x_1, x_3 labeled 0, x_2 labeled 1,

if x_1, x_2, x_3 all labeled 0,

choose a, b, c to form the curve below:

choose a, b, c to form the curve below:



\Rightarrow Three points can be scattered.

$$\Rightarrow VC \geq 3$$

② If we have four points x_1, x_2, x_3, x_4 and $x_1 < x_2 < x_3 < x_4$

and if they are labeled 1, 0, 1, 0

<> if at least two of x_1, x_2, x_3, x_4 are equal

\Rightarrow at least two equal points have different labels

\Rightarrow It is obvious that these four points cannot be shattered since same points give same $\text{sgn}(ax^2+bx+c)$

<> if $x_1 < x_2 < x_3 < x_4$ as below:

$$\begin{array}{cccc} | & 0 & | & 0 \\ \hline x_1 & x_2 & x_3 & x_4 \end{array}$$

$$\left\{ \begin{array}{l} ax_1^2 + bx_1 + c > 0 \\ ax_2^2 + bx_2 + c \leq 0 \\ ax_3^2 + bx_3 + c > 0 \\ ax_4^2 + bx_4 + c \leq 0 \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} ax_1^2 + bx_1 + c > ax_2^2 + bx_2 + c \\ ax_3^2 + bx_3 + c > ax_2^2 + bx_2 + c \\ ax_3^2 + bx_3 + c > ax_4^2 + bx_4 + c \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} (x_1 - x_2)(a(x_1 + x_2) + b) > 0 \\ (x_3 - x_2)(a(x_3 + x_2) + b) > 0 \\ (x_3 - x_4)(a(x_3 + x_4) + b) > 0 \end{array} \right.$$

$$\& \left\{ \begin{array}{l} x_1 - x_2 < 0 \Rightarrow a(x_1 + x_2) + b < 0 \\ x_3 - x_2 > 0 \quad a(x_3 + x_2) + b > 0 \\ x_3 - x_4 < 0 \quad a(x_3 + x_4) + b < 0 \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} a(x_3 - x_1) > 0 \\ a(x_2 - x_4) > 0 \end{array} \right.$$

$$\& \left\{ \begin{array}{l} x_3 - x_1 > 0 \Rightarrow a > 0 \\ x_2 - x_4 < 0 \quad a < 0 \end{array} \right.$$

\Rightarrow No hypothesis in the set can separate these points.

$$\Rightarrow VC < 4$$

\Rightarrow VC dimension of H is 3.

3.(a).

i	Label	Hypothesis 1 (1st iteration)					Hypothesis 2 (2nd iteration)											
		D ₀	f ₁ = [x > 2]	f _r = [y > 6]	h ₁ = [f ₁]	D ₁	f ₁ = [x > 10]	f _r = [y > 11]	h ₂ = [f ₂]	(12)	(13)	(14)	(15)	(16)	(17)	(18)	(19)	(20)
1	-	$\frac{1}{10}$	-	+	-	$\frac{1}{16}$	-	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-
2	-	$\frac{1}{10}$	-	-	-	$\frac{1}{16}$	-	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-
3	+	$\frac{1}{10}$	+	+	+	$\frac{1}{16}$	-	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-
4	-	$\frac{1}{10}$	-	-	-	$\frac{1}{16}$	-	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-
5	-	$\frac{1}{10}$	-	+	-	$\frac{1}{16}$	-	-	-	-	$\frac{1}{10}$	-	+	+	+	-	-	+
6	-	$\frac{1}{10}$	+	+	+	$\frac{1}{4}$	-	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-
7	+	$\frac{1}{10}$	+	+	+	$\frac{1}{16}$	+	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-
8	-	$\frac{1}{10}$	-	-	-	$\frac{1}{16}$	-	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-
9	+	$\frac{1}{10}$	-	+	-	$\frac{1}{4}$	-	-	-	-	$\frac{1}{10}$	-	+	+	+	-	-	+
10	+	$\frac{1}{10}$	+	+	+	$\frac{1}{16}$	-	-	-	-	$\frac{1}{10}$	-	-	-	-	-	-	-

$$(a). D_0 = \frac{1}{10}$$

$$(b). E(f_1) = \frac{1}{2} - \frac{1}{2}(8 \times \frac{1}{10} - 2 \times \frac{1}{10}) = 0.2$$

$$\mathcal{E}(f_2) = \frac{1}{2} - \frac{1}{2}(7 \times \frac{1}{10} - 3 \times \frac{1}{10}) = 0.3$$

$$\varepsilon(f_1) < \varepsilon(f_2) \Rightarrow h_1 = f_1$$

$$(C). \quad \varepsilon_0 = \varepsilon(f_1) = 0.2$$

$$\alpha_0 = \frac{1}{2} \log_2 \left(\frac{1 - \varepsilon_0}{\varepsilon_0} \right) = 1$$

$$Z_0 = D_0 \times 2^{-\frac{d_0}{\lambda}} \times 8 + D_0 \times 2^{\frac{d_0}{\lambda}} \times 2 = 0.8$$

$$D_1(i) = \begin{cases} \frac{D_0 \times 2^{-\alpha_0}}{Z_0} &= \frac{0.05}{0.8} = \frac{1}{16} \quad \text{label}_i = h_1(x_i) \\ \frac{D_0 \times 2^{\alpha_0}}{Z_0} &= \frac{0.2}{0.8} = \frac{1}{4} \quad \text{label}_i \neq h_1(x_i) \end{cases}$$

$\bar{t}=1, 2, \dots, 10$

$$E(f_1) = \frac{1}{2} - \frac{1}{2}(\frac{1}{16} \times 6 + \frac{1}{4} \times 1 - \frac{1}{16} \times 2 - \frac{1}{4} \times 1) = 0.375$$

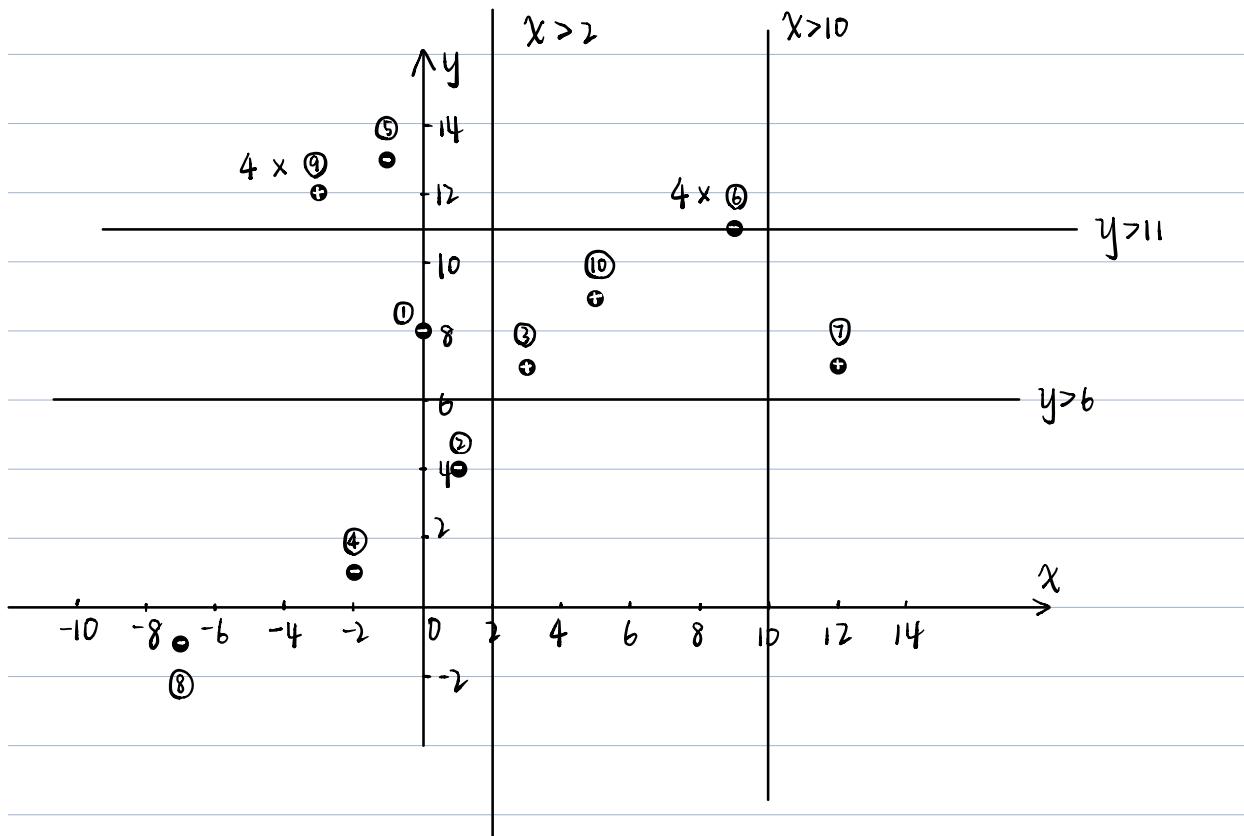
$$\mathcal{E}(f_2) = \frac{1}{2} - \frac{1}{2} (\frac{1}{16} \times 4 + \frac{1}{4} \times 2 - \frac{1}{16} \times 4) = 0.25$$

$$\mathcal{E}(f_1) > \mathcal{E}(f_2) \Rightarrow h_2 = f_2$$

(d). $\mathcal{E}_1 = \mathcal{E}(f_2) = 0.25$

$$\alpha_1 = \frac{1}{2} \log_2 \left(\frac{1 - \mathcal{E}_1}{\mathcal{E}_1} \right) = 0.8$$

$$\Rightarrow H_{\text{final}} = \text{sign}(\alpha_0 h_1 + \alpha_1 h_2) = \text{sign}([x > 2] + 0.8 [y > 11])$$



4.(a). i. One vs. All:

learns K classifiers

All vs. All:

learns $C_k^2 = \frac{K(K-1)}{2}$ classifiers

ii. One vs. All:

uses m examples to learn each classifier.

$\frac{m}{K}$ positive examples & $(m - \frac{m}{K})$ negative examples.

All vs. All:

uses $\frac{2m}{K}$ examples to learn each classifier.

$\frac{m}{K}$ positive examples & $\frac{m}{K}$ negative examples.

iii. One vs. All:

Suppose after the learning process, K classifiers give K vectors respectively: $\vec{w}_1, \vec{w}_2, \dots, \vec{w}_K$, and each classifier classifies one example \vec{x} as positive or negative based on $\text{sign}(\vec{w}_i^\top \vec{x})$, $i=1, 2, \dots, K$ (like Perceptron algorithm)

The final class label is the label corresponding to the vector that gives the maximum $\vec{w}_i^\top \vec{x}$, $i=1, 2, \dots, K$

$$\Rightarrow \text{label} = \arg \max_{y \in \{1, 2, \dots, K\}} \vec{w}_y^\top \vec{x}$$

All vs. All:

After applying all the $\frac{K(K-1)}{2}$ classifiers to example \vec{x} and getting different labels, the final class label is the result from majority vote of these labels.

iv. One vs. All:

There are k classifiers and each classifier uses m examples to learn: $k \cdot m$
 $\Rightarrow O(mk)$

All vs. All:

There are $\frac{k(k-1)}{2}$ classifiers and each classifier uses $\frac{2m}{k}$ examples to learn: $\frac{k(k-1)}{2} \cdot \frac{2m}{k} = m(k-1)$
 $\Rightarrow O(mk)$

(b). I do not have a preference for the two schemes because both of them are able to learn m examples with k labels, and they have the same computational complexity for Perceptron classifier.

(c). Using a KERNEL PERCEPTRON changes the analysis above.

The computational complexity of the KERNEL PERCEPTRON is $O(n^2)$, where n is the number of examples used during training.

\Rightarrow One vs. All:

There are k classifiers and each classifier uses m examples to learn: $k \cdot O(m^2)$
 $\Rightarrow O(m^2k)$

All vs. All:

There are $\frac{k(k-1)}{2}$ classifiers and each classifier uses $\frac{2m}{k}$ examples to learn: $\frac{k(k-1)}{2} \cdot O\left(\frac{2m}{k}\right)^2 = O(m^2)$
 $\Rightarrow O(m^2)$

I would prefer All vs. All scheme when using a KERNEL PERCEPTRON because

it has better computational complexity.

(d). One vs. All:

There are k classifiers and each classifier uses m examples with d dimensionality to learn: $k \cdot O(dm^2)$
 $\Rightarrow O(m^2 dk)$

All vs. All:

There are $\frac{k(k-1)}{2}$ classifiers and each classifier uses $\frac{2m}{k}$ examples with d dimensionality to learn: $\frac{k(k-1)}{2} \cdot O(\frac{2m}{k} d)$
 $\Rightarrow O(m^2 d)$

All vs. All training paradigm is the most efficient.

(e). One vs. All:

There are k classifiers and each classifier uses m examples with d dimensionality to learn: $k \cdot O(dm^2)$
 $\Rightarrow O(md^2 k)$

All vs. All:

There are $\frac{k(k-1)}{2}$ classifiers and each classifier uses $\frac{2m}{k}$ examples with d dimensionality to learn: $\frac{k(k-1)}{2} \cdot O(d^2 \frac{2m}{k})$
 $\Rightarrow O(md^2 k)$

One vs. All and All vs. All are the same efficient.

(f). Counting:

Counting runs all the $\frac{k(k-1)}{2}$ classifiers on one example, and then does the majority vote. $O(k) \Rightarrow O(k^2 + k)$

⇒ The overall evaluation time complexity per example is $O(K^2)$.

Knockout:

There are K classes. Since Knockout compares two classes at a time, and if one loses, it will never be considered again.

⇒ Knockout runs $(K-1)$ classifiers to exclude $(K-1)$ classes and remain one class.

⇒ The overall evaluation time complexity per example is $O(K)$.