FIRST NAME: _MohammadReza_     LAST NAME: _Ebrahimi_

STUDENT NUMBER: _____

# ECE 421S/ECE1513S — Introduction to Machine Learning
## Final Examination

**April 17th, 2019**
**6:30 p.m. − 9:00 p.m.**

Instructors: Ashish Khisti and Ben Liang and Amir Ashouri

## Instructions

- Please read the following instructions carefully.
- You have 2 hour 30 minutes to complete the exam.
- Please make sure that you have a complete exam booklet.
- Please answer *all* questions. Read each question carefully.
- The value of each question is indicated. Allocate your time wisely!
- $\mathcal{N}(\mu, \sigma^2)$ denotes a Gaussian density function with mean $\mu$ and variance $\sigma^2$.
- No additional pages will be collected beyond this answer book.
- This examination is closed-book; One 8.5 × 11 aid-sheet is permitted. A non-programmable calculator is also allowed.
- Good luck!

1. (20 MARKS) Consider linear regression with training data obtained from a noisy target function $y = f(\mathbf{x}) = \mathbf{c}^T\mathbf{x} + \epsilon$, where $y \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^{d+1}$ (with bias term $x_0 = 1$ as usual), $\mathbf{c} \in \mathbb{R}^{d+1}$, and **the noise term $\epsilon$ is independently generated with zero mean and variance $\sigma^2$**. Suppose we have a set of $N$ such training examples, $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_N, y_N)\}$, $i.e.$, $y_n = \mathbf{c}^T\mathbf{x}_n + \epsilon_n$ for $n = 1, 2, \ldots, N$.

Let $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$, $\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix}$, and $X = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}$. We assume the columns of $X$ are linearly independent, so that $X^TX$ is invertible. Let $\mathbf{w}_{\text{LS}}$ be the linear least-squares estimator, $i.e.$, the in-sample error $E_{\text{in}}(\mathbf{w}) = \frac{1}{N}\|X\mathbf{w} - \mathbf{y}\|^2$ is minimized when $\mathbf{w} = \mathbf{w}_{\text{LS}}$.

5 marks

(a) Show that $\mathbf{w}_{\text{LS}} = X^\dagger \mathbf{y}$, where $X^\dagger$ is the pseudo-inverse of $X$, $i.e.$, $X^\dagger = (X^TX)^{-1}X^T$. You may use without proof the fact that the gradient of $E_{\text{in}}(\mathbf{w})$ with respect to $\mathbf{w}$ is $\frac{2}{N}X^T(X\mathbf{w} - \mathbf{y})$.

$$\nabla_W \, E_{in}(w) \Big|_{W=W_{LS}} = 0$$

$$\Rightarrow \frac{2}{N} X^T (X W_{LS} - y) = 0 \quad \textcircled{1}$$

$$\Rightarrow X^TX \, W_{LS} - X^T y = 0$$

$$\Rightarrow W_{LS} = (X^TX)^{-1} X^T y \quad \textcircled{4}$$

$$= X^\dagger y$$

Scanned by CamScanner

5 marks

(b) Find a simplified expression for $E_{\text{in}}(\mathbf{w}_{\text{LS}})$ as a function of $N$, $X$, $X^{\dagger}$, and $\epsilon$. You may use without proof the fact that $(I - XX^{\dagger})^k = I - XX^{\dagger}$ for any positive integer $k$.

$$E_{in}(W_{LS}) = \frac{1}{N} \| X w_{LS} - y \|^2$$

$$= \frac{1}{N} \| XX^{\dagger}y - y \|^2$$

$$= \frac{1}{N} \| (XX^{\dagger} - I) y \|^2$$

$$= \frac{1}{N} \| (XX^{\dagger} - I)(Xc + \epsilon) \|^2$$

$$= \frac{1}{N} \| \underbrace{(XX^{\dagger} - I)Xc}_{0} + (XX^{\dagger} - I)\epsilon \|^2 \quad \textcircled{2}$$

$$= \boxed{\frac{1}{N} \| (XX^{\dagger} - I)\epsilon \|^2} \quad \textcircled{3}$$

**Note:**
$$= \frac{1}{N} \| (I - XX^{\dagger})\epsilon \|^2$$

$$= \frac{1}{N} \epsilon^T \underbrace{(XX^{\dagger} - I)^T (XX^{\dagger} - I)}_{symmetric} \epsilon$$

$$= \frac{1}{N} \epsilon^T (XX^{\dagger} - I)^2 \epsilon$$

$$= \boxed{\frac{1}{N} \epsilon^T (I - XX^{\dagger}) \epsilon} \quad \longleftarrow \text{This is acceptable too.}$$

**5 marks**

(c) Observe the result in Part (b) for the extreme case of $d = 0$, i.e., $X = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$. Find the expected

in-sample error $\mathbb{E}_\epsilon(E_{\text{in}}(\mathbf{w}_{\text{LS}}))$ as a function of $N$ and $\sigma$. (A side comment that can be safely ignored: a similar conclusion can be drawn for general $d$ values, but you do not need to prove that.)

$$E_{in}(W_{LS}) = \frac{1}{N} e^T (I - XX^\dagger) e$$

$$= \frac{1}{N} tr\left( e^T (I - XX^\dagger) e \right)$$

$$= \frac{1}{N} tr\left( (I - XX^\dagger) e e^T \right)$$

$$\Rightarrow \mathbb{E}_\epsilon \left[ E_{in}(W_{LS}) \right] = \frac{1}{N} tr\left( (I - XX^\dagger) \underbrace{\mathbb{E}[e e^T]}_{\sigma^2 I} \right)$$

$$= \frac{\sigma^2}{N} tr\left( I - XX^\dagger \right)$$

$$tr(XX^\dagger) = tr(X^\dagger X) = tr(I_{d+1}) = d+1$$

$$\Rightarrow \mathbb{E}_\epsilon \left[ E_{in}(W_{LS}) \right] = \frac{\sigma^2}{N} (N - d - 1) = \sigma^2 \left( 1 - \frac{d+1}{N} \right)$$

solution for general $d$ values.

Scanned by CamScanner

5 marks

(c) Observe the result in Part (b) for the extreme case of $d = 0$, *i.e.*, $X = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$. Find the expected

in-sample error $\mathbb{E}_\epsilon[E_{in}(\mathbf{w}_{LS})]$ as a function of $N$ and $\sigma$. (Side comment: you should see that $\mathbb{E}_\epsilon[E_{in}(\mathbf{w}_{LS})]$ strictly increases in $N$ and converges to $\sigma^2$. A similar conclusion holds for general $d$ values, but you do not need to prove that.)

$$E_{in}(W_{LS}) = \frac{1}{N}\|(XX^\dagger - I)\epsilon\|^2$$

$$X^\dagger = (X^TX)^{-1}X^T = (N)^{-1}X^T = \frac{1}{N}[1,\cdots,1]$$

$$\Rightarrow E_{in}(W_{LS}) = \frac{1}{N}\|(\frac{1}{N}XX^T - I)\epsilon\|^2 = \frac{1}{N^3}\left\|\begin{bmatrix} \epsilon_1 + \cdots + \epsilon_N - N\epsilon_1 \\ \epsilon_1 + \cdots + \epsilon_N - N\epsilon_2 \\ \vdots \\ \epsilon_1 + \cdots + \epsilon_N - N\epsilon_N \end{bmatrix}\right\|^2$$

$$= \frac{1}{N^3}\left[(\epsilon_1 + \cdots + \epsilon_N - N\epsilon_1)^2 + \cdots + (\epsilon_1 + \cdots + \epsilon_N - N\epsilon_N)^2\right]$$

$$\Rightarrow \mathbb{E}_\epsilon(E_{in}(W_{LS})) = \frac{1}{N^3}\left[\text{var}\left((1-N)\epsilon_1 + \epsilon_2 + \cdots + \epsilon_N\right) + \cdots + \text{var}\left(\epsilon_1 + \cdots + (1-N)\epsilon_N\right)\right]$$

$$= \frac{1}{N^3}\left[N\left((1-N)^2 + (N-1)\sigma^2\right)\right] = \sigma^2(1-\frac{1}{N}) \quad \text{(5)}$$

*other solutions w/ same answer are ok*

5 marks

(d) For any **out-of-sample** data point $\mathbf{x}$, the expected out-of-sample error is $\mathbb{E}_\epsilon[(f(\mathbf{x}) - \mathbf{w}_{LS}^T\mathbf{x})^2]$. Show that it is at least $\sigma^2$.

$$\mathbb{E}_\epsilon\left((f(x) - W_{LS}^Tx)^2\right) = \mathbb{E}_\epsilon\left((c^Tx + \epsilon - W_{LS}^Tx)^2\right)$$

$$= \mathbb{E}_\epsilon\left[((c - W_{LS})^Tx + \epsilon)^2\right]$$

$$= \left[(c - W_{LS})^Tx\right]^2 + \mathbb{E}(\epsilon^2)$$

$$= \underbrace{\left[(c - W_{LS})^Tx\right]^2}_{\geqslant 0} + \sigma^2 \geqslant \sigma^2 \quad \text{(5)}$$

**2.** (20 MARKS) Consider a binary classification problem in two dimensions. Let $\mathcal{H}_0$ denote the hypothesis class of all linear classifiers in two dimensions passing **through the origin**. Thus any $h \in \mathcal{H}_0$ can be expressed as:

$$h(\mathbf{x}) = \text{sign}(w_1 x_1 + w_2 x_2)$$

where $\mathbf{x} = (x_1, x_2)$ is the data point and $\mathbf{w} = (w_1, w_2) \in \mathbb{R}^2$.
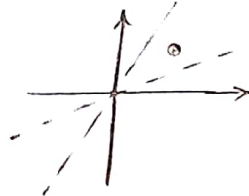
Furthermore let $g_1(\mathbf{x})$ and $g_2(\mathbf{x})$ be arbitrary classification functions such that $g_i : \mathbb{R}^2 \to \{-1, +1\}$ and assume that neither $g_1(\cdot)$ nor $g_2(\cdot)$ are in $\mathcal{H}_0$.

6 marks      (a) Let $m_{\mathcal{H}_0}(N)$ denote the growth function of the hypothesis class $\mathcal{H}_0$. Compute the growth function for $N = 1, 2, 3, 4$. Provide a brief justification for each computation.
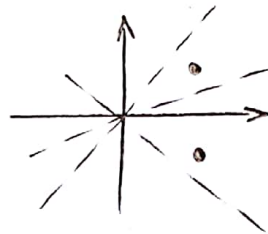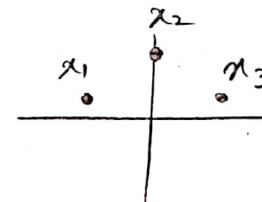
① $- m_{\mathcal{H}_0}(1) = 2$
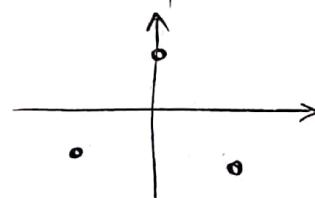
① $- m_{\mathcal{H}_0}(2) = 4$

② $- m_{\mathcal{H}_0}(3) = 6$

① if the points are on one side of a line

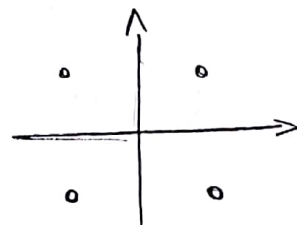② otherwise

$x_1$   $x_2$   $x_3$    cannot generate

$-+-/+-+$

cannot generate

$---/+++$

② $- m_{\mathcal{H}_0}(4) = 8$

cannot select only one point. $\longrightarrow$ $4 \times 2$ possibilities cannot be generated.

if no justification provided $-2$

**4 marks**

(b) Let us define a new hypothesis class $\mathcal{M} = \mathcal{H}_0 \cup \{g_1, g_2\}$. Let $m_{\mathcal{M}}(N)$ denote the growth function of $\mathcal{M}$. What is the maximum possible value $m_{\mathcal{M}}(N)$ for $N = 1, 2, 3, 4$ over all choices of $g_1(\cdot)$ and $g_2(\cdot)$.

$g_1$ and $g_2$ can at most add one element to $\mathcal{M}(\underline{x}_1, \cdots, \underline{x}_N)$

each, as long as $|\mathcal{M}(\underline{x}_1, \cdots, \underline{x}_N)| \leq 2^N$

$$\Rightarrow \max_{g_1, g_2} m_{\mathcal{M}}(N) = \begin{cases} 2 & N = 1 \\ 4 & N = 2 \\ 8 & N = 3 \\ 10 & N = 4 \end{cases}$$

**4 marks**

(c) Suppose we train the hypothesis class $\mathcal{M}$ defined in part (b) over a training set $\mathcal{D} = \{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N)\}$ (where $\mathbf{x}_i \overset{iid}{\sim} p_{\mathbf{x}}(\cdot)$ and $y_i = f(\mathbf{x}_i)$) and produce an output hypothesis $m(\cdot) \in \mathcal{M}$ that minimizes the training error. Let $E_{\text{in}}(m)$ denote the average classification error of $m(\cdot)$ on the training set $\mathcal{D}$ and $E_{\text{out}}(m)$ denote the test error for this hypothesis. Assume that $g_1(\cdot)$ and $g_2(\cdot)$ are **fixed** hypothesis selected before observing $\mathcal{D}$. Find a relation between $E_{\text{in}}(m)$ and $E_{\text{out}}(m)$ in the following sense:

With Probability at-least $1 - \delta$, $E_{\text{out}}(m) \leq E_{\text{in}}(m) + \Omega$, where you should provide a bound on $\Omega$ using the VC dimension theory.

VC bound: $\Omega = \sqrt{\dfrac{8}{N} \log\left( \dfrac{4 m_{\mathcal{M}}(2N)}{\delta} \right)}$

$m_{\mathcal{M}}(N) \leq N^{d_{vc}(\mathcal{M})} + 1$

part b $\longrightarrow$ $\boxed{d_{vc}(\mathcal{M}) = 3}$ ②

$$\Rightarrow \Omega \leq \sqrt{\dfrac{8}{N} \log\left( \dfrac{4((2N)^3 + 1)}{\delta} \right)}$$

②

Scanned by CamScanner

**6 marks**

(d) In this part, suppose that the hypothesis $g_1(\mathbf{x}) = 1$ for all $\mathbf{x} \in \mathbb{R}^2$. However $g_2(\cdot)$ is selected **after** observing the dataset $\mathcal{D}$ (stated in part (c)) as follows:

Let $\mathcal{F}$ denote the class of all hypothesis where the decision boundaries are horizontal in the $x_1 - x_2$ plane i.e.,
$$f(\mathbf{x}) = \text{sign}(a \cdot x_2 - c)$$
where $c \in \mathbb{R}$ and $a \in \{-1, +1\}$. We select $g_2(\cdot)$ by finding $f \in \mathcal{F}$ that minimizes the average classification error over $\mathcal{D}$, i.e., $E_{\text{in}}(f)$

The hypothesis $g_2(\cdot)$, thus selected, is included in $\mathcal{M}$ along with $g_1(\cdot)$. We train the hypothesis class $\mathcal{M}$ on the **same** dataset $\mathcal{D}$ (used to select $g_2(\cdot)$), and output $m \in \mathcal{M}$ that minimizes the training error.

Find a relation between $E_{\text{in}}(m)$ and $E_{\text{out}}(m)$ in the following sense: with probability at-least $1 - \delta$, $E_{\text{out}}(m) \leq E_{\text{in}}(m) + \Omega$, where you should provide a bound on $\Omega$ using the VC dimension theory.

This is equivalent to use hypothesis class $\mathcal{H}_1 = \mathcal{H}_0 \cup \mathcal{F}$ ①

$$m_{\mathcal{H}_1}(3) = 8$$
$$m_{\mathcal{H}_1}(4) \leq 10 \longrightarrow \text{from part } b \quad\Bigg\} \Rightarrow \boxed{d_{vc}(\mathcal{H}_1) = 3} \quad ④$$

$$\Rightarrow \quad \Omega \leq \sqrt{\frac{8}{N} \log\left(\frac{4\left((2N)^3 + 1\right)}{\delta}\right)} \quad ①$$

if used upper bound on $d_{vc}$:

$$d_{vc}(\mathcal{H}_0 \cup \mathcal{F}) \leq 1 + d_{vc}(\mathcal{H}_0) + d_{vc}(\mathcal{F}) = 5$$

assign half mark. 3/6

**total/6**

3. (20 MARKS) Consider a target function $f(x) = x^2$ where the domain of input $x$ is the interval $[-1, 1]$. The training data set contains only one example: $\mathcal{D} = \{(U, U^2)\}$, where $U$ is sampled uniformly from $[-1, 1]$. Our hypothesis set $\mathcal{H}$ consists of all lines of the form $h(x) = ax$, for $a \in \mathbb{R}$. Let $g^{\mathcal{D}}(x)$ be the output hypothesis given training data set $\mathcal{D}$.

5 marks      (a) Show that the average hypothesis, *i.e.*, the expectation of $g^{\mathcal{D}}(x)$ over $\mathcal{D}$, is $\bar{g}(x) = 0$.

$$g^{\mathcal{D}}(U) = U^2 \Rightarrow a^{\mathcal{D}} U = U^2 \Rightarrow a^{\mathcal{D}} = U$$

$$\Rightarrow \boxed{g^{\mathcal{D}}(x) = U x} \quad (2)$$

$$\bar{g}(x) = \mathbb{E}_U(U x) = \underbrace{\mathbb{E}_U(U)}_{0} x = 0 \Rightarrow \boxed{\bar{g}(x) = 0} \quad (3)$$

5 marks      (b) Find bias$(x)$ and var$(x)$.

$$\text{bias}(x) = \left(\bar{g}(x) - f(x)\right)^2 = (0 - x^2)^2 = \underline{x^4} \quad (2)$$

$$\text{var}(x) = \mathbb{E}_{\mathcal{D}}\left[\left(g^{\mathcal{D}}(x) - \bar{g}(x)\right)^2\right]$$

$$= \mathbb{E}_U\left[(U - \bar{U})^2 x^2\right] = \text{var}(U) \, x^2$$

$$= \underline{\frac{1}{3} x^2} \quad (3)$$

Scanned by CamScanner

5 marks      (c) Find the expected out-of-sample error $\mathbb{E}_D[E_{out}(g^D)]$.

Solution 1

$$\mathbb{E}_D\left[E_{out}(g^D)\right] = \mathbb{E}_x\left[bias(x) + var(x)\right] \;①$$

$$= \mathbb{E}_x\left(\tfrac{1}{3}x^2 + x^4\right)$$

$$= \int_{-1} \tfrac{1}{2}\left(\tfrac{1}{3}x^2 + x^4\right)dx = \tfrac{1}{9} + \tfrac{1}{5} = \tfrac{14}{45} \;③$$

Solution 2

$$\mathbb{E}_D\left[E_{out}(g^D)\right] = \mathbb{E}_u\left[\mathbb{E}_x\left((g^D(x) - f(x))^2\right)\right]$$

$$= \mathbb{E}_{u,x}\left[(ux - x^2)^2\right] = \cdots$$

acceptable if answer is $\tfrac{14}{45}$

5 marks      (d) Suppose we actually have two data examples: $\mathcal{D} = \{(x_1, y_1), (x_2, y_2)\}$, where $x_1$ and $x_2$ are independently and uniformly sampled from $[-1, 1]$. However, we use the leave-one-out cross validation method, so that the training data set always contains only one example. What is the expectation of the cross validation estimation $E_{CV}$ over $\mathcal{D}$?

Solution 1

$$\mathbb{E}_D\left[E_{cv}\right] = \mathbb{E}_D\left[\tfrac{1}{2}(e_1 + e_2)\right] = \mathbb{E}_D\left[e_1\right]$$

$$e_1 = \left(f(x_1) - \bar{g}_1(x_1)\right)^2$$

$$= \left(x_1^2 - x_2 x_1\right)^2 = x_1^4 + x_2^2 x_1^2 - 2x_2 x_1^3 \;②$$

$$\mathbb{E}_D\left[e_1\right] = \mathbb{E}_{x_1, x_2}\left[x_1^4 + x_2^2 x_1^2 - 2x_2 x_1^3\right]$$

$$= \mathbb{E}(x_1^4) + \mathbb{E}(x_1^2)^2 - 2\mathbb{E}(x_2)\mathbb{E}(x_1^3) \quad \overset{0}{\diagdown}$$

$$= \tfrac{1}{5} + \left(\tfrac{1}{3}\right)^2 = \tfrac{14}{45} \;③$$

Sol 2   $E_{cv}$ is an unbiased estimator of average out of sample error.

total/10

$$\Rightarrow part\ (c) \longrightarrow \mathbb{E}^D(E_{cv}) = \mathbb{E}^D[E_{out}] = \tfrac{14}{45}$$

4. (20 MARKS) In ancient times, there was a village surrounded by hundreds of lakes. Each lake was either poisonous (P) or healthy (H). Anyone who ate fish from a poisonous lake would die immediately while anyone who ate fish from a healthy lake would survive. All fish looked identical and tasted identical and there was no way of knowing whether a lake was poisonous or healthy.

Fortunately a famous chemist visited the village and was told of this dilemma. She suggested using the pH level of water to determine whether a given lake was poisonous or healthy. She hypothesized that lakes with poisonous fish would have higher pH value than healthy lakes. Accordingly she visited each lake and collected the pH value of the water in each lake. The data set is denoted by:

$$\mathcal{D} = \{l_1, l_2, \ldots, l_N\}$$

where $l_i$ denotes the pH level of lake $i \in \{1, 2, \ldots, N\}$. We assume that $l_1 \leq l_2 \leq l_3 \leq \ldots \leq l_N$.

You are hired as a machine learning scientist to help determine the probability that a lake is poisonous given its pH value. In order to do this you propose to use a Gaussian Mixture Model (GMM) as follows:

- Pr(lake is poisonous) = $p_1$
- Pr(lake is healthy) = $p_2$
- $f(\text{pH} = l \mid \text{lake is poisonous}) \sim \mathcal{N}(\mu_1, \sigma_1^2)$
- $f(\text{pH} = l \mid \text{lake is healthy}) \sim \mathcal{N}(\mu_2, \sigma_2^2)$
- $\boxed{\mu_1 \leq \mu_2}$ ✗ $\mu_1 \gg \mu_2$

Here $f(\cdot)$ denotes the conditional density function for the pH value. You decide to use the EM algorithm to train the above GMM on the dataset $\mathcal{D}$.

5 marks

(a) Using the above GMM, provide an expression of the probability that a lake is poisonous given that its pH level is measured to be $l$.

$$Pr(P \mid l) = \frac{f(l \mid P)\, P_1}{f(l \mid P)\, P_1 + f(l \mid H)\, P_2} \quad \textcircled{5}$$

or in terms of $\mathcal{N}(\mu_i, \sigma_i^2)$

Scanned by CamScanner

5 marks

**(b)** Write down the pseudocode for the EM algorithm with hard decisions that finds the parameters of the GMM. Assume that we initialize the algorithm in such a way that $B_1^0 = \{l_1, l_2, \ldots, l_K\}$ denotes one cluster of lakes and $B_2^0 = \{l_{K+1}, \ldots, l_N\}$ denotes it's complement.

- Initialize $B_1^{(new)}$, $B_2^{(new)}$ by $B_1^0$, $B_2^0$

- $B_1 \leftarrow B_1^{(new)}$, $B_2 \leftarrow B_2^{(new)}$

① $\mu_1 \leftarrow \dfrac{1}{|B_1|} \sum_{l_i \in B_1} l_i$ , $\mu_2 \leftarrow \dfrac{1}{|B_2|} \sum_{l_i \in B_2} l_i$

① $\sigma_1^2 \leftarrow \dfrac{1}{|B_1|} \sum_{l_i \in B_1} (l_i - \mu_1)^2$, $\sigma_2^2 \leftarrow \dfrac{1}{|B_2|} \sum_{l_i \in B_2} (l_i - \mu_2)^2$

① $P_1 \leftarrow |B_1|/N$ , $P_2 \leftarrow |B_2|/N$

- $B_1^{(new)}$ , $B_2^{(new)} \leftarrow \phi$

- for each $l$ in $B_1 \cup B_2$:

① $\begin{cases} - \text{if } Pr(P|l) \geqslant \frac{1}{2} : & B_1^{(new)} \leftarrow B_1^{(new)} \cup l \\ - \text{else}: & B_2^{(new)} \leftarrow B_2^{(new)} \cup l \end{cases}$

- run untill $[B_1^{(new)}, B_2^{(new)}] = [B_1, B_2]$ ①

5 marks

**(c)** Write down the pseudocode for the EM algorithm with soft decisions that finds the parameters of the GMM. Do state the initialization of your algorithm explicitly.

- Initialize $\mu_i \leftarrow mean(B_i^0)$, $\sigma_i^2 \leftarrow var(B_i^0)$, $P_i \leftarrow \dfrac{|B_i|}{N}$

⟶ // E-step:

- for $i$ in 1 to N:

① $r_i^{(1)} \leftarrow Pr(P|l_i)$, $r_i^{(2)} \leftarrow 1 - Pr(P|l_i)$

// M-step:

- for $c$ in 1 to 2:

① $\mu_c = \dfrac{\sum_{i=1}^{N} r_i^{(c)} l_i}{\sum_{i=1}^{N} r_i^{(c)}}$ , ① $\sigma_c^2 = \dfrac{\sum_{i=1}^{N} (l_i - \mu_c)^2 r_i^{(c)}}{\sum_{i=1}^{N} r_i^{(c)}}$ , ① $P_c = \dfrac{\sum r_i^{(c)}}{N}$

- run untill log-likelihood change $< e$

(-1) each

if init. or stopping not specified

total/10

Scanned by CamScanner

5 marks      **(d)** Suppose that $N = 5$ and we observe $\mathcal{D} = \{1, 2, 3, 6, 10\}$. Let $\mathcal{B}_1 = \{1, 2, 3\}$ and $\mathcal{B}_2 = \{6, 10\}$. Execute the hard decision EM algorithm and compute the parameters of the GMM.

$\underline{\text{Iteration 1}}$ :      $B_1 = \{1, 2, 3\}$   $B_2 = \{6, 10\}$

$\text{(2)}$   M $\begin{cases} \mu_1 = 2 & \mu_2 = 8 \\ \sigma_1^2 = 0.667 & \sigma_2^2 = 4 \\ P_1 = 3/5 & P_2 = 2/5 \end{cases}$

$\text{(3)}$   E $\begin{cases} Pr(1 \in B_1 \mid \ell = 1) = 0.998 \longrightarrow 1 \in B_1^{(new)} \\ Pr(2 \in B_1 \mid \ell = 2) = 0.997 \longrightarrow 2 \in B_1^{(new)} \\ Pr(3 \in B_1 \mid \ell = 3) = 0.975 \longrightarrow 3 \in B_1^{(new)} \\ Pr(6 \in B_1 \mid \ell = 6) = 0 \longrightarrow 6 \in B_2^{(new)} \\ Pr(10 \in B_1 \mid \ell = 10) = 0 \longrightarrow 10 \in B_2^{(new)} \end{cases}$

$\Rightarrow B_1^{(new)} = \{1, 2, 3\}$      $B_2^{(new)} = \{6, 10\}$
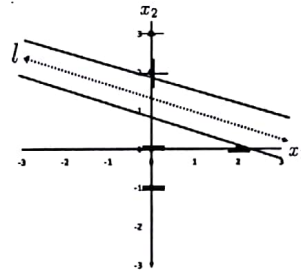
$B_1^{(new)} = B_1 \Rightarrow$ stop.

Scanned by CamScanner

1. (20 MARKS) Consider a dataset with 5 data points in $\mathbb{R}^2$

$X_1 = (0,0)$    $y_1 = -1$
$X_1 = (2,0)$    $y_2 = -1$
$X_1 = (0,-1)$   $y_2 = -1$
$X_1 = (0,2)$    $y_3 = +1$
$X_1 = (0,3)$    $y_4 = +1$

Figures below correctly classify the given data using a hyperplane ($l$)



(a)                (b)

3 marks     (a) Given the provided classifiers and the 5 data points above, choose the one that optimally does a maximum margin classification using SVM.

(a) has a larger margin.

(3)

5 marks     (b) Using the given data points, formulate a Q-P optimization by writing down the objective function and the constraints.

$$\begin{cases} \min\limits_{w,b} \quad \frac{1}{2} w^T w \\ s.t. \quad y_n (w^T x_n + b) \geq 1 \end{cases}$$

$$\equiv \begin{cases} \min\limits_{w_1, w_2, b} \quad \frac{1}{2}(w_1^2 + w_2^2) \longrightarrow \text{(-2) if wrong} \\ \\ s.t. \quad -b \geq 1 \qquad (1) \quad \bigcirc \\ \qquad -2w_1 - b \geq 1 \quad (2) \quad \bigcirc \\ \qquad -w_2 - b \geq 1 \quad (3) \quad \bigcirc \\ \qquad 2w_2 + b \geq 1 \quad (4) \quad \bigcirc \\ \qquad 3w_2 + b \geq 1 \quad (5) \quad \bigcirc \end{cases}$$

total/10

7 marks      (c) Show/argue that the solution you chose in (a) is the optimal margin for the given points.

combining (1) and (4) $\longrightarrow$ $W_2 \geqslant 1$ ②

combining (1) and (2) $\longrightarrow$ $W_1 \leq 0$ ②

$\Rightarrow$ $W_2^* = 1$, $W_1^* = 0$, $b^* = -1$ $\Rightarrow$ max-margin hyperplane:

②      ① $x_2 - 1 = 0$

5 marks      (d) Find the support vectors and the $E_{CV}$ bound for the following data.

① ① ①

support vectors: $x_1$ $x_2$ $x_4$

$E_{CV} \leq \dfrac{\text{\# of support vectors}}{N} = 3/5$ ②

total/10