



INDIAN INSTITUTE OF TECHNOLOGY
KHARAGPUR

Stamp / Signature of the Invigilator

EXAMINATION (End Semester)

SEMESTER (Spring 2024-2025)

Roll Number

Section

Name

Subject Number

C

S

6

0

0

2

0

Subject Name

FOUNDATIONS OF ALGORITHM DESIGN AND MACHINE LEARNING

Department / Center of the Student

Additional sheets

Important Instructions and Guidelines for Students

1. You must occupy your seat as per the Examination Schedule/Sitting Plan.
2. Do not keep mobile phones or any similar electronic gadgets with you even in the switched off mode.
3. Loose papers, class notes, books or any such materials must not be in your possession, even if they are irrelevant to the subject you are taking examination.
4. Data book, codes, graph papers, relevant standard tables/charts or any other materials are allowed only when instructed by the paper-setter.
5. Use of instrument box, pencil box and non-programmable calculator is allowed during the examination. However, exchange of these items or any other papers (including question papers) is not permitted.
6. Write on both sides of the answer script and do not tear off any page. **Use last page(s) of the answer script for rough work.** Report to the invigilator if the answer script has torn or distorted page(s).
7. It is your responsibility to ensure that you have signed the Attendance Sheet. Keep your Admit Card/Identity Card on the desk for checking by the invigilator.
8. You may leave the examination hall for wash room or for drinking water for a very short period. Record your absence from the Examination Hall in the register provided. Smoking and the consumption of any kind of beverages are strictly prohibited inside the Examination Hall.
9. Do not leave the Examination Hall without submitting your answer script to the invigilator. **In any case, you are not allowed to take away the answer script with you.** After the completion of the examination, do not leave the seat until the invigilators collect all the answer scripts.
10. During the examination, either inside or outside the Examination Hall, gathering information from any kind of sources or exchanging information with others or any such attempt will be treated as 'unfair means'. Do not adopt unfair means and do not indulge in unseemly behavior.

Violation of any of the above instructions may lead to severe punishment.

Signature of the Student

To be filled in by the examiner

Question Number

1

2

3

4

5

6

7

8

9

10

Total

Marks Obtained

Marks obtained (in words)

Signature of the Examiner

Signature of the Scrutineer

Indian Institute of Technology Kharagpur
Department of Computer Science and Engineering

Foundations of Algorithm Design and Machine Learning (CS60020)

Spring 2024-2025

April-2025

End-Semester Examination

Maximum Marks: 100

Instructions:

- Write your answers in the question paper itself. Be brief and precise. Answer all questions.
 - Write the answers only in the respective spaces provided. The last two blank pages may be used for rough work or leftover answers.
 - In case you may need more space/pages, please ask for additional sheets in the exam hall and attach the same with this booklet while submitting.
 - If you use any algorithm / result / formula covered in the class, just mention it, do not elaborate (unless the same thing has been explicitly asked to answer in the question).
-

Q1. [Decision Tree Learning]

8 marks

The table (shown right) contains data samples of six patients examined in a hospital. Use entropy based information gain measures to construct a minimal decision tree that can predict whether or not a patient is likely to have a heart attack. Show each step of your computation.

Patient IDs	Attributes				Heart Attack
	Gender	Smoker	Exercise	ChestPain	
1	Female	No	Regular	Yes	YES
2	Male	No	Never	Yes	YES
3	Male	Yes	Never	No	YES
4	Female	No	Never	No	NO
5	Male	Yes	Regular	Yes	YES
6	Male	No	Regular	No	NO

Solution: [Helper Data: $\log_2 3 = 1.585$]

The entropy of the training examples (given sample set S) is,

$$Entropy(S) = -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) = \log_2 3 - \frac{2}{3} = 0.9183$$

The entropy for the attributes, Gender, Smoker, Exercise, and ChestPain are (respectively),

$$\begin{aligned} Entropy(\text{Gender}) &= \frac{4}{6} \left[-\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \left(\frac{1}{4}\right) \right] + \frac{2}{6} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right] = \frac{5}{3} - \frac{1}{2} \log_2 3 = 0.8742 \\ Entropy(\text{Smoker}) &= \frac{2}{6} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2}\right) - \frac{0}{2} \log_2 \left(\frac{0}{2}\right) \right] + \frac{4}{6} \left[-\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) \right] = \frac{2}{3} = 0.6667 \\ Entropy(\text{Exercise}) &= \frac{3}{6} \left[-\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) \right] + \frac{3}{6} \left[-\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) \right] = \log_2 3 - \frac{2}{3} = 0.9183 \\ Entropy(\text{ChestPain}) &= \frac{3}{6} \left[-\frac{3}{3} \log_2 \left(\frac{3}{3}\right) - \frac{0}{3} \log_2 \left(\frac{0}{3}\right) \right] + \frac{3}{6} \left[-\frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right) \right] = \frac{1}{2} \log_2 3 - \frac{1}{3} = 0.4591 \end{aligned}$$

So, the information gain (which is $Entropy(S) - Entropy(\text{Attribute})$) is highest for ChestPain.

For ChestPain = YES, all *three* samples result HeartAttack = YES.

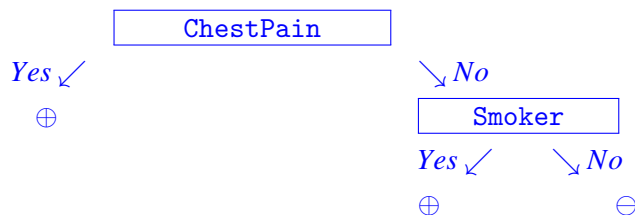
For ChestPain = NO, *two* samples result HeartAttack = YES and *one* sample results HeartAttack = NO. So, the measured entropy for ChestPain = NO is,

$$Entropy(S|_{\text{ChestPain=NO}}) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) = \log_2 3 - \frac{2}{3} = 0.9183$$

The entropy for the attributes, Gender, Smoker, and Exercise are (respectively),

$$\begin{aligned} Entropy(\text{Gender}) &= \frac{2}{3} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right] + \frac{1}{3} \left[-\frac{0}{1} \log_2 \left(\frac{0}{1}\right) - \frac{1}{1} \log_2 \left(\frac{1}{1}\right) \right] = \frac{2}{3} = 0.6667 \\ Entropy(\text{Smoker}) &= \frac{1}{3} \left[-\frac{1}{1} \log_2 \left(\frac{1}{1}\right) - \frac{0}{1} \log_2 \left(\frac{0}{1}\right) \right] + \frac{2}{3} \left[-\frac{0}{2} \log_2 \left(\frac{0}{2}\right) - \frac{2}{2} \log_2 \left(\frac{2}{2}\right) \right] = \frac{2}{3} = 0 \\ Entropy(\text{Exercise}) &= \frac{1}{3} \left[-\frac{0}{1} \log_2 \left(\frac{0}{1}\right) - \frac{1}{1} \log_2 \left(\frac{1}{1}\right) \right] + \frac{2}{3} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right] = \frac{2}{3} = 0.6667 \end{aligned}$$

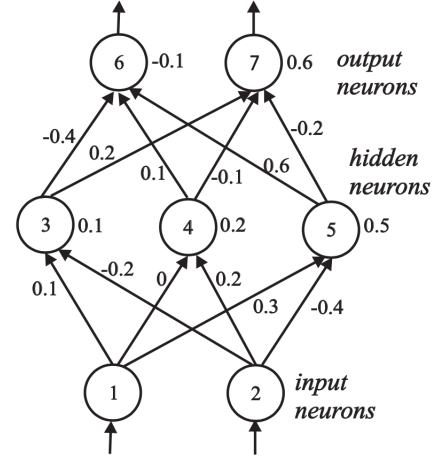
So, the information gain (which is $Entropy(S|_{\text{ChestPain=NO}}) - Entropy(\text{Attribute})$) is highest for Smoker. This finally completes the classification of all examples perfectly. The final decision tree is:



Q2. [Artificial Neural Networks (ANN): Backpropagation]

18 marks

Consider the 3-layer ANN (having one input, one output and one hidden layer) given in the figure (right) which is being trained to distinguish between nails (output encoding as 10) and screws (output encoding as 01). Let the learning rate be $\eta = 0.1$ and the initial weights (w_{ij} from a node i to another node j) are mentioned in the directed edges of the figure. Also, the bias (w_{0j} for a node j) is specified beside each node (neuron) directly. Assume that, each node (neuron), n_j ($3 \leq j \leq 7$), applies the default sigmoid activation function (i.e., $o_j = \sigma(s_j) = \frac{1}{1+\exp(-s_j)}$) over the weighted input sum (i.e., $s_j = w_{0j} + \sum_i w_{ij}o_i$).



Suppose you train this ANN with only *one* example:

$T = \{(0.6, 0.1), \text{nail}\}$, which indicates that when the inputs are $o_1 = 0.6$ and $o_2 = 0.1$, the true outcomes become $d_6 = 1$ and $d_7 = 0$ (target class being nail). Answer the following questions.

- (a) Allowing a *forward pass* with the training example T , compute the values of the outputs of all nodes/neurons. In particular, show your calculation in details for every s_j and o_j ($3 \leq j \leq 7$). (5)

Solution:

$$\begin{aligned}
 s_3 &= 0.1 + (0.1 \times 0.6 - 0.2 \times 0.1) = 0.14 \\
 o_3 &= \frac{1}{1 + \exp(-0.14)} = 0.535 \\
 s_4 &= 0.2 + (0 \times 0.6 + 0.2 \times 0.1) = 0.22 \\
 o_4 &= \frac{1}{1 + \exp(-0.22)} = 0.555 \\
 s_5 &= 0.5 + (0.3 \times 0.6 - 0.4 \times 0.1) = 0.64 \\
 o_5 &= \frac{1}{1 + \exp(-0.64)} = 0.655 \\
 s_6 &= -0.1 + (-0.4 \times 0.535 + 0.1 \times 0.555 + 0.6 \times 0.655) = 0.135 \\
 o_6 &= \frac{1}{1 + \exp(-0.135)} = 0.534 \\
 s_7 &= 0.6 + (0.2 \times 0.535 - 0.1 \times 0.555 - 0.2 \times 0.655) = 0.521 \\
 o_7 &= \frac{1}{1 + \exp(-0.521)} = 0.627
 \end{aligned}$$

- (c) Present the (delta) weight update rule/equation followed during backpropagation of ANNs. (2)

Solution:

$$\begin{aligned}
 w_{ij}^{new} &\leftarrow w_{ij} + \Delta w_{ij} \text{ (weight update) and } \Delta w_{ij} = \eta \times \delta_j \times o_i \\
 \text{where, (error) } \delta_j &= \begin{cases} (d_j - o_j) \times o_j \times (1 - o_j), & \text{for output layer nodes} \\ \sum_k (w_{jk} \times \delta_k) \times o_j \times (1 - o_j) & \text{for hidden layer nodes} \end{cases}
 \end{aligned}$$

-
- (c) Allowing a *backward pass* with the training example T , compute the updates for each weight in the network (as per your equation given in Part (b)). In particular, show your calculation in details for every δ_j and Δw_{ij} (where, $0 \leq i \leq 5$ and $3 \leq j \leq 7$). (11)

Solution:

$$\delta_7 = (0 - 0.627) \times 0.627 \times (1 - 0.627) = -0.147$$

$$\Delta w_{07} = 0.1 \times (-0.147) \times 1 = -0.015$$

$$\Delta w_{37} = 0.1 \times (-0.147) \times 0.535 = -0.008$$

$$\Delta w_{47} = 0.1 \times (-0.147) \times 0.555 = -0.008$$

$$\Delta w_{57} = 0.1 \times (-0.147) \times 0.655 = -0.010$$

$$\delta_6 = (1 - 0.534) \times 0.534 \times (1 - 0.534) = 0.116$$

$$\Delta w_{06} = 0.1 \times 0.116 \times 1 = 0.012$$

$$\Delta w_{36} = 0.1 \times 0.116 \times 0.535 = 0.006$$

$$\Delta w_{46} = 0.1 \times 0.116 \times 0.555 = 0.006$$

$$\Delta w_{56} = 0.1 \times 0.116 \times 0.655 = 0.008$$

$$\delta_5 = (0.6 \times 0.116 - 0.2 \times (-0.147)) \times 0.655 \times (1 - 0.655) = 0.0227$$

$$\Delta w_{05} = 0.1 \times 0.0227 \times 1 = 0.0023$$

$$\Delta w_{15} = 0.1 \times 0.0227 \times 0.6 = 0.0014$$

$$\Delta w_{25} = 0.1 \times 0.0227 \times 0.1 = 0.0002$$

$$\delta_4 = (0.1 \times 0.116 - 0.1 \times (-0.147)) \times 0.555 \times (1 - 0.555) = 0.0065$$

$$\Delta w_{04} = 0.1 \times 0.0065 \times 1 = 0.00065$$

$$\Delta w_{14} = 0.1 \times 0.0065 \times 0.6 = 0.00039$$

$$\Delta w_{24} = 0.1 \times 0.0065 \times 0.1 = 0.00007$$

$$\delta_3 = (-0.4 \times 0.116 + 0.2 \times (-0.147)) \times 0.535 \times (1 - 0.535) = -0.0189$$

$$\Delta w_{03} = 0.1 \times (-0.0189) \times 1 = -0.00189$$

$$\Delta w_{13} = 0.1 \times (-0.0189) \times 0.6 = -0.00113$$

$$\Delta w_{23} = 0.1 \times (-0.0189) \times 0.1 = -0.00019$$

Q3. [Support Vector Machines (SVM)]**10 marks**

Consider a set of 2-dimensional training data points (x_1, x_2) belonging to two classes ' \oplus ' and ' \ominus ', respectively, as shown below.

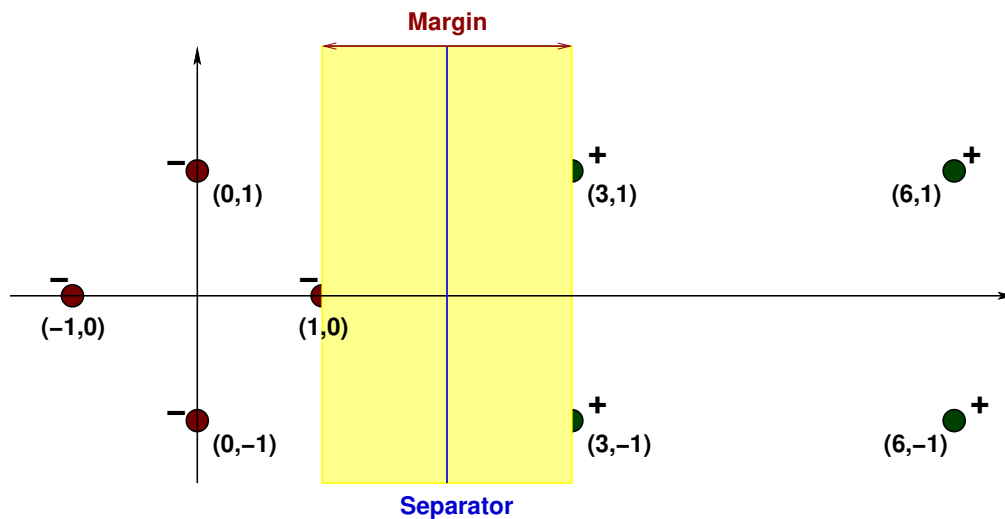
- Class ' \oplus ': $(3, 1)$; $(3, -1)$; $(6, 1)$; $(6, -1)$
- Class ' \ominus ': $(1, 0)$; $(0, 1)$; $(0, -1)$; $(-1, 0)$

We design a linear hard-margin SVM to classify these linearly separable points. Answer the following.

- (a) Pictorially (graphically) represent the constellation of data points and the optimal separating hyperplane. Write the equation of the optimal separator and mention the width of the margin (figuring it out manually from the diagram/graph you have shown). (2)

Solution:

The constellation of data points and the optimal separator (with margin) is presented below.



SVM tries to maximize the margin between two classes of data points. Therefore, the optimal decision boundary crosses the point $(2, 0)$ and is parallel to vertical axis. Thus, the equation of optimal separator is given as, $x_1 - 2 = 0$, having the width of the margin = 2-units.

- (b) Which data points are the support vectors here? (2)

Solution:

Support vectors are $(3, 1)$, $(3, -1)$ and $(1, 0)$. These three points have minimum perpendicular distance from the separator line (Euclidean distance of 1 unit).

-
- (c) What weight vector and threshold (bias) value are being learnt using hard-margin SVM training algorithm with these eight training points? Show the detailed calculations. (4)

Solution:

Let the weight vector learnt be of the form $w = [w_1, w_2]^T$ and threshold/bias is b .

From the three support vectors, $(3, 1)$, $(3, -1)$ and $(1, 0)$ (which are the closest points from the separating line), we get,

$$3w_1 + w_2 + b = +1$$

$$3w_1 - w_2 + b = +1$$

$$w_1 + b = -1$$

Solving above equations, we get, $w_1 = 1$, $w_2 = 0$, and $b = -2$.

- (d) Using the learnt weights and threshold values (in part (c)), what is the margin you get for the optimal classifier? Derive mathematically. (2)

Solution:

$$\text{Margin} = \frac{2}{\|w\|} = \frac{2}{\sqrt{w_1^2 + w_2^2}} = 2.$$

Q4. [Classifier Evaluation]**9 marks**

You wrote a spam filtering program by yourself and now you are testing your program on 100 emails among which you already knew that 20% emails are spams. However, upon running your program on those 100 email corpus, it predicted $\frac{1}{2}$ of the 'spam' emails as non-spam. Answer the following:

- (a) In order to push the $Accuracy \geq 75\%$, how many 'non-spam' emails at most (maximum) you can afford to mis-predict as spams? (3)

Solution:

Suppose, the spam filtering program can afford to mis-predict at most M 'non-spam' emails as spams. As per the problem, among 100 total test emails, 20 are actual spams, and hence $TP = FN = 10$. So, we have, $FP = M$ and $TN = 80 - M$.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} = \frac{10 + (80 - M)}{100} \geq \frac{3}{4} \implies M \leq 15$$

So, this spam filtering program can afford to mis-predict at most 15 'non-spam' emails as spams.

- (b) With the derived setup in Part (a), i.e., when your $Accuracy$ is exactly 75%, present the confusion matrix (in tabular form). (3)

Solution:

		(Actual)	
		Spam Emails	Non-Spam Emails
(Predicted)	Spam Emails	10 (TP)	15 (FP)
	Non-Spam Emails	10 (FN)	65 (TN)

Alternative Approach:

		(Actual)	
		Non-Spam Emails	Spam Emails
(Predicted)	Non-Spam Emails	65 (TP)	10 (FP)
	Spam Emails	15 (FN)	10 (TN)

- (c) As per your confusion matrix that you presented in Part (b), calculate $Precision$, $Recall$ and F_1 -score of your spam filtering program. (3)

Solution:

$$\begin{aligned}
 Precision &= \frac{TP}{TP + FP} = \frac{10}{10 + 15} = 0.4 \\
 Recall &= \frac{TP}{TP + FN} = \frac{10}{10 + 10} = 0.5 \\
 F_1\text{-score} &= 2 \times \frac{Precision \times Recall}{Precision + Recall} = 2 \times \frac{0.4 \times 0.5}{0.4 + 0.5} \approx 0.44
 \end{aligned}$$

Alternative Approach:

$$\begin{aligned}
 Precision &= \frac{TP}{TP + FP} = \frac{65}{65 + 10} \approx 0.87 \\
 Recall &= \frac{TP}{TP + FN} = \frac{65}{65 + 15} \approx 0.81 \\
 F_1\text{-score} &= 2 \times \frac{Precision \times Recall}{Precision + Recall} \approx 2 \times \frac{0.87 \times 0.81}{0.87 + 0.81} \approx 0.84
 \end{aligned}$$

Q5. [Dimensionality Reduction: Principal Component Analysis (PCA)]**8 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.

Solution:

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

The covariance matrix can be constructed as:

$$\begin{aligned} \text{CoVar}(x,x) &= \text{Var}(x) = \frac{[(4-3)^2 + (2-3)^2 + (5-3)^2 + (1-3)^2]}{4} = \frac{5}{2} \\ \text{CoVar}(x,y) &= \text{CoVar}(y,x) = \frac{[(4-3) \times (1-2) + (2-3) \times (3-2) + (5-3) \times (4-2) + (1-3) \times (0-2)]}{4} = \frac{3}{2} \\ \text{CoVar}(y,y) &= \text{Var}(y) = \frac{[(1-2)^2 + (3-2)^2 + (4-2)^2 + (0-2)^2]}{4} = \frac{5}{2} \end{aligned}$$

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

To compute eigenvalues, we make $|\text{CoVar} - \lambda I| = 0$, which gives:

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

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

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

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

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

(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).

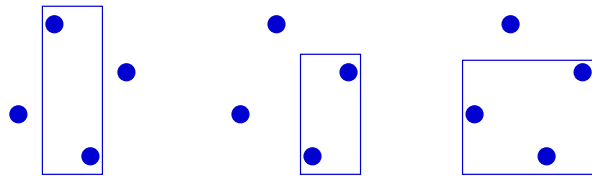
- (a) What the VC-dimension of axis-aligned rectangles in a 2-dimensional plane? Justify / Prove. (5)

Solution:

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

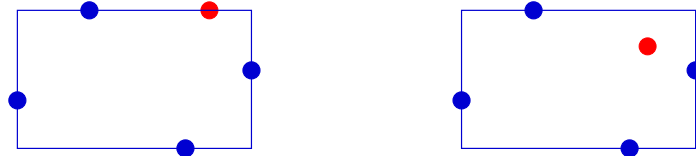
- There exist 4 points that can be shattered. Hence, $d_{VC} \geq 4$.

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.

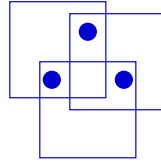
(b) What is the VC-dimension of axis-aligned squares in a 2-dimensional plane? Justify / Prove. (5)

Solution:

The VC-dimension of axis-aligned squares is 3. We prove $d_{VC} = 3$ as follows:

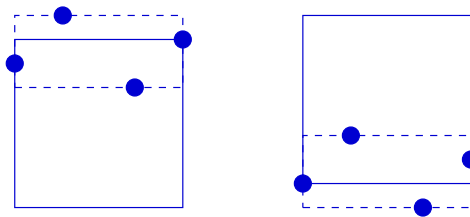
- There exist 3 points that can be shattered. Hence, $d_{VC} \geq 3$.

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 an axis-aligned square.

- (c) Let the VC-dimensions of two hypothesis classes, H_1 and H_2 , be $VCDim(H_1) = d_1$ and $VCDim(H_2) = d_2$. Prove that, the VC-Dimension of the union of these hypothesis, i.e. $H = H_1 \cup H_2$, will be at most $(d_1 + d_2 + 1)$, i.e. $VCDim(H) \leq VCDim(H_1) + VCDim(H_2) + 1$. (5)

Solution:

By the definition of growth function on any N points for a hypothesis class \mathcal{H} , we know that,

$$m_{\mathcal{H}}(N) \leq \sum_{i=0}^{d_{VC}} \binom{N}{i}, \text{ where } VCDim(\mathcal{H}) = d_{VC} \text{ is the VC-dimension of } \mathcal{H}.$$

Let the growth functions on any N points of the hypothesis classes, H_1 , H_2 and H , are denoted by $m_{H_1}(N)$, $m_{H_2}(N)$ and $m_H(N)$, respectively. Since we have $H = H_1 \cup H_2$, we can write

$$m_H(N) \leq m_{H_1}(N) + m_{H_2}(N).$$

Taking $N = d_1 + d_2 + 2$, we can proceed as follows:

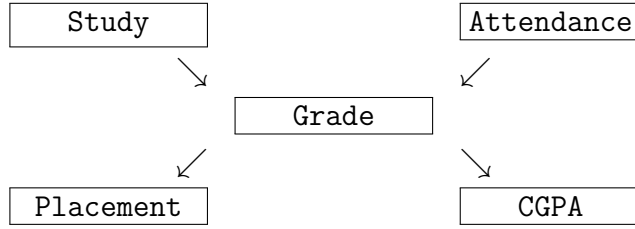
$$\begin{aligned} m_H(N) \leq m_{H_1}(N) + m_{H_2}(N) &\leq \sum_{i=0}^{d_1} \binom{N}{i} + \sum_{i=0}^{d_2} \binom{N}{i} \\ &= \sum_{i=0}^{d_1} \binom{d_1 + d_2 + 2}{i} + \sum_{i=0}^{d_2} \binom{d_1 + d_2 + 2}{i} \\ &= \sum_{i=0}^{d_1} \binom{d_1 + d_2 + 2}{i} + \sum_{i=0}^{d_2} \binom{d_1 + d_2 + 2}{d_1 + d_2 + 2 - i} \\ &= \sum_{i=0}^{d_1} \binom{d_1 + d_2 + 2}{i} + \sum_{i=d_1+2}^{d_1+d_2+2} \binom{d_1 + d_2 + 2}{i} \\ &= \sum_{i=0}^{d_1+d_2+2} \binom{d_1 + d_2 + 2}{i} - \binom{d_1 + d_2 + 2}{d_1 + 1} \\ &= 2^{d_1+d_2+2} - \binom{d_1 + d_2 + 2}{d_1 + 1} \\ &< 2^{d_1+d_2+2} = 2^N \\ \implies m_H(d_1 + d_2 + 2) &< 2^{d_1+d_2+2} \implies (d_1 + d_2 + 2) \text{ is a break point of } H \\ &\implies VCDim(H) \leq d_1 + d_2 + 1 \end{aligned}$$

Therefore, $VCDim(H) \leq VCDim(H_1) + VCDim(H_2) + 1$.

[Proved]

Q7. [Bayesian Learning: Expectation-Maximization (EM) Algorithm]**16 marks**

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



You 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 . You have 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 you need to estimate for the given Bayes Network?

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

(1)

Solution:

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

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

-
- (b) Now, you need to simulate the first E-step of the EM algorithm. Before you start, you 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: (5)

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

(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:

$$\begin{aligned} \therefore \mathbb{E}(p_7 = 1 \mid s_7, a_7, g_7, c_7, \theta) &= \frac{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \theta)}{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \theta) + \mathbb{P}(p_7 = 0, s_7, a_7, g_7, c_7 \mid \theta)} \\ &= \frac{\theta(p_7 = 1 \mid g_7) \cdot \theta(g_7 \mid s_7, a_7) \cdot \theta(s_7) \cdot \theta(a_7)}{\theta(p_7 = 1 \mid g_7) \cdot \theta(g_7 \mid s_7, a_7) \cdot \theta(s_7) \cdot \theta(a_7) + \theta(p_7 = 0 \mid g_7) \cdot \theta(g_7 \mid s_7, a_7) \cdot \theta(s_7) \cdot \theta(a_7)} \\ &= \frac{0.5 \times 0.5 \times 0.5 \times 0.5}{2 \times 0.5 \times 0.5 \times 0.5 \times 0.5} = 0.5 = \mathbb{E}(p_7 = 1 \mid g_7 = 1, \theta(p \mid 1)) \end{aligned}$$

$$\begin{aligned} \therefore \mathbb{E}(c_8 = 1 \mid s_8, a_8, g_8, p_8, \theta) &= \frac{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \theta)}{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \theta) + \mathbb{P}(c_8 = 0, s_8, a_8, g_8, p_8 \mid \theta)} \\ &= \frac{\theta(c_8 = 1 \mid g_8) \cdot \theta(g_8 \mid s_8, a_8) \cdot \theta(s_8) \cdot \theta(a_8)}{\theta(c_8 = 1 \mid g_8) \cdot \theta(g_8 \mid s_8, a_8) \cdot \theta(s_8) \cdot \theta(a_8) + \theta(c_8 = 0 \mid g_8) \cdot \theta(g_8 \mid s_8, a_8) \cdot \theta(s_8) \cdot \theta(a_8)} \\ &= \frac{0.5 \times 0.5 \times 0.5 \times 0.5}{2 \times 0.5 \times 0.5 \times 0.5 \times 0.5} = 0.5 = \mathbb{E}(c_8 = 1 \mid g_8 = 1, \theta(c \mid 1)) \end{aligned}$$

-
- (c) Now, you need 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)) after this M-step? (5)
(Note that, you can use the expected count only when the variable is unobserved in an example)

Solution:

10 parameters will get the updated values as follows:

$$\begin{aligned}
 \theta(s) &= \mathbb{P}(S=1) = \frac{\#\{S=1\}}{\#K} = \frac{4}{8} = 0.5 \\
 \theta(a) &= \mathbb{P}(A=1) = \frac{\#\{A=1\}}{\#K} = \frac{4}{8} = 0.5 \\
 \theta(g|00) &= \mathbb{P}(G=1 | S=0, A=0) = \frac{\#\{G=1, S=0, A=0\}}{\#\{S=0, A=0\}} = \frac{0}{3} = 0.0 \\
 \theta(g|01) &= \mathbb{P}(G=1 | S=0, A=1) = \frac{\#\{G=1, S=0, A=1\}}{\#\{S=0, A=1\}} = \frac{1}{1} = 1.0 \\
 \theta(g|10) &= \mathbb{P}(G=1 | S=1, A=0) = \frac{\#\{G=1, S=1, A=0\}}{\#\{S=1, A=0\}} = \frac{1}{1} = 1.0 \\
 \theta(g|11) &= \mathbb{P}(G=1 | S=1, A=1) = \frac{\#\{G=1, S=1, A=1\}}{\#\{S=1, A=1\}} = \frac{3}{3} = 1.0 \\
 \theta(p|0) &= \mathbb{P}(P=1 | G=0) = \frac{\#\{G=0\} \cdot \mathbb{E}[P=1]}{\#\{G=0\}} = \frac{(1 \times 1.0 + 2 \times 0.0)}{3} = 0.33 \\
 \theta(p|1) &= \mathbb{P}(P=1 | G=1) = \frac{\#\{G=1\} \cdot \mathbb{E}[P=1]}{\#\{G=1\}} = \frac{(4 \times 1.0 + 1 \times 0.5)}{5} = 0.9 \\
 \theta(c|0) &= \mathbb{P}(C=1 | G=0) = \frac{\#\{G=0\} \cdot \mathbb{E}[C=1]}{\#\{G=0\}} = \frac{(1 \times 1.0 + 2 \times 0.0)}{3} = 0.33 \\
 \theta(c|1) &= \mathbb{P}(C=1 | G=1) = \frac{\#\{S=1\} \cdot \mathbb{E}[C=1]}{\#\{G=1\}} = \frac{(3 \times 1.0 + 1 \times 0.0 + 1 \times 0.0)}{5} = 0.7
 \end{aligned}$$

(d) Lastly, you again simulate the second E-step of the EM algorithm. What are the following expectation values that will get calculated in this E-step? In particular, calculate the following: (5)

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

Solution:

$$\begin{aligned}
 &\therefore \mathbb{E}(p_7 = 1 \mid s_7, a_7, g_7, c_7, \theta) \\
 &= \frac{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \theta)}{\mathbb{P}(p_7 = 1, s_7, a_7, g_7, c_7 \mid \theta) + \mathbb{P}(p_7 = 0, s_7, a_7, g_7, c_7 \mid \theta)} \\
 &= \frac{\theta(p_7 = 1 \mid g_7) \cdot \theta(g_7 \mid s_7, a_7) \cdot \theta(s_7) \cdot \theta(a_7)}{\theta(p_7 = 1 \mid g_7) \cdot \theta(g_7 \mid s_7, a_7) \cdot \theta(s_7) \cdot \theta(a_7) + \theta(p_7 = 0 \mid g_7) \cdot \theta(g_7 \mid s_7, a_7) \cdot \theta(s_7) \cdot \theta(a_7)} \\
 &= \frac{0.9 \times 1.0 \times 0.5 \times 0.5}{1.0 \times 1.0 \times 0.5 \times 0.5} = 0.9 = \mathbb{E}(p_7 = 1 \mid g_7 = 1, \theta(p \mid 1))
 \end{aligned}$$

$$\begin{aligned}
 &\therefore \mathbb{E}(c_8 = 1 \mid s_8, a_8, g_8, p_8, \theta) \\
 &= \frac{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \theta)}{\mathbb{P}(c_8 = 1, s_8, a_8, g_8, p_8 \mid \theta) + \mathbb{P}(c_8 = 0, s_8, a_8, g_8, p_8 \mid \theta)} \\
 &= \frac{\theta(c_8 = 1 \mid g_8) \cdot \theta(g_8 \mid s_8, a_8) \cdot \theta(s_8) \cdot \theta(a_8)}{\theta(c_8 = 1 \mid g_8) \cdot \theta(g_8 \mid s_8, a_8) \cdot \theta(s_8) \cdot \theta(a_8) + \theta(c_8 = 0 \mid g_8) \cdot \theta(g_8 \mid s_8, a_8) \cdot \theta(s_8) \cdot \theta(a_8)} \\
 &= \frac{0.7 \times 1.0 \times 0.5 \times 0.5}{1.0 \times 1.0 \times 0.5 \times 0.5} = 0.7 = \mathbb{E}(c_8 = 1 \mid g_8 = 1, \theta(c \mid 1))
 \end{aligned}$$

Q8. [Unsupervised Learning: K-Means Clustering]

16 marks

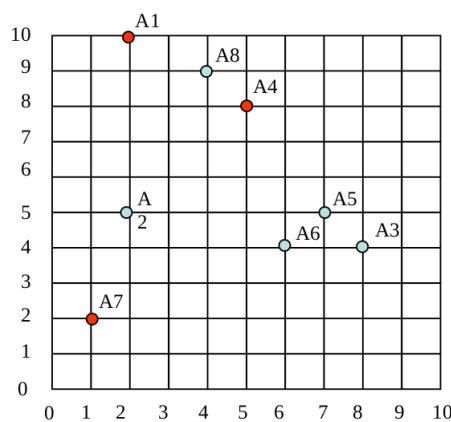
Suppose the following dataset (consisting of (x,y) -coordinates of eight data points in a 2-dimensional plane) is given: $(1,2)$, $(2,5)$, $(2,10)$, $(4,9)$, $(5,8)$, $(6,4)$, $(7,5)$, and $(8,4)$. You need to run K -Means algorithm (till termination) with $K = 3$ to cluster these points. Assume that, Euclidean distance measure is used as the distance computing function for the dataset. Answer the following.

- (a) Assuming the initial centroids as $(1,2)$, $(2,10)$, and $(5,8)$, show in details the execution of K -Means algorithm (till termination) with $K = 3$. In particular, indicate the set of points that come under each cluster after every iteration and also compute their centroid to be used for the next iteration. Indicate when and how you decided to terminate/stop. (8)

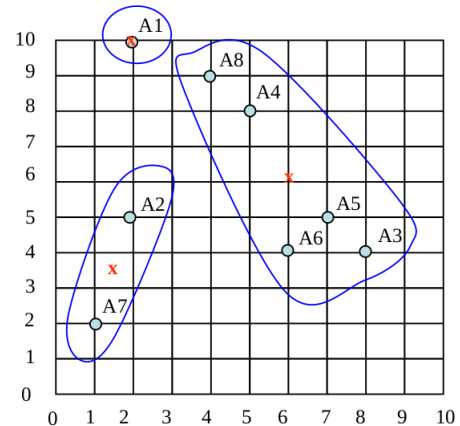
Solution:

Iteration	Output	Cluster-1	Cluster-2	Cluster-3
0 (init.)	Points	—	—	—
	Centroid	$(1,2)$	$(2,10)$	$(5,8)$
1 (cont.)	Points	$(1,2)$, $(2,5)$	$(2,10)$	$(4,9)$, $(5,8)$, $(6,4)$, $(7,5)$, $(8,4)$
	New Centroid	$(1.5, 3.5)$	$(2, 10)$	$(6, 6)$
2 (cont.)	Points	$(1,2)$, $(2,5)$	$(2,10)$, $(4,9)$	$(5,8)$, $(6,4)$, $(7,5)$, $(8,4)$
	New Centroid	$(1.5, 3.5)$	$(3, 9.5)$	$(6.5, 5.25)$
3 (cont.)	Points	$(1,2)$, $(2,5)$	$(2,10)$, $(4,9)$, $(5,8)$	$(6,4)$, $(7,5)$, $(8,4)$
	New Centroid	$(1.5, 3.5)$	$(3.67, 9)$	$(7, 4.33)$
4 (stop)	Points	$(1,2)$, $(2,5)$	$(2,10)$, $(4,9)$, $(5,8)$	$(6,4)$, $(7,5)$, $(8,4)$
	New Centroid	$(1.5, 3.5)$	$(3.67, 9)$	$(7, 4.33)$

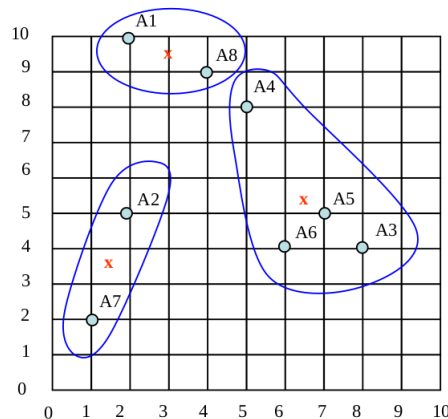
Decision to terminate: When the set of points inside clusters remain unchanged across iterations.



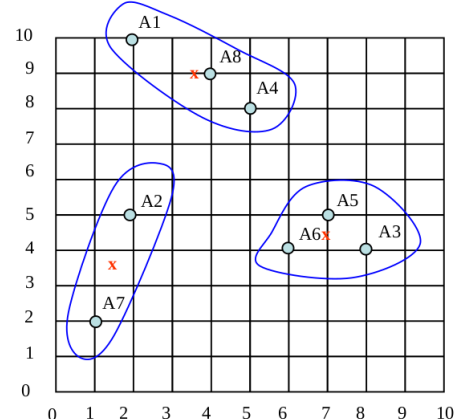
Iteration-0



Iteration-1



Iteration-2



Iteration-3,4

- (b) Upon termination, compute the average silhouette coefficient (SC) of the overall clustering only for the three clusters formed. Show your calculations in details. (8)

Note: For your convenience, pairwise distances between points are shown in the following table.

Pairwise Distance		Data Points							
		(1, 2)	(2, 5)	(2, 10)	(4, 9)	(5, 8)	(6, 4)	(7, 5)	(8, 4)
Data Points	(1, 2)	0.000							
	(2, 5)	3.162	0.000						
	(2, 10)	8.062	5.000	0.000					
	(4, 9)	7.616	4.472	2.236	0.000				
	(5, 8)	7.211	4.243	3.606	1.414	0.000			
	(6, 4)	5.385	4.123	7.211	5.385	4.123	0.000		
	(7, 5)	6.708	5.000	7.071	5.000	3.606	1.414	0.000	
	(8, 4)	7.280	6.083	8.485	6.403	5.000	2.000	1.414	0.000

Solution:

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

$$P_1 (1, 2) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{3.162}{1}\right)}{\left(\frac{8.062+7.616+7.211+5.385+6.708+7.280}{6}\right)} = 0.551$$

$$P_2 (2, 5) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{3.162}{1}\right)}{\left(\frac{5.000+4.472+4.243+4.123+5.000+6.083}{6}\right)} = 0.344$$

$$P_3 (2, 10) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{2.236+3.606}{2}\right)}{\left(\frac{8.062+5.000+7.211+7.071+8.485}{5}\right)} = 0.592$$

$$P_4 (4, 9) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{2.236+1.414}{2}\right)}{\left(\frac{7.616+4.472+5.385+5.000+6.403}{5}\right)} = 0.684$$

$$P_5 (5, 8) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{3.606+1.414}{2}\right)}{\left(\frac{7.211+4.243+4.123+3.606+5.000}{5}\right)} = 0.481$$

$$P_6 (6, 4) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{1.414+2.000}{2}\right)}{\left(\frac{5.385+4.123+7.211+5.385+4.123}{5}\right)} = 0.675$$

$$P_7 (7, 5) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{1.414+1.414}{2}\right)}{\left(\frac{6.708+5.000+7.071+5.000+3.606}{5}\right)} = 0.742$$

$$P_8 (8, 4) : SC = 1 - \frac{a}{b} = 1 - \frac{\left(\frac{2+1.414}{2}\right)}{\left(\frac{7.280+6.083+8.485+6.403+5.000}{5}\right)} = 0.743$$

$$\text{Cluster-1} : \left(\{P_1, P_2\}\right) \text{ Average-SC} = \frac{0.551 + 0.344}{2} = 0.448$$

$$\text{Cluster-2} : \left(\{P_3, P_4, P_5\}\right) \text{ Average-SC} = \frac{0.592 + 0.684 + 0.481}{3} = 0.586$$

$$\text{Cluster-3} : \left(\{P_6, P_7, P_8\}\right) \text{ Average-SC} = \frac{0.675 + 0.742 + 0.743}{3} = 0.720$$

$$\text{Overall} : \text{Average-SC} = \frac{0.448 + 0.586 + 0.720}{3} = 0.585$$

