Name: Yongqian Li
Student Id # : 004997466

## Submission instructions

- Submit your solutions electronically on the course Gradescope site as PDF files.

- If you plan to typeset your solutions, please use the LaTeX solution template. If you must submit scanned handwritten solutions, please use a black pen on blank white paper and a high-quality scanner app.

# 1 PAC Learning [25 pts]

*all examples*

(a). Iterate in $S_{train}$, for each data point $(x_i, y_i)$ in the $m$ examples

If its label $y_i = 1$, the hypothesis is $h_{x_i}$. ~~If none of~~ ~~the~~ If after iterations no $y_i = 1$ is found, the hypothesis is $h^-$. The training error of this algorithm is always $0$ since the prediction is consistent with ~~the~~ $S_{train}$. *all*
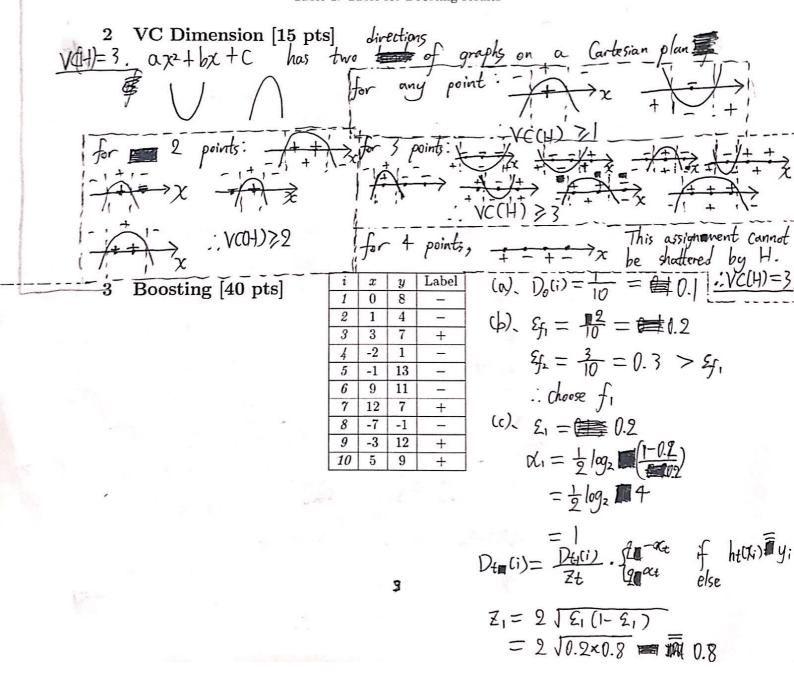
(b). The hypothesis $A(S_{train})$ is always consistent with $S_{train}$, and it makes at most $1$ mistake in $D$ since any $h^* \in H_{singleton}$ is either $h^-$ or has only one $x \in X$ s.t. $h^*(x) = 1$, and mistakes are only made when $A(S_{train})$ is $h^-$ with $h^* \neq h^-$ (Otherwise we would have predicted $h^*$ correctly).

$\therefore$ Loss is either $0$ or $\frac{1}{N}$.

If the true $h^*$ is $h_{x^*}$, we will make mistakes only if $x^*$ does not appear in $S_{train}$. $\therefore P(loss > 0) = (1 - \frac{1}{N})^m$

$\because P(loss > 0) \geq P(loss > \varepsilon)$ $\therefore$ If $\sigma \geq P(loss > 0)$, then $\sigma \geq P(loss > \varepsilon)$

When $loss > \varepsilon$, $loss = \frac{1}{N} \Rightarrow \frac{1}{N} > \varepsilon$ $\therefore (1 - \frac{1}{N})^m < (1 - \varepsilon)^m$

$\therefore$ let $\sigma \geq (1 - \varepsilon)^m > (1 - \frac{1}{N})^m \geq P(loss > \varepsilon)$,

$\sigma \geq (1 - \varepsilon)^m$ $e^{-\varepsilon m} \leq \sigma$

if $-\varepsilon m \leq \log \sigma$

$m \geq -\frac{\log \sigma}{\varepsilon} = \frac{\log \frac{1}{\sigma}}{\varepsilon}$

$\therefore$ If $m \geq \dfrac{\log \frac{1}{\sigma}}{\varepsilon}$,

$\sigma \geq (1 - \varepsilon)^m > (1 - \frac{1}{N})^m$

$\therefore$ Loss is either $0$ or $\frac{1}{N}$ $\qquad \therefore P(\text{loss} > \varepsilon) = P(\text{loss} > 0) = (1 - \frac{1}{N})^m$

$\therefore \quad \sigma \geq (1 - \varepsilon)^m > P(\text{loss} > \varepsilon)$

$\therefore$ when $m \geq \dfrac{\log \frac{1}{\sigma}}{\varepsilon}$, $\qquad \sigma \geq P(l_{0,1}(A(S_{train})) > \varepsilon)$

| | | Hypothesis 1 (1st iteration) | | | | Hypothesis 2 (2nd iteration) | | | |
|---|---|---|---|---|---|---|---|---|---|
| $i$ | Label | $D_0$ | $f_1 \equiv$ $[x > 9]$ | $f_2 \equiv$ $[y > 6]$ | $h_1 \equiv$ $[x > 9]$ | $D_1$ | $f_1 \equiv$ $[x > 9]$ | $f_2 \equiv$ $[y > 11]$ | $h_2 \equiv$ $[y > 11]$ |
| (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) | (9) | (10) |
| 1 | − | 0.1 | − | + | − | 0.0625 | − | − | − |
| 2 | − | 0.1 | − | − | − | 0.0625 | − | − | − |
| 3 | + | 0.1 | + | + | + | 0.0625 | − | − | − |
| 4 | − | 0.1 | − | − | − | 0.0625 | − | − | − |
| 5 | − | 0.1 | − | + | − | 0.0625 | − | + | + |
| 6 | − | 0.1 | + | + | + | 0.25 | − | − | − |
| 7 | + | 0.1 | + | + | + | 0.0625 | + | − | − |
| 8 | − | 0.1 | − | − | − | 0.0625 | − | − | − |
| 9 | + | 0.1 | − | + | − | 0.25 | − | + | + |
| 10 | + | 0.1 | + | + | + | 0.0625 | − | − | − |

Table 1: Table for Boosting results

## 2 VC Dimension [15 pts]

$VC(H) = 3$. $ax^2 + bx + c$ has two directions of graphs on a Cartesian plane



for any point:

$VC(H) \geq 1$

for 2 points:

$\therefore VC(H) \geq 2$

for 3 points:

$\therefore VC(H) \geq 3$

for 4 points:

This assignment cannot be shattered by H.

$\therefore VC(H) = 3$

## 3 Boosting [40 pts]

| $i$ | $x$ | $y$ | Label |
|---|---|---|---|
| 1 | 0 | 8 | − |
| 2 | 1 | 4 | − |
| 3 | 3 | 7 | + |
| 4 | -2 | 1 | − |
| 5 | -1 | 13 | − |
| 6 | 9 | 11 | − |
| 7 | 12 | 7 | + |
| 8 | -7 | -1 | − |
| 9 | -3 | 12 | + |
| 10 | 5 | 9 | + |

(a). $D_0(i) = \dfrac{1}{10} = 0.1$

(b). $\varepsilon_{f_1} = \dfrac{2}{10} = 0.2$

$\varepsilon_{f_2} = \dfrac{3}{10} = 0.3 > \varepsilon_{f_1}$

$\therefore$ choose $f_1$

(c). $\varepsilon_1 = 0.2$

$\alpha_1 = \dfrac{1}{2} \log_2 \left( \dfrac{1-0.2}{0.2} \right)$

$= \dfrac{1}{2} \log_2 4$

$= 1$

$D_{t+1}(i) = \dfrac{D_t(i)}{z_t} \cdot \begin{cases} e^{-\alpha_t} \\ e^{\alpha_t} \end{cases}$ if $h_t(x_i) = y_i$ else

$z_1 = 2\sqrt{\varepsilon_1(1-\varepsilon_1)}$

$= 2\sqrt{0.2 \times 0.8} = 0.8$

$\because D_0(i) = 0.1$

~~$\therefore$ when $h_1(x_i) = y_i$, $D_2(i) = 2^{-1}$~~

$\therefore$ when $h_1(x_i) = y_i$, $D_1(i) = \dfrac{0.1 \cdot \frac{1}{2}}{0.8} = \dfrac{0.1}{0.8 \times 2} \approx$ ~~0.45985~~ ~~0.16~~ $0.0625$

when $h_1(x_i) \neq y_i$, $D_1(i) = \dfrac{0.1 \cdot 2}{0.8} \approx$ ~~0.25~~ $0.25$

(d).

~~$H_{final} = sgn(\sum \alpha_t h_t(x_i))$~~

$\varepsilon_2 = 0.0625 \times 4 = 0.25$

$\alpha_2 = \frac{1}{2} \log_2\left(\frac{1-0.25}{0.25}\right) = \frac{1}{2}\log_2 3 \approx 0.792$

$H_{final}(x) = sgn\left(\sum_t \alpha_t h_t(x_i)\right)$

$\qquad = sgn\left(h_1(x) + 0.792\, h_2(x)\right)$

$\qquad = sgn\left([x > 2] + 0.792\,[y > 11]\right)$

# 4 Multi-class classification [60 pts]

Consider a multi-class classification problem with $k$ class labels $\{1, 2, \ldots k\}$. Assume that we are given $m$ examples, labeled with one of the $k$ class labels. Assume, for simplicity, that we have $m/k$ examples of each type.

Assume that you have a learning algorithm $L$ that can be used to learn Boolean functions. (E.g., think about $L$ as the Perceptron algorithm). We would like to explore several ways to develop learning algorithms for the multi-class classification problem.

There are two schemes to use the algorithm $L$ on the given data set, and produce a multi-class classification:

- **One vs. All:** For every label $i \in [1, k]$, a classifier is learned over the following data set: the examples labeled with the label $i$ are considered "positive", and examples labeled with any other class $j \in [1, k], j \neq i$ are considered "negative".

- **All vs. All:** For every pair of labels $(i, j)$, a classifier is learned over the following data set: the examples labeled with one class $i \in [1, k]$ are considered "positive", and those labeled with the other class $j \in [1, k], j \neq i$ are considered "negative".

(a). i. for One v.s All, $k$ classifiers are learned;

for All vs. All, $\dfrac{k(k-1)}{2}$ classifiers are learned.

ii. for One vs. All, $m$ examples;

for All vs. All, $\dfrac{2m}{k}$ examples.

iii. for One vs. All, the label $y_i = \arg\max_{y \in \{1,2,\cdots,k\}} w_y^T x_i$

for All vs. All, each label gets $(k-1)$ votes,

there are 2 methods:

(1). Tournament: start with any pair of labels, compare $\overset{<a,\ b>}{W_j^T x_i}\ \text{and}\ W$ $W_a^T x_i$ and $W_b^T x_i$

the winner (with larger $W^T x_i$) compare with another label; continue with winners until ~~is~~ the final winner is found

(2). Majority: compare each pair of labels $<a,b>$, the larger one of $W_a^T x_i$ and $W_b^T x_i$ indicates the winner. Then choose the label which wins more often than any other label.

~~would choose All vs. All because for One against All, each class needs to be separable from the others, which is a very strong assumption. Each comparation of One vs. All involves a class with $\frac{m}{k}$ examples and the other class with $\frac{(k-1)m}{k}$ examples, so they may be very unbalanced.~~

iv. The computational complexity for One vs. ~~All~~ All is $k \cdot m \Rightarrow O(km)$; that of All vs. All is $\frac{k(k-1)}{2} \cdot \frac{m}{k} \Rightarrow O(km)$. ~~d is the dimension of each example.~~

(b). I would choose ~~All~~ One vs. All because it ~~is more efficient~~ and only need to store $k$ classifiers instead of $\frac{k(k-1)}{2}$. It is also easier to implement, and performs as good as ~~All~~ All vs. All in this case.

(c). ~~Let d be the dimension of each example.~~ Kernal Perceptron calculates each classifier in $O(\underset{n^2}{\boxed{\phantom{x}}})$. $\therefore$ One vs. All: $O(km^2)$
All vs. All: $\frac{k(k+1)}{2} \cdot (\frac{m}{k})^2 \Rightarrow O(m^2)$
$\therefore$ Now All vs. All is more efficient, so I ~~prefer~~ would prefer All vs. All.
(since it needs $K(x_i, x_j)$ between each pair)

(d). One vs. All: $\underset{k \cdot O(dm^2)}{\cancel{k \cdot dm}} \Rightarrow O(kdm^2)$
All vs. All: $\frac{k(k-1)}{2} \cdot O(d \frac{m^2}{k^2}) \Rightarrow O(dm^2)$
$\therefore$ All vs. All is ~~more~~ more efficient.

(e). One vs. All: $k \cdot O(d^2 m) \Rightarrow O(kd^2 m)$
All vs. All: $\frac{k(k-1)}{2} \cdot O(d^2 \frac{m}{k}) \Rightarrow O(kd^2 m)$
$\therefore$ They are the same efficient.

(f). Counting: for each example, test all $\frac{k(k-1)}{2}$ classifiers
$\therefore \frac{k(k-1)}{2} \cdot d \Rightarrow O(k^2 d)$

Knockout: for each example, eliminate 1 classifier at each step until the winner is found. $\therefore (k-1)$ steps.
$(k-1) \cdot d \Rightarrow O(kd)$