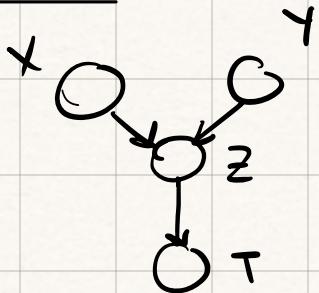


1) DGM



Factorization is $p(x_v) = \prod_i p(x_i | x_{\neg i})$

$$\Rightarrow p(x_v) = p(x)p(y)p(z|x,y)p(t|z)$$

Also, $X \perp\!\!\!\perp Y | T \Leftrightarrow X$ and Y are d-separated by T .

We have seen that sets X and Y are d-separated by set D if

1) $d \in D$ and (v_{i-1}, d, v_{i+1}) is not a V-structure

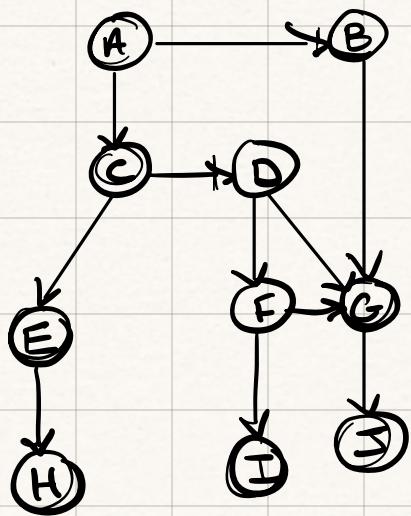
or 2) $d \notin D$ and (v_{i-1}, d, v_{i+1}) is a V-structure but no descendants of d is in D .

1) is not respected because $X-Y-Z$ is a V-structure.

2) is not respected, because although $X-Y-Z$ is a V-structure there is a descendant of $d \in Z$ in D (t is T).

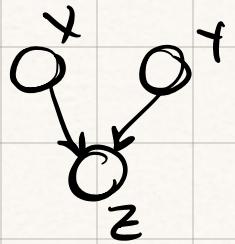
So $X \perp\!\!\!\perp Y | T$ is not true in general.

2) D-separation in DGM



- a) $C \perp\!\!\!\perp B$ No
- b) $C \perp\!\!\!\perp B | A$ Yes
- c) $C \perp\!\!\!\perp B | A, J$ No
- d) $C \perp\!\!\!\perp B | A, D, J$ Yes
- e) $C \perp\!\!\!\perp G$ No
- f) $C \perp\!\!\!\perp G | B$ No
- g) $C \perp\!\!\!\perp G | B, D$ Yes
- h) $C \perp\!\!\!\perp G | B, D, H$ Yes
- i) $C \perp\!\!\!\perp G | B, D, H, E$ Yes
- j) $B \perp\!\!\!\perp I | J$ No

3. Positive interactions in V structure



$$\alpha := P(X=1), \beta := P(X=1 | Z=1)$$

$$\gamma := P(X=1 | Z=1, Y=1).$$

a)

(i) Want $\gamma < \alpha$.

With

$$\begin{array}{|c|} \hline P(X=1) \\ \hline 0.5 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline P(Y=1) \\ \hline 0.5 \\ \hline \end{array}$$

$$\text{then } \beta = P(X=1 | Z=1) = \frac{P(X=1, Z=1)}{P(Z=1)} \quad (\text{product rule, } P(X, Y) = P(X)P(Y))$$

$$= \frac{\sum_{i=0,1} P(X=1, Y=i, Z=1)}{\sum_{i,j=0,1} P(X=i, Y=j, Z=1)} \quad (\text{marginalization})$$

$$= \frac{\sum_{i=0,1} P(X=i) P(Y=i) P(Z=1 | X=i, Y=i)}{\sum_{i,j=0,1} P(X=i) P(Y=j) P(Z=1 | X=i, Y=j)} \quad (\text{graph factorization})$$

$$= \frac{0.25(P(Z=1 | X=1, Y=0) + P(Z=1 | X=1, Y=1))}{\sum_{j=0,1} 0.5 P(Y=j) P(Z=1 | X=0, Y=j) + 0.5 P(Y=j) P(Z=1 | X=1, Y=j)}.$$

$$\sum_{j=0,1} 0.5 P(Y=j) P(Z=1 | X=0, Y=j) + 0.5 P(Y=j) P(Z=1 | X=1, Y=j).$$

(*) .

$$\Rightarrow \beta = \frac{P(Z=1 | X=1, Y=0) + P(Z=1 | X=1, Y=1)}{P(Z=1 | X=0, Y=0) + P(Z=1 | X=0, Y=1) + P(Z=1 | X=1, Y=0) + P(Z=1 | X=1, Y=1)}$$

IP	X	Y	$P(Z=1 X,Y)$
	1	1	0.1
	1	0	0.1
	0	1	1
	0	0	1

Then, $\beta = \frac{0.1+0.1}{0.1+0.1+1+1} = 0.091$

Also, $\gamma = P(X=1|Z=1, Y=1) = \frac{P(X=1, Y=1, Z=1)}{P(Z=1, Y=1)}$ (product rule)

$$= \frac{P(X=1)P(Y=1)P(Z=1|X=1, Y=1)}{\sum_{i=0,1} P(X=i, Y=1, Z=1)} \quad (\text{graph fact } \& \text{ marginalization})$$

$$= \frac{0.25P(Z=1|X=1, Y=1)}{P(X=0)P(Y=1)P(Z=1|X=0, Y=1) + P(X=1)P(Y=1)P(Z=1|X=1, Y=1)}$$

(graph factorization)

$\Rightarrow \gamma = \frac{P(Z=1|X=1, Y=1)}{P(Z=1|X=0, Y=1) + P(Z=1|X=1, Y=1)}$ (gx)

In this case

$$\gamma = \frac{0.1}{1+0.1} = 0.091$$

Z	β	γ
0.5	0.091	0.091

and $\gamma < \beta$.

i) $\alpha < \gamma < \beta$. Again let me set

$$\begin{array}{|c|} \hline P(X=1) \\ 0.5 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline P(Y=1) \\ 0.5 \\ \hline \end{array}$$

Equations (x) for β & (y) for γ are valid again.

$$\beta = \frac{P(Z=1 | X=1, Y=0) + P(Z=1 | X=1, Y=1)}{P(Z=1 | X=0, Y=0) + P(Z=1 | X=1, Y=0) + P(Z=1 | X=0, Y=1) + P(Z=1 | X=1, Y=1)}$$

$$\gamma = \frac{P(Z=1 | X=1, Y=1)}{P(Z=1 | X=0, Y=1) + P(Z=1 | X=1, Y=1)}$$

With.

X	Y	$P(Z=1 X, Y)$
1	1	1
1	0	1
0	1	0.2
0	0	0.1

I have $\beta = \frac{1+1}{1+1+0.2+0.1} \approx 0.87$

$$\gamma = \frac{1}{1+0.2} \approx 0.83$$

So $\begin{array}{ccc} \alpha & \beta & \gamma \\ 0.5 & 0.87 & 0.83 \end{array}$ and $\alpha < \gamma < \beta$.

iii) $\beta < \alpha < \gamma$

X	Y	$P(Z=1 X, Y)$
1	1	1
1	0	0.1
0	1	0.2
0	0	1

$$\beta = \frac{P(Z=1|X=0, Y=0) + P(Z=1|X=1, Y=1)}{P(Z=1|X=0, Y=0) + P(Z=1|X=0, Y=1) + P(Z=1|X=1, Y=0) + P(Z=1|X=1, Y=1)}$$

$$\gamma = \frac{P(Z=1|X=1, Y=1)}{P(Z=1|X=0, Y=1) + P(Z=1|X=1, Y=1)}$$

$$\beta = \frac{1+0.1}{1+0.1+0.2+1} \approx 0.478$$

$$S_0 \quad \gamma = \frac{1}{1+0.2} \approx 0.833$$

α	β	γ
0.5	0.478	0.833

and $\beta < \alpha < \gamma$.

b).

i) Say $Z = \text{late}$? with $Z=1$ is Yes & $Z=0$ is No

$X = \text{Has watch?}$ with $X=1$ Yes & $X=0$ No

and $Y = \text{Ate cereal?}$ with $Y=1$ Yes & $Y=0$ No

Conditionals were

X	Y	$P(Z=1 X, Y)$	
1	1	0.1	has watch, ate cereal
1	0	0.1	has watch, didn't eat cereal
0	1	1	forgot watch, ate cereal
0	0	1	forgot watch, didn't eat cereal

$P(X=1) = 0.5$, but we see that we've forgotten his watch, the probability of being late is very high ($P(Z=1|X=0, Y=i) = 1$). So it makes sense that

given the fact that he is late, the probability that he has his watch is smaller, i.e.

$$P(X=1|Y=1, Z=1) < P(X=1).$$

b)

X	Y	$P(Z=1 X,Y)$
1	1	1
1	0	1
0	1	0.2
0	0	0.1

Say $Z=\text{late}$? with $Z=1$ is Yes & $Z=0$ is No

$X=\text{Forget watch}$? with $X=1$ Yes & $X=0$ No

$Y=\text{Ate cereal}$? with $Y=1$ Yes & $Y=0$ No

Because it is certain that he is late when he forgets his watch, it is clear that the probability that he forget his watch is bigger if we knew that he is late. That is

$$P(X=1|Z=1), P(X=1|Z=1, Y=1) > P(X=1).$$

$$\beta, \gamma > \alpha.$$

Also, we he has his watch, he is more likely to be late.

when $Y=1$, so

$$P(X=1|Z=1, Y=1) < P(X=1|Z=1)$$

$$\gamma \qquad \qquad \qquad \beta$$

$$\Rightarrow \alpha < \gamma < \beta.$$

c) $X \setminus Y$		$P(Z=1 X, Y)$
1	1	1
1	0	0.1
0	1	0.2
0	0	1

We knew that $Z=1$ is certain to happen if $X=Y=1$ or $X=Y=0$ at the same time, and unlikely in the other cases.

So given $Y=1$ & $Y=0$, it is much more probable that $X=1$. This is why

$$\gamma > \alpha, \beta.$$

If we just knew that $Z=1$, then the probability that $X=1$ is smaller than just $Y=1$

$$\text{So } \beta < \alpha < \gamma.$$

4) Equivalence of directed tree w/ undirected tree.

Let $G(V, E)$ be a directed tree.

Since G is a directed tree, every node in it has at most one parent.

Now let \bar{G} be the neutralized graph of G . Then it has the same node set $J = V$ & an edge set

$$\bar{E} = \left\{ \{i, j\} : (i, j) \in E \right\} \cup \left\{ \{k, l\} : k, l \in T_i, k \neq l, \forall i \in V \right\}$$

same edges as E
but undirected.

but since any node in V has at most one parent, $\Pi_i : i \leq 1$, and we cannot find two nodes $K \neq L$ in Π_i s.t. $K \neq L$, so $\{\{K, L\} : K, L \in \Pi_i, K \neq L, \forall i \in V\} = \{\emptyset\}$ (empty)

and $\bar{E} = E$ but undirected edges.

Now, any $p \in \mathcal{L}(G)$ factorizes as

$$p(X_V) = \prod_{i=1}^N p(X_i | X_{\Pi_i}) \quad \forall i \in V$$

but since $|\Pi_i| \leq 1$, this is a product of factors $p(X_i | X_K)$ $\forall i \in V$ with $K \in \Pi_i$ is at most one node. I.e. we just have factors between a parent & its children

Also, any $p' \in \mathcal{L}(\bar{G})$ factorizes as

$$p'(X_{\bar{V}}) = \frac{1}{Z} \prod_C \phi_c(X_c)$$

but since $\bar{E} = E$ but undirected, the cliques C in \bar{G} link just two nodes, so it is a product of factors $\phi_c(X_1, X_2)$. where $X_1 \neq X_2$ are parent & children.

These are the same conditional independence statements as $\mathcal{L}(G)$. So if we have some $p \in \mathcal{L}(G)$ then we have $p \in \mathcal{L}(\bar{G})$, and vice versa.

$$\Rightarrow \mathcal{L}(G) = \mathcal{L}(\bar{G}).$$

5.) Hammersley-Clifford Counter example

$$G: X_1 - X_2 - X_3 - X_4 - X_1$$

Probabilities that have $p=1/8$:

$$(0,0,0,0), (1,0,0,0), (1,1,0,0), (1,1,1,0), (0,0,0,1), (0,0,1,1) \\ (0,1,1,1), (1,1,1,1)$$

all others have zero.

Assuming we can show that Global Markov property is satisfied,

Let's show that p cannot factorize according to G :

Usual factorization in a UGM:

$$p(X) = \frac{1}{2} \prod_c \psi_c(x_c) = \frac{1}{2} (\psi_{12}(X_1, X_2) \psi_{23}(X_2, X_3) \psi_{34}(X_3, X_4) \psi_{41}(X_4, X_1))$$

We know

$$p(0,0,0,0) \sim \psi_{12}(0,0) \psi_{23}(0,0) \psi_{34}(0,0) \psi_{41}(0,0) = \frac{1}{8}$$

and

$$p(0,0,1,1) \sim \psi_{12}(0,0) \psi_{23}(0,1) \psi_{34}(1,1) \psi_{41}(1,0) = \frac{1}{8}.$$

This means, in particular, that $\psi_{12}(0,0)$, $\psi_{23}(0,1)$ and $\psi_{41}(0,0)$ cannot be zero (otherwise $p \neq 1/8$, but zero).

Now, consider

$$p(0,0,1,0) \sim \psi_{12}(0,0) \psi_{23}(0,1) \psi_{34}(1,0) \psi_{41}(0,0) = 0.$$

\uparrow \uparrow \uparrow
 Net zero Net zero Net zero

so $\psi_{34}(1,0)$ has to be zero.

$$\text{but consider } p(1,1,1,0) \sim \psi_{12}(1,1) \psi_{23}(1,1) \psi_{34}(1,0) \psi_{41}(0,1) = \frac{1}{8}.$$

but if $\psi_{34}(1,0)$ is zero, $p(1,1,1,0)$ is zero... Contradiction !!

p cannot factorize according to $G \Rightarrow p \notin \mathcal{L}(G)$, which contradicts the Hammersley-Clifford theorem.

(6) Bizarre conditional independence statement.

Attempt ...

$$\begin{aligned}
 p(x,y,z) &= p(z)p(x,y|z) \\
 &= p(z) \underbrace{p(x|z)p(y|z)}_{\text{(product rule)}} \\
 &= \frac{p(y,z)}{p(z)} \quad \text{by chain rule} \\
 &= p(x|z)p(y|z). \quad (\text{i})
 \end{aligned}$$

$$\text{but also } p(x,y,z) = p(z) \underbrace{p(x|z)p(y|z)}$$

$$\begin{aligned}
 &\frac{p(x,z)}{p(z)} \\
 &= p(y|z)p(x,z) \quad (\text{ii})
 \end{aligned}$$

$$(\text{i}) = (\text{ii})$$

$$p(x|z)p(y,z) = p(y|z)p(x,z)$$

Assume $x \perp\!\!\!\perp z$ but $y \not\perp\!\!\!\perp z$

$$\begin{aligned}
 \text{then } p(x|z) &= p(x) \Rightarrow p(x)p(y,z) = p(y|z)p(x)p(z) \\
 &\Rightarrow p(y,z) = p(y|z).
 \end{aligned}$$

which is not true in general. Similarly, if

$y \perp\!\!\!\perp z$ but $x \not\perp\!\!\!\perp z$

$$\text{I get } p(x,z) = p(x|z)$$

which is also not true in general.

So no, it is not true that $(X \amalg Z)$ or $(Y \amalg Z)$ in general.

This is the answer to b.

$$7.) p(x_i | e) = \sum_{k=1}^K \pi_k N(x_i | \vec{\mu}_k, \sigma_k^2 I).$$

Let $l_c(\theta; D_c)$ be the complete log likelihood, where

D_c is our "dataset" $D_c = \{(x_i, z_i) : i = 1, \dots, N\}$.

$$\begin{aligned} l_c(\theta; D_c) &= \sum_{i=1}^N \log p(x_i, z_i; \theta) \\ &= \sum_{i=1}^N \log \prod_{k=1}^K (\pi_k N(x_i | \vec{\mu}_k, \sigma_k^2 I))^{z_i^k} \\ &= \sum_{i=1}^N \sum_{k=1}^K z_i^k \log \pi_k N(x_i | \vec{\mu}_k, \sigma_k^2 I) \end{aligned}$$

E-step:

$$\mathbb{E}_{\theta^{(t)}}[l_c(\theta; D_c)] = \sum_{i=1}^N \sum_{k=1}^K z_i^k \log \pi_k N(x_i | \vec{\mu}_k, \sigma_k^2 I)$$

$$\boxed{\mathbb{E}_{\theta^{(t)}}[l_c(\theta; D_c)] = \sum_{i=1}^N \sum_{k=1}^K z_i^k \log \pi_k N(x_i | \vec{\mu}_k, \sigma_k^2 I)}$$

and define $\tau_i^k = p(z_i^k = 1 | x, e)$

$$\text{Bayes rule: } \tau_i^{k(t)} = \frac{p(x_i | z_i^k = 1, \theta_K^{(t)}) p(z_i^k = 1 | \pi_i)}{p(x_i | \theta_K^{(t)})}$$

$$\Rightarrow \tau_i^{k(t)} = \frac{\pi_k N(x_i | \vec{\mu}_k, \sigma_k^2 I)}{\sum_j \pi_j N(x_i | \vec{\mu}_j, \sigma_j^2 I)}$$

M-step is MLE for θ using $\tau_i^{k(t)}$ instead of $\sum \tau_i^k \gamma_{\theta}^{(t)}$
on $\sum \tau_i \gamma_{\theta}^{(t)}$

MLE for μ_k : take derivative w.r.t. $\vec{\mu}_j$ & set to 0.

$$\frac{\partial}{\partial \vec{\mu}_j} \sum \tau_i^{k(t)} (\log \pi_k + \log \frac{1}{(2\pi)^{d/2} |\sigma_k^2 I|^{1/2}} \exp(-\frac{1}{2} (x_i - \vec{\mu}_k)^T \sigma_k^{-2} I (x_i - \vec{\mu}_k))) = 0$$

$$\Rightarrow \frac{\partial}{\partial \vec{\mu}_j} \sum_{k,i} \tau_i^{k(t)} \left(-\frac{1}{2} (\log 2\pi^d + \log |\sigma_k^2 I|) - \frac{1}{2} (x_i - \vec{\mu}_k)^T \sigma_k^{-2} I (x_i - \vec{\mu}_k) \right) = 0.$$

$$\Rightarrow \frac{\partial}{\partial \vec{\mu}_j} \sum_{k,i} \tau_i^{k(t)} \left(-\frac{1}{2} (x_i - \vec{\mu}_k)^T \sigma_k^{-2} I (x_i - \vec{\mu}_k) \right) = 0.$$

Use $\frac{\partial}{\partial x} x^T A x = (A + A^T)x$

$$\Rightarrow \sum_{i=1}^N \tau_i^{k(t)} \underbrace{\left(-\frac{1}{2} (\sigma_j^{-2} I + (\sigma_j^{-2} I)^T) (x_i - \vec{\mu}_j) \right)}_{\sigma_j^{-2} I \text{ symmetric}} = 0.$$

$$\Rightarrow (\sigma_j^{-2} I)^T = \sigma_j^{-2} I$$

$$\Rightarrow \sum_{i=1}^N \tau_i^{k(t)} \cancel{\sigma_j^{-2} I} x_i = \sum_{i=1}^N \tau_i^{k(t)} \cancel{\sigma_j^{-2} I} \vec{\mu}_j$$

$$\Rightarrow \boxed{\vec{\mu}_j^{(t+1)} = \frac{\sum_{i=1}^N \tau_i^{k(t)} x_i}{\sum_{i=1}^N \tau_i^{k(t)}}}$$

Now, for σ_j^2 (dropping terms that do not depend on σ_k^2)

$$\frac{\partial}{\partial \sigma_j^2} \mathcal{L}(\theta, \Sigma) \Big|_{\theta^{(t)}} = \frac{\partial}{\partial \sigma_j^2} \sum_{i,k} \tau_i^{k(t)} \left(\frac{1}{2} \log |\sigma_k^{-2} I| - \frac{1}{2} (\mathbf{x}_i - \vec{\mu}_k)^T \right.$$

$$\left. \sigma_k^{-2} I (\mathbf{x}_i - \vec{\mu}_k) \right) = 0$$

$$\Rightarrow -\frac{1}{2} \sum_{i,k} \tau_i^{k(t)} \left(\frac{\partial}{\partial \sigma_j^2} \log |\sigma_k^{-2} I| - \frac{\partial}{\partial \sigma_j^2} \underbrace{(\mathbf{x}_i - \vec{\mu}_k)^T \sigma_k^{-2} I (\mathbf{x}_i - \vec{\mu}_k)}_{\det \sigma_k^{-2} I = \prod_{i=1}^d \sigma_k^{-2} I} \right) = 0$$

(*)

The eigenvalues of $\sigma_k^{-2} I$. I is the $d \times d$ identity matrix,

so it has d eigenvalues of 1 so $\det \sigma_k^{-2} I = \prod_{i=1}^d \sigma_k^{-2} = (\sigma_k^{-2})^d$

$$\Rightarrow \frac{\partial}{\partial \sigma_j^2} \log |\sigma_k^{-2} I| = \frac{\partial}{\partial \sigma_j^2} \log (\sigma_k^{-2})^d = d \frac{\partial}{\partial \sigma_j^2} \log \sigma_k^{-2} = d \sigma_j^2$$

$(*) = (\mathbf{x}_i - \vec{\mu}_k)^T \sigma_k^{-2} I (\mathbf{x}_i - \vec{\mu}_k)$ is a scalar so it is equal to its trace

$$= \text{Tr} (\mathbf{x}_i - \vec{\mu}_k)^T \sigma_k^{-2} I (\mathbf{x}_i - \vec{\mu}_k)$$

Cyclic property of trace two times:

$$= \text{Tr} \sigma_k^{-2} I (\mathbf{x}_i - \vec{\mu}_k) (\mathbf{x}_i - \vec{\mu}_k)^T$$

$$= \text{Tr} \underbrace{\sigma_k^{-2} I (\mathbf{x}_i - \vec{\mu}_k) (\mathbf{x}_i - \vec{\mu}_k)^T}_{\equiv \tilde{\Sigma}_k}$$

$$= \text{Tr} \sigma_k^{-2} I \tilde{\Sigma}_k$$

$$\text{and } \frac{\partial}{\partial \sigma_j^2} \sum_k \text{Tr} \sigma_k^{-2} I \tilde{\Sigma}_k = \text{Tr} I \tilde{\Sigma}_j = \sum_j = (\mathbf{x}_i - \vec{\mu}_j) (\mathbf{x}_i - \vec{\mu}_j)^T$$

$$\text{So I have } \sum_{i=1}^N \tau_i^{j(t)} \left(d \sigma_j^2 - (\mathbf{x}_i - \vec{\mu}_j) (\mathbf{x}_i - \vec{\mu}_j)^T \right) = 0.$$

$$\Rightarrow d \sum_{i=1}^N \tau_i^{j(+)} \sigma_j^i = \sum_{i=1}^N \tau_i^{j(+)} (x_i - \vec{\mu}_j) (x_i - \vec{\mu}_j)^T$$

$$\Rightarrow \sigma_j^{Z(t+1)} = \frac{\sum_{i=1}^N \tau_i^{j(+)} (x_i - \vec{\mu}_j^{(t+1)}) (x_i - \vec{\mu}_j^{(t+1)})^T}{d \sum_{i=1}^N \tau_i^{j(+)} }$$

Finally for π_k (dropping all terms that don't depend on π_k):

$$\frac{\partial}{\partial \pi_j} \mathcal{L}_c(\theta, \sigma) \Big|_{j+1} = \frac{\partial}{\partial \pi_j} \sum_{i,k} \tau_i^{k(+)} \log \pi_k$$

With the constraint that $\sum_{k=1}^K \pi_k = 1$

$$\Rightarrow \underbrace{\sum_{k=1}^K \pi_k - 1}_{= g(\pi_k)} = 0 \Leftrightarrow g(\pi_k) = 0$$

Build the Lagrangian:

$$\mathcal{L}(\pi_k, \lambda) = \sum_{i,k} \tau_i^{k(+)} \log \pi_k + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

Maximize w.r.t. π_j

$$\frac{\partial \mathcal{L}}{\partial \pi_j} = \sum_{i=1}^N \tau_i^{j(+)} \cdot \frac{1}{\pi_j} + \lambda = 0.$$

$$\Rightarrow \pi_j^{(t+1)} = - \sum_{i=1}^N \frac{\tau_i^{j(+)}}{\lambda}$$

But need $g(\pi_j^{(t+1)}) = 0 \Rightarrow - \sum_{j=1}^K \sum_{i=1}^N \frac{\tau_i^{j(+)}}{\lambda} - 1 = 0.$

$$\Rightarrow - \sum_{i=1}^N \sum_{j=1}^n C_i^{j+1} = \lambda$$

$$\Rightarrow - \sum_{i=1}^N 1 = \lambda \Rightarrow \boxed{\lambda = -N}$$

and

$$\boxed{\pi_j^{(t+1)} = \frac{1}{N} \sum_i C_i^{k(t)}}$$

=

$\vec{\mu}_j^{(t+1)}$, $\sigma_j^2 {}^{(t+1)} I$ and $\pi_j^{(t+1)}$ are the M-step update parameters for a Gaussian mixture model.