

EEE - 8803

JEET KIRAN PAWANI
GTID: 903397407

PGM ASSIGNMENT - 4

1] Naive Bayes model with class C ①

Discrete evidence variables x_1, \dots, x_n

$$\text{CPD is } P(C=c) = \theta_c \quad P(x_i = a_i | C=c) = \theta_{a_i/c} \quad i=1, \dots, n$$

$$a_i \in \text{Val}(x_i) \quad c \in \text{Val}(C)$$

Dataset $D = \{x^{[1]}, \dots, x^{[m]}\}$ where $x^{[m]}$ is a complete assignment of evidence variables x_1, \dots, x_n .

If parameters are assigned uniformly,

$$\theta_c^0 = \frac{1}{|\text{Val}(C)|} \quad \theta_{a_i/c}^0 = \frac{1}{|\text{Val}(x_i)|} \quad \text{for all } a_i \in$$

TP - EM algo converges in one iteration.

Now lets start with step 1

$$\theta_c^1 = \frac{1}{m} \sum P(C=c | x^{(m)}, \theta^0) \rightarrow ②$$

Each step of EM is given by $j(y|i) = \rho^{(y|x^{(i)}, \theta^{t-1})}$

Then new parameter is calculated by

$$q^{t(y)} = \frac{1}{n} \sum_{i=1}^n j(y|i) \rightarrow ①$$

(1) being used to derive ②

$$\begin{aligned} P(C=c | x^{(m)}, \theta^0) &= \frac{P(C=c, x^{(m)}, \theta^0)}{P(x^{(m)}, \theta^0)} \\ &= \frac{P(x^{(m)} | C=c, \theta^0) P(C=c, \theta^0)}{P(x^{(m)}) P(\theta^0)} \\ &= \frac{P(x^{(m)} | C=c, \theta^0) P(C=c | \theta^0)}{P(x^{(m)}) P(\theta^0)} \end{aligned}$$

$$P(c/x(m), \theta^0) = \frac{P(x(m)/c, \theta^0) P(c/\theta^0)}{P(x(m)/\theta^0)}$$

$$\Rightarrow \theta_c' = \frac{1}{m} \sum_m P(c/x(m), \theta^0)$$

$$= \frac{1}{m} \sum_m \frac{P(x(m)/c, \theta^0) P(c/\theta^0)}{P(x(m)/\theta^0)}$$

$$= \frac{1}{m} \sum_m \left[\frac{P(x(m)/c, \theta^0) P(c/\theta^0)}{\sum_{c' \in \text{Val}(c)} P(x(m)/c', \theta^0) P(c'/\theta^0)} \right]$$

we know $x(m)$ ranges from $\{x(1), \dots, x(n)\}$
writing it to product form

$$= \frac{1}{m} \sum_m \left[\prod_{c' \in \text{Val}(c)} P(c/\theta^0) \prod_{i=1}^n P(x_i(m)/c, \theta^0) \right]$$

$$\text{Given } P(c=c) = \theta_c \quad P(x_i=a/c=c) = \theta_{xi}/c$$

$$\Rightarrow P(c/\theta_c) = \frac{1}{\text{Val}(c)}, \quad P(x_i=a/c=c, \theta^0) = \frac{1}{\text{Val}(x_i)}$$

$$\theta_c' = \frac{1}{m} \sum_m \left[\frac{\frac{1}{\text{Val}(c)}}{\sum_{c' \in \text{Val}(c)} \frac{1}{\text{Val}(c)}} \prod_{i=1}^n \frac{1}{\text{Val}(x_i)} \right]$$

$\prod_{i=1}^n \frac{1}{\text{Val}(x_i)}$ is independent of m

$$\text{also } \sum_{c \in \text{Val}(c)} \frac{1}{\text{Val}(c)} = 1$$

$$\begin{aligned}\theta_c' &= \frac{1}{m} \sum_{i=1}^m \left[\frac{1}{|\text{Val}(c)|} \right] \\ &= \frac{1}{m} \times \frac{m}{|\text{Val}(c)|} = \underline{\underline{\frac{1}{|\text{Val}(c)|}}} \\ \boxed{\theta_c' = \frac{1}{|\text{Val}(c)|}}\end{aligned}$$

we know $q_j(x_i | y) = \frac{\sum_{c: x_i \in c} d(y/c)}{\sum_{c: x_i \in c} 1}$

Using this as we have $x_i \in c$ so indicator function we get.

$$\theta_{x_i/c} = \frac{\sum_{c: x_i \in c} P(c/x(m), \theta^0) \cdot 1(x_i(m) = x)}{\sum_{c: x_i \in c} P(c/x(m), \theta^0)}$$

we know $P(c/x(m), \theta^0) = \frac{1}{|\text{Val}(c)|}$

$$= \frac{\sum_{c: x_i \in c} \frac{1}{|\text{Val}(c)|} \cdot 1(x_i(m) = x)}{\sum_{c: x_i \in c} \frac{1}{|\text{Val}(c)|}} \Rightarrow \frac{m}{|\text{Val}(c)|} \text{ as } c \text{ circles of } m$$

$$= \frac{\sum_{c: x_i \in c} \frac{1}{|\text{Val}(c)|} \cdot 1(x_i(m) = x)}{\frac{m}{|\text{Val}(c)|}}$$

$\frac{m}{|\text{Val}(c)|} \rightarrow$ this is indicator function when $x_i(m) = x$

\Rightarrow this can be written as - as we know all are equally distributed

$$\theta_{x_i/c} = \cancel{\frac{\sum_{c: x_i \in c} 1(x_i(m) = x)}{m}}$$

$\sum_m 1(x_i(m) = z)$ is number of times $x_i(m) = z$ occurs in data.

This is for the step 1 now we try to generalize wrong induction

$$\text{To prove } \theta_{c=c}^t = \frac{1}{|\text{Val}(c)|}$$

$$\theta_{x_i=l}^t = \frac{\sum 1(x_i(m) = z)}{|\cancel{\text{Val}(l)}|}$$

Now we try to prove θ_c^{t+1} as $\frac{1}{\text{Val}(c)}$

Similar to previous steps

$$\theta_{c=c}^{t+1} = \frac{1}{m} \sum P(c/x_i(m), \theta^t)$$

From (3) we can write this as

$$= \frac{1}{m} \sum_m \frac{P(x_i(m)/c, \theta^t) P(c/\theta^t)}{P(x_i(m)/\theta^t)} \rightarrow \text{adding } c'$$

$$= \frac{1}{m} \sum_m \frac{P(x_i(m)/c, \theta^t) P(c/\theta^t)}{\sum_{c' \in \text{Val}(c)} P(x_i(m)/c', \theta^t) P(c'/\theta^t)}$$

$x_i(m)$ writing as prod of $\prod_{i=1}^n x_i(m)$

$$\theta_{c=c}^{t+1} = \frac{1}{m} \sum_m \frac{P(c/\theta^t) \prod_{i=1}^n P(x_i(m)/c, \theta^t)}{\sum_{c' \in \text{Val}(c)} P(c'/\theta^t) \prod_{i=1}^n P(x_i(m)/c', \theta^t)}$$

$$P(x_i(m)/c, \theta^t) = P(x_i(m)/c', \theta^t) = \frac{1}{|\text{Log}(x_i)|}$$

$$P(c/\theta^t) = P(c'/\theta^t) = \frac{1}{|\text{Val}(c)|}$$

$$\theta_{c=c}^{t+1} = \frac{1}{m} \sum_m \frac{\frac{1}{\text{val}(c)}}{\prod_{i=1}^m \frac{1}{\text{val}(x_i)}}$$

$\sum_{i \in \text{Val}(c)} \frac{1}{\text{val}(x_i)}$

(3)

$$= \frac{1}{m} \sum_m \frac{1}{|\text{val}(c)|} = \frac{1}{m} \frac{m}{|\text{val}(c)|} = \underline{\underline{\frac{1}{|\text{val}(c)|}}}$$

$\theta_{c=c}^{t+1} = \frac{1}{|\text{val}(c)|}$

$$\theta_{x_i=c | c=c}^{t+1} = \frac{\sum_m \frac{1}{\text{val}(c)} \frac{1}{P(c|x_i(m), \theta^t)}}{\sum_m P(c|x_i(m), \theta^t)}$$

(from prev.)

$$= \sum_m \frac{1}{|\text{val}(c)|} \frac{1}{P(c|x_i(m), \theta^t)}$$

$\theta_{x_i=c | c=c}^{t+1} = \sum_m \frac{1}{m} \frac{1}{|\text{val}(c)|} \frac{1}{P(c|x_i(m), \theta^t)}$

Thus, by induction we proved that when uniformly initialized at θ^0 step, θ^m also converges in 1 step.

$$\theta_{c=c}^t = \underline{\underline{\frac{1}{|\text{val}(c)|}}}$$

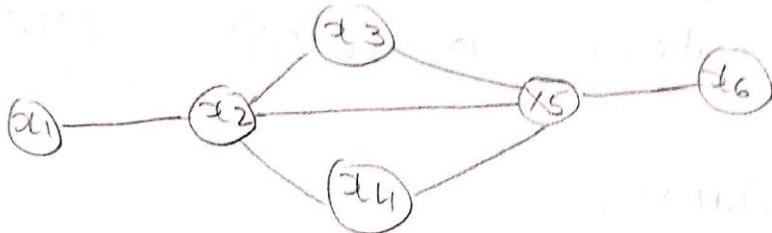
$$\theta_{x_i=c | c=c}^t = \underline{\underline{\frac{\sum_m 1_{\{x_i(m)=c\}}}{m}}}$$

2] Dataset $\{x^n \mid n=1 \dots N\}$

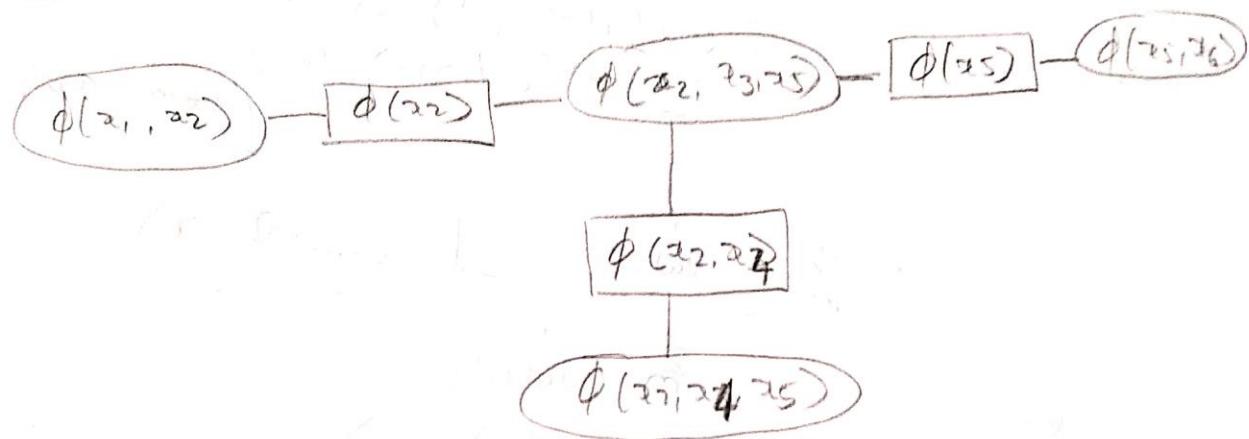
$$x = (x_1, \dots, x_6)$$

empirical distribution $E(x)$

$$p(x_1, \dots, x_6) = \frac{1}{2} \phi_1(x_1, x_2) \phi_2(x_2, x_3, x_5) \phi_3(x_2, x_4, x_5) \phi_4(x_5, x_6)$$



Junction tree



join the above junction tree

$$p(x_1, \dots, x_6) = \frac{\phi(x_1, x_2) \phi(x_2, x_3, x_5) \phi(x_2, x_4, x_5) \phi(x_5, x_6)}{\phi(x_2)}$$

$$\phi(x_1, x_2) = \frac{\phi(x_1/x_2)}{\phi(x_2)} \quad \frac{p(x_2/x_4, x_5)}{p(x_2, x_4)} = \frac{p(x_5/x_2, x_4)}{p(x_2, x_4)}$$

$$\frac{p(x_5, x_6)}{p(x_5)} = p(x_6/x_5)$$

$$p(x_1 \dots x_6) = p(x_1/x_2)p(x_2/x_3, x_5)p(x_4/x_2, x_5) \rightarrow \textcircled{1} P(x_6/x_5) \quad (4)$$

Since we have identified clique potentials, the normalization constant $(2) = 1$

As we are forming Bayesian networks or clusters the clique potentials are independent

Taking log likelihood -

$$L = \sum_n \log (P(x_1^n/x_2)) + \log P(x_2^n, x_3^n, x_5^n) \\ + \log P(x_4^n/x_2^n, x_5^n) + \log P(x_6^n/x_5^n)$$

from \textcircled{1}

$$\phi(x_1, x_2) = p(x_1/x_2)$$

$$\phi(x_2, x_3, x_5) = p(x_2, x_3, x_5)$$

$$\phi(x_2, x_4, x_5) = p(x_4/x_2, x_5)$$

$$\phi(x_5, x_6) = p(x_6/x_5)$$

Since all terms in dataset are independent as clique potentials are independent. The max. likelihood solution in this Bayesian network can be simply setting each factor to its empirical distributions.

$$\phi(x_1, x_2) = \epsilon(x_1/x_2)$$

$$\phi(x_2, x_3, x_5) = \epsilon(x_2, x_3, x_5)$$

$$\phi(x_2, x_4, x_5) = \epsilon(x_4/x_2, x_5)$$

$$\phi(x_5, x_6) = \underline{\epsilon(x_6/x_5)}$$

3] N observations $x = x^1, \dots, x^N$
 the log likelihood of belief network to generate x is

$$\log p(x) = \sum_{n=1}^N \sum_{i=1}^K \log p(x_i^n | pa(x_i^n))$$

$\theta_{s^i}(t) = p(x_i = s | pa(x_i) = t)$
 representing prob that var x_i is in state s given
 parents of x_i are of states t .

Using Lagrangian

$$L = \sum_{n=1}^N \sum_{i=1}^K \log p(x_i^n | pa(x_i^n)) + \sum_{i=1}^K \sum_{t \in S} \lambda_{ti} (1 - \sum_s \theta_{s^i}(t))$$

To show $\theta_{s^i}(t) = \frac{\sum_{n=1}^N \mathbb{I}[x_i^n = s] \mathbb{I}[pa(x_i^n) = t]}{\sum_{n=1}^N \sum_s \mathbb{I}[x_i^n = s] \mathbb{I}[pa(x_i^n) = t]}$

\Rightarrow We can write $\theta_{s^i}(t)$ as indicator function

form
 i.e., $\theta_{s^i}(t) = p(x_i | pa(x_i)) \mathbb{I}[x_i = s] \mathbb{I}[pa(x_i) = t]$

$$\log p(x) = \sum_{n=1}^N \sum_{i=1}^K \log \left(\frac{x_i^n}{\sum_s p(x_i^n | pa(x_i))} \right)$$

Since we are using a discrete variable, the normalization

constant $= 1$

$$\Rightarrow \sum_s \theta_{s^i}(t) = \sum_s p(x_i = s | pa(x_i) = t) = 1$$

$\rightarrow \textcircled{1}$

The log likelihood function in the form of Oseen (5) be written as $\approx \log(\alpha_i^n / \text{ratio}^n)$

$$L = \sum_{i=1}^K \sum_{n=1}^N \underbrace{\mathbb{I}(x_i^n = s) \mathbb{I}(\text{pa}(x_i^n) = t^i)}_{+ \sum_{i=1}^K \sum_{t^i} \lambda^i(t^i) \left(1 - \sum_s \theta_s^i(t^i)\right)} \log \theta_s^i(t^i)$$

Now we differentiate w.r.t $\theta_s^i(t^i)$ and set the eq values to 0.

$$\frac{\partial L}{\partial \theta_s^i(t^i)} = 0$$

$$\rightarrow \frac{\partial L}{\partial \theta_s^i(t^i)} = \sum_{n=1}^N \mathbb{I}(x_i^n = s) \mathbb{I}(\text{pa}(x_i^n) = t^i) \frac{1}{\theta_s^i(t^i)} - \lambda^i(t^i) = 0$$

$$\Rightarrow \sum_{n=1}^N \mathbb{I}(x_i^n = s) \mathbb{I}(\text{pa}(x_i^n) = t^i) \frac{1}{\theta_s^i(t^i)} - \lambda^i(t^i) = 0$$

$$\Rightarrow \lambda^i(t^i) = \sum_{n=1}^N \mathbb{I}(x_i^n = s) \mathbb{I}(\text{pa}(x_i^n) = t^i) \frac{1}{\theta_s^i(t^i)}$$

or rearranging

$$\theta_s^i(t^i) = \frac{1}{\lambda^i(t^i)} \sum_{n=1}^N \mathbb{I}(x_i^n = s) \mathbb{I}(\text{pa}(x_i^n) = t^i)$$

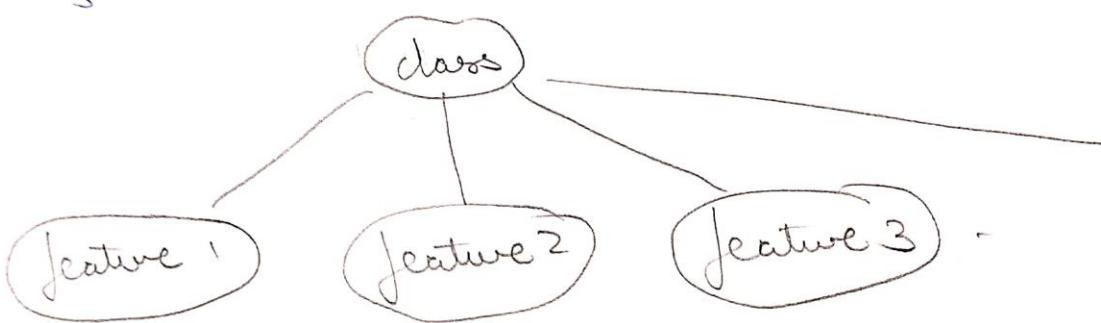
here $\lambda^i(t^i)$ is the normalizing constant

$$\Rightarrow \lambda^i(t^i) = 1 \Rightarrow \lambda^i(t^i) = \sum_s \theta_s^i(t^i)$$

$$\begin{aligned}
 \rightarrow \theta_{s^j}(t^s) &= \frac{\sum_{n=1}^N \mathbb{I}(\alpha_j^n = s) \mathbb{I}(\text{pa}(\alpha_j^n) = t^j)}{\sum_s \theta_{s^j}(t^s)} \\
 &= \frac{\sum_{n=1}^N \mathbb{I}(\alpha_j^n = s) \mathbb{I}(\text{pa}(\alpha_j^n) = t^j)}{\sum_s \mathbb{P}(\alpha_j = s | \text{pa}(\alpha_j) = t^j)} \\
 &= \frac{\sum_{n=1}^N \mathbb{I}(\alpha_j^n = s) \mathbb{I}(\text{pa}(\alpha_j^n) = t^j)}{\sum_s \mathbb{I}(\alpha_j = s) \mathbb{P}(\text{pa}(\alpha_j) = t^j)} \\
 &\quad \text{written in form of indicator functions} \\
 &\quad \text{taking sum all all } \alpha_j
 \end{aligned}$$

$$\begin{aligned}
 \theta_{s^j}(t^s) &= \frac{\sum_{n=1}^N \mathbb{I}(\alpha_j^n = s) \mathbb{I}(\text{pa}(\alpha_j^n) = t^j)}{\sum_{n=1}^N \sum_s \mathbb{I}(\alpha_j^n = s) \mathbb{I}(\text{pa}(\alpha_j^n) = t^j)}
 \end{aligned}$$

4] Graphical Bayesian model for Naive Bayes setup



$$p(\text{class}, \text{features}_1, \dots, n) = p(\text{class}) \prod_{i=1}^n p(\text{feature}_i | \text{class})$$

Input Test Data

`xd =`

```
1      0      0      1      1      1      1      0
```

Probability of Input test data being about politics is

`ProbOfPolitics =`

```
0.8306
```

`fx >>`

First I load all the Sequences for Politics and Sports.

Find mean of each row.

Then load the test datapoint.

Then I use the formula `probPolitics * (meanPolitics.^testdatapoints) .* (1-meanPolitics).^(1-testdatapoints)`

Probability of Sequence being Politics is = 0.8306

5)

```
Probability of a drum unit problem given the evidence is:  
First value for No and second value for yes using Joint probability found by me  
3.810115e-01 No  
6.189885e-01 Yes  
First value for No and second value for yes using multpot function  
3.810115e-01 No  
6.189885e-01 Yes
```

In this Question, I first make the network in the form of Brml array and then keep random probabilities to all. I have set a tolerance value of 0.0001 and number of iterations of EM to 30. It converges after almost 15-17 iterations.

Now we have new potentials of Each data. I have self found joint probability and then even used the multipot function of brml library. Both answers are matching.

The condition given is quality of paper is bad , but no bad smell or wrinkled paper. I have set that according, and used the condpot function of brml library and obtained the values of whether it is a drum problem or not.

Probability that it is drum problem is 0.6189885

Probability that it is not a drum problem 0.3810115

6] Structure learning of tree augmented Naive Bayes (7)
 Class variable C and feature variables X_1, \dots, X_n

a) Class variable C is connected to X_i 's and also X_i 's have one additional parent. As per the hint given by TM, this can be maximized by maximizing the BIC score.

$$\Delta(G_i) = \text{Score}_{\text{BIC}}(G_i : D) - \text{Score}_{\text{BIC}}(G'_i : D)$$

G_i is naive bayes graph where we have been given one parent as C and one additional parent say X_j where j is anything except the node C itself.

This can be maximized by

$$\Delta(G_i) = \text{Score}_{\text{BIC}}(G_i : D) - \text{Score}(G'_i : D)$$

G'_i is naive bayes graph with only C as parent

G_i is a graph with parents C and X_j as stated above

Let's define weight $w_{j \rightarrow i}$ as weight of the edge which connects X_j to X_i .

$$w_{j \rightarrow i} = \text{FanScore}(X_i | X_j, C : D) - \text{FanScore}(X_i | C : D)$$

We know Fan Score $(X_i | C : D)$

$$= m \left[\hat{\pi}_p(X_i | C) - \hat{\pi}(X_i) \right] - \frac{\log m}{2} (\text{Val}(X) - 1) \text{Val}(C)$$

On substituting we get

$$\begin{aligned}
 w_{j \rightarrow i} &= m \left[\hat{I}_P(x_i; x_j, c) - \frac{\hat{H}(x_j)}{2} \right] \\
 &\quad - \frac{\log m}{2} (\text{Val}(x_i) - 1) (\text{Val}(x_j) - 1) |c| \\
 &\quad - m \left(\hat{I}_P(x_i; c) - \frac{\hat{H}(x_i)}{2} \right) \\
 &\quad + \frac{\log m}{2} (\text{Val}(x_i) - 1) |c| \\
 &= m \left[\hat{I}_P(x_i; x_j, c) - \hat{I}_P(x_i; c) \right] \\
 &\quad - \frac{\log m}{2} (x_i \cdot 2) + \frac{\log m}{2} (x_j \cdot 2)
 \end{aligned}$$

as all x_i, x_j, c are binary (2 possibilities)

$$\Rightarrow w_{j \rightarrow i} = m \left[\hat{I}_P(x_i; x_j, c) - \hat{I}_P(x_i; c) \right] - \frac{\log m}{2}$$

Now, we need to prove $w_{j \rightarrow i} = w_{i \rightarrow j}$

$$\begin{aligned}
 &\hat{I}_P(x_i; x_j, c) - \hat{I}_P(x_i; c) \\
 &= H(x_i) + H(x_j, c) - H(x_i, x_j, c) \\
 &\quad - H(x_i) - H(c) + H(x_i, c) \\
 &\text{adding } H(x_j) \\
 &= H(x_j) + H(x_j, c) - H(x_i, x_j, c) \\
 &\quad - H(x_j) - H(c) + H(x_i, c) \\
 &= I(x_j; x_i, c) - \hat{I}(x_j; c)
 \end{aligned}$$

$$\Rightarrow \boxed{w_{i \rightarrow j} = w_{j \rightarrow i}}$$

Thus, an undirected graph with x_i nodes and edge weights connecting to other x_j node can be made.

we can now choose x_i to be root and by using maximum spanning tree algo, the graph can be connected to a directed graph and then finally a edge can be added from $(\text{to } x_i)$.

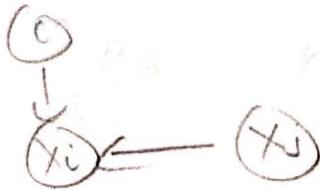
b) we must allow more flexibility for each x_i
we need to decide whether to add c as parent or not. So, we must consider all incoming connections to node x_i .

These 3 possibilities can be

a) only c as parent



b) parent x_j as parent



c) only x_j as parent



maximize

$$\Delta G_i = \text{Score with parents} - \text{Score without parents}$$

Case① only C as parent

$$\Delta(G) = \text{fanScore}(x_i | C:D) - \text{fanScore}(x_i : D)$$

$$= m \left(\hat{I}_p(x_i; C) - \frac{\hat{I}(x_i)}{D} \right) - \frac{\log(m)}{2} + \frac{\log(m)}{2}$$

$$\boxed{\Delta(G) = m \left(\hat{I}_p(x_i; C) \right) - \frac{\log(m)}{2}}$$

Case②

x_j and C as parents

$$\Delta(G) = \text{fanScore}(x_i | x_j, C:D) - \text{fanScore}(x_i : D)$$

$$= m \left[\hat{I}_p(x_i; x_j, C) \right] - 2\log(m) + \frac{\log(m)}{2}$$

$$\boxed{\Delta(G) = m \left[\hat{I}_p(x_i; x_j, C) \right] - 1.5 \log(m)}$$

Case③ only C as parent

$$\Delta G = \text{fanScore}(x_i | C:D) - \text{fanScore}(x_i : D)$$

$$= m \left(\hat{I}_p(x_i; C) \right) - \log(m) + \frac{\log(m)}{2}$$

$$\boxed{\Delta G = m \left(\hat{I}_p(x_i; C) \right) - \frac{\log(m)}{2}}$$

So, now we consider a graph which connects node i to x_i and all $x_j \rightarrow x_i$ for all nodes x_i of graph. (2)

→ We have a fully connected directed graph which if x_i taken as root, has x_j as parent as well as all x_j 's are parents.

weights can be set as

$$w_{c \rightarrow i} = \min \Delta(G)_c$$

$$w_{j \rightarrow i} = \begin{cases} \text{maximum of } (\Delta(G_j) \text{ when } x_j \text{ and } i \text{ are parents}) \\ \text{or } \Delta(G_j) \text{ when only } x_j \text{ is parent} \end{cases}$$

→ This is clearly seen as a directed graph.

→ This is clearly seen as a directed graph.
so we can run the maximum weighted
directed spanning tree algo to get a
Bayesian network. Due to regularization in BIC
scoring, there are possibility of negative scores
as well. We need to 0 them out. Thus a
graph with maximum difference in BIC scores
is created as this is a sum of weights
of the graph.

Thus after finding the max. directed spanning graph F , the parent set of each feature x_i in the resulting network can be defined as

$$\text{Parent in graph w.r.t } x_i = \begin{cases} \{x_j\} & \text{if } \leftarrow x_i \in F \\ \{x_j, y\} & \text{if } x_j \rightarrow x_i \text{ in } F \text{ and } \Delta_{x_j, y} > \Delta_{x_j, c} \\ \{x_j, c\} & \text{if } x_j \rightarrow x_i \text{ in } F \text{ and } \Delta_{x_j, c} > \Delta_{x_j, y} \\ \emptyset & \text{none of the above are satisfied.} \end{cases}$$

Scores

7 and 8

(10)

① sequence v_1, \dots, v_T & s_1, \dots, s_T is generated by marker chain. Single chain of length T , we have

$$p(v_1, \dots, v_T) = p^{(v_1)} \prod_{c=1}^{T-1} p(v_{t+1} | v_c)$$

The sequence of visible variables is given by

$$v = (v_1, \dots, v_T)$$

For a single marker chain labelled by h ,

$$p(v|h) = p(v_1|h) \prod_{c=1}^{T-1} p(v_{t+1}|v_t, h) \rightarrow ①$$

In total, there are a set of H such marker chains ($h=1 \dots T$)

$$\text{Distri of visible variables is } p(v) = \sum_{h=1}^H p(v|h) \rightarrow ②$$

set of latent variables v^n , $n=1 \dots N$

" independently drawn.

" energy function of $p(v^n)$ is given as:

The energy function

$$\sum_n \log p(v^n) \geq \text{potl}(h|v^n)$$

from ②

$$= \sum_n \log p(v^n|h) p(h) \geq \text{potl}(h|v^n)$$

from ①

$$\sum_n \langle \log p(v^n) \rangle_{\text{old}(h|v^n)} = \sum_{t=1}^{T-1} \sum_h \langle \log (v_{t+1}^n / v_t^n, h) \rangle_{\text{old}(h|v^n)}$$

↑ in log π changes to Σ

$$+ \sum_n \langle \log p(v_t^n | h) \rangle_{\text{old}(h|v^n)}$$

$$+ \sum_{t=1}^{T-1} \sum_n \langle \log p(h) \rangle_{\text{old}(h|v^n)}$$

This needs to be optimized and is possible by differentiating with variations

$$p(v_{t+1} = i | v_t = j, h) = \text{let this be } \Theta_{ij}^h$$

i = new state
 j = old state

Initial distribution is given by $p(v_0 | h)$
and $p(h)$ is prior, so these become 0 or give no contributions to variations

The contribution of transitions of energy with language term is

$$L = \sum_n \sum_{t=1}^{T-1} \sum_{i,j,h} p_{\text{old}}(h|v^n) \prod_{\substack{(v_{t+1} = i, v_t = j) \\ \text{states}}} \log \Theta_{ij}^h$$

$$+ \sum_n \sum_{h,j} \sum_i (1 - \Theta_{ij}^h)$$

calculated by Σ language

$$\Rightarrow \sum_n \sum_{t=1}^{T-1} p^{\text{old}}(h|v^n) \mathbb{I}[v_t^n = j, v_{t+1}^n = i] \frac{\partial}{\partial_{ij}} - \lambda_j = 0$$

$$d_j = \frac{\sum_{n=1}^T \sum_{t=1}^{T-n} P_{\text{old}}(h|v^n) \prod_{k=t+1}^n [v_k^n = j, v_{k+1}^n = i]}{\sum_{i,j} d_{ij}}$$

On rearranging

$$\alpha_{ij} = \frac{\sum_{t=1}^{T-1} p^{\text{old}}(h|v^n) \prod_{t'=j}^{i-1} (v_{t'} = j, v_{t'+1} = i)}{\sum_{k=1}^n \sum_{t=1}^{T-1} p^{\text{old}}(h|v^n) \prod_{t'=j}^{i-1} (v_{t'} = k, v_{t'+1} = i)}$$

we have used normalization, as this is discrete,

$$\Rightarrow \sum_i \alpha_{ij} = 1$$

$$\Rightarrow \boxed{O_{ij}} = \frac{\sum_{t=1}^{T-1} p^{\text{ord}}(h|v^n) \mathbb{I}(v_t^n=j, v_{t+1}^n=i)}{\sum_i \sum_n \sum_{t=1}^{T-1} p^{\text{ord}}(h|v^n) \mathbb{I}(v_t^n=j, v_{t+1}^n=i)}$$

This gives us the weighted numbers of j to i transitions across all the sequences.

variations across
Also, initial destruction plus initial energy and initial destruction D_i
is given by

By following previous steps

$$\boxed{p_i = \frac{\sum_{v^n} p^{\text{old}}(h|v^n) \prod [u_i^n = i]}{\sum_{i^n} \sum_{v^n} p^{\text{old}}(h|v^n) \prod [u_i^n = i]}}$$

The optimal prior $p(h)$ is given by

$$\boxed{p(h) = \frac{\sum_n p^{\text{old}}(h|v^n)}{\sum_h \sum_n p^{\text{old}}(h|v^n)}}$$

8)

Sequence belongs to Group 1

Sequence is: CATAGGCATTCTATGTGCTG
Sequence is: CCAGTTACGGACGCCGAAAG
Sequence is: CGGCCGCGCCTCCGGGAACG
Sequence is: ACATGAACTACATAGTATAA
Sequence is: GTTGGTCAGCACACGGACTG
Sequence is: CACTACGGCTACCTGGGCAA
Sequence is: CGGTCCGTCCGAGGCACTCG
Sequence is: CACCATCACCCCTTGCTAAGG
Sequence is: CAAATGCCTCACGCGTCTCA
Sequence is: GCCAAGCAGGGTCTCAACTT
Sequence is: CATGGACTGCTCCACAAAGG

Sequence belongs to Group 2

Sequence is: TGGAACCTTAAAAAAAAAAAAA
Sequence is: GTCTCCTGCCCTCTCTGAAC
Sequence is: GTGCCTGGACCTGAAAAGCC
Sequence is: AAAGTGCTCTGAAAACTCAC
Sequence is: CCTCCCCCTCCCCTTCCTGC
Sequence is: TAAGTGTCCCTGCTCCTAA
Sequence is: AAAGAACTCCCCTCCCTGCC
Sequence is: AAAAAAACGAAAAACCTAAG
Sequence is: GCGTAAAAAAAGTCCTGGGT

Sequence is:
CATAGGCATTCTATGTGCTG
Sequence belongs to Group 1
Sequence is:
CCAGTTACGGACGCCGAAAG
Sequence belongs to Group 1
Sequence is:
TGGAACCTTAAAAAAA
Sequence belongs to Group 2
Sequence is:
GTCTCCTGCCCTCTCTGAAC
Sequence belongs to Group 2
Sequence is:
GTGCCTGGACCTGAAAAGCC
Sequence belongs to Group 2 |
Sequence is:
CGGCCGCGCCTCGGGAACG
Sequence belongs to Group 1
Sequence is:
AAAGTGCTCTGAAAATCAC
Sequence belongs to Group 2
Sequence is:
ACATGAACTACATAGTATAA
Sequence belongs to Group 1
Sequence is:
GTTGGTCAGCACACGGACTG
Sequence belongs to Group 1
Sequence is:
CCTCCCCCTCCCCCTTCCTGC
Sequence belongs to Group 2
Sequence is:
CACTACGGCTACCTGGGCAA
Sequence belongs to Group 1
Sequence is:
CGGTCCGTCCGAGGGACTCG
Sequence belongs to Group 1
Sequence is:
TAAGTGTCTCTGCTCCTAA
Sequence belongs to Group 2
Sequence is:
CACCATCACCTTGCTAAGG
Sequence belongs to Group 1

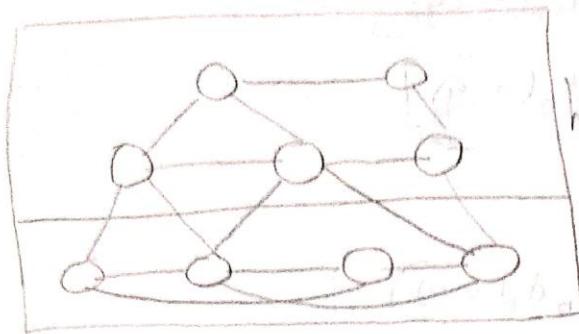
```
Sequence is:  
GCCAAGCAGGGTCTCAACTT  
Sequence belongs to Group 1  
Sequence is:  
CATGGACTGCTCCACAAAGG  
Sequence belongs to Group 1  
Sequence is:  
AAAAAAACGAAAAACCTAAG  
Sequence belongs to Group 2  
Sequence is:  
GCGTAAAAAAAGTCCTGGGT  
Sequence belongs to Group 2  
Log Likelihood is
```

```
loglikelihood =  
  
-483.6352
```

In this question I have just loaded the sequences.mat file. The visible nodes, A, C,G , T are set as 1,2 3 4 respectively. I then call the mixMMarkov function with iterations of 50.

Then in this function we get a value phgv if that is above 0.5 then the sequence belongs to cluster 2 else it belongs to cluster one. Results are attached as above.

9



6

hidden layer

output (visible) layer

5 hidden units, 4 visible units

$$P(x, y) = \frac{1}{Z} \exp \left(\sum_k \theta_k \phi_k(x, y) \right)$$

$\phi_k(x, y)$ are pairwise potentials for Boltzmann machines
 θ_k are weights of BMs

$$\text{To show } \frac{\partial \log P(x)}{\partial \theta_k} = \sum_y \phi_k(x, y) P(y|x) - \sum_{x,y} \phi_k(x, y) P(x, y)$$

\Rightarrow we know Z is normalization constant

$$\Rightarrow P(x, y) = \exp \left(\sum_k \theta_k \phi_k(x, y) \right)$$

$$\Rightarrow P(x, y) = \frac{\exp \left(\sum_k \theta_k \phi_k(x, y) \right)}{\sum_{x,y} \exp \left(\sum_k \theta_k \phi_k(x, y) \right)}$$

$$\Rightarrow Z = \sum_{x,y} \exp \left(\sum_k \theta_k \phi_k(x, y) \right) \quad P(x, y)$$

\Rightarrow to get $P(x)$ we have to marginalize $P(x, y)$

$$P(x) = \sum_y P(x, y)$$

$$P(x) = \sum_y \frac{\exp \left(\sum_k \theta_k \phi_k(x, y) \right)}{\sum_{x,y} \exp \left(\sum_k \theta_k \phi_k(x, y) \right)}$$

$$P(x) = \frac{\sum_y \exp\left(\sum_k \theta_k \phi_k(x, y)\right)}{\sum_{x,y} \exp\left(\sum_k \theta_k \phi_k(x, y)\right)}$$

taking log

$$\log P(x) = \log \left[\frac{\sum_y \exp\left(\sum_k \theta_k \phi_k(x, y)\right)}{\sum_{x,y} \exp\left(\sum_k \theta_k \phi_k(x, y)\right)} \right]$$

$$= \log \left[\frac{\sum_y \exp\left(\sum_k \theta_k \phi_k(x, y)\right)}{\sum_{x,y} \exp\left(\sum_k \theta_k \phi_k(x, y)\right)} \right] - \log \left[\sum_{x,y} \exp\left(\sum_k \theta_k \phi_k(x, y)\right) \right]$$

differentiating w.r.t θ_k

$$\begin{aligned} \frac{\partial \log P(x)}{\partial \theta_k} &= \frac{\partial}{\partial \theta_k} \left[\log \left(\frac{\sum_y \exp\left(\sum_k \theta_k \phi_k(x, y)\right)}{\sum_{x,y} \exp\left(\sum_k \theta_k \phi_k(x, y)\right)} \right) \right] \\ &= \sum_y 1 \cdot \frac{\exp\left(\sum_k \theta_k \phi_k(x, y)\right) \times \phi_k(x, y)}{\sum_{x,y} \exp\left(\sum_k \theta_k \phi_k(x, y)\right)} \rightarrow P(x, y) \\ &\quad - \sum_y \frac{\exp\left(\sum_k \theta_k \phi_k(x, y)\right)}{\sum_{x,y} \exp\left(\sum_k \theta_k \phi_k(x, y)\right)} \times \frac{\partial \phi_k(x, y)}{\partial \theta_k} \rightarrow P(x, y) \end{aligned}$$

$$\boxed{\frac{\partial \ln P}{\partial x} = \left(\frac{\partial \ln P}{\partial x} \right)_{\theta_k}}$$

$$P(x, y) = \frac{1}{2} \exp\left(\sum_k \theta_k \phi_k(x, y)\right)$$

$$P(x) = \frac{1}{2} \sum_y \exp\left(\sum_k \theta_k \phi_k(x, y)\right)$$

$$\frac{\partial \log P(x)}{\partial \theta_k} = \sum_y \left(\frac{P(x, y)}{P(x)} \right) \phi_k(x, y) - \sum_{x,y} P(x, y) \phi_k(x, y)$$

$$\boxed{\frac{\partial \log P(x)}{\partial \theta_k} = \sum_y P(y/x) \phi_k(x, y) - \sum_{x,y} P(x, y) \phi_k(x, y)}$$

$$\boxed{\frac{\partial \log P(x)}{\partial \theta_k} = \sum_y P(y/x) \phi_k(x, y) - \sum_{x,y} P(x, y) \phi_k(x, y)}$$