

Gibbs and Metropolis sampling (MCMC methods) and relations of Gibbs to EM

Lecture Outline

1. Gibbs

- the algorithm
- a bivariate example
- an elementary convergence proof for a (discrete) bivariate case
- more than two variables
- a counter example.

2. *EM* – again

- *EM* as a maximization/maximization method
- Gibbs as a variation of Generalized *EM*

3. Generating a Random Variable.

- Continuous r.v.s and an exact method based on transforming the cdf.
- The “accept/reject” algorithm.
- The Metropolis Algorithm

Gibbs Sampling

We have a joint density

$$f(x, y_1, \dots, y_k)$$

and we are interested, say, in some features of the marginal density

$$f(x) = \iint \dots \int f(x, y_1, \dots, y_k) dy_1, dy_2, \dots, dy_k.$$

For instance, suppose that we are interested in the average

$$E[X] = \int x f(x) dx.$$

If we can sample from the marginal distribution, then

$$\lim_{m \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = E[X]$$

without using $f(x)$ explicitly in integration. Similar reasoning applies to any other characteristic of the statistical model, i.e., of the *population*.

The Gibbs Algorithm for computing this average.

Assume we can sample the $k+1$ -many univariate conditional densities:

$$\begin{aligned} &f(X \mid y_1, \dots, y_k) \\ &f(Y_1 \mid x, y_2, \dots, y_k) \\ &f(Y_2 \mid x, y_1, y_3, \dots, y_k) \\ &\dots \\ &f(Y_k \mid x, y_1, y_3, \dots, y_{k-1}). \end{aligned}$$

Choose, arbitrarily, k initial values: $Y_1 = y_1^0, Y_2 = y_2^0, \dots, Y_k = y_k^0$.

Create:

- x^1 by a draw from $f(X \mid y_1^0, \dots, y_k^0)$
- y_1^1 by a draw from $f(Y_1 \mid x^1, y_2^0, \dots, y_k^0)$
- y_2^1 by a draw from $f(Y_2 \mid x^1, y_1^1, y_3^0, \dots, y_k^0)$
- \dots
- y_k^1 by a draw from $f(Y_k \mid x^1, y_1^1, \dots, y_{k-1}^1)$.

This constitutes one Gibbs “pass” through the $k+1$ conditional distributions,

yielding values: $(x^1, y_1^1, y_2^1, \dots, y_k^1)$.

Iterate the sampling to form the second “pass”

$$(x^2, y_1^2, y_2^2, \dots, y_k^2).$$

Theorem: (under general conditions)

The distribution of x^n converges to $F(x)$ as $n \rightarrow \infty$.

Thus, we may take the last n X -values after many Gibbs passes:

$$\frac{1}{n} \sum_{i=m}^{m+n} X^i \approx E[X]$$

or take just the last value, $x_i^{n_i}$ of n -many sequences of Gibbs passes

$$(i = 1, \dots, n) \quad \frac{1}{n} \sum_{i=1}^n X_i^{n_i} \approx E[X]$$

to solve for the average, $= \int x f(x) dx$.

A bivariate example of the Gibbs Sampler.

Example: Let X and Y have similar truncated conditional exponential distributions:

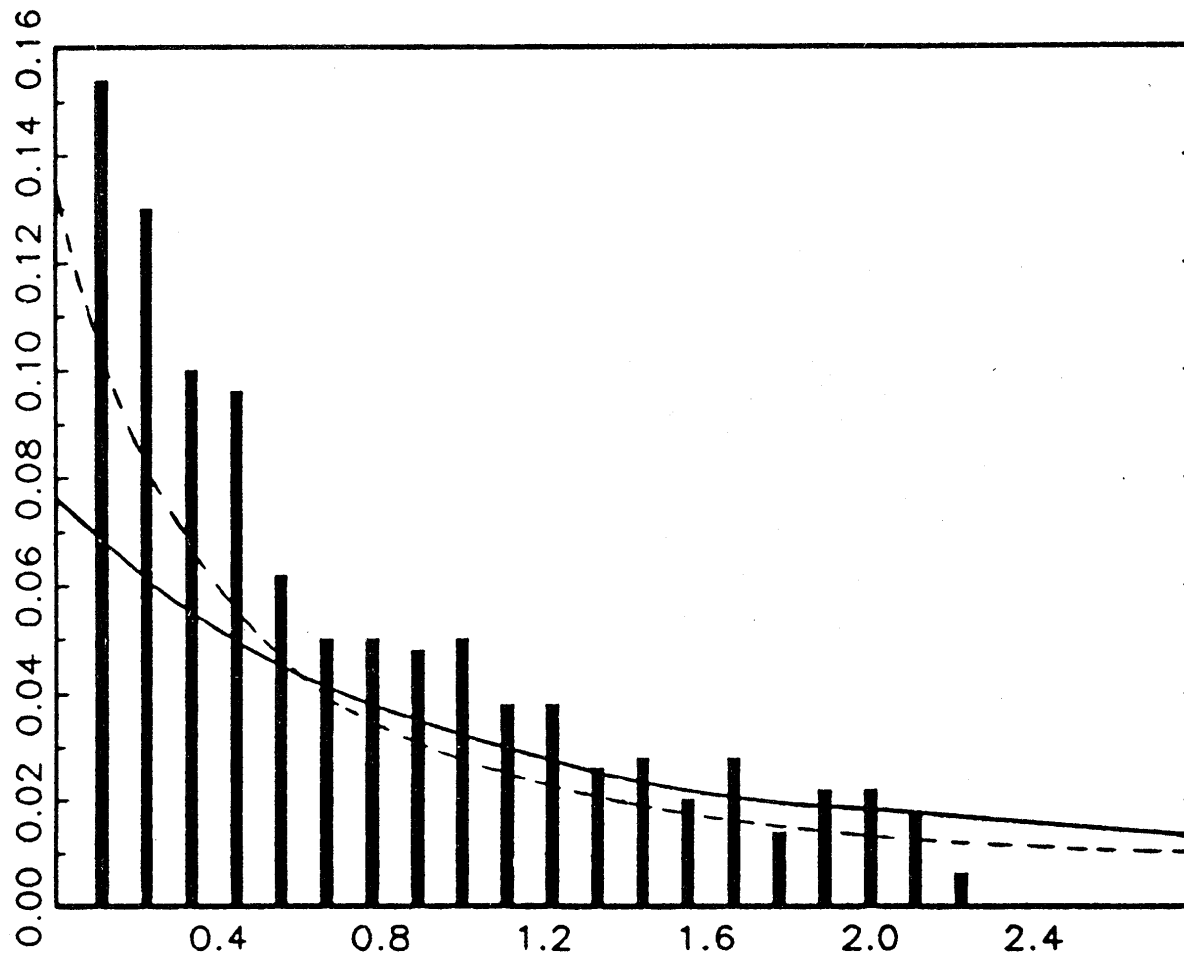
$$f(X|y) \propto ye^{-yx} \text{ for } 0 < X < \mathbf{b}$$

$$f(Y|x) \propto xe^{-xy} \text{ for } 0 < Y < \mathbf{b}$$

where \mathbf{b} is a known, positive constant.

Though it is not convenient to calculate, the marginal density $f(X)$ is readily simulated by Gibbs sampling from these (truncated) exponentials.

Below is a histogram for X , $\mathbf{b} = 5.0$, using a sample of 500 terminal observations with 15 Gibbs' passes per trial, $x_i^{n_i}$ ($i = 1, \dots, 500, n_i = 15$) (from Casella and George, 1992).



Histogram for X , $b = 5.0$, using a sample of 500 terminal observations with 15 Gibbs' passes per trial,

$x_i^{n_i}$ ($i = 1, \dots, 500, n_i = 15$). Taken from (Casella and George, 1992).

Here is an alternative way to compute the marginal $f(X)$ using the same Gibbs Sampler.

Recall the law of conditional expectations (assuming $E[X]$ exists):

$$E[E[X | Y]] = E[X]$$

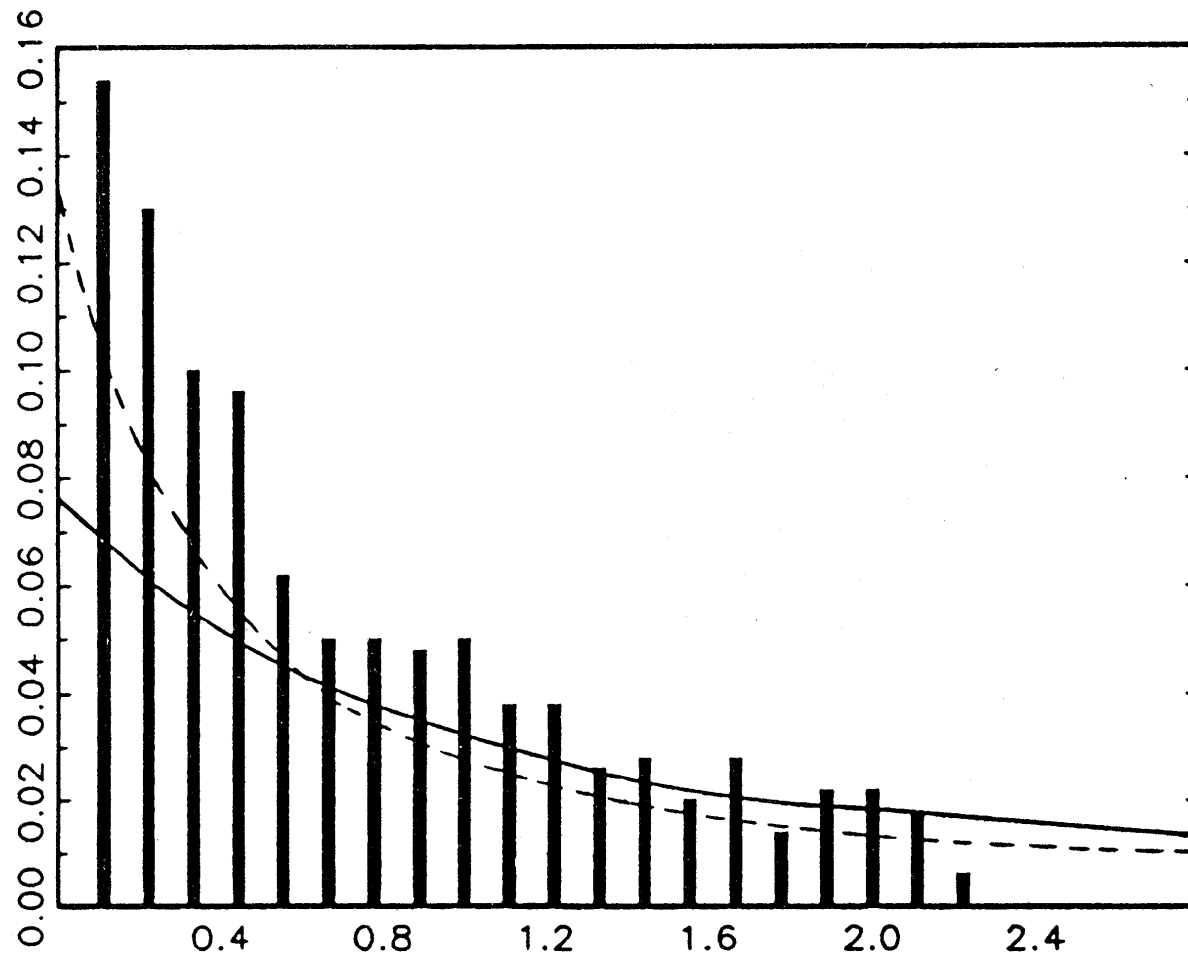
Thus $E[f(x|Y)] = \int f(x | y)f(y)dy = f(x)$.

Now, use the fact that the Gibbs sampler gives us a simulation of the marginal density $f(Y)$ using the penultimate values (for Y) in each Gibbs' pass, above: $y_i^{n_i-1}$ ($i = 1, \dots, 500; n_i = 15$).

Calculate $f(x | y_i^{n_i-1})$, which by assumption is feasible.

Then note that:

$$f(x) \approx \frac{1}{n} \sum_{i=1}^n f(x | y_i^{n_i-1})$$



The **solid line** graphs the alternative Gibbs Sampler estimate of the marginal $f(x)$ from the same sequence of 500 Gibbs' passes, using $\int f(x | y)f(y)dy = f(x)$. The **dashed-line** is the exact solution. Taken from (Casella and George, 1992).

An elementary proof of convergence in the case of 2 x 2 Bernoulli data

Let (X, Y) be a bivariate variable, marginally, each is Bernoulli

$$\begin{array}{c} X \\ 0 \quad 1 \\ Y^0 \left[\begin{array}{cc} p_1 & p_2 \\ p_3 & p_4 \end{array} \right] \\ 1 \end{array}$$

where $p_i \geq 0$, $\sum p_i = 1$, marginally

$$\mathbf{P}(X=0) = p_1 + p_3 \quad \text{and} \quad \mathbf{P}(X=1) = p_2 + p_4$$

$$\mathbf{P}(Y=0) = p_1 + p_2 \quad \text{and} \quad \mathbf{P}(Y=1) = p_3 + p_4.$$

The conditional probabilities $\mathbf{P}(X|y)$ and $\mathbf{P}(Y|x)$ are evident:

$\mathbf{P}(Y|x)$:

$$\begin{array}{c}
 \mathbf{Y} \\
 \begin{array}{cc}
 0 & 1 \\
 \left[\begin{array}{cc}
 \frac{p_1}{p_1+p_3} & \frac{p_2}{p_2+p_4} \\
 \frac{p_3}{p_1+p_3} & \frac{p_4}{p_2+p_4}
 \end{array} \right]
 \end{array}
 \end{array}$$

$\mathbf{P}(X|y)$:

$$\begin{array}{c}
 \mathbf{Y} \\
 \begin{array}{cc}
 0 & 1 \\
 \left[\begin{array}{cc}
 \frac{p_1}{p_1+p_2} & \frac{p_2}{p_1+p_2} \\
 \frac{p_3}{p_3+p_4} & \frac{p_4}{p_3+p_4}
 \end{array} \right]
 \end{array}
 \end{array}$$

Suppose (for illustration) that we want to generate the marginal distribution of X by the Gibbs Sampler, using the sequence of iterations of draws between the two conditional probabilities $\mathbf{P}(X|y)$ and $\mathbf{P}(Y|x)$.

That is, we are interested in the sequence $\langle x^i : i = 1, \dots \rangle$ created from the starting value $y^0 = 0$ or $y^0 = 1$.

Note that:

$$\begin{aligned} \mathbf{P}(X^n = 0 \mid x^i : i = 1, \dots, n-1) &= \mathbf{P}(X^n = 0 \mid x^{n-1}) \textit{ the Markov property} \\ &= \mathbf{P}(X^n = 0 \mid y^{n-1} = 0) \mathbf{P}(Y^{n-1} = 0 \mid x^{n-1}) + \mathbf{P}(X^n = 0 \mid y^{n-1} = 1) \mathbf{P}(Y^{n-1} = 1 \mid x^{n-1}) \end{aligned}$$

Thus, we have the four (positive) transition probabilities:

$$\mathbf{P}(X^n = j | x^{n-1} = i) = p_{ij} > 0, \text{ with } \sum_i \sum_j p_{ij} = 1 \quad (i, j = 0, 1).$$

With the transition probabilities positive, it is an (old) ergodic theorem that, $\mathbf{P}(X^n)$ converges to a (unique) *stationary* distribution, independent of the starting value (y^0).

Next, we confirm the easy fact that the marginal distribution $\mathbf{P}(X)$ is that same distinguished *stationary* point of this Markov process.

$$\begin{aligned}
& \mathbf{P}(X^n = 0) \\
&= \mathbf{P}(X^n = 0 \mid x^{n-1} = 0) \mathbf{P}(X^{n-1} = 0) + \mathbf{P}(X^n = 0 \mid x^{n-1} = 1) \mathbf{P}(X^{n-1} = 1) \\
&= \mathbf{P}(X^n=0 \mid y^{n-1}=0) \mathbf{P}(Y^{n-1}=0 \mid x^{n-1} = 0) \mathbf{P}(X^{n-1} = 0) \\
&\quad + \mathbf{P}(X^n=0 \mid y^{n-1}=1) \mathbf{P}(Y^{n-1}=1 \mid x^{n-1} = 0) \mathbf{P}(X^{n-1} = 0) \\
&\quad + \mathbf{P}(X^n=0 \mid y^{n-1}=0) \mathbf{P}(Y^{n-1}=0 \mid x^{n-1} = 1) \mathbf{P}(X^{n-1} = 1) \\
&\quad + \mathbf{P}(X^n=0 \mid y^{n-1}=1) \mathbf{P}(Y^{n-1}=1 \mid x^{n-1} = 1) \mathbf{P}(X^{n-1} = 1) \\
&= \mathbf{E}_{\mathbf{P}} [\mathbf{E}_{\mathbf{P}} [X^n=0 \mid X^{n-1}]] \\
&= \mathbf{E}_{\mathbf{P}} [X^n=0] \\
&= \mathbf{P}(X^n = 0) .
\end{aligned}$$

The *Ergodic* Theorem:

Definitions:

- A *Markov chain*, X_0, X_1, \dots satisfies

$$\mathbf{P}(X_n | x_i: i = 1, \dots, n-1) = \mathbf{P}(X_n | x_{n-1})$$

- The distribution $F(x)$, with density $f(x)$, for a Markov chain is *stationary* (or *invariant*) if

$$\int_{\mathbf{A}} f(x) dx = \int \mathbf{P}(X_n \in \mathbf{A} | x_{n-1}) f(x) dx.$$

- The Markov chain is *irreducible* if each set with positive \mathbf{P} -probability is visited at some point (almost surely).

- An irreducible Markov chain is *recurrent* if, for each set \mathbf{A} having positive \mathbf{P} -probability, with positive \mathbf{P} -probability the chain visits \mathbf{A} infinitely often.
- A Markov chain is *periodic* if for some integer $k > 1$, there is a partition into k sets $\{\mathbf{A}_1, \dots, \mathbf{A}_k\}$ such that $\mathbf{P}(X_{n+1} \in \mathbf{A}_{j+1} \mid x_n \in \mathbf{A}_j) = 1$ for all $j = 1, \dots, k-1 \pmod{k}$. That is, the chain cycles through the partition. Otherwise, the chain is *aperiodic*.

Theorem: If the Markov chain X_0, X_1, \dots is irreducible with an invariant probability distribution $F(x)$ then:

1. the Markov chain is recurrent
2. F is the unique invariant distribution

If the chain is aperiodic, then for F -almost all x_0 , both

$$3. \lim_{n \rightarrow \infty} \sup_{\mathbf{A}} | \mathbf{P}(X_n \in \mathbf{A} | X_0 = x_0) - \int_{\mathbf{A}} \mathbf{f}(x) dx | = 0$$

And for any function \mathbf{h} with $\int \mathbf{h}(x) dx < \infty$,

$$4. \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=0}^{n-1} \mathbf{h}(X_i) = \int \mathbf{h}(x) \mathbf{f}(x) dx \quad (= \mathbf{E}_{\mathbf{F}}[\mathbf{h}(x)]),$$

That is, the *time average* of $\mathbf{h}(X)$ equals its *state-average*, a.e. F .

A (now-familiar) puzzle.

Example (continued): Let X and Y have similar conditional exponential distributions:

$$f(X|y) \propto ye^{-yx} \text{ for } 0 < X$$

$$f(Y|x) \propto xe^{-xy} \text{ for } 0 < Y$$

To solve for the marginal density $f(X)$ use Gibbs sampling from these exponential distributions. The resulting sequence does *not* converge!

Question: Why does this happen?

Answer: (Hint: Recall HW #1, problem 2.) Let θ be the statistical parameter for X with $f(X|\theta)$ the exponential model. What “prior” density for θ yields the *posterior* $f(\theta|x) \propto xe^{-x\theta}$?

Then, what is the “prior” expectation for X ?

Remark: Note that $W = X\theta$ is pivotal. What is its distribution?

More on this puzzle:

The conjugate prior for the parameter θ in the exponential distribution is the Gamma $\Gamma(\alpha, \beta)$.

$$f(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta} \quad \text{for } \theta, \alpha, \beta > 0,$$

Then the posterior for θ based on $\mathbf{x} = (x_1, \dots, x_n)$, n iid observations from the exponential distribution is

$$f(\theta|\mathbf{x}) \text{ is Gamma } \Gamma(\alpha', \beta')$$

where $\alpha' = \alpha+n$ and $\beta' = \beta + \sum x_i$.

Let $n=1$, and consider the limiting distribution as $\alpha, \beta \rightarrow 0$.

This produces the “posterior” density $f(\theta|x) \propto x e^{-x\theta}$, which is mimicked in Bayes theorem by the improper “prior” density

$f(\theta) \propto 1/\theta$. But then $E_{\mathbf{F}}(\theta)$ does not exist!

Part 2 EM – again

- **EM as a maximization/maximization method**
- **Gibbs as a variation of Generalized EM**

EM as a maximization/maximization method.

Recall:

$L(\theta ; \mathbf{x})$ is the likelihood function for θ with respect to the incomplete data \mathbf{x} .

$L(\theta ; (\mathbf{x}, \mathbf{z}))$ is the likelihood for θ with respect to the complete data (\mathbf{x}, \mathbf{z}) .

And $L(\theta ; \mathbf{z} | \mathbf{x})$ is a *conditional likelihood* for θ with respect to \mathbf{z} , given \mathbf{x} ;

which is based on $h(\mathbf{z} | \mathbf{x}, \theta)$: the conditional density for the data \mathbf{z} , given (\mathbf{x}, θ) .

Then as
$$f(\mathbf{X} | \theta) = f(\mathbf{X}, \mathbf{Z} | \theta) / h(\mathbf{Z} | \mathbf{x}, \theta)$$

we have
$$\log L(\theta ; \mathbf{x}) = \log L(\theta ; (\mathbf{x}, \mathbf{z})) - \log L(\theta ; \mathbf{z} | \mathbf{x}) \quad (*)$$

As below, we use the EM algorithm to compute the mle

$$\hat{\theta} = \operatorname{argmax}_{\Theta} L(\theta ; \mathbf{x})$$

With $\hat{\theta}_0$ an arbitrary choice, define

$$(E\text{-step}) \quad Q(\theta | \mathbf{x}, \hat{\theta}_0) = \int_{\mathcal{Z}} [\log \mathbf{L}(\theta ; \mathbf{x}, \mathbf{z})] \mathbf{h}(\mathbf{z} | \mathbf{x}, \hat{\theta}_0) d\mathbf{z}$$

and

$$H(\theta | \mathbf{x}, \hat{\theta}_0) = \int_{\mathcal{Z}} [\log \mathbf{L}(\theta ; \mathbf{z} | \mathbf{x})] \mathbf{h}(\mathbf{z} | \mathbf{x}, \hat{\theta}_0) d\mathbf{z}.$$

then $\log \mathbf{L}(\theta ; \mathbf{x}) = Q(\theta | \mathbf{x}, \hat{\theta}_0) - H(\theta | \mathbf{x}, \hat{\theta}_0)$,

as we have integrated-out \mathbf{z} from (*) using the conditional density $\mathbf{h}(\mathbf{z} | \mathbf{x}, \hat{\theta}_0)$.

The *EM algorithm* is an iteration of

- i. the *E*-step: determine the integral $Q(\theta | \mathbf{x}, \hat{\theta}_j)$,
- ii. the *M*-step: define $\hat{\theta}_{j+1}$ as $\mathit{argmax}_{\Theta} Q(\theta | \mathbf{x}, \hat{\theta}_j)$.

Continue until there is convergence of the $\hat{\theta}_j$.

Now, for a *Generalized EM* algorithm.

Let be $P(\mathbf{Z})$ any distribution over the augmented data \mathbf{Z} , with density $p(\mathbf{z})$
Define the function F by:

$$\begin{aligned} F(\theta, P(\mathbf{Z})) &= \int_{\mathbf{Z}} [\log \mathbf{L}(\theta; \mathbf{x}, \mathbf{z})] p(\mathbf{z}) d\mathbf{z} - \int_{\mathbf{Z}} \log p(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\ &= \mathbf{E}_P [\log \mathbf{L}(\theta; \mathbf{x}, \mathbf{z})] - \mathbf{E}_P [\log p(\mathbf{z})] \end{aligned}$$

When $p(\mathbf{Z}) = h(\mathbf{Z} | \mathbf{x}, \hat{\theta}_0)$ from above, then $F(\theta, P(\mathbf{Z})) = \log \mathbf{L}(\theta; \mathbf{x})$.

Claim: For a fixed (arbitrary) value $\theta = \hat{\theta}_0$, $F(\hat{\theta}_0, P(\mathbf{Z}))$ is maximized over distributions $P(\mathbf{Z})$ by choosing $p(\mathbf{Z}) = h(\mathbf{Z} | \mathbf{x}, \hat{\theta}_0)$.

Thus, the *EM* algorithm is a sequence of *M-M* steps: the old *E*-step now is a max over the second term in $F(\hat{\theta}_0, P(\mathbf{Z}))$, given the first term. The second step remains (as in *EM*) a max over θ for a fixed second term, which does not involve θ

Suppose that the augmented data \mathbf{Z} are multidimensional.

Consider the *GEM* approach and, instead of maximizing the choice of $P(\mathbf{Z})$ over all of the augmented data – instead of the old *E*-step – instead maximize over only *one* coordinate of \mathbf{Z} at a time, alternating with the (old) *M*-step.

This gives us the following link with the Gibbs algorithm: Instead of maximizing at each of these two steps, use the conditional distributions, we sample from them!

Part 3) Generating a Random Variable

3.1) Continuous r.v.'s – an Exact Method using transformation of the CDF

- Let Y be a continuous r.v. with **cdf** $F_Y(\bullet)$. Then the range of $F_Y(\bullet)$ is $(0, 1)$, and as a r.v. F_Y it is distributed $U \sim \text{Uniform}(0,1)$. Thus the *inverse* transformation $F_Y^{-1}(U)$ gives us the desired distribution for Y .

Examples:

- If $Y \sim \text{Exponential}(\lambda)$ then $F_Y^{-1}(U) = -\lambda \ln(1-U)$ is the desired Exponential.

And from known relationships between the Exponential distribution and other members of the Exponential Family, we may proceed as follows.

Let U_j be *iid* Uniform(0,1), so that $Y_j = -\lambda \ln(U_j)$ are *iid* Exponential(λ)

- $Z = -2 \sum_{j=1}^n \ln(U_j) \sim \chi^2_{2n}$ a Chi-squared distribution on $2n$ degrees of freedom

Note only even integer values possible here, alas!

- $Z = -\beta \sum_{j=1}^a \ln(U_j) \sim \text{Gamma } \Gamma(a, \beta)$ – with integer values only for a .

- $Z = \frac{\sum_{j=1}^a \ln(U_j)}{\sum_{j=1}^{a+b} \ln(U_j)} \sim \text{Beta}(a,b)$ – with integer values only for a .

3.2) The “Accept/Reject” algorithm for approximations using pdf’s.

Suppose we want to generate $Y \sim \text{Beta}(a,b)$, for non-integer values of a and b , say $a = 2.7$ and $b = 6.3$.

Let (U, V) be independent Uniform(0, 1) random variables. Let $c \geq \max_y f_Y(y)$.
Now calculate $P(Y \leq y)$ as follows:

$$\begin{aligned} P(V \leq y, U \leq (1/c) f_Y(V)) &= \int_0^y \int_0^1 f_Y(v)^{1/c} du dv \\ &= (1/c) \int_0^y f_Y(v) dv \\ &= (1/c) P(Y \leq y). \end{aligned}$$

So: (i) generate independent (U, V) Uniform(0,1)

(ii) If $U < (1/c)f_Y(V)$, set $Y = V$, otherwise, return to step (i).

Note: The waiting time for one value of Y with this algorithm is c , so we want c small. Thus, choose $c = \max_y f_Y(y)$. But we waste generated values of U, V whenever $U \geq (1/c)f_Y(V)$, so we want to choose a better approximation distribution for V than the uniform.

Let $Y \sim f_Y(y)$ and $V \sim f_V(v)$.

- Assume that these two have common support, i.e., the smallest closed sets of measure one are the same.
- Also, assume that $\mathbf{M} = \sup_y [f_Y(y) / f_V(y)]$ exists, i.e., $\mathbf{M} < \infty$.

Then generate the *r.v.* $Y \sim f_Y(y)$ using

$U \sim \text{Uniform}(0,1)$ and $V \sim f_V(v)$, with (U, V) independent, as follows:

- (i) Generate values (u, v) .
- (ii) If $u < (1/\mathbf{M}) f_Y(v) / f_V(v)$ then set $y = v$.
If not, return to step (i) and redraw (u, v) .

Proof of correctness for the accept/reject algorithm:

The generated r.v. Y has a *cdf*

$$\begin{aligned} P(Y \leq y) &= P(V \leq y \mid \text{stop}) \\ &= P(V \leq y \mid U < (1/M) f_Y(v) / f_V(y)) \\ &= \frac{P(V \leq y, U < (1/M) f_Y(V) / f_V(V))}{P(U < (1/M) f_Y(V) / f_V(V))} \\ &= \frac{\int_{-\infty}^y \int_0^{(1/M) f_Y(v) / f_V(v)} du f_V(v) dv}{\int_{-\infty}^{\infty} \int_0^{(1/M) f_Y(v) / f_V(v)} du f_V(v) dv} \\ &= \int_{-\infty}^y f_Y(v) dv. \end{aligned}$$

Example: Generate $Y \sim \text{Beta}(2.7, 6.3)$.

Let $V \sim \text{Beta}(2, 6)$. Then $M = 1.67$ and we may proceed with the algorithm.

3.3) Metropolis algorithm for “heavy-tailed” target densities.

As before, let $Y \sim f_Y(y)$, $V \sim f_V(v)$, $U \sim \text{Uniform}(0,1)$, with (U, V) independent.

Assume only that Y and V have a common support.

Metropolis Algorithm:

Step₀: Generate v_0 and set $z_0 = v_0$. For $i = 1, \dots,$

Step_i: Generate (u_i, v_i)

Define
$$\rho_i = \min \left\{ \frac{f_Y(v_i)}{f_V(v_i)} \times \frac{f_V(z_{i-1})}{f_Y(z_{i-1})}, 1 \right\}$$

Let
$$z_i = \begin{cases} v_i & \text{if } u_i \leq \rho_i \\ z_{i-1} & \text{if } u_i > \rho_i. \end{cases}$$

Then, as $i \rightarrow \infty$, the r.v. Z_i converges in distribution to the random variable Y .

Additional References

Casella, G. and George, E. (1992) “Explaining the Gibbs Sampler,”
Amer. Statistician **46**, 167-174.

Flury, B. and Zoppe, A. (2000) “Exercises in EM,” *Amer. Staisticalian* **54**,
207-209.

Hastie, T., Tibshirani, R, and Friedman, J. *The Elements of Statistical Learning*. New York: Spring-Verlag, 2001, sections 8.5-8.6.

Tierney, L. (1994) “Markov chains for exploring posterior distributions”
(with discussion) *Annals of Statistics* **22**, 1701-1762.