# Indian Institute of Technology Kharagpur
## Department of Computer Science and Engineering

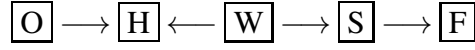| | | |
|---|---|---|
| **Machine Learning (CS60050)** | End-Semester Examination | **Spring Semester, 2022-2023** |
| **Date:** 18-Apr-2023 (AN) | Answer _all_ questions. | **Maximum Marks:** 100 |

**Q1.** **[ Bayesian Networks ]** | 11 marks |

Let us learn some aspects about our life through inference in the following Bayesian Network shown.

$$O \longrightarrow H \longleftarrow W \longrightarrow S \longrightarrow F$$

The variables of our interest are as follows:

    **O**: being Optimistic,      **W**: Working hard,      **H**: being Happy,

    **S**: founding a Start-up company,      **F**: being Famous.

The conditional probability tables for the model are given as:

$$\mathbb{P}(O=true) = 0.5 \quad \mathbb{P}(W=true) = 0.4$$
$$\mathbb{P}(H=true \mid O=true,\ W=true) = 0.9 \quad \mathbb{P}(H=true \mid O=true,\ W=false) = 0.7$$
$$\mathbb{P}(H=true \mid O=false,\ W=true) = 0.5 \quad \mathbb{P}(H=true \mid O=false,\ W=false) = 0.2$$
$$\mathbb{P}(S=true \mid W=true) = 0.6 \quad \mathbb{P}(S=true \mid W=false) = 0.2$$
$$\mathbb{P}(F=true \mid S=true) = 0.4 \quad \mathbb{P}(F=true \mid S=false) = 0.1$$

Compute the following probabilities. Show your calculations in details.

**(a)** $\mathbb{P}(H=false \mid O=false,\ W=true,\ S=true,\ F=true) = ?$ **(3)**
**Solution:**

$$\mathbb{P}(H=f \mid O=f,W=t,S=t,F=t) = \frac{\mathbb{P}(H=f,O=f,W=t,S=t,F=t)}{\mathbb{P}(O=f,W=t,S=t,F=t)}$$

$$= \frac{\mathbb{P}(H=f,O=f,W=t,S=t,F=t)}{\sum\limits_{h\in\{t,f\}} \mathbb{P}(H=h,O=f,W=t,S=t,F=t)}$$

$$= \frac{0.5 \times 0.5 \times 0.4 \times 0.6 \times 0.4}{(0.5 \times 0.5 \times 0.4 \times 0.6 \times 0.4) + (0.5 \times 0.5 \times 0.4 \times 0.6 \times 0.4)} = 0.5$$

**(b)** $\mathbb{P}(H=true \mid S=true,\ F=true) = ?$ **(4)**
**Solution:**

$$\mathbb{P}(H=t \mid S=t,F=t) = \frac{\mathbb{P}(H=t,S=t,F=t)}{\mathbb{P}(S=t,F=t)}$$

$$= \frac{\sum\limits_{o\in\{t,f\}}\sum\limits_{w\in\{t,f\}} \mathbb{P}(H=t,S=t,F=t,O=o,W=w)}{\sum\limits_{h\in\{t,f\}}\sum\limits_{o\in\{t,f\}}\sum\limits_{w\in\{t,f\}} \mathbb{P}(H=h,S=t,F=t,O=o,W=w)}$$

$$= \frac{\mathbb{P}(F=t \mid S=t) \sum\limits_{w\in\{t,f\}} \mathbb{P}(W=w)\mathbb{P}(S=t \mid W=w) \sum\limits_{o\in\{t,f\}} \mathbb{P}(O=o)\mathbb{P}(H=t \mid O=o,W=w)}{\mathbb{P}(F=t \mid S=t) \sum\limits_{w\in\{t,f\}} \mathbb{P}(W=w)\mathbb{P}(S=t \mid W=w) \sum\limits_{o\in\{t,f\}} \mathbb{P}(O=o) \sum\limits_{h\in\{t,f\}} \mathbb{P}(H=h \mid O=o,W=w)}$$

$$= \frac{0.4 \times (0.4 \times 0.6 \times (0.5 \times 0.9 + 0.5 \times 0.5) + 0.6 \times 0.2 \times (0.5 \times 0.7 + 0.5 \times 0.2))}{0.4 \times (0.4 \times 0.6 \times 1 + 0.6 \times 0.2 \times 1)} = 0.6167$$
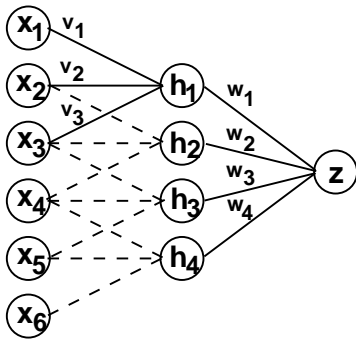
**(c)** $\mathbb{P}(F = true \mid H = true) = ?$ **(4)**

**Solution:**

$$\mathbb{P}(F = t \mid H = t) = \frac{\mathbb{P}(F = t, H = t)}{\mathbb{P}(H = t)}$$

$$= \frac{\displaystyle\sum_{o \in \{t,f\}} \mathbb{P}(O = o) \sum_{w \in \{t,f\}} \mathbb{P}(W = w)\mathbb{P}(H = t \mid O = o, W = w) \sum_{s \in \{t,f\}} \mathbb{P}(S = s \mid W = w)\mathbb{P}(F = t \mid S = s)}{\displaystyle\sum_{o \in \{t,f\}} \mathbb{P}(O = o) \sum_{w \in \{t,f\}} \mathbb{P}(W = w)\mathbb{P}(H = t \mid O = o, W = w) \sum_{s \in \{t,f\}} \mathbb{P}(S = s \mid W = w) \sum_{g \in \{t,f\}} \mathbb{P}(F = g \mid S = s)}$$

$$= \frac{0.5 \times (0.4 \times 0.9 \times (0.6 \times 0.4 + 0.4 \times 0.1) + 0.6 \times 0.7 \times (0.2 \times 0.4 + 0.8 \times 0.1) + 0.5 \times (0.4 \times 0.5 \times (0.6 \times 0.4 + 0.4 \times 0.1) + 0.6 \times 0.2 \times (0.2 \times 0.4 + 0.8 \times 0.1))}{0.5 \times (0.4 \times 0.9 \times 1 + 0.6 \times 0.7 \times 1) + 0.5 \times (0.4 \times 0.5 \times 1 + 0.6 \times 0.2 \times 1)}$$

$$= \quad 0.2211$$

---

**Q2.    [ Artificial Neural Networks ]**     <span style="border:1px solid">**5 marks**</span>

Consider the following convolutional neural network architecture.



In the first layer, we have a one-dimensional convolution with a single filter of size 3, such that $h_i = s\left( \sum_{j=1}^{3} v_j.x_{i+j-1} \right)$. The second layer is fully connected, such that $z = \sum_{i=1}^{4} w_i.h_i$. The hidden units' activation function $s(x)$ is the logistic (sigmoid) function of the form $s(x) = \frac{1}{1+e^{-x}}$. The output unit is linear (no activation function). We perform gradient descent on the loss function, $\mathcal{L} = (y - z)^2$, where $y$ is the training label for $x$.

Compute the following.

**(a)** What will be the expression for $\frac{\delta \mathcal{L}}{\delta w_i}$? **(2)**

**Solution:**

$$\mathcal{L} = (y - z)^2 \implies \frac{\delta \mathcal{L}}{\delta z} = -2(y - z)$$

$$z = \sum_{i=1}^{4} w_i.h_i \implies \frac{\delta z}{\delta w_i} = h_i$$

$$\therefore \quad \frac{\delta \mathcal{L}}{\delta w_i} \quad = \quad \frac{\delta \mathcal{L}}{\delta z} \cdot \frac{\delta z}{\delta w_i} \quad = \quad -2(y - z).h_i$$

**(b)** What will be the expression for $\frac{\delta \mathcal{L}}{\delta v_j}$? **(3)**

**Solution:**

$$h_i = s\left( \sum_{j=1}^{3} v_j.x_{i+j-1} \right) \implies \frac{\delta h_i}{\delta v_j} = h_i.\left(1 - h_i\right).x_{i+j-1} \qquad \dots \text{since, } s'(x) = s(x)\left(1 - s(x)\right)$$

$$z = \sum_{i=1}^{4} w_i.h_i \implies \frac{\delta z}{\delta h_i} = w_i$$

$$\therefore \quad \frac{\delta \mathcal{L}}{\delta v_j} = \frac{\delta \mathcal{L}}{\delta z} \cdot \frac{\delta z}{\delta v_j} = \frac{\delta \mathcal{L}}{\delta z} \cdot \sum_{i=1}^{4} \frac{\delta z}{\delta h_i} \cdot \frac{\delta h_i}{\delta v_j} \quad = \quad -2(y - z).\sum_{i=1}^{4} w_i.h_i.\left(1 - h_i\right).x_{i+j-1}$$

| Instance | True Class | $\mathbb{P}(A,\ldots,Z, M_1)$ | $\mathbb{P}(A,\ldots,Z, M_2)$ |
|---|---|---|---|
| 1 | + | 0.73 | 0.61 |
| 2 | + | 0.69 | 0.03 |
| 3 | − | 0.44 | 0.68 |
| 4 | − | 0.55 | 0.31 |
| 5 | + | 0.67 | 0.45 |
| 6 | + | 0.47 | 0.09 |
| 7 | − | 0.08 | 0.38 |
| 8 | − | 0.15 | 0.05 |
| 9 | + | 0.45 | 0.01 |
| 10 | − | 0.35 | 0.04 |

## Q3.  [ Classifier Evaluation ]  5 marks

You are asked to evaluate the performance of two classification models, $M_1$ and $M_2$. The test set you have chosen contains 26 binary attributes, labeled as $A$ through $Z$.

The above table shows the posterior probabilities obtained by applying the models to the test set. (Only the posterior probabilities for the positive class are shown). As this is a two-class problem, $\mathbb{P}(-) = 1 - \mathbb{P}(+)$ and $P(- \mid A,\ldots,Z) = 1 - \mathbb{P}(+ \mid A,\ldots,Z)$. Assume that, we are mostly interested in detecting instances from the positive class.

For both models, $M_1$ and $M_2$, suppose you choose the cutoff threshold to be $t = 0.5$. In other words, any test instances whose posterior probability is greater than $t$ will be classified as a positive example. Compute the precision, recall, and F-measure for both models at this threshold value.  (5)

**Solution:**

When $t = 0.5$, the confusion matrix for $M_1$ and $M_2$ are shown in the tables.

| $M_1$ | | Prediction | |
|---|---|---|---|
| | | + | − |
| **Actual** | + | 3 | 2 |
| | − | 1 | 4 |

| $M_2$ | | Prediction | |
|---|---|---|---|
| | | + | − |
| **Actual** | + | 1 | 4 |
| | − | 1 | 4 |

For $M_1$:

Precision $= \frac{3}{4} = 0.75.$  Recall $= \frac{3}{5} = 0.6.$  F-score $= \frac{2\times 0.75\times 0.6}{0.75+0.6} = \frac{2}{3} = 0.67.$

For $M_2$:

Precision $= \frac{1}{2} = 0.5.$  Recall $= \frac{1}{5} = 0.2.$  F-score $= \frac{2\times 0.5\times 0.2}{0.5+0.2} = \frac{2}{7} = 0.29.$

## Q4.  [ Computational Learning Theory ]  12 marks

Answer the following questions.

**(a)** Let the growth function $m_H(N)$ for some hypothesis set, $H$ ($N$ = number of training examples), be $m_H(N) = \frac{1}{2}N^2 + \frac{1}{2}N + 1$. Determine the Generalization Bound ($\Omega$) for $E_{out}$ with at least 95% probability (confidence) when the number of training examples are 1000.  (2)

**Solution:**

Given that, $1 - \delta = 0.95$, $N = 10^3$, and $m_H(N) = \frac{1}{2}N^2 + \frac{1}{2}N + 1$.

We know that, $E_{out} \leq E_{in} + \Omega$, where Generalization Bound, $\Omega = \sqrt{\frac{8}{N} . \ln\left(\frac{4.m_H(2N)}{\delta}\right)}$.

Therefore, $\Omega = \sqrt{\frac{8}{1000} \ln\left(\frac{4\times 2001001)}{0.05}\right)} = 0.389.$

**(b)** Consider the feature transform $\mathbf{z} = [L_0(x)\ \ L_1(x)\ \ L_2(x)]^T$ with Legendre polynomials and the linear model $h(x) = \mathbf{w}^T.\mathbf{z}$. For the regularized hypothesis with $\mathbf{w} = [-1\ \ +2\ \ -1]^T$, what is $h(x)$ explicitly as a function of $x$?  (2)

**Solution:**

$$L_{(x)} = 1, \quad L_1(x) = x, \quad L_2(x) = \frac{1}{2}(3x^2 - 1)$$

We may write,

$$
\begin{aligned}
h(x) &= [-1 + 2 - 1].[L_0(x) \ L_1(x) \ L_2(x)]^T \\
&= -L_0(x) + 2L_1(x) - L_2(x) \\
&= -1 + 2x - \frac{1}{2}(3x^2 - 1) \quad = -\frac{3}{2}x^2 + 2x - \frac{1}{2}
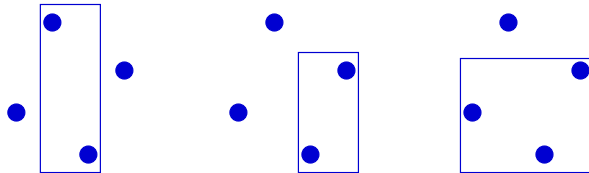\end{aligned}
$$

**(d)** What the VC-dimension of axis-aligned rectangles in a 2-dimensional plane? Derive / Prove.  **(4)**

**Solution:**

The VC-dimension of axis-aligned rectangles is 4. We prove $d_{VC} = 4$ as follows:

- There exist 4 points that can be shattered. Hence, $d_{VC} \geq 4$.
  *Proof*: It is clear that capturing just 1 point and all 4 points are both trivial, because a bounding rectangle can cover them easily. The figure below shows how we can capture a general constellation of 2 points and 3 points.



- No set of 5 points can be shattered. Hence, $d_{VC} < 5$.
  *Proof*: Suppose we have 5 points. A shattering must allow us to select all 5 points and allow us to select 4 points without the 5-th.



  Our minimum enclosing axis-aligned rectangle that allows us to select all five points is defined by only four points – one for each edge. So, it is clear that the fifth point must lie either on an edge or on the inside of the rectangle. This prevents us from selecting four points without the fifth, thereby disallowing the possibility to realize all dichotomies for general constellations of 5 points.
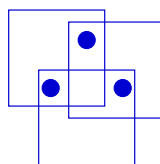
**(e)** What is the VC-dimension of axis-aligned squares in a 2-dimensional plane? Derive / Prove.  **(4)**

**Solution:**

The VC-dimension of axis-aligned squares is 3. We prove $d_{VC} = 3$ as follows:
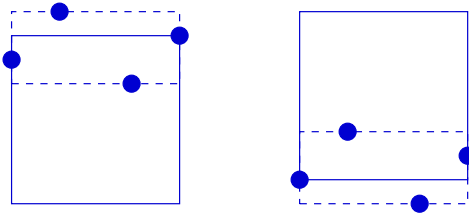
- There exist 3 points that can be shattered. Hence, $d_{VC} \geq 3$.
  *Proof*: Again, 1 point and 3 points are trivial, because a bounding square can cover them easily. The figure below shows how we can capture a general constellation of 2 points.

- No set of 4 points can be shattered. Hence, $d_{VC} < 4$.
  *Proof*: Suppose we have four points arranged such that they define a rectangle. Now, suppose we want to select two points ($A$ and $C$, in this case).



The minimum enclosing square for $A$ and $C$ must contain either $B$ or $D$ – so we cannot capture just two points with a axis-aligned square.

---

## Q5. [ Bias and Variance ]  6 marks

For $z \in \mathbb{R}$, you are trying to estimate a true function $g(z) = 2z^2$ with *linear (least-squares) regression*, where the regression function is a line $h(z) = wz$ that goes through the origin and $w \in \mathbb{R}$. Each sample point $x \in \mathbb{R}$ is drawn from the *uniform distribution* on $[-1, 1]$ and has a corresponding label $y = g(x) \in \mathbb{R}$. There is no noise in the labels. We train the model with *just one sample point*! Call it $x$, and assume $x \neq 0$. We want to apply the bias-variance decomposition to this model.

What is the <u>bias</u> and <u>variance</u> of your model $h(z)$ as a function of a test point $z \in \mathbb{R}$? Your final bias and variance both should not include an x; work out the expectations. **(3+3)**

(*Hint: start by working out the value of the least-squares weight w.*)

**Solution:**

The least-squares solution for linear regression is $w = X^{\dagger}y$, where $X$ is the $1 \times 1$ matrix $[x]$.

Hence $X^{\dagger} = (X^T X)^{-1}.X^T = \frac{x}{x^2} = \frac{1}{x}$.

Then, $w = X^{\dagger}y = \frac{1}{x}(2x^2) = 2x$, and $h(z) = wz = 2xz$.

$$\therefore \; Bias[h(z)] = \mathbb{E}[h(z)] - g(z) = \mathbb{E}[2xz] - 2z^2 = 2z.\mathbb{E}[x] - 2z^2 = 2z.\int_{-1}^{1} x\frac{1}{2}dx - 2z^2 = -2z^2$$

(However, it may seem obvious that $\mathbb{E}[x] = 0$ for a uniformly distributed $x \in [-1, 1]$; hence the integral is not required.)

$$\therefore \; Var[h(z)] = Var[2xz] = \mathbb{E}[4x^2z^2] - \mathbb{E}[2xz]^2 = \int_{-1}^{1} 4x^2z^2\frac{1}{2}dx - 4z^2.\mathbb{E}[x]^2 = \frac{2}{3}x^3z^2\Big|_{-1}^{1} - 0 = \frac{4}{3}z^2$$

<u>Alternative Approach (Variance Computation):</u>

$$Var[h(z)] = Var[2xz] = 4z^2.Var[x] = 4z^2.\mathbb{E}[(x - \mathbb{E}[x])^2] = 4z^2.\int_{-1}^{1} x^2\frac{1}{2}dx = \frac{2}{3}x^3z^2\Big|_{-1}^{1} - 0 = \frac{4}{3}z^2$$

---

## Q6. [ Unsupervised Learning ]  16 marks

Suppose, six points ($P_1$, $P_2$, $P_3$, $P_4$, $P_5$ and $P_6$) are provided in a 2-dimensional plane. The Euclidean distance between a pair of these points are provided in the table below.

If you use *Hierarchical Agglomerative Clustering* technique to form the single-link dendrogram, initially each point will form separate clusters, denoted as, $\{P_1\}$, $\{P_2\}$, $\{P_3\}$, $\{P_4\}$, $\{P_5\}$ and $\{P_6\}$. Then, at the first (bottom-most grouping) phase, the algorithm selects $\{P_1\}$ and $\{P_2\}$ clusters to merge and

|      | $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_5$ | $P_6$ |
|------|-------|-------|-------|-------|-------|-------|
| $P_1$ | 0.00 |      |      |      |      |      |
| $P_2$ | 0.12 | 0.00 |      |      |      |      |
| $P_3$ | 0.51 | 0.25 | 0.00 |      |      |      |
| $P_4$ | 0.84 | 0.16 | 0.14 | 0.00 |      |      |
| $P_5$ | 0.28 | 0.77 | 0.70 | 0.45 | 0.00 |      |
| $P_6$ | 0.34 | 0.61 | 0.93 | 0.20 | 0.67 | 0.00 |

form new cluster $\{P_1, P_2\}$, as the distance considered for grouping here was, $dist(P_1, P_2) = 0.12$ (the minimum among all pairs), for both *single-linkage* and *complete-linkage* variants.

Now, you need to complete the rest of the phases mentioning the next new cluster formed and the distance considered that time for both *single-linkage* and *complete-linkage* variations in the following.

**(a)** Using ***Single** Linkage Hierarchical Agglomerative Clustering* technique to form the single-link dendrogram, complete the remaining phases (missing entries) in the following table. **(4)**
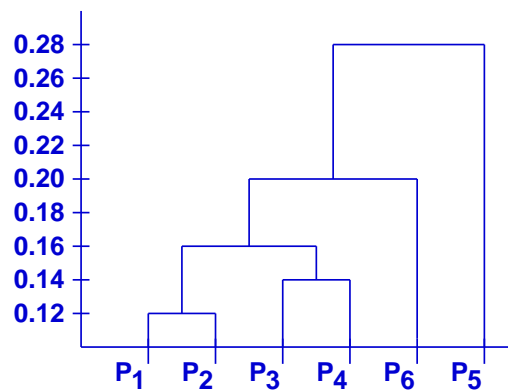
| Phase $\rightarrow$ | 1st | 2nd | 3rd | 4th | 5th |
|---------------------|-----|-----|-----|-----|-----|
| **New Cluster Formed** | $\{P_1, P_2\}$ |  |  |  |  |
| **Distance Considered** | 0.12 |  |  |  |  |

Show final result of hierarchical clustering with *single linkage* by drawing a dendrogram. **(1)**

**Solution:**

| Phase $\rightarrow$ | 1st | 2nd | 3rd | 4th | 5th |
|---------------------|-----|-----|-----|-----|-----|
| **New Cluster Formed** | $\{P_1, P_2\}$ | $\{P_3, P_4\}$ | $\{P_1, P_2, P_3, P_4\}$ | $\{P_1, P_2, P_3, P_4, P_6\}$ | $\{P_1, P_2, P_3, P_4, P_5, P_6\}$ |
| **Distance Considered** | 0.12 | 0.14 | 0.16 | 0.20 | 0.28 |

Note: At every phase, the clustering is here based on choosing the minimum among the *minimum* distances between a pair of existing clusters.



**Dendrogram for Single–Linkage Hierarchical Clustering Process**

**(b)** Using ***Complete** Linkage Hierarchical Agglomerative Clustering* technique to form the complete-link dendrogram, complete the remaining phases (missing entries) in the above table. **(4)**

| Phase $\rightarrow$ | 1st | 2nd | 3rd | 4th | 5th |
|---------------------|-----|-----|-----|-----|-----|
| **New Cluster Formed** | $\{P_1, P_2\}$ |  |  |  |  |
| **Distance Considered** | 0.12 |  |  |  |  |

Show final result of hierarchical clustering with *complete linkage* by drawing a dendrogram. **(1)**

| Phase → | 1st | 2nd | 3rd | 4th | 5th |
|---|---|---|---|---|---|
| **New Cluster Formed** | $\{P_1, P_2\}$ | $\{P_3, P_4\}$ | $\{P_1, P_2, P_6\}$ | $\{P_3, P_4, P_5\}$ | $\{P_1, P_2, P_3, P_4, P_5, P_6\}$ |
| **Distance Considered** | 0.12 | 0.14 | 0.61 | 0.70 | 0.93 |

**Solution:**

Note: At every phase, the clustering is here based on choosing the minimum among the *maximum* distances between a pair of existing clusters.



**Dendrogram for Complete–Linkage
Hierarchical Clustering Process**

(c) Suppose, for both the above variants (single and complete linkage) of hierarchical clustering, we stop after 4th phase. Compute the average silhouette coefficient (SC) of the overall clustering for both these cases. **(3+3)**

**Solution:**

Let $a$ indicate the average distance of a point to other points within its cluster, and $b$ indicate the minimum of the average distance of a point to points in another cluster.

Single-Linkage Hierarchical Clustering Case:

Here after 4th phase, the two clusters formed are: $\{P_1, P_2, P_3, P_4, P_6\}$ and $\{P_5\}$.

The silhouette coefficient (SC) for each of the points are computed as:

$$P_1 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.12 + 0.51 + 0.84 + 0.34}{4}\right)}{0.28} = -0.62$$

$$P_2 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.12 + 0.25 + 0.16 + 0.61}{4}\right)}{0.77} = 0.63$$

$$P_3 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.51 + 0.25 + 0.14 + 0.93}{4}\right)}{0.70} = 0.35$$

$$P_4 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.84 + 0.16 + 0.14 + 0.20}{4}\right)}{0.45} = 0.26$$

$$P_5 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{0}{\left(\frac{0.28 + 0.77 + 0.70 + 0.45 + 0.67}{5}\right)} = 1$$

$$P_6 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.34 + 0.61 + 0.93 + 0.20}{4}\right)}{0.67} = 0.22$$

$$\text{Cluster-1} \quad : \quad \left(\{P_1, P_2, P_3, P_4, P_6\}\right) \quad \text{Average-SC} = \frac{-0.62 + 0.63 + 0.35 + 0.26 + 0.22}{5} = 0.17$$

$$\text{Cluster-2} \quad : \quad \left(\{P_5\}\right) \quad \text{Average-SC} = \frac{1}{1} = 1$$

$$\text{Overall} \quad : \quad \text{Average-SC} = \frac{0.17 + 1}{2} = 0.585$$

Complete-Linkage Hierarchical Clustering Case:

Here after 4th phase, the two clusters formed are: $\{P_1, P_2, P_6\}$ and $\{P_3, P_4, P_5\}$.

The silhouette coefficient (SC) for each of the points are computed as:

$$P_1 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.12+0.34}{2}\right)}{\left(\frac{0.51+0.84+0.28}{3}\right)} = 0.58$$

$$P_2 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.12+0.61}{2}\right)}{\left(\frac{0.25+0.16+0.77}{3}\right)} = 0.07$$

$$P_3 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.14+0.70}{2}\right)}{\left(\frac{0.51+0.25+0.93}{3}\right)} = 0.25$$

$$P_4 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.14+0.45}{2}\right)}{\left(\frac{0.84+0.16+0.20}{3}\right)} = 0.26$$

$$P_5 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.70+0.45}{2}\right)}{\left(\frac{0.28+0.77+0.67}{3}\right)} = -0.003$$

$$P_6 \quad : \quad SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{0.34+0.61}{2}\right)}{\left(\frac{0.93+0.20+0.67}{3}\right)} = 0.21$$

$$\text{Cluster-1} \quad : \quad \left(\{P_1, P_2, P_6\}\right) \quad \text{Average-SC} \; = \frac{0.58+0.07+0.21}{3} = 0.29$$

$$\text{Cluster-2} \quad : \quad \left(\{P_3, P_4, P_5\}\right) \quad \text{Average-SC} \; = \frac{0.25+0.26-0.003}{3} = 0.17$$

$$\text{Overall} \quad : \quad \text{Average-SC} \; = \frac{0.29+0.17}{2} = 0.23$$

---

**Q7.** **[ Ensemble Learning ]** $\boxed{\textbf{10 marks}}$

In this problem, we study how boosting algorithm performs on a very simple classification problem. We are given with four training points, $P_1$, $P_2$, $P_3$, $P_4$, in a 1-dimensional line ($x$-valued) having their respective values as $x = 1$, $x = 2$, $x = 3$, $x = 4$ and their corresponding 2-class ($+/-$) labels as $-$, $+$, $-$, $+$, respectively.

We shall use decision stumps as our weak learner / hypothesis. Decision stump classifier chooses a constant value $c$ and classifies all points where $x \geq c$ as one class and other points where $x < c$ as the other class. In our given example, let us chose one such decision stump as follows: $x \geq 3$ region is classified as '+' zone and $x < 3$ region is classified as '−' zone.

Answer the following questions.

**(a)** What is the initial weight assigned to each data point? **(1)**

**Solution:**
Since $Weight(P_1) = Weight(P_2) = Weight(P_3) = Weight(P_4)$ and
$$Weight(P_1) + Weight(P_2) + Weight(P_3) + Weight(P_4) = 1,$$
$\therefore Weight(P_1) = Weight(P_2) = Weight(P_3) = Weight(P_4) = \frac{1}{4}.$

**(b)** How many different decision stumps are possible for the data points given? **(1)**

**Solution:**
(4 separators/stumps) x (2 different class organizations for each)
$\quad$ = 8 different decision stumps are possible.

**(c)** Which data point(s) will have weights increased after the boosting process as per the decision stump considered in the problem? **(1)**

**Solution:**

Since the given decision stump will misclassify $P_2$ and $P_3$ as '+' and '−', respectively, so only these two data points ($P_2$ and $P_3$) *may* have their weights increased after boosting.

**(d)** What will be weights of all the data points after boosting is performed? Show your approach. **(4)**

**Solution:**

$$\varepsilon_t = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

$$\alpha_t = \frac{1}{2}\ln\left(\frac{1-\frac{1}{2}}{\frac{1}{2}}\right) = \frac{1}{2}\ln 1 = 0$$

Also, normalization factor, $Z = 2\sqrt{\varepsilon_t(1-\varepsilon_t)} = 2\sqrt{\frac{1}{2}\times\frac{1}{2}} = 1$.

So, for correctly classified data-points ($P_1$ and $P_4$), the weight will become:

$$Weight(P_1) = Weight(P_4) = \frac{W_t \cdot e^{-\alpha_t}}{Z} = \frac{1}{4}$$

So, for wrongly classified data-points ($P_2$ and $P_3$), the weight will become:

$$Weight(P_2) = Weight(P_3) = \frac{W_t \cdot e^{\alpha_t}}{Z} = \frac{1}{4}$$

**(e)** Indicate whether the following statements are <u>true</u> / <u>false</u>. Give a brief justification.

(i) We cannot perfectly classify all the training examples given in this problem by only applying boosting algorithm (AdaBoost). **(1.5)**

(ii) The training error of boosting classifier (combination of all the weak classifier) monotonically decreases as the number of iterations in the boosting algorithm increases. **(1.5)**

**Solution:**

(i) <u>True</u>, since the data is not linearly separable.

(ii) <u>False</u>, since boosting minimizes loss function: $\sum\limits_{i=1}^{m} e^{-y_i \cdot f(x_i)}$,

which does not necessary mean that training error monotonically decrease.

---

**Q8.** **[ Principal Component Analysis ]** **5 marks**

Given the $(x,y)$-coordinates of four data points in two-dimensional space: $(4,1)$, $(2,3)$, $(5,4)$ and $(1,0)$, calculate the first principal component. Show your calculations in details. **(5)**

**Solution:**

The mean of the given data points is: $\left(\frac{4+2+5+1}{4}, \frac{1+3+4+0}{4}\right) = (3,2)$.

The covariance matix can be constructed as:

$$CoVar(x,x) = Var(x) = \frac{[(4-3)^2 + (2-3)^2 + (5-3)^2 + (1-3)^2]}{4} = \frac{5}{2}$$

$$CoVar(x,y) = CoVar(y,x) = \frac{[(4-3)\times(1-2) + (2-3)\times(3-2) + (5-3)\times(4-2) + (1-3)\times(0-2)]}{4} = \frac{3}{2}$$

$$CoVar(y,y) = Var(y) = \frac{[(1-2)^2 + (3-2)^2 + (4-2)^2 + (0-2)^2]}{4} = \frac{5}{2}$$

$$\therefore CoVar = \begin{bmatrix} CoVar(x,x) & CoVar(x,y) \\ CoVar(y,x) & CoVar(y,y) \end{bmatrix} = \begin{bmatrix} \frac{5}{2} & \frac{3}{2} \\ \frac{3}{2} & \frac{5}{2} \end{bmatrix}.$$

To compute eigenvalues, we make $\left| CoVar - \lambda I \right| = 0$, which gives:

$$(\frac{5}{2} - \lambda)^2 - \frac{9}{4} = 0 \quad \implies \quad \lambda = 4, 1$$

The corresponding eigenvector with respect to the highest eigenvalue is the principal component, which is computed as,

$$\begin{bmatrix} \frac{5}{2} & \frac{3}{2} \\ \frac{3}{2} & \frac{5}{2} \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = 4. \begin{bmatrix} x \\ y \end{bmatrix} \quad \implies \quad [x, y]^T = \left[ \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right]^T.$$

Alternative Approach:

Since the mean of the given data points, $X = \begin{bmatrix} 4 & 1 \\ 2 & 3 \\ 5 & 4 \\ 1 & 0 \end{bmatrix}$ is $(3, 2)$, we can center the given points with

respect to mean as, $\hat{X} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \\ 2 & 2 \\ -2 & -2 \end{bmatrix}$.

Now, $\hat{X}^T . \hat{X} = \begin{bmatrix} 10 & 6 \\ 6 & 10 \end{bmatrix}$.

(Divide by 4 if you want the sample covariance matrix, but we do not care about the magnitude here.)

Its eigenvectors are $\left[ \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right]^T$ for eigenvalue 16 and $\left[ \frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \right]^T$ for eigenvalue 4. The former eigenvector is chosen to be the principal component (as the corresponding eigenvalue is the highest).

---

**Q9.** **[ Kernel Functions ]** **5 marks**

Answer the following.

**(a)** Let $k_1$ and $k_2$ be (valid) kernels; that is, $k_1(\mathbf{x}, \mathbf{y}) = \Phi_1(\mathbf{x})^T . \Phi_1(\mathbf{y})$ and $k_2(\mathbf{x}, \mathbf{y}) = \Phi_2(\mathbf{x})^T . \Phi_2(\mathbf{y})$. Show that $k = k_1 + k_2$ is a valid kernel by explicitly constructing a corresponding feature mapping $\Phi(\mathbf{z})$. **(2)**

**Solution:**

Note that,

$k(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y}) + k_2(\mathbf{x}, \mathbf{y}) = \Phi_1(\mathbf{x})^T . \Phi_1(\mathbf{y}) + \Phi_2(\mathbf{x})^T . \Phi_2(\mathbf{y}) = \left[ \Phi_1(\mathbf{x})^T \ \ \Phi_2(\mathbf{x})^T \right] . \left[ \Phi_1(\mathbf{y}) \ \ \Phi_2(\mathbf{y}) \right]^T$

If we let $\Phi(\mathbf{z}) = \left[ \Phi_1(\mathbf{x}) \ \ \Phi_2(\mathbf{x}) \right]^T$, then we have $k(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{z})^T . \Phi(\mathbf{z})$.

Therefore, $k = k_1 + k_2$ is a valid kernel.

**(b)** The polynomial kernel is defined to be

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T . \mathbf{y} + c)^d,$$

where $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, and $c \geq 0$. When we take $d = 2$, this kernel is called the *quadratic kernel*. Find the feature mapping $\Phi(\mathbf{z})$ that corresponds to the quadratic kernel. **(3)**

**Solution:**

First we expand the dot product inside, and square the entire sum. We will get a sum of the squares of the components and a sum of the cross products.

$$(\mathbf{x}^T . \mathbf{y} + c)^d = \left( c + \sum_{i=1}^{n} x_i y_i \right)^2$$

$$= c^2 + \sum_{i=1}^{n} x_i^2 y_i^2 + \sum_{i=2}^{n} \sum_{j=1}^{i-1} 2 x_i y_i x_j y_j + \sum_{i=1}^{n} 2 x_i y_i c \quad = \Phi(\mathbf{x})^T . \Phi(\mathbf{y})$$

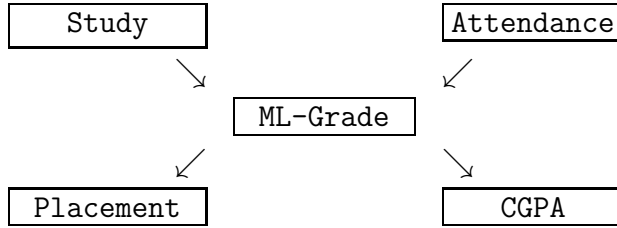Pulling this sum into a dot product of x components and y components, we have,

$$\Phi(\mathbf{x}) \;=\; \left[c, x_1^2, \ldots, x_n^2, \sqrt{2}x_1x_2, \ldots, \sqrt{2}x_1x_n, \sqrt{2}x_2x_3, \ldots, \sqrt{2}x_2x_n, \ldots, \sqrt{2}x_{n-1}x_n, \sqrt{2c}x_1, \ldots, \sqrt{2c}x_n\right]$$

$$\Phi(\mathbf{y}) \;=\; \left[c, y_1^2, \ldots, y_n^2, \sqrt{2}y_1y_2, \ldots, \sqrt{2}y_1y_n, \sqrt{2}y_2y_3, \ldots, \sqrt{2}y_2y_n, \ldots, \sqrt{2}y_{n-1}y_n, \sqrt{2c}y_1, \ldots, \sqrt{2c}y_n\right]$$

In this feature mapping, we have $c$, the squared components of the vector, $\sqrt{2}$ multiplied by all of the cross terms, and $\sqrt{2c}$ multiplied by all of the components.

---

**Q10.** **[ Expectation-Maximization Algorithm]** <span>15 marks</span>

Consider the Bayes Network structure shown below. From the figure below, we abbreviate as follows:
**S** = Study well, **A** = high Attendance, **G** = good ML-Grade, **P** = better Placement, **C** = high CGPA.



We are given the following $K = 8$ training examples as shown below, where only two examples contain unobserved values (marked with ?), namely, $p_7$ and $c_8$. We like to simulate a few steps of the simplified EM algorithm by hand.

| K | S | A | G | P | C |
|---|---|---|---|---|---|
| $k=1$ | 1 | 0 | 1 | 1 | 1 |
| $k=2$ | 0 | 1 | 1 | 1 | 0 |
| $k=3$ | 1 | 1 | 1 | 1 | 1 |
| $k=4$ | 0 | 0 | 0 | 0 | 1 |
| $k=5$ | 0 | 0 | 0 | 1 | 0 |
| $k=6$ | 0 | 0 | 0 | 0 | 0 |
| $k=7$ | 1 | 1 | 1 | ? | 1 |
| $k=8$ | 1 | 1 | 1 | 1 | ? |

Notation: Here, $s_k$, $a_k$, $g_k$, $p_k$, and $c_k$ indicate the values of **S**, **A**, **G**, **P**, and **C**, respectively, as seen in the $k$-th example/row. For example, $s_1 = 1$, $a_1 = 0$, $g_1 = 1$, $p_1 = 1$, and $c_1 = 1$.

Answer the following questions:

**(a)** Given that *all variables are Boolean*, how many basic parameters we need to estimate for the given Bayes Network?

For example, one parameter will be $\theta(g \mid 11)$, which stands for $\mathbb{P}(G = 1 \mid S = 1, A = 1)$. **(2)**

**Solution:**

We need to estimate 10 parameters, which are given as follows:

$$\theta(s) \;=\; \mathbb{P}(S=1) \qquad\qquad \theta(a) \;=\; \mathbb{P}(A=1)$$
$$\theta(g \mid 00) \;=\; \mathbb{P}(G=1 \mid S=0,\, A=0) \qquad \theta(g \mid 01) \;=\; \mathbb{P}(G=1 \mid S=0,\, A=1)$$
$$\theta(g \mid 10) \;=\; \mathbb{P}(G=1 \mid S=1,\, A=0) \qquad \theta(g \mid 11) \;=\; \mathbb{P}(G=1 \mid S=1,\, A=1)$$
$$\theta(p \mid 0) \;=\; \mathbb{P}(P=1 \mid G=0) \qquad\qquad \theta(p \mid 1) \;=\; \mathbb{P}(P=1 \mid G=1)$$
$$\theta(c \mid 0) \;=\; \mathbb{P}(C=1 \mid G=0) \qquad\qquad \theta(c \mid 1) \;=\; \mathbb{P}(C=1 \mid G=1)$$

**(b)** Now, we like to simulate the first E-step of the EM algorithm. Before we start, we initialize all the parameters as 0.5, and then proceed to execute the E-step. What are the following expectation values that will get calculated in this E-step? In particular, calculate the following: **(2+2)**

- $\mathbb{E}(p_7 = 1 \mid s_7, a_7, g_7, c_7, \boldsymbol{\theta}) =?$
- $\mathbb{E}(c_8 = 1 \mid s_8, a_8, g_8, p_8, \boldsymbol{\theta}) =?$

*(Note that, only two examples (k=7 and k=8) contains unobserved variables, where $p_7 =?$, but $s_7 = a_7 = g_7 = c_7 = 1$; and $c_8 =?$, but $s_8 = a_8 = g_8 = p_8 = 1$, respectively.)*

**Solution:**

$$\therefore \quad \mathbb{E}(p_7 = 1 \mid s_7, a_7, g_7, c_7, \boldsymbol{\theta})$$

$$= \frac{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \boldsymbol{\theta})}{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \boldsymbol{\theta}) + \mathbb{P}(p_7 = 0, s_7, a_7, g_7, c_7 \mid \boldsymbol{\theta})}$$

$$= \frac{\boldsymbol{\theta}(p_7 = 1 \mid g_7).\boldsymbol{\theta}(g_7 \mid s_7, a_7).\boldsymbol{\theta}(s_7).\boldsymbol{\theta}(a_7)}{\boldsymbol{\theta}(p_7 = 1 \mid g_7).\boldsymbol{\theta}(g_7 \mid s_7, a_7).\boldsymbol{\theta}(s_7).\boldsymbol{\theta}(a_7) + \boldsymbol{\theta}(p_7 = 0 \mid g_7).\boldsymbol{\theta}(g_7 \mid s_7, a_7).\boldsymbol{\theta}(s_7).\boldsymbol{\theta}(a_7)}$$

$$= \frac{0.5 \times 0.5 \times 0.5 \times 0.5}{2 \times 0.5 \times 0.5 \times 0.5 \times 0.5} \quad = \quad 0.5 \quad = \quad \mathbb{E}(p_7 = 1 \mid g_7 = 1, \boldsymbol{\theta}(p \mid 1))$$

$$\therefore \quad \mathbb{E}(c_8 = 1 \mid s_8, a_8, g_8, p_8, \boldsymbol{\theta})$$

$$= \frac{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \boldsymbol{\theta})}{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \boldsymbol{\theta}) + \mathbb{P}(c_8 = 0, s_8, a_8, g_8, p_8 \mid \boldsymbol{\theta})}$$

$$= \frac{\boldsymbol{\theta}(c_8 = 1 \mid g_8).\boldsymbol{\theta}(g_8 \mid s_8, a_8).\boldsymbol{\theta}(s_8).\boldsymbol{\theta}(a_8)}{\boldsymbol{\theta}(c_8 = 1 \mid g_8).\boldsymbol{\theta}(g_8 \mid s_8, a_8).\boldsymbol{\theta}(s_8).\boldsymbol{\theta}(a_8) + \boldsymbol{\theta}(c_8 = 0 \mid g_8).\boldsymbol{\theta}(g_8 \mid s_8, a_8).\boldsymbol{\theta}(s_8).\boldsymbol{\theta}(a_8)}$$

$$= \frac{0.5 \times 0.5 \times 0.5 \times 0.5}{2 \times 0.5 \times 0.5 \times 0.5 \times 0.5} \quad = \quad 0.5 \quad = \quad \mathbb{E}(c_8 = 1 \mid g_7 = 1, \boldsymbol{\theta}(c \mid 1))$$

**(c)** Now, we like to simulate the first M-step of the EM algorithm. What will be the estimated values of all the model parameters (which you identified in part (a)) that we obtain in this M-step? **(5)**

*(Note that, we use the expected count only when the variable is unobserved in an example)*

**Solution:**

10 parameters will get the updated values as follows:

$$\boldsymbol{\theta}(s) \quad = \quad \mathbb{P}(S = 1) = \frac{\#\{S = 1\}}{\#K} = \frac{4}{8} = 0.5$$

$$\boldsymbol{\theta}(a) \quad = \quad \mathbb{P}(A = 1) = \frac{\#\{A = 1\}}{\#K} = \frac{4}{8} = 0.5$$

$$\boldsymbol{\theta}(g \mid 00) \quad = \quad \mathbb{P}(G = 1 \mid S = 0, A = 0) = \frac{\#\{G = 1, S = 0, A = 0\}}{\#\{S = 0, A = 0\}} = \frac{0}{3} = 0.0$$

$$\boldsymbol{\theta}(g \mid 01) \quad = \quad \mathbb{P}(G = 1 \mid S = 0, A = 1) = \frac{\#\{G = 1, S = 0, A = 1\}}{\#\{S = 0, A = 1\}} = \frac{1}{1} = 1.0$$

$$\boldsymbol{\theta}(g \mid 10) \quad = \quad \mathbb{P}(G = 1 \mid S = 1, A = 0) = \frac{\#\{G = 1, S = 1, A = 0\}}{\#\{S = 1, A = 0\}} = \frac{1}{1} = 1.0$$

$$\boldsymbol{\theta}(g \mid 11) \quad = \quad \mathbb{P}(G = 1 \mid S = 1, A = 1) = \frac{\#\{G = 1, S = 1, A = 1\}}{\#\{S = 1, A = 1\}} = \frac{3}{3} = 1.0$$

$$\boldsymbol{\theta}(p \mid 0) \quad = \quad \mathbb{P}(P = 1 \mid G = 0) = \frac{\#\{G = 0\}.\mathbb{E}[P = 1]}{\#\{G = 0\}} = \frac{(1 \times 1.0 + 2 \times 0.0)}{3} = 0.33$$

$$\boldsymbol{\theta}(p \mid 1) \quad = \quad \mathbb{P}(P = 1 \mid G = 1) = \frac{\#\{G = 1\}.\mathbb{E}[P = 1]}{\#\{G = 1\}} = \frac{(4 \times 1.0 + 1 \times 0.5)}{5} = 0.9$$

$$\boldsymbol{\theta}(c \mid 0) \quad = \quad \mathbb{P}(C = 1 \mid G = 0) = \frac{\#\{G = 0\}.\mathbb{E}[C = 1]}{\#\{G = 0\}} = \frac{(1 \times 1.0 + 2 \times 0.0)}{3} = 0.33$$

$$\boldsymbol{\theta}(c \mid 1) \quad = \quad \mathbb{P}(C = 1 \mid G = 1) = \frac{\#\{S = 1\}.\mathbb{E}[C = 1]}{\#\{G = 1\}} = \frac{(3 \times 1.0 + 1 \times 0.0 + 1 \times 0.0)}{5} = 0.7$$

**(d)** Last, let us (again) simulate the second E-step of the EM algorithm. What a re the following expectation values that will get calculated in this E-step? In particular, calculate the following: **(2+2)**

- $\mathbb{E}(p_7 = 1 \mid s_7, a_7, g_7, c_7, \theta) = ?$
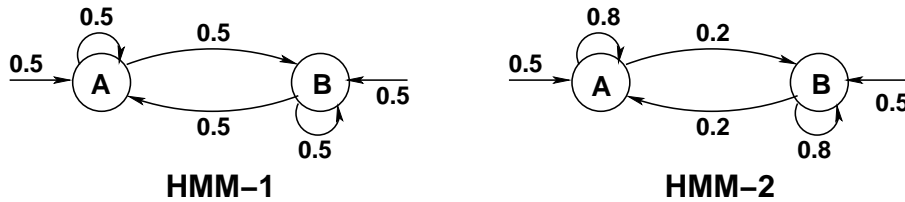- $\mathbb{E}(c_8 = 1 \mid s_8, a_8, g_8, p_8, \theta) = ?$

**Solution:**

$$\therefore \quad \mathbb{E}(p_7 = 1 \mid s_7, a_7, g_7, c_7, \theta)$$

$$= \frac{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \theta)}{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \theta) + \mathbb{P}(p_7 = 0, s_7, a_7, g_7, c_7 \mid \theta)}$$

$$= \frac{\theta(p_7 = 1 \mid g_7).\theta(g_7 \mid s_7, a_7).\theta(s_7).\theta(a_7)}{\theta(p_7 = 1 \mid g_7).\theta(g_7 \mid s_7, a_7).\theta(s_7).\theta(a_7) + \theta(p_7 = 0 \mid g_7).\theta(g_7 \mid s_7, a_7).\theta(s_7).\theta(a_7)}$$

$$= \frac{0.9 \times 1.0 \times 0.5 \times 0.5}{1.0 \times 1.0 \times 0.5 \times 0.5} \quad = \quad 0.9 \quad = \quad \mathbb{E}(p_7 = 1 \mid g_7 = 1, \theta(p \mid 1))$$

$$\therefore \quad \mathbb{E}(c_8 = 1 \mid s_8, a_8, g_8, p_8, \theta)$$

$$= \frac{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \theta)}{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \theta) + \mathbb{P}(c_8 = 0, s_8, a_8, g_8, p_8 \mid \theta)}$$

$$= \frac{\theta(c_8 = 1 \mid g_8).\theta(g_8 \mid s_8, a_8).\theta(s_8).\theta(a_8)}{\theta(c_8 = 1 \mid g_8).\theta(g_8 \mid s_8, a_8).\theta(s_8).\theta(a_8) + \theta(c_8 = 0 \mid g_8).\theta(g_8 \mid s_8, a_8).\theta(s_8).\theta(a_8)}$$

$$= \frac{0.7 \times 1.0 \times 0.5 \times 0.5}{1.0 \times 1.0 \times 0.5 \times 0.5} \quad = \quad 0.7 \quad = \quad \mathbb{E}(c_8 = 1 \mid g_7 = 1, \theta(c \mid 1))$$

---

**Q11.** **[ Hidden Markov Models ]** $\boxed{\text{10 marks}}$

The following figure above presents two HMMs. States are represented by circles and transitions by directed edges. In both, emissions are deterministic and listed inside the states (either *A* or *B*).



**HMM–1**          **HMM–2**

Transition probabilities and starting probabilities are listed next to the relevant edges. For example, in HMM-1 we have a probability of 0.5 to start with the state that emits *A* and a probability of 0.5 to transition to the state that emits *B* if we are now in the state that emits *A*.

Notation: In the questions below, $O_{100} = A$ means that the 100-th symbol emitted by an HMM is *A*.

Answer the following.

**(a)** Calculate $\mathbb{P}(O_{100} = A, \ O_{101} = A, \ O_{102} = A)$ for HMM-1 and HMM-2, respectively. **(2+2)**

**Solution:**

Note that, for HMM-1,

$$\mathbb{P}(O_{100} = A, O_{101} = A, O_{102} = A)$$
$$= \mathbb{P}(O_{100} = A, O_{101} = A, O_{102} = A, S_{100} = A, S_{101} = A, S_{102} = A),$$

since if we are not always in state *A* we will not be able to emit *A*.

Given the Markov property, this can be written as:

$$\mathbb{P}(O_{100}=A, O_{101}=A, O_{102}=A, S_{100}=A, S_{101}=A, S_{102}=A)$$
$$= \mathbb{P}(O_{100}=A \mid S_{100}=A).\mathbb{P}(S_{100}=A).$$
$$\mathbb{P}(O_{101}=A \mid S_{101}=A).\mathbb{P}(S_{101}=A \mid S_{100}=A).$$
$$\mathbb{P}(O_{102}=A \mid S_{102}=A).\mathbb{P}(S_{102}=A \mid S_{101}=A)$$
$$= 1 \times 0.5 \times 1 \times 0.5 \times 1 \times 0.5 \quad = 0.125$$

[ Here, since the model is fully symmetric, $\mathbb{P}(S_{100}=A)=0.5$. ]

Folowing similar lines in HMM-2, $\mathbb{P}(O_{100}=A, O_{101}=A, O_{102}=A)$ can be written as:

$$\mathbb{P}(O_{100}=A, O_{101}=A, O_{102}=A, S_{100}=A, S_{101}=A, S_{102}=A)$$
$$= \mathbb{P}(O_{100}=A \mid S_{100}=A).\mathbb{P}(S_{100}=A).$$
$$\mathbb{P}(O_{101}=A \mid S_{101}=A).\mathbb{P}(S_{101}=A \mid S_{100}=A).$$
$$\mathbb{P}(O_{102}=A \mid S_{102}=A).\mathbb{P}(S_{102}=A \mid S_{101}=A)$$
$$= 1 \times 0.5 \times 1 \times 0.8 \times 1 \times 0.8 \quad = 0.32$$

**(b)** Calculate $\mathbb{P}(O_{100}=A,\ O_{101}=B,\ O_{102}=A,\ O_{103}=B)$ for HMM-1 and HMM-2. respectively. **(2+2)**

**Solution:**

For HMM-1, this can be expressed as:

$$\mathbb{P}(O_{100}=A,\ O_{101}=B,\ O_{102}=A,\ O_{103}=B)$$
$$= \mathbb{P}(O_{100}=A, O_{101}=B, O_{102}=A, O_{103}=B, S_{100}=A, S_{101}=B, S_{102}=A, S_{103}=B)$$
$$= \mathbb{P}(O_{100}=A \mid S_{100}=A).\mathbb{P}(S_{100}=A).$$
$$\mathbb{P}(O_{101}=B \mid S_{101}=B).\mathbb{P}(S_{101}=B \mid S_{100}=A).$$
$$\mathbb{P}(O_{102}=A \mid S_{102}=A).\mathbb{P}(S_{102}=A \mid S_{101}=B).$$
$$\mathbb{P}(O_{103}=B \mid S_{103}=B).\mathbb{P}(S_{103}=B \mid S_{102}=A)$$
$$= 1 \times 0.5 \times 1 \times 0.5 \times 1 \times 0.5 \times 1 \times 0.5 \quad = 0.0625$$

Similarly, for HMM-2,

$$\mathbb{P}(O_{100}=A,\ O_{101}=B,\ O_{102}=A,\ O_{103}=B)$$
$$= \mathbb{P}(O_{100}=A, O_{101}=B, O_{102}=A, O_{103}=B, S_{100}=A, S_{101}=B, S_{102}=A, S_{103}=B)$$
$$= \mathbb{P}(O_{100}=A \mid S_{100}=A).\mathbb{P}(S_{100}=A).$$
$$\mathbb{P}(O_{101}=B \mid S_{101}=B).\mathbb{P}(S_{101}=B \mid S_{100}=A).$$
$$\mathbb{P}(O_{102}=A \mid S_{102}=A).\mathbb{P}(S_{102}=A \mid S_{101}=B).$$
$$\mathbb{P}(O_{103}=B \mid S_{103}=B).\mathbb{P}(S_{103}=B \mid S_{102}=A)$$
$$= 1 \times 0.5 \times 1 \times 0.2 \times 1 \times 0.2 \times 1 \times 0.2 \quad = 0.004$$

**(c)** Assume you are told that a casino has been using one of the two HMMs to generate streams of letters. You are also told that among the first 1000 letters emitted, 500 are $A$s and 500 are $B$s. Can you tell which of the HMMs is being used by this casino? Explain. **(2)**

**Solution:**

While we saw in the previous part (b) that it is much more less likely to switch between $A$ and $B$ in HMM-2, this is only true if we switch at every step. However, when aggregating over 1000 steps, since the two HMMs are both symmetric, *both are likely to generate the same number of As and Bs*. So, the casino may have been using any of these HMMs.

--- **— END —** ---