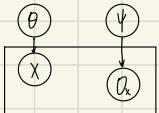


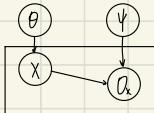


Likelihood of Data and Observation Models

Let $Y = \{Y_1, \dots, Y_n\}$ be some set of random variables, and $\alpha = \{\alpha_1, \dots, \alpha_m\}$ be their observability variable. The observability model is a joint distribution $P_{\text{miss}}(X, \alpha) = P(X) P_{\text{miss}}(\alpha | X)$, where $P(X)$ is parameterized by parameter θ and $P_{\text{miss}}(\alpha | X)$ is parameterized by ψ .



random missing values



deliberate missy values

Define a new set of variables $Y = \{Y_1, \dots, Y_n\}$, $\text{Val}(Y_i) = \text{Val}(x_i) \cup \{?\}$

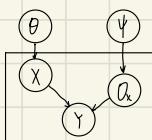
The actual observation is \bar{Y} , which is a deterministic function of X and α .

$$Y_i = \begin{cases} x_i & \alpha_i = 0 \\ ? & \alpha_i = 1 \end{cases}$$

"?" represents a missing value

e.g. $X_i \sim \text{Bernoulli}(\theta)$ $\alpha_i \sim \text{Bernoulli}(t)$

$$P(Y_i=1) = \theta t \quad P(Y_i=0) = (1-\theta)t \quad P(Y_i=?) = (1-t)$$



Given $Y=?$: α must be 0

Given $Y=0$: α must be 0, $X=Y$

$$L(\theta, t; D) = \prod \frac{\theta^{M_{10}} (1-\theta)^{M_{00}}}{(1-t)^{M_{10}} + t^{M_{00}}} \cdot \theta^{M_{01}} (1-\theta)^{M_{11}}$$

$$\max_{\theta, t} L(\theta, t; D) \Rightarrow \left\{ \begin{array}{l} \max_t \frac{\theta^{M_{10}} (1-\theta)^{M_{00}}}{(1-t)^{M_{10}} + t^{M_{00}}} \\ \max_\theta \theta^{M_{01}} (1-\theta)^{M_{11}} \end{array} \right.$$

$$\theta^* = \frac{M_{10} + M_{01}}{M_{10} + M_{01} + M_{00}} \quad \left[\text{separable} \right]$$

$$\theta^* = \frac{M_{01}}{M_{10} + M_{01}}$$

$(X \perp \alpha | Y)$ even though knowing Y activates the V-structure (context-specific independence)

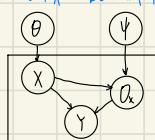
The independence reveals in the likelihood function

e.g. $X \sim \text{Bernoulli}(\theta)$ $\alpha_i | x_i \sim \text{Bernoulli}(\psi_{i|x_i})$ $\alpha_i | x_i \sim \text{Bernoulli}(\psi_{\alpha|x_i})$

$$P(Y_i=1) = \theta \cdot \psi_{i|x_i}$$

$$P(Y_i=0) = (1-\theta) \cdot \psi_{i|x_i}$$

$$P(Y_i=?) = \theta \cdot (1-\psi_{i|x_i}) \cdot (1-\theta) \cdot (1-\psi_{i|x_i})$$



knowing $Y=0$, we can infer X from α

$$\therefore X \perp \alpha | Y$$

$$L(\theta, \psi_{i|x_i}, \psi_{\alpha|x_i}; D) = \theta^{M_{10}} \cdot \psi_{i|x_i}^{M_{00}} \cdot (1-\theta)^{M_{01}} \cdot \psi_{\alpha|x_i}^{M_{11}} \cdot \left[\theta \cdot (1-\psi_{i|x_i}) \cdot (1-\theta) \cdot (1-\psi_{i|x_i}) \right]^{M_{00}}$$

Decomposing of Observation Mechanism

A missing data model P_{miss} is Missing Completely at Random (MCAR) if $P_{\text{miss}} \models (X \perp \alpha)$

In the case of MCAR, the likelihood function decomposes as a product

MCAR is sufficient but not necessary for decomposition of likelihood function

g. always observe the first toss. decide whether or not hide the 2nd toss based on the first toss

$$P(X=Y_1, Y=Y_1) = (\neg \theta_{11}) \cdot (\neg \theta_{12}) \cdot \psi(\theta_{11}|X=Y_1)$$

$$P(X=Y_1, Y=Y_2) = (\neg \theta_{11}) \cdot \theta_{12} \cdot \psi(\theta_{12}|X=Y_1)$$

$$P(X=Y_2, Y=?) = (\neg \theta_{12}) \cdot [\neg \psi(\theta_{12}|X=Y_1)]$$

:

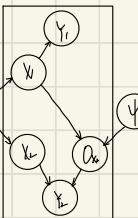
$$[(\theta_{11}, D) = [(\neg \theta_{11}) (\neg \theta_{12}) \psi(\theta_{12}|X=Y_1)]^{MCY_1, Y_1}] [(\neg \theta_{11}) \theta_{12} \psi(\theta_{12}|X=Y_1)]^{MCY_1, Y_2}] [(\neg \theta_{12}) [\neg \psi(\theta_{12}|X=Y_1)]]^{MCY_2, Y_1}] \\ [\theta_{11} (\neg \theta_{12}) \psi(\theta_{12}|X=Y_1)]^{MCY_1, Y_2}] [\theta_{11} \theta_{12} \psi(\theta_{12}|X=Y_1)]^{MCY_1, Y_1}] [\theta_{12} [\neg \psi(\theta_{12}|X=Y_1)]]^{MCY_2, Y_2}]$$

$$= (\neg \theta_{11})^{MCY_1} \cdot \theta_{12}^{MCY_1} \cdot$$

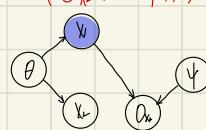
$$(\neg \theta_{11})^{MCY_2} \cdot \theta_{12}^{MCY_2}$$

$$\psi(\theta_{12}|X=Y_1)^{MCY_1, Y_1} + \psi(\theta_{12}|X=Y_2)^{MCY_2, Y_1}$$

$$\psi(\theta_{12}|X=Y_1)^{MCY_1, Y_2} + \psi(\theta_{12}|X=Y_2)^{MCY_2, Y_2}$$



$$(O_1 \perp X_2 | X_1)$$



the conditional independence helps to decouple $P(X)$ from $P(O_i|X)$

let y be a tuple of observations. Variables X are partitioned into \geq sets.

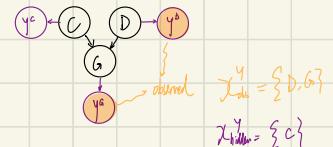
$$X_{\text{obs}}^y = \{X_i : Y_i \neq ?\} \quad (\text{observed vars}) \quad \text{and} \quad X_{\text{hidden}}^y = \{X_i : Y_i = ?\} \quad (\text{unobserved vars})$$

A missing data Model P_{missing} is missing at random (MAR) if

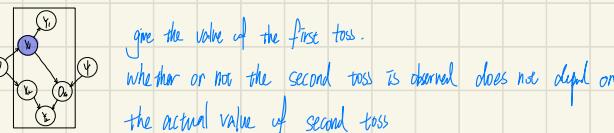
\forall observations y with $P_{\text{missing}}(y) > 0$, $\nexists X_{\text{hidden}}^y \in \text{Val}(X_{\text{hidden}}^y)$, $P_{\text{missing}} \models (O_i \perp X_{\text{hidden}}^y | X_{\text{obs}}^y)$

where O_i are specific values of observation variables given y

(given observed vars, whether or not other vars are observed does not depend on the actual value of those vars)



eg. give the value of the first toss.



whether or not the second toss is observed does not depend on the actual value of second toss

$$P_{\text{missing}}(X_{\text{hidden}}^y | X_{\text{obs}}^y, O_i) = P_{\text{missing}}(X_{\text{hidden}}^y | X_{\text{obs}}^y)$$

given observed vars, the observation pattern does not give additional informations about unobserved vars

$$P_{\text{missing}}(y) = \sum_{X_{\text{hidden}}^y} P(X_{\text{obs}}^y, X_{\text{hidden}}^y) \cdot P_{\text{missing}}(O_i | X_{\text{obs}}^y, X_{\text{hidden}}^y)$$

$$= \sum_{X_{\text{hidden}}^y} P(X_{\text{obs}}^y, X_{\text{hidden}}^y) \cdot P_{\text{missing}}(O_i | X_{\text{obs}}^y)$$

$$= P_{\text{missing}}(O_i | X_{\text{obs}}^y) \cdot P(X_{\text{obs}}^y)$$

depends only on ψ depends only on θ

The Likelihood Function

Given a BN G over a set of vars \mathcal{L} , each instance has a different set of observed vars

$D[m]$ and $O[m]$ are observed vars and values of m th instance

$H[m]$ is the missing vars in m th instance

$$l(\theta; D) = \log L(\theta; D) = \sum_{m=1}^M \log P(O[m] | \theta)$$



$$\text{fully observed: } L(\theta; D) = \theta_X^{ME[X]} (\perp \theta_X) \theta_Y^{ME[Y]} (\perp \theta_Y) \theta_{XY}^{ME[X,Y]} (\perp \theta_{XY}) \theta_{Y|X}^{ME[Y|X]} (\perp \theta_{Y|X})$$

$L(\theta; D)$ is log-concave, has a close-form for global maxi

$f(\theta; D)$ can be decomposed into function of each θ

partially observed: $D = \{(x?; y?), (x=x_0; y=y_0); \dots; (x=x_M; y=y_M)\}$

$$L(D, D[\cdot]) = P(X=x^*, Y=y^*) + P(X=x^*, Y=y^*)$$

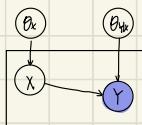
$$= (\perp \theta_X) \theta_{Y|X}^{ME[Y|X]} + \theta_X \theta_{Y|X}^{ME[Y|X]}$$

$$L(D, D[\cdot, \cdot]) = \theta_X^{ME[X]} (\perp \theta_X) \theta_{Y|X}^{ME[Y|X]} (\perp \theta_{Y|X}) \theta_{Y|X}^{ME[Y|X]} (\perp \theta_{Y|X})$$

$$f(\theta; D) = L(D, D[\cdot]) + f(\theta; D[\cdot, \cdot])$$

$$= \underbrace{\log [(\perp \theta_X) \theta_{Y|X}^{ME[Y|X]} + \theta_X \theta_{Y|X}^{ME[Y|X]}]}_{\text{can not be decomposed; and not necessarily log-concave}} + f(\theta; D[\cdot, \cdot])$$

each partially observed data introduce a log-plus term.



when Y is observed and X is not

$\theta_X \rightarrow X \rightarrow Y \leftarrow \theta_{Y|X}$ is an active trail

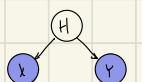
$$\theta_X \perp \theta_{Y|X}$$

$$P(\theta|Y) = \sum_{\text{active}} P(\theta|X|Y) = \sum_X \frac{P(\theta) f(X|Y|\theta)}{p(Y)} = \frac{P(\theta)}{p(Y)} \cdot \sum_X P(X|Y|\theta) = \underbrace{\frac{P(\theta)}{p(Y)}}_{\text{can not be decomposed as } f_1(\theta_X) \cdot f_2(\theta_{Y|X})} \cdot L(\theta, Y)$$

lose the property of parameter independence

if θ_X is known, $P(\theta_X|Y) = C \cdot \log [(\perp \theta_X) \theta_{Y|X}^{ME[Y|X]} + \theta_X \theta_{Y|X}^{ME[Y|X]}]$ the local decomposability is lost

g-



$X \nmid \text{observed}$

$$p(X|Y) = \sum_{\text{active}} p(X) p(X|H) p(Y|H)$$

$$L(\theta; D) = \prod_{(x,y) \in \text{active}} \left[\sum_{\text{active}} p(X) p(X|H) p(Y|H) \right]^{ME[Y|X]}$$

cannot be decomposed as $p(X|H)$ and $p(Y|H)$ (as in MLE)

In the presence of partially observed data, we lose all important properties of likelihood function

log-concavity, close-form solution, decomposability of $L(\theta; D)$

• Identifiability

ability to identify uniquely a model from data

e.g. randomly choose one from 2 biased coins, then toss

$(\theta_{H1}, \theta_{H2}, \theta_{T1})$

$$P(\theta_{H1}) = P(X^1) \stackrel{\text{def}}{=} \cdot (1 + p(X^1))^{-M+1}$$

$$P(\theta_{H2}) = \theta_{H1} \cdot \theta_{H2} + (1 - \theta_{H1}) \cdot \theta_{H2}$$

$P(X^1)$ is a weighted average of θ_{H1} and θ_{H2}

different choices of $(\theta_{H1}, \theta_{H2}, \theta_{T1})$ can give the same likelihood

Given the observation (no matter how many instances), we cannot hope to recover a unique set of parameters

Suppose we have a parametric model with $\theta \in \Theta$ that defines a distribution $p(x|\theta)$

A choice of θ is identifiable if there is no $\theta' \neq \theta$ s.t. $p(x|\theta) = p(x|\theta')$

A model is identifiable if θ is identifiable $\forall \theta \in \Theta$

• Gradient Ascent

Let B be a Bayesian Network over \mathcal{X} that induces $p(x)$

let o be a tuple of observations for some variables

$$\frac{\partial p(o)}{\partial p(x|u)} = \frac{1}{p(x|u)} p(x|u, o) \quad \text{if } p(x|u) > 0, \quad x \in \text{Var}(B) \cup \text{Val}(B)$$

Consider a full assignment ξ .

$$p(\xi) = \prod_{x \in \mathcal{X}} p(\xi(x) | \xi(x_<))$$

$$\frac{\partial p(\xi)}{\partial p(x|u)} = \sum_{x \in \mathcal{X}} \frac{\partial p(\xi(x) | \xi(x_<))}{\partial p(x|u)} = \frac{p(\xi, \theta)}{p(x|u, \theta)} \quad \text{if } \xi(x, p(x)) = (x, u)$$

otherwise

Consider a partial assignment o

$$p(o) = \sum_{\xi(o) = o} p(\xi)$$

$$\begin{aligned} \frac{\partial p(o)}{\partial p(x|u)} &= \sum_{\xi(o) = o} \frac{\partial p(\xi, \theta)}{\partial p(x|u)} = \underbrace{\sum_{\xi(o) = o, \xi(x) = x, \xi(u) = u} \frac{1}{p(x|u, \theta)} p(\xi, \theta)}_{p(x|u) \neq 0} \\ &= \frac{1}{p(x|u)} p(o, x, u) \end{aligned}$$

Let G be a Bayesian network over \mathcal{X} . $D = \{o_1, \dots, o_m\}$ be a partially observed dataset

$$\frac{\partial p(\theta; D)}{\partial p(x|u)} = \frac{1}{p(x|u)} \sum_{m=1}^M p(x, u | o_m, \theta)$$

$$\frac{\partial p(\theta; D)}{\partial p(x|u)} = \frac{1}{p(x|u)} \cdot \sum_{m=1}^M \frac{\partial p(o_m; \theta)}{\partial p(x|u)} = \sum_{m=1}^M \frac{1}{p(o_m; \theta)} \cdot \frac{\partial p(o_m; \theta)}{\partial p(x|u)}$$

$$\begin{aligned} &= \frac{1}{p(x|u)} \sum_{m=1}^M p(x, u | o_m; \theta) \\ &= \frac{1}{p(x|u)} \sum_{m=1}^M p(x|u | o_m; \theta) \end{aligned}$$

Def compute gradient (θ , ϕ , D):

for each $\Theta \rightarrow \Theta \in G$:

for each $x_i, u_i \in \text{Val}(X, U)$
 $\bar{M}(x_i, u_i) \geq 0$

} can use similar data structure as Factor

for $m = 1 \dots M$:

run clique tree calibration on $\langle G, \theta \rangle$ with evidence $O[m]$ // $\sum_m P(X, U | O[m], \theta)$ for each x, u

for each $\Theta \rightarrow \Theta \in G$:

for each $x_i, u_i \in \text{Val}(X, U)$

$$\bar{M}(x_i, u_i) += P(x_i, u_i | O[m])$$

for each $\Theta \rightarrow \Theta \in G$:

for each $x_i, u_i \in \text{Val}(X, U)$

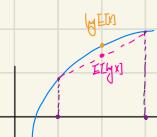
$$\nabla P(x_i | u_i) = \frac{1}{\bar{M}(x_i, u_i)} \bar{M}(x_i, u_i)$$

return $\{\nabla P(x_i | u_i)\}$

EM Algorithm



$$\begin{aligned}
 \sum_m \log P(Y[m]; \theta) &= \sum_m \log \sum_{\text{nonempty } \Omega_m} P(Y[m], X[m]; \theta) \\
 &= \sum_m \log \sum_{\text{nonempty } \Omega_m} \Omega_m(X[m]) \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \\
 &= \sum_m \log \mathbb{E}_{\Omega_m \sim \Omega} \left[\frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \right] \quad \text{construct a probability } \Omega_m(x) \quad \forall m = 1 \dots M \\
 &\geq \sum_m \mathbb{E}_{\Omega_m \sim \Omega} \left[\log \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \right] \quad (\log) \text{ is concave, Jensen's Inequality} \\
 &= \sum_m \mathbb{E}_{\Omega_m \sim \Omega} \left[\Omega_m(X[m]) \log \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \right]
 \end{aligned}$$



$$L(\theta; D) \geq \sum_m \sum_{\text{nonempty } \Omega_m} \Omega_m(X[m]) \log \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \quad \text{a lower bound for } L(\theta; D)$$

$L(\theta; D)$ is concave when $P(X, Y; \theta)$ is log-concave in θ :
 $\sum_m \sum_{\text{nonempty } \Omega_m} \Omega_m(X[m]) \log \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])}$

$$L(\theta; D) = \sum_m \log \mathbb{E}_{\Omega_m \sim \Omega} \left[\frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \right] = \sum_m \mathbb{E}_{\Omega_m \sim \Omega} \left[\log \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \right] \quad \text{if } \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])} \text{ is a constant } \forall m = 1 \dots M$$

$$\Omega_m(X[m]) = \frac{1}{z_m} \sum_{y \in \text{Val}(Y)} P(Y[m], X[m]; \theta)$$

$$\sum_{y \in \text{Val}(Y)} P(Y[m], X[m]; \theta) = \frac{1}{z_m} \sum_{y \in \text{Val}(Y)} P(Y[m], X[m]; \theta) = \frac{1}{z_m} \cdot P(Y[m]; \theta) = z_m = P(Y[m]; \theta)$$

$$\Omega_m(X[m]) = \frac{P(Y[m], X[m]; \theta)}{P(Y[m]; \theta)} = P(X[m] | Y[m]; \theta) + X[m] \text{ grad}(x)$$

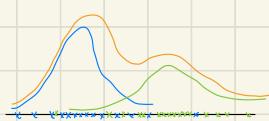
then maximize the lower bound.

$$\theta := \arg \max_{\theta} \sum_m \sum_{\text{nonempty } \Omega_m} \Omega_m(X[m]) \log \frac{P(Y[m], X[m]; \theta)}{\Omega_m(X[m])}$$

g. mixture of Gaussians

$X \sim \text{Bernoulli}(p_x)$

$(Y|X=j) \sim N(\mu_j, \Sigma_j)$



$$\begin{aligned} l(b) &= \sum_{m=1}^M \log P(Y[m]; b) = \sum_{m=1}^M \log \sum_{j=1}^K Q_m(X[m], Y[m]; b) = \sum_{m=1}^M \log \left(\sum_{j=1}^K Q_m(X[m]) \frac{P(X[m], Y[m]; b)}{Q_m(X[m])} \right) = \sum_{m=1}^M \log \mathbb{E}_{X[m] \sim p_x} \left[\frac{P(X[m], Y[m]; b)}{Q_m(X[m])} \right] \\ &> \sum_{m=1}^M \mathbb{E}_{X[m] \sim p_x} \left[\log \frac{P(X[m], Y[m]; b)}{Q_m(X[m])} \right] = \sum_{m=1}^M \sum_{j=1}^K Q_m(X[m]) \log \frac{P(X[m], Y[m]; b)}{Q_m(X[m])} \end{aligned}$$

E-Step:

$$Q_m(X[m]) = P(X[m] | Y[m], b) = \frac{P(X[m], b) \cdot P(Y[m] | X[m], b)}{P(Y[m], b)}$$

$$Q_m(X) = \frac{1}{2} \left(p_x \cdot N(Y[m]; \mu_i, \Sigma_i) \right) = \frac{p_x}{2} \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (Y[m] - \mu_i)^T \Sigma_i^{-1} (Y[m] - \mu_i) \right)$$

$$Q_m(X^*) = \frac{1}{2} \left((1-p_x) \cdot N(Y[m]; \mu_0, \Sigma_0) \right) = \frac{1-p_x}{2} \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_0|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (Y[m] - \mu_0)^T \Sigma_0^{-1} (Y[m] - \mu_0) \right)$$

M-Step:

$$\begin{aligned} &\max_b \sum_{m=1}^M \sum_{j=1}^K Q_m(X[m]) \cdot \log \frac{P(X[m], Y[m]; b)}{Q_m(X[m])} \\ &= \max_b \sum_{m=1}^M \sum_{j=1}^K Q_m(X[m]) \cdot \left[\log \frac{p_x}{(2\pi)^{\frac{D}{2}} |\Sigma_j|^{\frac{1}{2}}} - \frac{1}{2} (Y[m] - \mu_j)^T \Sigma_j^{-1} (Y[m] - \mu_j) \right] \end{aligned}$$

MAX μ :

$$\begin{aligned} \nabla_{\mu_j} &= \sum_{m=1}^M Q_m(j) \cdot -\frac{1}{2} \cdot 2 \sum_j (Y[m] - \mu_j) \cdot \\ &= \sum_j \left[\sum_{m=1}^M Q_m(j) Y[m] - \sum_{m=1}^M Q_m(j) \mu_j \right] \rightarrow \\ \mu_j &= \frac{\sum_{m=1}^M Q_m(j) Y[m]}{\sum_{m=1}^M Q_m(j)} \end{aligned}$$

MAX Σ : let $A_j = \Sigma_j^{-1}$

$$\max_{A_j} \sum_{m=1}^M Q_m(j) \left[-\frac{1}{2} \log |A_j| - \frac{1}{2} (Y[m] - \mu_j)^T A_j (Y[m] - \mu_j) \right]$$

$$= \max_{A_j} \sum_{m=1}^M Q_m(j) \left[\log \det(A_j) - \text{tr}(A_j, (Y[m] - \mu_j)(Y[m] - \mu_j)^T) \right]$$

$$= \max_{A_j} \log \det(A_j) \sum_{m=1}^M Q_m(j) - \text{tr}(A_j, \frac{1}{m} \sum_{m=1}^M Q_m(j) (Y[m] - \mu_j)(Y[m] - \mu_j)^T)$$

$$\nabla_{A_j} = A_j^{-1} \sum_{m=1}^M Q_m(j) - \frac{1}{m} \sum_{m=1}^M Q_m(j) (Y[m] - \mu_j)(Y[m] - \mu_j)^T \rightarrow$$

$$A_j^{-1} = \frac{1}{m} \sum_{m=1}^M Q_m(j) (Y[m] - \mu_j)(Y[m] - \mu_j)^T$$

MAX θ

$$\max_{\theta} \sum_{m=1}^M \sum_{j=1}^K Q_m(j) \log \frac{\theta_j \cdot N(Y[m], \mu_j, \Sigma_j)}{Q_m(j)}$$

↓

$$\min_{\theta} \sum_{m=1}^M \sum_{j=1}^K Q_m(j) \log \theta_j$$

$$\text{s.t. } \sum_j \theta_j = 1$$

$$\log \det(X + \lambda I) = \log \det \left[X^{\frac{1}{2}} (I + X^{\frac{1}{2}} \alpha X^{\frac{1}{2}}) X^{\frac{1}{2}} \right]$$

$$= \log \det \left[X (I + X^{\frac{1}{2}} \alpha X^{\frac{1}{2}}) \right]$$

$$= \log \det(X) + \log \det(I + X^{\frac{1}{2}} \alpha X^{\frac{1}{2}})$$

λ_i is the i th eigenvalue
of $X^{\frac{1}{2}} \alpha X^{\frac{1}{2}}$

$$\approx \log \det(X) + \sum_i \lambda_i$$

$$= \log \det(X) + \text{tr}(X^{\frac{1}{2}} \alpha X^{\frac{1}{2}})$$

$$= \log \det(X) + \text{tr}(\alpha X, X^{\frac{1}{2}})$$

$$f_i(\theta, \mathbf{v}) = -\sum_{m=1}^M \sum_{j \in \text{faulty}} (\alpha_m(j)) \log p_j + V \left(\sum_j p_j - 1 \right)$$

$$\text{Therefore } \hat{\theta}_j = -\frac{1}{V} \sum_{m=1}^M \sum_{j \in \text{faulty}} (\alpha_m(j)) + V : \approx 0$$

$$\hat{\theta}_j = \frac{\sum_m (\alpha_m(j))}{V}$$

$$g(\mathbf{V}) = \frac{1}{V} \sum_{m=1}^M f_i(\theta, \mathbf{V}) = -\frac{1}{V} \sum_{m=1}^M \sum_{j \in \text{faulty}} (\alpha_m(j)) \log \frac{\sum_{m=1}^M \alpha_m(j)}{V} + \sum_j \sum_{m=1}^M \alpha_m(j) - V$$

$$= \frac{1}{V} \sum_{m=1}^M \sum_{j \in \text{faulty}} (\alpha_m(j)) \log V - V + C$$

$$\text{Therefore } \hat{V} = \frac{1}{V} \sum_{m=1}^M \sum_{j \in \text{faulty}} (\alpha_m(j)) - 1 : \approx 0$$

$$V^* = \sum_{m=1}^M \sum_j (\alpha_m(j)) = \sum_{m=1}^M 1 = M$$

$$\hat{V}^* = \frac{\sum_m (\alpha_m(j))}{M}$$

eg. Factor Analysis

$$Z \sim N(0, I)$$

$$X/Z \sim N(\mu + \lambda Z, \Sigma)$$

$$f(\theta) = \prod_{m=1}^M \log P(X[m]; \theta)$$

$$= \prod_{m=1}^M \log \int_{Z[m]}^{+\infty} P(X[m], Z[m]; \theta) dZ[m]$$

$$= \prod_{m=1}^M \log \int_{Z[m]}^{+\infty} Q_m(Z[m]) \cdot \frac{P(X[m], Z[m]; \theta)}{Q_m(Z[m])} dZ[m] = \prod_{m=1}^M \log E_{Z[m] \sim Q_m} \left[\frac{P(X[m], Z[m]; \theta)}{Q_m(Z[m])} \right]$$

$$\geq \prod_{m=1}^M E_{Z[m] \sim Q_m} \left[\log \frac{P(X[m], Z[m]; \theta)}{Q_m(Z[m])} \right]$$

$$= \prod_{m=1}^M \int_{Z[m]}^{+\infty} Q_m(Z[m]) \cdot \log \frac{P(X[m], Z[m]; \theta)}{Q_m(Z[m])} dZ[m]$$

$$Q_m(z[m]) = \frac{1}{z} P(X[m], z[m]; \theta) \quad \forall z[m] \in \text{val}(z)$$

$$\int_{Z[m]}^{+\infty} Q_m(Z[m]) dZ[m] = \frac{1}{z} \int_{Z[m]}^{+\infty} P(X[m], Z[m]; \theta) dZ[m] = \frac{1}{z} P(X[m]; \theta) \approx 1$$

$$Q_m(z[m]) = \frac{P(X[m], z[m]; \theta)}{P(X[m]; \theta)} = P(z[m] | X[m]; \theta)$$

E-Step:

$$P(Z[m], X[m]; \theta) = P(Z[m]) \cdot P(X[m] | Z[m], \theta)$$

$$= \frac{1}{(2\pi)^{\frac{M}{2}}} \exp\left(-\frac{1}{2} Z[m]^T Z[m]\right) \cdot \frac{1}{(2\pi)^{\frac{M}{2}} |E|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} [X[m] - (U + \lambda Z[m])]^T E^{-1} [X[m] - (U + \lambda Z[m])] \right)$$

$$= \frac{1}{z} \exp \left\{ -\frac{1}{2} \left[Z[m]^T Z[m] + X[m]^T E^{-1} X[m] - 2 X[m]^T E^{-1} (U + \lambda Z[m]) + (U + \lambda Z[m])^T E^{-1} (U + \lambda Z[m]) \right] \right\}$$

$$= \frac{1}{z} \exp \left\{ -\frac{1}{2} \left[Z[m]^T Z[m] + X[m]^T E^{-1} X[m] - 2 X[m]^T E^{-1} U - 2 X[m]^T E^{-1} \lambda Z[m] + U^T E^{-1} U + 2 U^T E^{-1} \lambda Z[m] + (\lambda Z[m])^T E^{-1} (\lambda Z[m]) \right] \right\}$$

$$= \frac{1}{z} \exp \left\{ -\frac{1}{2} \left[Z[m]^T \begin{bmatrix} I + \lambda^T E^{-1} \lambda & -\lambda^T E^{-1} \\ -E^{-1} \lambda & E^{-1} \end{bmatrix} \begin{bmatrix} Z[m] \\ X[m] \end{bmatrix} \right] \right\}$$

$$\begin{bmatrix} Z[m] \\ X[m] \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ u \end{bmatrix}, \Sigma \right)$$

$$\Sigma = \begin{bmatrix} I + \lambda^T E^{-1} \lambda & -\lambda^T E^{-1} \\ -E^{-1} \lambda & E^{-1} \end{bmatrix} = \begin{bmatrix} I & \lambda^T \\ -\lambda & E^{-1} \end{bmatrix}$$

$$\begin{bmatrix} Z_{11} & Z_{12} \\ Z_{21} & Z_{22} \end{bmatrix}^{-1} = \begin{bmatrix} (Z_{11} - Z_{12} Z_{21}^{-1} Z_{12})^{-1} & - (Z_{11} - Z_{12} Z_{21}^{-1} Z_{12})^{-1} Z_{12} Z_{21}^{-1} \\ - Z_{21} Z_{21}^{-1} (Z_{11} - Z_{12} Z_{21}^{-1} Z_{12})^{-1} & Z_{21}^{-1} Z_{21}^{-1} (Z_{11} - Z_{12} Z_{21}^{-1} Z_{12})^{-1} Z_{12} Z_{21}^{-1} \end{bmatrix}$$

by block elimination

$$(Z^{[m]} | X^{[m]}) \sim N(-\lambda(E + \lambda I)^{-1}(X^{[m]} - u), I - \lambda^T(E + \lambda I)^{-1}\lambda)$$

$$Q_m(Z^{[m]}) = P(Z^{[m]} | X^{[m]}) = N\left(Z^{[m]}, -\lambda(E + \lambda I)^{-1}(X^{[m]} - u), I - \lambda^T(E + \lambda I)^{-1}\lambda\right)$$

$E[Z^{[m]} | X^{[m]}]$ $\text{cov}(Z^{[m]} | X^{[m]})$

M-step:

$$\begin{aligned} & \max_{\lambda, u, E} \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \log \frac{P(X^{[m]}, Z^{[m]}, \theta)}{Q_m(Z^{[m]})} dZ^{[m]} \\ &= \max_{\lambda, u, E} \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \left[\log P(Z^{[m]}) + \log P(X^{[m]} | Z^{[m]}, \theta) - \log Q_m(Z^{[m]}) \right] dZ^{[m]} \\ &= \max_{\lambda, u, E} \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \log P(X^{[m]} | Z^{[m]}, \theta) dZ^{[m]} \\ &= \max_{\lambda, u, E} \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \left[\log \frac{1}{(2\pi)^{\frac{M}{2}} |\lambda|^{1/2}} - \frac{1}{2} (\lambda^{[m]} - (u + \lambda Z^{[m]}))^T \lambda^{-1} (\lambda^{[m]} - (u + \lambda Z^{[m]})) \right] dZ^{[m]} \end{aligned}$$

MAX λ :

$$\begin{aligned} \nabla_\lambda &= \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \lambda^{-1} (\lambda^{[m]} - (u + \lambda Z^{[m]})) Z^{[m] T} dZ^{[m]} \\ &= \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \lambda^{-1} (\lambda^{[m]} - u) Z^{[m] T} dZ^{[m]} - \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \cdot \lambda^{-1} \lambda Z^{[m]} Z^{[m] T} dZ^{[m]} \\ &= \sum_{m=1}^M \lambda^{-1} (\lambda^{[m]} - u) E_{Z^{[m]} | X^{[m]}} [Z^{[m]}]^T - \sum_{m=1}^M \lambda^{-1} \lambda E_{Z^{[m]} | X^{[m]}} [Z^{[m]} Z^{[m]}]^T \Rightarrow \\ &\lambda = \left(\sum_{m=1}^M (\lambda^{[m]} - u) E_{Z^{[m]} | X^{[m]}} [Z^{[m]}]^T \right) \left(\sum_{m=1}^M E_{Z^{[m]} | X^{[m]}} [Z^{[m]} Z^{[m]}]^T \right)^{-1} \end{aligned}$$

MAX u :

$$\begin{aligned} \nabla_u &= \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \lambda^{-1} (\lambda^{[m]} - u - \lambda Z^{[m]}) dZ^{[m]} \\ &= \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) (\lambda^{[m]} - u) dZ^{[m]} - \sum_{m=1}^M \lambda^{-1} \lambda \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) Z^{[m]} dZ^{[m]} \\ &= \lambda^{-1} \sum_{m=1}^M (\lambda^{[m]} - u) - \lambda^{-1} \lambda \sum_{m=1}^M E_{Z^{[m]} | X^{[m]}} [Z^{[m]}] \\ &= \lambda^{-1} \sum_{m=1}^M (\lambda^{[m]} - u) - \lambda^{-1} \lambda \cdot \sum_{m=1}^M \lambda^{-1} (\lambda^{[m]} - u - \lambda Z^{[m]})^T (\lambda^{[m]} - u) \\ &= (\lambda^{-1} + \lambda^{-1} \lambda \lambda^T (\lambda + \lambda I)^{-1}) \sum_{m=1}^M (\lambda^{[m]} - u) \Rightarrow \\ u &= \frac{1}{M} \sum_{m=1}^M \lambda^{[m]} \end{aligned}$$

MAX E :

$$\begin{aligned} & \max_E \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \cdot \left\{ -\frac{1}{2} (\lambda^{[m]} - (u + \lambda Z^{[m]}))^T \lambda^{-1} (\lambda^{[m]} - (u + \lambda Z^{[m]})) \right\} dZ^{[m]} \\ &= \min_E \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \left\{ \log \det(\lambda) + \text{tr}(\lambda^{-1}, (\lambda^{[m]} - (u + \lambda Z^{[m]})) (\lambda^{[m]} - (u + \lambda Z^{[m]}))^T) \right\} dZ^{[m]} \\ &\text{let } J = \lambda^{-1} \\ &= \min_J \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \left\{ -\log \det(J) + \text{tr}(J, (\lambda^{[m]} - (u + \lambda Z^{[m]})) (\lambda^{[m]} - (u + \lambda Z^{[m]}))^T) \right\} dZ^{[m]} \\ &\nabla_J = \sum_{m=1}^M \int_{-\infty}^{+\infty} Q_m(Z^{[m]}) \left\{ -J^{-1} + (\lambda^{[m]} - u - \lambda Z^{[m]})(\lambda^{[m]} - u - \lambda Z^{[m]})^T \right\} dZ^{[m]} \end{aligned}$$

$$\begin{aligned}
&= \sum_{m=1}^M -J^{-1} + \sum_{m=1}^M \sum_{z[m]}^{t+m} Q_m(z[m]) \lambda z[m] z[m]^T \rightarrow \lambda z[m] (X(m)u)^T + (X(m)u)(X(m)u)^T d(z[m]) \\
&= -M J^{-1} + \sum_{m=1}^M \lambda E_{z[m|z[m]]} [z[m] z[m]^T] \rightarrow \lambda E_{z[m|z[m]]} [z[m]] (X(m)u)^T + (X(m)u)(X(m)u)^T := \\
&\bar{J} = J^{-1} = \frac{1}{M} \sum_{m=1}^M \lambda E_{z[m|z[m]]} [z[m] z[m]^T] \rightarrow \lambda E_{z[m|z[m]]} [z[m]] (X(m)u)^T + (X(m)u)(X(m)u)^T
\end{aligned}$$

EM algorithm for Bayesian Network

let $o[m]$ be observed vars of math data and $z[m]$ be hidden vars of math data

$$\begin{aligned}
\text{(1)} &= \sum_{m=1}^M \log P(o[m]; \theta) \\
&= \sum_{m=1}^M \log \sum_{z[m|o[m], u]} P(o[m], z[m], u) \\
&= \sum_{m=1}^M \log \sum_{z[m]} Q_m(z[m]) \cdot \frac{P(o[m], z[m], u)}{Q_m(z[m])} \\
\text{①} &\geq \sum_{m=1}^M \sum_{z[m]} Q_m(z[m]) \log \frac{P(o[m], z[m], u)}{Q_m(z[m])} \\
&= \sum_{m=1}^M \sum_{z[m]} Q_m(z[m]) \cdot \left[\sum_{(o \rightarrow z) \in \theta} \log P((o[m], z[m]) \times o) | ((o[m], z[m]) \times u) \right] - \log(Q_m(z[m]))
\end{aligned}$$

E-Step:

to achieve equality at ① $\frac{P(o[m], z[m], u)}{Q_m(z[m])}$ is constant

$$Q_m(z[m]) = \frac{1}{Z} P(o[m], z[m], u)$$

$$\sum_{z[m]} Q_m(z[m]) = \frac{1}{Z} \sum_{z[m]} P(o[m], z[m], u) = \frac{1}{Z} P(o[m], u) = 1$$

$$Q_m(z[m]) = \frac{P(o[m], z[m], u)}{P(o[m], u)} = p(z[m] | o[m], u)$$

M-Step:

Find θ_{new} s.t $x \in \text{val}(x)$

$$\text{min. } -\sum_{m=1}^M \sum_{z[m]} Q_m(z[m]) \log P((o[m], z[m]) \times o) | ((o[m], z[m]) \times u)$$

s.t $\sum_{x \in \text{val}(x)} \theta_{\text{new}} = 1$

$$\begin{aligned}
\mathcal{L}(\theta_{\text{new}}, v) &= -\sum_{m=1}^M \sum_{z[m]} Q_m(z[m]) \log P((o[m], z[m]) \times o) | ((o[m], z[m]) \times u) + v \left(\sum_x \theta_{\text{new}} - 1 \right) \\
&= -\sum_{x \in \text{val}(x)} \log \theta_{\text{new}} \sum_{m=1}^M \sum_{z[m|o[m], u]} Q_m(z[m]) \cdot \mathbb{1}\{(o[m], z[m]) \times u = (x, u)\} + v \left(\sum_x \theta_{\text{new}} - 1 \right)
\end{aligned}$$

$$\nabla_{\theta_{\text{new}}} \mathcal{L} = -\frac{1}{\theta_{\text{new}}} \cdot \sum_{m=1}^M \sum_{z[m|o[m], u]} Q_m(z[m]) \mathbb{1}\{(o[m], z[m]) \times u = (x, u)\} + v := 0$$

$$\theta_{\text{new}}(v) = \frac{1}{V} \cdot \sum_{m=1}^M \sum_{z[m|o[m], u]} Q_m(z[m]) \mathbb{1}\{(o[m], z[m]) \times u = (x, u)\}$$

let $h(x, u) = \sum_{m=1}^M \sum_{z[m]} Q_m(z[m]) \mathbb{1}\{(o[m], z[m]) \times u = (x, u)\}$ is constant $\forall x, u$

$$\theta_{x|u}(y) = \frac{1}{V} h(x|u)$$

$$g(w) = \inf_{\theta_{x|u}} L(\theta_{x|u}, V)$$

$$= - \sum_m b_y \left(\frac{1}{V} h(x|u) \right) \cdot h(x|u) + V \left(\sum_m \frac{1}{V} h(x|u) - 1 \right)$$

$$= - \sum_m \left(b_y(h(x|u)) - b_V \right) h(x|u) + \sum_m h(x|u) - V$$

$$T_b g = \frac{1}{V} \cdot \sum_m h(x|u) - 1 := 0$$

$$V^* = \sum_m h(x|u)$$

$$\begin{aligned} \theta_{x|u}^* &= \frac{h(x|u)}{\sum_m h(x|u)} = \frac{\sum_m \sum_{m \in M} \alpha_m(z|m) \cdot \mathbb{1}_{\{(z|m), z|m\} \perp x|u\}}{\sum_m \sum_{m \in M} \alpha_m(z|m) \cdot \mathbb{1}_{\{(z|m), z|m\} \perp x|u\}} \\ &= \frac{\sum_m \sum_{m \in M} p(z|m | o|m; \theta) \cdot \mathbb{1}_{\{(z|m), z|m\} \perp x|u\}}{\sum_m \sum_{m \in M} p(z|m | o|m; \theta) \cdot \mathbb{1}_{\{(z|m), z|m\} \perp x|u\}} \\ &= \frac{\sum_m p(x|u | o|m; \theta)}{\sum_m p(x|u | o|m; \theta)} \end{aligned}$$

def compute_marginal_Q(G, θ, D): partially observed dataset:

for each $\emptyset \rightarrow X \in G$

for each $(x|u) \in \text{val}(x|u)$:

$$Q(x|u) = 0$$

for $m=1, \dots, M$,

run inference with evidence $o|m$

for each $\emptyset \rightarrow X \in G$

for each $(x|u) \in \text{val}(x|u)$:

$$Q(x|u) += p(x|u | o|m) \quad \leftarrow$$

return $\{Q(x|u) \mid \emptyset \rightarrow X \in G, X(u) \in \text{val}(x|u)\}$

def EM_Bayesian_Net(G, θ^{*}, D):

for $t=0, 1, 2, \dots$ until convergence

$$\{Q(x|u)\} = \text{Compute_Marginal}_Q(G, \theta^t, D)$$

for each $\emptyset \rightarrow X \in G$

for each $(x|u) \in \text{val}(x|u)$:

$$\theta_{x|u} = \frac{Q(x|u)}{\sum_X Q(x|u)} \quad \leftarrow$$

return θ^t