

$f(y|\theta)$: observed data density

$g(y, z|\theta)$: complete data density

$$h(z|y, \theta) \triangleq \frac{g(y, z|\theta)}{f(y|\theta)} \quad \text{missing data density}$$

$$f(y|\theta) = \frac{g(y, z|\theta)}{h(z|\theta, y)}$$

$$-\log f(y|\theta) = -\log g(y, z|\theta) - [-\log h(z|y, \theta)]$$

$$-\frac{\partial}{\partial \theta} \log f(y|\theta) = -\frac{\partial}{\partial \theta} \log g(y, z|\theta) - \left[-\frac{\partial}{\partial \theta} \log h(z|y, \theta) \right]$$

$$\underbrace{I_y(\theta)}_{\text{obs}} = \underbrace{I_{yz}(\theta)}_{\text{complete}} - \underbrace{I_{zy}(\theta)}_{\text{missing}}$$

How to compute $I_{yz}(\theta)$ and $I_{zy}(\theta)$?

$$\text{Let } S(y|\theta) = \frac{\partial}{\partial \theta} \log f(y|\theta)$$

$$s(y, z|\theta) = \frac{\partial}{\partial \theta} \log g(y, z|\theta)$$

Louis 1982 showed:

$$I_{z|y}(\theta) = E_{z|y} [S(y, z|\theta) S(y, z|\theta)^T] - S(y|\theta) S(y|\theta)^T$$

with expectation taken wrt $h(z|y, \theta)$

$$\Rightarrow I_Y(\theta) = I_{Y,z}(\theta) - \underbrace{E_{z|y} [S(y, z|\theta) S(y, z|\theta)^T]}_{\text{Computation only on complete data}} - \underbrace{S(y|\theta) S(y|\theta)^T}_{= 0 \text{ for } \theta = \hat{\theta}}$$

$$\text{Note } I_{Y,z}(\hat{\theta}) = -E \left[\frac{\partial^2}{\partial \theta^2} \log g(y, z|\theta) \mid \theta_{n-1}, y \right] \\ = -Q''(\hat{\theta} \mid \hat{\theta})$$

Meilijson (89) stated that if y 's are i.i.d, then

$$I_Y(\theta) = \text{Var}(S(y|\theta))$$

$$\frac{1}{n} \sum_{i=1}^n$$

$$S(y|\theta) = \sum_{i=1}^n S(y_i|\theta)$$

$$I_Y(\theta) = \text{Var}(S(Y|\theta))$$

$$= \frac{1}{n} \sum S(Y_i|\theta) S(Y_i|\theta)^T - \underbrace{\left[\frac{1}{n} \sum S(Y_i|\theta) \right]}_{\theta \text{ for } \theta = \hat{\theta}} \left[\frac{1}{n} \sum S(Y_i|\theta) \right]^T$$

How?

Lecture '82 shared

$$S(Y|\theta) = E_{Z|Y} [S(Y, Z|\theta)]$$

$$I_Y(\hat{\theta}) = \frac{1}{n} \sum E_{Z|Y} [S(Y_i, Z|\hat{\theta})] E_Z [S(Y_i, Z|\hat{\theta})]^T$$

Ex: Mixture models

$$\ell_i(\theta) = z_i \log \phi(Y_i|\mu_1, \sigma_1^2) + (1-z_i) \log \phi(Y_i|\mu_2, \sigma_2^2) + z_i \log \lambda + (1-z_i) \log (1-\lambda)$$

$$S(Y_i, Z_i|\theta) = \begin{pmatrix} \frac{Y_i - \mu_1}{\sigma_1^2} z_i \\ \frac{Y_i - \mu_2}{\sigma_2^2} (1-z_i) \\ \left(-\frac{1}{\sigma_1^2} + \frac{1}{\sigma_1^3} (Y_i - \mu_1) \right) z_i \\ \left(-\frac{1}{\sigma_2^2} + \frac{1}{\sigma_2^3} (Y_i - \mu_2) \right) (1-z_i) \\ \frac{z_i}{\lambda} + \frac{1-z_i}{1-\lambda} \end{pmatrix}$$

$$[X]_{\text{排}} [x|y]$$

It is MCMC,
but not Gibbs.
exactly

Data Augmentation (Tanner + Wong '87)

We have some data y_1, \dots, y_n , and we propose a model

$$y \sim p(y|\theta)$$

$$\theta \sim \pi(\theta)$$

So we have a likelihood $p(y|\theta)$ and a prior $\pi(\theta)$.

We want the posterior $p(\theta|y)$.

~~$P(4) = 2$ $P(4) = 2$ $P(4) = 2$~~

BUT because of some missing data z , it is difficult to evaluate $p(\theta|y)$. ~~that~~

Let (y, z) be the "complete data" and we want to know

$$p(\theta|y) = \int \underbrace{p(\theta|y, z)} \underbrace{p(z|y)} dz \quad (1)$$

"complete data posterior" predictive density for Z_i .

↳ easy to calculate

Notice also that we have

$$p(z|y) = \int p(z|y, \theta) p(\theta|y) d\theta \quad \text{easy to calculate} \quad (2)$$

→ easy to calculate

Substituting, ~~we get~~ (2) into (1), we get

$$p(\theta|y) = \int p(\theta|y, z) \left[\int p(z|y, \theta') p(\theta'|y) d\theta' \right] dz$$

$$= \iint p(\theta | y, z) p(z | y, \theta') dz \int p(\theta' | y) d\theta'$$

Let $K(\theta, \theta') = \int p(\theta | y, z) p(z | y, \theta') d\theta'$

$$p_{(1)}(\theta|y) = \int k(\theta, \theta') p_{(2)}(\theta'|y) d\theta'$$

(fixed point system) If $p(e'(y))$ is the ~~IF ① and ②~~ true posterior, then it will map to itself.

Let T be a functional which maps a function $g(\theta)$ to

$$Tg(\theta) = \int K(\theta, \theta') g(\theta') d\theta'$$

If we take some initial value $P_0(\theta|y)$, what happens when we iterate ~~$T_{P_0}(\theta|y)$~~ ~~$T_{P_0}(\theta|y)$~~ . $T_{P_0}(\theta|y)$
i.e.

$$\cancel{p_0(\theta)} \quad \cancel{g_{j+1}(\theta)} \quad \cancel{p_{j+1}} \quad p_1(\theta|y) = T p_0(\theta|y) = \int K(\theta, \theta') p_0(\theta'|y) d\theta'$$

Does the sequence $\{P_i(\theta|y)\}$ converge to anything?

~~(D plot) the true posterior is the~~

① $\{p_i(\theta|y)\} \rightarrow p(\theta|y)$ monotonically (always getting closer)
to true posterior

(2) $p(y)$ is the unique solution to system of integral equations

(3) $p_i(\theta|y) \rightarrow p(\theta|y)$ linearly.

DA algorithm

Roughly: At iteration i

① Sample $z_1 \rightarrow z_m \sim p(z|y)$

② Estimate $p_{i+1}(\theta|y) = \frac{1}{m} \sum_{i=1}^m p(\theta|y, z_i)$

Repeat until $\|p_{i+1}(\theta|y) - p_i(\theta|y)\| < \epsilon$. complete data posterior

① Pick some initial $p_0(\theta|y)$. At step i

②a Generate $\theta \sim p_i(\theta|y)$

②b Generate $z \sim p(z|\theta, \theta)$

Repeat ②a - ②b m times to get $z_1 \rightarrow z_m$.

③ Let $p_{i+1}(\theta|y) = \frac{1}{m} \sum p(\theta|y, z_i)$

mixture of conditional densities.

Monte Carlo estimate of

$$\int p(\theta|y, z) p(z|y) dz$$

DA not so straightforward when there are more than 2 "missing" components.

Use Gibbs sampling for more complicated problems

Ex. 3-stage hierarchical model

~~$$y|\alpha \sim N(\alpha, 1)$$~~

$$\alpha|\theta \sim N(\theta, 1)$$

$$\theta \sim \pi(\theta) = N(0, 1)$$

$$p(y, \alpha, \theta) = \underbrace{p(y|\alpha, \theta)}_{N(\alpha, 1)} \underbrace{p(\alpha|\theta)}_{N(\theta, 1)} \underbrace{\pi(\theta)}_{N(0, 1)}$$

~~$$p(\theta|y, \alpha)$$~~

$$p(\alpha|y, \theta) \propto p(y|\alpha, \theta) p(\alpha|\theta)$$

$$= N\left(y + \frac{1}{2}(\theta - y), \frac{1}{2}\right)$$

$$p(\theta|y, \alpha) \propto p(y|\alpha, \theta) p(\alpha|\theta) \pi(\theta)$$

$$N(\quad)$$

① find some $p_i(\theta|y)$ [Normal?]

②a sample $\theta \sim p_i(\theta|y)$

②b sample $\alpha \sim p(\alpha|y, \theta)$

Repeat m times to get $\alpha_1, \dots, \alpha_m$

③ Let $p_{\#}(\theta|y) = \frac{1}{m} \sum p_i(\theta|y, \alpha)$

~~...~~

	<u>EM</u>	<u>DA</u>
Paradigm Philosophy	Likelihood	Bayesian
Target	MLE	Posterior dist.
Method to handle z	Average (conditional)	sample (conditional)
Need $(z y)$	compute E	sample from
Std. errors	No, by default	Yes, have entire posterior (need prior)
Rate of convergence	linear	linear
Monotonicity	✓	✓

Both methods "impute" missing data via the specification of the complete data model.

Once $p(y, z)$ is specified, everything else is determined.