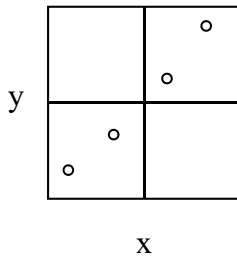


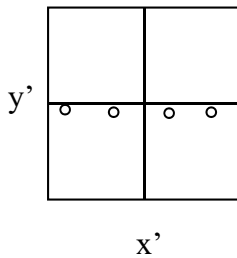
COMP5331 Knowledge Discovery in Databases (Fall Semester 2019)
Homework 2 Solution

Q1.(a) No.

In the original coordinate system, we have the following 4 data points. The threshold is set to 4. In the following, we could not find any clusters in any subspace. There is no grid unit involved in the original coordinate system.



However, in the new coordinate system, we will obtain the following.



We could find one grid unit on the 1-dimensional axis, y' , with the same threshold of 4. Thus, the total number of grid units involved in all clusters based on the new coordinate system is greater than that based on the original coordinate system

(b) When the size of the subspace is larger, the no. of grid units is larger. The probability that a grid unit contains a sufficient number of points is smaller.

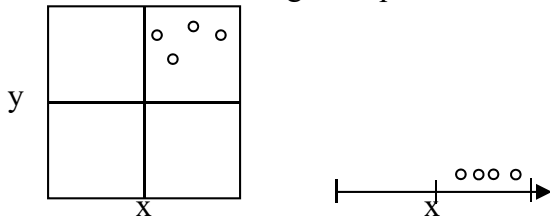
(c) (i) No.

Apriori-like algorithm is valid when the following property holds. If there exists a dense unit *w.r.t* a subspace S , then there exists a dense unit *w.r.t* a subspace S' where $S' \subseteq S$.

Now, we give a counter-example that the above property does not hold under Condition 1. Consider $T_2 = 4$ and $T_1 = 5$ (in count instead of density)

$\because T_2 < T_1 \therefore$ Condition 1 holds.

However, from the following example



we find that there exists a dense unit in subspace $\{x, y\}$ but not in subspace $\{x\}$.

(ii) No.

Let us use the example from Part (c)(i).

In this example, c is equal to 2.

Suppose that $T_2 = 4$. In this case, according to Condition 1, we have $T_1 = 2 T_2 = 2 \cdot 4 = 8$.

Q2.

(a)

$$\text{mean vector} = \begin{pmatrix} \frac{6+8+5+9}{4} \\ \frac{6+8+9+5}{4} \end{pmatrix} = \begin{pmatrix} 7 \\ 7 \end{pmatrix} \quad \begin{pmatrix} \frac{7+9+6+10}{4} \\ \frac{7+9+10+6}{4} \end{pmatrix} = \begin{pmatrix} 8 \\ 8 \end{pmatrix}$$

$$\text{For data } (7, 7), \text{ difference from mean vector} = \begin{pmatrix} 6-7 \\ 6-7 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

$$\text{For data } (9, 9), \text{ difference from mean vector} = \begin{pmatrix} 8-7 \\ 8-7 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\text{For data } (6, 10), \text{ difference from mean vector} = \begin{pmatrix} 5-7 \\ 9-7 \end{pmatrix} = \begin{pmatrix} -2 \\ 2 \end{pmatrix}$$

$$\text{For data } (10, 6), \text{ difference from mean vector} = \begin{pmatrix} 9-7 \\ 5-7 \end{pmatrix} = \begin{pmatrix} 2 \\ -2 \end{pmatrix}$$

$$Y = \begin{pmatrix} -1 & 1 & -2 & 2 \\ -1 & 1 & 2 & -2 \end{pmatrix}$$

$$\Sigma = \frac{1}{4} Y Y^T = \frac{1}{4} \begin{pmatrix} -1 & 1 & -2 & 2 \\ -1 & 1 & 2 & -2 \end{pmatrix} \begin{pmatrix} -1 & -1 \\ 1 & 1 \\ -2 & 2 \\ 2 & -2 \end{pmatrix}$$

$$= \frac{1}{4} \begin{pmatrix} 10 & -6 \\ -6 & 10 \end{pmatrix}$$

$$= \begin{pmatrix} \frac{5}{2} & -\frac{3}{2} \\ -\frac{3}{2} & \frac{5}{2} \end{pmatrix}$$

$$\begin{vmatrix} \frac{5}{2} - \lambda & -\frac{3}{2} \\ -\frac{3}{2} & \frac{5}{2} - \lambda \end{vmatrix} = 0 \implies \left(\frac{5}{2} - \lambda\right)^2 - \left(-\frac{3}{2}\right)^2 = 0 \implies \lambda = 4 \text{ or } \lambda = 1$$

when $\lambda = 4$,

$$\begin{pmatrix} \frac{5}{2} - 4 & -\frac{3}{2} \\ -\frac{3}{2} & \frac{5}{2} - 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \implies \begin{pmatrix} -\frac{3}{2} & -\frac{3}{2} \\ -\frac{3}{2} & -\frac{3}{2} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \implies \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \implies x_1 + x_2 = 0$$

$$\text{We choose the eigenvector of unit length: } \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} \end{pmatrix}.$$

When $\lambda = 1$,

$$\begin{pmatrix} \frac{5}{2}-1 & -\frac{3}{2} \\ -\frac{3}{2} & \frac{5}{2}-1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{pmatrix} \frac{3}{2} & -\frac{3}{2} \\ -\frac{3}{2} & \frac{3}{2} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow x_1 - x_2 = 0$$

We choose the eigenvector of unit length: $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{pmatrix}$.

Thus, $\Phi = \begin{pmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{2} \\ -\frac{\sqrt{2}}{2} & \frac{2}{2} \end{pmatrix}$, $Y = \Phi^T X = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{2} \\ \frac{2}{2} & \frac{2}{2} \end{pmatrix} X$.

For data (6, 6), $Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{2}{2} & \frac{2}{2} \end{pmatrix} \begin{pmatrix} 6 \\ 6 \end{pmatrix} = \begin{pmatrix} 0 \\ 8.49 \\ 9.90 \end{pmatrix}$

For data (8, 8), $Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{2}{2} & \frac{2}{2} \end{pmatrix} \begin{pmatrix} 8 \\ 8 \end{pmatrix} = \begin{pmatrix} 0 \\ 11.31 \\ 9.90 \end{pmatrix}$

For data (5, 9), $Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{2}{2} & \frac{2}{2} \end{pmatrix} \begin{pmatrix} 5 \\ 9 \end{pmatrix} = \begin{pmatrix} -2.83 \\ 9.90 \\ 9.90 \end{pmatrix}$

For data (9, 5), $Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{2}{2} & \frac{2}{2} \end{pmatrix} \begin{pmatrix} 9 \\ 5 \end{pmatrix} = \begin{pmatrix} 2.83 \\ 9.90 \\ 9.90 \end{pmatrix}$

Thus, (6, 6) is reduced to (8.49);
 (8, 8) is reduced to (11.31);
 (5, 9) is reduced to (9.90);
 (9, 5) is reduced to (9.90).

(b)(i)

$$\alpha = 2$$

$$\beta = -7$$

(ii)

Yes. It is possible.

We know that we could express

$$x_q = 2 \cdot x_p - 7 \quad \dots\dots\dots(*)$$

$$y_q = 2 \cdot y_p - 7 \quad \dots\dots\dots(**)$$

According to the steps in Part (a), when we substitute equations (*) and (**) into the derivation in Part (a), it is easy to verify that we could obtain the same matrix Φ .

$$\Phi = \begin{pmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix}$$

$$\text{Thus, we have } Y = \Phi^T X = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix} X.$$

$$\text{For data } (5, 5), Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix} \begin{pmatrix} 5 \\ 5 \end{pmatrix} = \begin{pmatrix} 0 \\ 7.07 \end{pmatrix}$$

$$\text{For data } (9, 9), Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix} \begin{pmatrix} 9 \\ 9 \end{pmatrix} = \begin{pmatrix} 0 \\ 12.73 \end{pmatrix}$$

$$\text{For data } (3, 11), Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix} \begin{pmatrix} 3 \\ 11 \end{pmatrix} = \begin{pmatrix} -5.66 \\ 9.90 \end{pmatrix}$$

$$\text{For data } (11, 3), Y = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix} \begin{pmatrix} 11 \\ 3 \end{pmatrix} = \begin{pmatrix} 5.66 \\ 9.90 \end{pmatrix}$$

Thus, (5, 5) is reduced to (7.07);
 (9, 9) is reduced to (12.73);
 (3, 11) is reduced to (9.90);
 (11, 3) is reduced to (9.90).

(c)
 Yes. It is possible.

According to the steps in Part (a), we could calculate the same matrix Σ based on these 8 data points. The remaining steps in Part (a) are the same.

The answer of this part is as follows.

(6, 6) (two occurrences) is reduced to (8.49);
 (8, 8) (two occurrences) is reduced to (11.31);
 (5, 9) (two occurrences) is reduced to (9.90);
 (9, 5) (two occurrences) is reduced to (9.90).

$$\begin{aligned}
\text{(d)(i)} \quad & P(X=1|Y=1)=0 \quad P(X=1|Y=2)=1 \quad P(X=1|Y=3)=0.5 \\
& P(X=2|Y=1)=0.8 \quad P(X=2|Y=2)=0 \quad P(X=2|Y=3)=0.5 \\
& P(X=3|Y=1)=0.2 \quad P(X=3|Y=2)=0 \quad P(X=3|Y=3)=0 \\
& H(X|Y=1) = -0.8\log 0.8 - 0.2\log 0.2 = 0.7219, \quad H(X|Y=2)=0, \quad H(X|Y=3)=1
\end{aligned}$$

$$\begin{aligned}
\therefore H(X|Y) &= \frac{5}{8}H(X|Y=1) + \frac{1}{8}H(X|Y=2) + \frac{1}{4}H(X|Y=3) \\
&= \frac{5}{8} \times 0.7219 + \frac{1}{8} \times 0 + \frac{1}{4} \times 1 \\
&= 0.7012
\end{aligned}$$

(ii)

$$\begin{aligned}
H(X|Y) &= -\frac{1}{2}\log 0.8 - \frac{1}{8}\log 0.2 - \frac{1}{8}\log 1 - \frac{1}{8}\log \frac{1}{2} - \frac{1}{8}\log \frac{1}{2} \\
&= 0.7012
\end{aligned}$$

Q3

(a)(i)

$$\text{Info}(T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

For attribute NoOfPapers,

$$\text{Info}(T_{\text{enough}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(T_{\text{many}}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

$$\text{Info}(T_{\text{few}}) = -0 \log 0 - 1 \log 1 = 0$$

$$\text{Info}(\text{NoOfPapers}) = 1/8 \text{Info}(T_{\text{enough}}) + 1/2 \text{Info}(T_{\text{many}}) + 3/8 \text{Info}(T_{\text{few}}) = 0.405$$

$$\text{SplitInfo}(\text{NoOfPapers}) = -\frac{1}{8} \log \frac{1}{8} - \frac{1}{2} \log \frac{1}{2} - \frac{3}{8} \log \frac{3}{8} = 1.4056$$

$$\text{Gain}(\text{NoOfPapers}, T) = \frac{1 - 0.405}{1.4056} = 0.4233$$

For attribute Age,

$$\text{Info}(T_{\text{young}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{\text{old}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(\text{Age}, T) = \frac{1}{2} \text{Info}(T_{\text{young}}) + \frac{1}{2} \text{Info}(T_{\text{old}}) = 1$$

$$\text{SplitInfo}(\text{Age}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Gain}(\text{Age}, T) = \frac{1 - 1}{1} = 0$$

For attribute Major,

$$\text{Info}(T_{\text{ComputerScience}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{\text{CivilEngineering}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

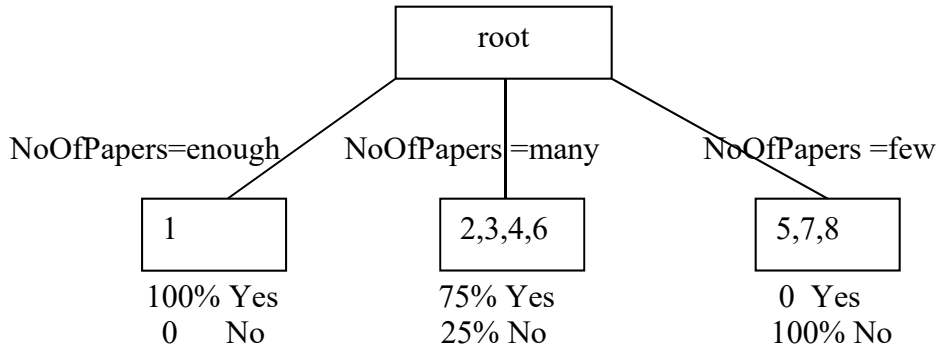
$$\text{Info}(T_{\text{DataScience}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(\text{Major}, T) = 1/2 \text{Info}(T_{\text{ComputerScience}}) + 1/4 \text{Info}(T_{\text{CivilEngineering}}) + 1/4 \text{Info}(T_{\text{DataScience}}) = 1$$

$$\text{SplitInfo}(\text{Major}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} = 1.5$$

$$\text{Gain}(\text{Major}, T) = \frac{1 - 1}{1.5} = 0$$

We choose attribute **NoOfPapers** for Splitting:



Consider the node for “NoOfPapers=many”

$$\text{Info}(T) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

For attribute Age,

$$\text{Info}(T_{\text{young}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{\text{old}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(\text{Age}, T) = \frac{1}{2} \text{Info}(T_{\text{young}}) + \frac{1}{2} \text{Info}(T_{\text{old}}) = 0.5$$

$$\text{SplitInfo}(\text{Age}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Gain}(\text{Age}, T) = \frac{0.8113 - 0.5}{1} = 0.3113$$

For attribute Major,

$$\text{Info}(T_{\text{ComputerScience}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{\text{CivilEngineering}}) = -1 \log 1 - 0 \log 0 = 0$$

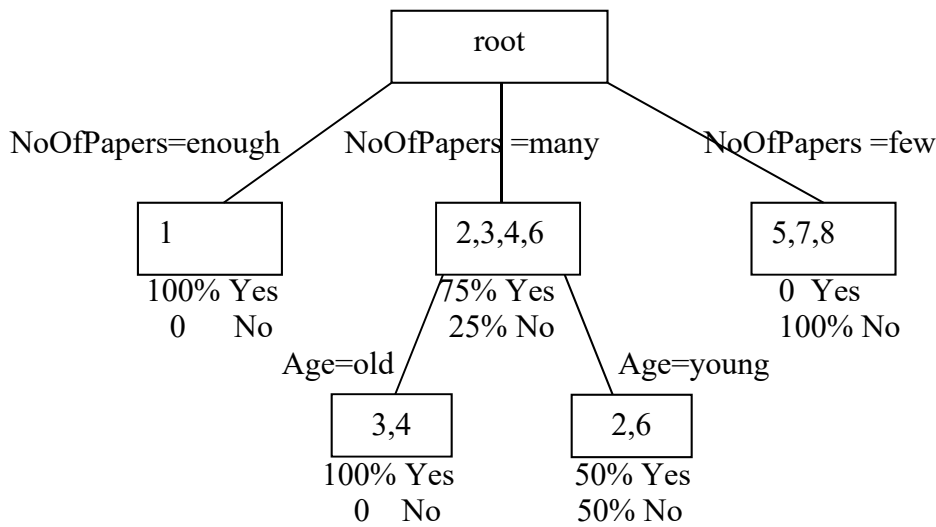
$$\text{Info}(T_{\text{DataScience}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(\text{Major}) = \frac{1}{2} \text{Info}(T_{\text{ComputerScience}}) + \frac{1}{4} \text{Info}(T_{\text{CivilEngineering}}) + \frac{1}{4} \text{Info}(T_{\text{DataScience}}) = 0.5$$

$$\text{SplitInfo}(\text{Major}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} = 1.5$$

$$\text{Gain}(\text{Major}, T) = \frac{0.8113 - 0.5}{1.5} = 0.2075$$

We choose attribute **Age** for Splitting:



(ii) It is 50% probability that this PhD student will become a professor.

(b)(i) Person A said this because his claim is based on this “instance” of the “training set” only and based on some specific parameter values for this decision tree.

(ii) Person B said this because he thought that the result is just based on this “instance” of the training set. It is possible that other attributes could be found in X based on other instances of the training set.

(iii) Person A is more unreasonable since the “instance” of the training set is just a “small” snapshot of the whole population.

(c)

Difference:

The definition of the gain used in C4.5 is different from that used in ID3.

The gain used in C4.5 is equal to the gain used in ID3 divided by SplitInfo.

The reason why there is a difference is described as follows.

In ID3, there is a higher tendency to choose an attribute containing more values (e.g., attribute identifier and attribute HKID). Thus, splitInfo in C4.5 is used to penalize an attribute containing more values. If this value is larger, the penalty is larger.

Q4

$$\begin{aligned}
 \text{(a) } P(\text{SIR} = \text{Yes}) &= \sum_{x \in \{\text{Yes}, \text{No}\}} \sum_{y \in \{\text{Yes}, \text{No}\}} P(\text{SIR} = \text{Yes} | \text{AP} = x, \text{P} = y) P(\text{AP} = x, \text{P} = y) \\
 &= 0.7 \times 0.3 \times 0.6 + 0.45 \times 0.3 \times 0.4 + 0.55 \times 0.7 \times 0.6 + 0.2 \times 0.7 \times 0.4 \\
 &= 0.467
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{SIR} = \text{Yes} | \text{AP} = \text{Yes}, \text{P} = \text{Yes}, \text{WBC} = \text{Low}) \\
 &= \frac{P(\text{WBC} = \text{Low} | \text{AP} = \text{Yes}, \text{P} = \text{Yes}, \text{SIR} = \text{Yes})}{P(\text{WBC} = \text{Low} | \text{AP} = \text{Yes}, \text{P} = \text{Yes})} P(\text{SIR} = \text{Yes} | \text{AP} = \text{Yes}, \text{P} = \text{Yes}) \\
 &= \frac{P(\text{WBC} = \text{Low} | \text{SIR} = \text{Yes}) P(\text{SIR} = \text{Yes} | \text{AP} = \text{Yes}, \text{P} = \text{Yes})}{\sum_{x \in \{\text{Yes}, \text{No}\}} P(\text{WBC} = \text{Low} | \text{SIR} = x) P(\text{SIR} = x | \text{AP} = \text{Yes}, \text{P} = \text{Yes})} \\
 &= \frac{0.4 \times 0.7}{0.4 \times 0.7 + 0.7 \times 0.3} \\
 &= 0.5714
 \end{aligned}$$

$$P(\text{SIR} = \text{No} | \text{AP} = \text{Yes}, \text{P} = \text{Yes}, \text{WBC} = \text{Low}) = 1 - 0.5714 = 0.4286$$

Since $P(\text{SIR} = \text{Yes} | \text{AP} = \text{Yes}, \text{P} = \text{Yes}, \text{WBC} = \text{Low}) > P(\text{SIR} = \text{No} | \text{AP} = \text{Yes}, \text{P} = \text{Yes}, \text{WBC} = \text{Low})$, it is more likely that the person has systemic inflammation reaction.

(a) Disadvantages:

The Bayesian Belief network classifier requires a predefined knowledge about the network.

The Bayesian Belief Network classifier cannot work directly when the network contains cycles.

Q5

(a)(i)

<p>Consider $t = 1$.</p> $f_1 = \sigma(W_f [x_1, y_0] + b_f)$ $= \sigma\left(\begin{pmatrix} 0.8 \\ 0.4 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.6 \\ 0 \end{pmatrix} + 0.2\right)$ $= \sigma(0.8 \cdot 0.3 + 0.4 \cdot 0.6 + 0.1 \cdot 0 + 0.2)$ $= \sigma(0.68)$ $= 0.6637$ $i_1 = \sigma(W_i [x_1, y_0] + b_i)$ $= \sigma\left(\begin{pmatrix} 0.9 \\ 0.8 \\ 0.7 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.6 \\ 0 \end{pmatrix} + 0.5\right)$ $= \sigma(0.9 \cdot 0.3 + 0.8 \cdot 0.6 + 0.7 \cdot 0 + 0.5)$ $= \sigma(1.25)$ $= 0.7773$ $a_1 = \tanh(W_a [x_1, y_0] + b_a)$ $= \tanh\left(\begin{pmatrix} 0.4 \\ 0.2 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.6 \\ 0 \end{pmatrix} + 0.3\right)$ $= \tanh(0.4 \cdot 0.3 + 0.2 \cdot 0.6 + 0.1 \cdot 0 + 0.3)$ $= \tanh(0.54)$ $= 0.4930$ $s_1 = f_1 \cdot s_0 + i_1 \cdot a_1$ $= 0.6637 \cdot 0 + 0.7773 \cdot 0.4930$ $= 0.3832$ $o_1 = \sigma(W_o [x_1, y_0] + b_o)$ $= \sigma\left(\begin{pmatrix} 0.6 \\ 0.4 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.6 \\ 0 \end{pmatrix} + 0.2\right)$ $= \sigma(0.6 \cdot 0.3 + 0.4 \cdot 0.6 + 0.1 \cdot 0 + 0.2)$ $= \sigma(0.62)$ $= 0.6502$ $y_3 = o_1 \cdot \tanh(s_1)$ $= 0.6502 \cdot \tanh(0.3832)$ $= 0.2376$	<p>Consider $t = 2$.</p> $f_2 = \sigma(W_f [x_2, y_1] + b_f)$ $= \sigma\left(\begin{pmatrix} 0.8 \\ 0.4 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.1 \\ 1.0 \\ 0.2376 \end{pmatrix} + 0.2\right)$ $= \sigma(0.8 \cdot 0.1 + 0.4 \cdot 1.0 + 0.1 \cdot 0.2376 + 0.2)$ $= \sigma(0.7038)$ $= 0.6690$ $i_2 = \sigma(W_i [x_2, y_1] + b_i)$ $= \sigma\left(\begin{pmatrix} 0.9 \\ 0.8 \\ 0.7 \end{pmatrix} \begin{pmatrix} 0.1 \\ 1.0 \\ 0.2376 \end{pmatrix} + 0.5\right)$ $= \sigma(0.9 \cdot 0.1 + 0.8 \cdot 1.0 + 0.7 \cdot 0.2376 + 0.5)$ $= \sigma(1.5564)$ $= 0.8258$ $a_2 = \tanh(W_a [x_2, y_1] + b_a)$ $= \tanh\left(\begin{pmatrix} 0.4 \\ 0.2 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.1 \\ 1.0 \\ 0.2376 \end{pmatrix} + 0.3\right)$ $= \tanh(0.4 \cdot 0.1 + 0.2 \cdot 1.0 + 0.1 \cdot 0.2376 + 0.3)$ $= \tanh(0.5638)$ $= 0.5108$ $s_2 = f_2 \cdot s_1 + i_2 \cdot a_2$ $= 0.6690 \cdot 0.3832 + 0.8258 \cdot 0.5108$ $= 0.6782$ $o_2 = \sigma(W_o [x_2, y_1] + b_o)$ $= \sigma\left(\begin{pmatrix} 0.6 \\ 0.4 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.1 \\ 1.0 \\ 0.2376 \end{pmatrix} + 0.2\right)$ $= \sigma(0.6 \cdot 0.1 + 0.4 \cdot 1.0 + 0.1 \cdot 0.2376 + 0.2)$ $= \sigma(0.6838)$ $= 0.6646$ $y_2 = o_2 \cdot \tanh(s_2)$ $= 0.6646 \cdot \tanh(0.6782)$ $= 0.3923$
--	---

(ii)

<p>When $t = 1$,</p> $\text{error} = y_1 - y$ $= 0.2376 - 0.2$ $= 0.0376$	<p>When $t = 2$,</p> $\text{error} = y_2 - y$ $= 0.3923 - 0.4$ $= -0.0077$
--	---

(b) (i)

<p>Consider $t = 1$.</p> $r_1 = \sigma(W_r [x_1, y_0] + b_r)$	<p>Consider $t = 2$.</p> $r_2 = \sigma(W_r [x_2, y_1] + b_r)$
--	--

$= \sigma\left(\begin{pmatrix} 0.3 \\ 0.2 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.6 \\ 0 \end{pmatrix} + 0.5\right)$ $= \sigma(0.3 \cdot 0.3 + 0.2 \cdot 0.6 + 0.1 \cdot 0 + 0.5)$ $= \sigma(0.71)$ $= 0.6704$ $a_1 = \tanh(W_a [x_1, r_1 y_0] + b_a)$ $= \tanh\left(\begin{pmatrix} 0.4 \\ 0.3 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.6 \\ 0.6704 \cdot 0 \end{pmatrix} + 0.1\right)$ $= \tanh(0.4 \cdot 0.3 + 0.3 \cdot 0.6 + 0.1 \cdot 0 + 0.1)$ $= \tanh(0.4)$ $= 0.3799$ $u_1 = \sigma(W_u [x_1, y_0] + b_u)$ $= \sigma\left(\begin{pmatrix} 0.4 \\ 0.2 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.6 \\ 0 \end{pmatrix} + 0.1\right)$ $= \sigma(0.4 \cdot 0.3 + 0.2 \cdot 0.6 + 0.1 \cdot 0 + 0.1)$ $= \sigma(0.34)$ $= 0.5842$ $y_1 = (1 - u_1) \cdot y_0 + u_1 \cdot a_1$ $= (1 - 0.5842) \cdot 0 + 0.5842 \cdot 0.3799$ $= 0.2220$	$= \sigma\left(\begin{pmatrix} 0.3 \\ 0.2 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.1 \\ 1.0 \\ 0.2220 \end{pmatrix} + 0.5\right)$ $= \sigma(0.3 \cdot 0.1 + 0.2 \cdot 1.0 + 0.1 \cdot 0.2220 + 0.5)$ $= \sigma(0.7522)$ $= 0.6797$ $a_2 = \tanh(W_a [x_2, r_2 y_1] + b_a)$ $= \tanh\left(\begin{pmatrix} 0.4 \\ 0.3 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.1 \\ 1.0 \\ 0.6797 \cdot 0.2220 \end{pmatrix} + 0.1\right)$ $= \tanh(0.4 \cdot 0.1 + 0.3 \cdot 1.0 + 0.1 \cdot 0.1509 + 0.1)$ $= \tanh(0.4551)$ $= 0.4261$ $u_2 = \sigma(W_u [x_2, y_1] + b_u)$ $= \sigma\left(\begin{pmatrix} 0.4 \\ 0.2 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.1 \\ 1.0 \\ 0.2220 \end{pmatrix} + 0.1\right)$ $= \sigma(0.4 \cdot 0.1 + 0.2 \cdot 1.0 + 0.1 \cdot 0.2220 + 0.1)$ $= \sigma(0.3622)$ $= 0.5896$ $y_2 = (1 - u_2) \cdot y_1 + u_2 \cdot a_2$ $= (1 - 0.5896) \cdot 0.2220 + 0.5896 \cdot 0.4261$ $= 0.3423$
--	--

(ii)

<p>When $t = 1$,</p> $\text{error} = y_1 - y$ $= 0.2220 - 0.2$ $= 0.0220$	<p>When $t = 2$,</p> $\text{error} = y_2 - y$ $= 0.3423 - 0.4$ $= -0.0577$
--	---

(c) The neural network has an assumption that records in the training set are “independent”. In some cases, records in the training set are related/correlated to (or dependent on) other records in the training set. Thus, the neural network could not capture this “dependent” scenario well if the training set has dependent records.