# 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
$$f(x, y_1, \ldots, y_k)$$
and we are interested, say, in some features of the marginal density
$$f(x) = \iint \ldots \int f(x, y_1, \ldots, y_k) \, dy_1, dy_2, \ldots, dy_k.$$

*F*or 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
$$lim_{m \to \infty} \frac{1}{n} \sum_{i=1}^{n} X_i = E[X]$$
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:*

$$f(X \mid y_1, \ldots, y_k)$$
$$f(Y_1 \mid x, y_2, \ldots, y_k)$$
$$f(Y_2 \mid x, y_1, y_3, \ldots, y_k)$$
$$\ldots$$
$$f(Y_k \mid x, y_1, y_3, \ldots, y_{k-1}).$$

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

Create:      $x^1$ by a draw from $f(X \mid y_1^0, \ldots, y_k^0)$

$y_1^1$ by a draw from $f(Y_1 \mid x^1, y_2^0, \ldots, y_k^0)$

$y_2^1$ by a draw from $f(Y_2 \mid x^1, y_1^1, y_3^0 \ldots, y_k^0)$

$\ldots$

$y_k^1$ by a draw from $f(Y_k \mid x^1, y_1^1, \ldots, y_{k-1}^1).$

This constitutes one Gibbs "pass" through the k+1 conditional distributions,

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

Iterate the sampling to form the second "pass"

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

*Theorem*: (under general conditions)

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

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

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

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

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

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

A bivariate example of the Gibbs Sampler.

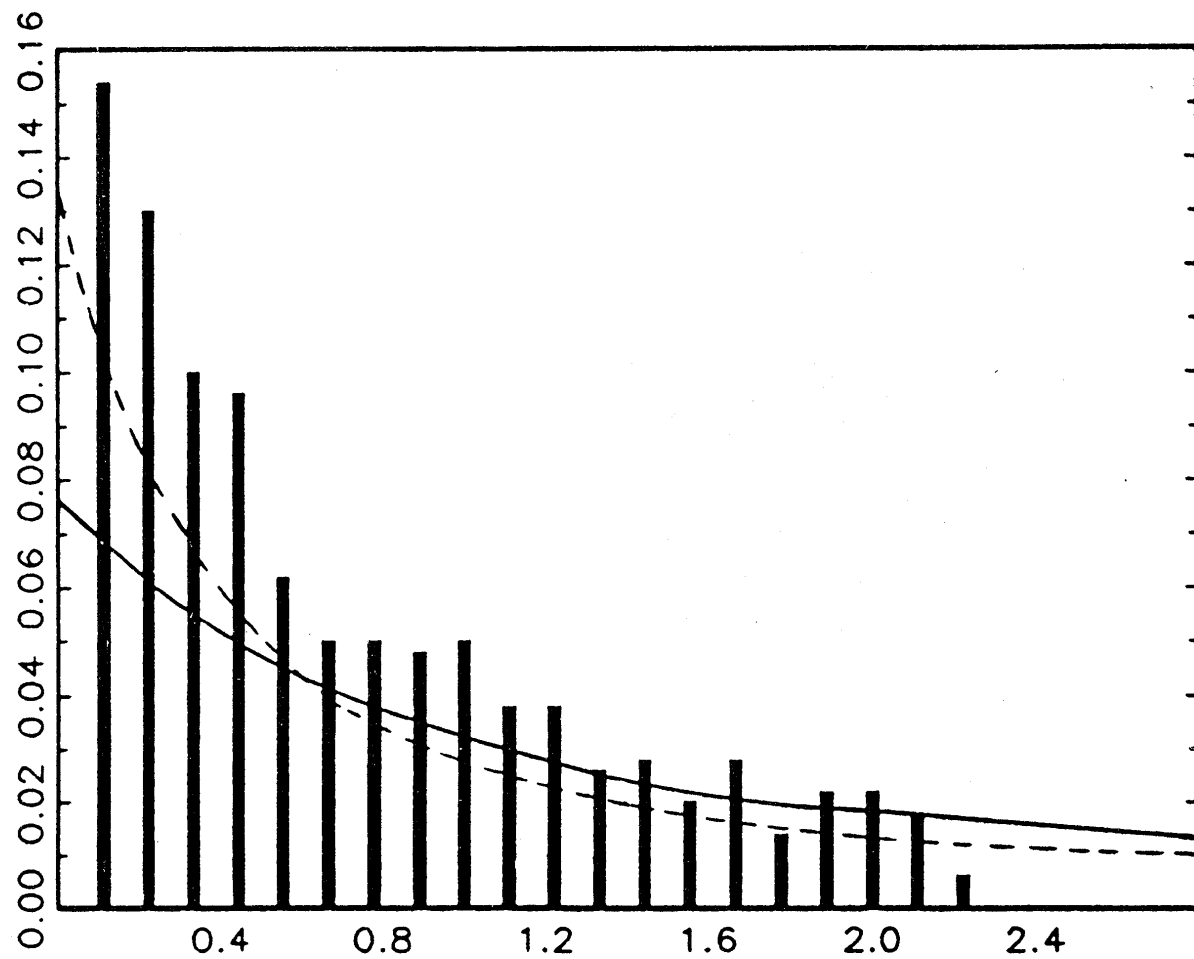*Example*: Let $X$ and $Y$ have similar truncated conditional exponential distributions:

$$f(X \mid y) \propto y e^{-yx} \text{ for } 0 < X < \boldsymbol{b}$$

$$f(Y \mid x) \propto x e^{-xy} \text{ for } 0 < Y < \boldsymbol{b}$$

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

Histogram for $X$, **b** = 5.0, using a sample of 500 terminal observations with 15 Gibbs' passes per trial, $x_i^{n_i}$ ($i = 1,\ldots, 500$, $n_i = 15$). 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):
$$E[\ E[X\,|\,Y]\ ] = E[X]$$
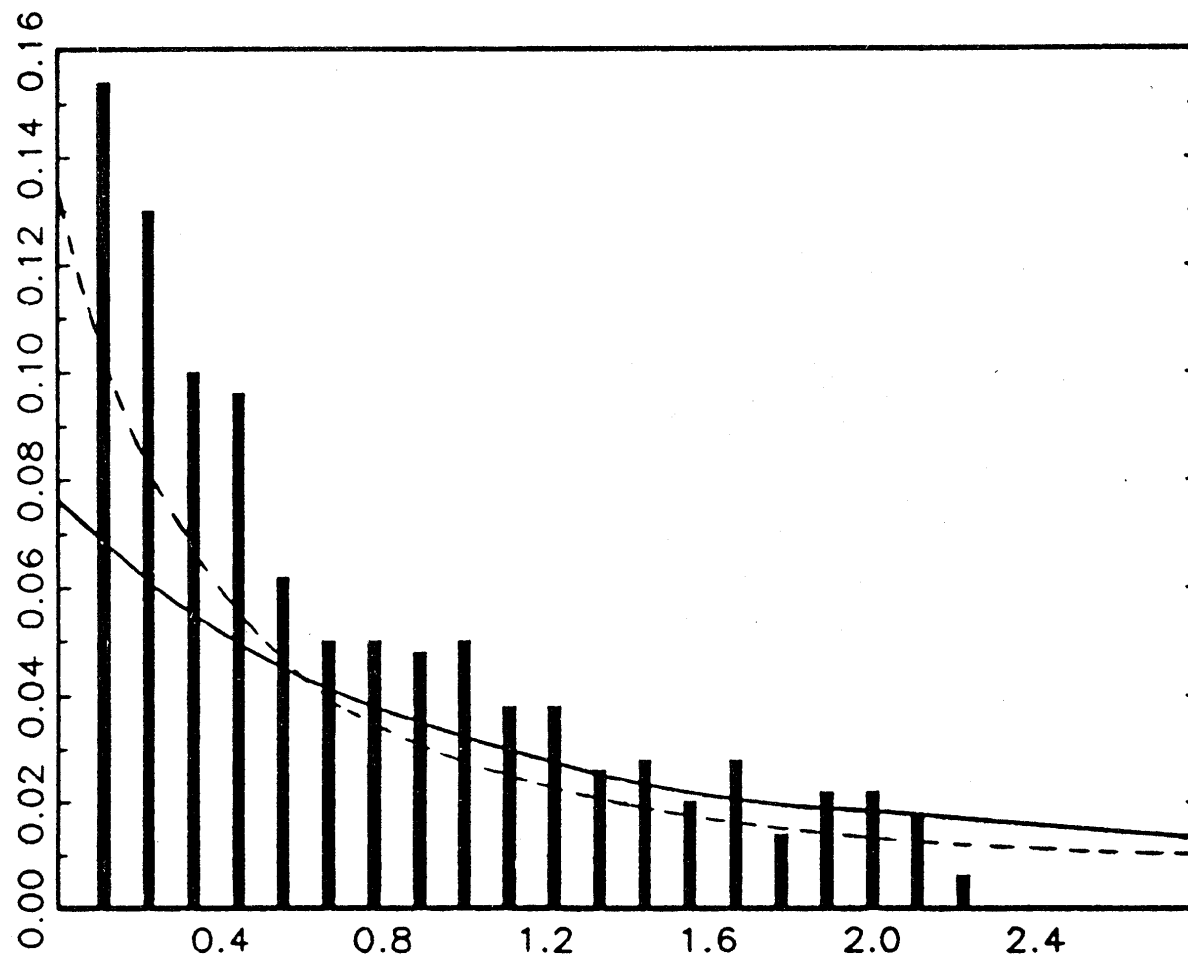
Thus
$$E[f(x|Y)] = \int f(x\,|\,y)f(y)dy = f(x).$$

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: $\quad y_i^{n_i-1}$ (i = 1, …500; $n_i$ = 15).

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

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

The **solid line** graphs the alternative Gibbs Sampler estimate of the marginal $f(x)$ from eth same sequence of 500 Gibbs' passes, using $\int f(x \mid y) f(y) dy = f(x)$. 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

$$
\begin{array}{cc}
 & X \\
 & \begin{array}{cc} 0 & 1 \end{array} \\
Y \begin{array}{c} 0 \\ 1 \end{array} & \begin{bmatrix} p_1 & p_2 \\ p_3 & p_4 \end{bmatrix}
\end{array}
$$

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

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

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

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

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

$$
\begin{array}{cc}
 & X \\
 & \begin{array}{cc} 0 & \quad 1 \end{array} \\
Y \begin{array}{c} 0 \\ 1 \end{array} & \left[ \begin{array}{cc} \dfrac{p_1}{p_1+p_3} & \dfrac{p_2}{p_2+p_4} \\[2ex] \dfrac{p_3}{p_1+p_3} & \dfrac{p_4}{p_2+p_4} \end{array} \right]
\end{array}
$$

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

$$
\begin{array}{cc}
 & X \\
 & \begin{array}{cc} 0 & \quad 1 \end{array} \\
Y \begin{array}{c} 0 \\ 1 \end{array} & \left[ \begin{array}{cc} \dfrac{p_1}{p_1+p_2} & \dfrac{p_2}{p_1+p_2} \\[2ex] \dfrac{p_3}{p_3+p_4} & \dfrac{p_4}{p_3+p_4} \end{array} \right]
\end{array}
$$

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 probabilites $\mathbf{P}(X|y)$ and $\mathbf{P}(Y|x)$.

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

Note that:

$$\mathbf{P}(X^n = 0 \,|\, x^i : i = 1, \ldots, n\text{-}1) = \mathbf{P}(X^n = 0 \,|\, x^{n-1}) \quad \textit{the Markov property}$$

$$= \mathbf{P}(X^n = 0 \,|\, y^{n-1} = 0) \, \mathbf{P}(Y^{n-1} = 0 \,|\, x^{n-1}) \; + \; \mathbf{P}(X^n = 0 \,|\, y^{n-1} = 1) \, \mathbf{P}(Y^{n-1} = 1 \,|\, x^{n-1})$$

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

$$\mathbf{P}(X^n = \mathrm{j} \mid x^{n-1} = i) = p_{ij} > 0, \text{ with } \Sigma_i \, \Sigma_j \, p_{ij} = 1 \quad (i, j = 0, 1).$$

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.

$$\mathbf{P}(X^n = 0)$$

$$= \quad \mathbf{P}(X^n = 0 \mid x^{n-1} = 0)\, \mathbf{P}(X^{n-1} = 0) \;+\; \mathbf{P}(X^n = 0 \mid x^{n-1} = 1)\, \mathbf{P}(X^{n-1} = 1)$$

$$= \quad \mathbf{P}(X^n{=}0 \mid y^{n-1}{=}0)\, \mathbf{P}(Y^{n-1}{=}0 \mid x^{n-1} = 0)\, \mathbf{P}(X^{n-1} = 0)$$

$$+ \quad \mathbf{P}(X^n{=}0 \mid y^{n-1}{=}1)\, \mathbf{P}(Y^{n-1}{=}1 \mid x^{n-1} = 0)\, \mathbf{P}(X^{n-1} = 0)$$

$$+ \quad \mathbf{P}(X^n{=}0 \mid y^{n-1}{=}0)\, \mathbf{P}(Y^{n-1}{=}0 \mid x^{n-1} = 1)\, \mathbf{P}(X^{n-1} = 1)$$

$$+ \quad \mathbf{P}(X^n{=}0 \mid y^{n-1}{=}1)\, \mathbf{P}(Y^{n-1}{=}1 \mid x^{n-1} = 1)\, \mathbf{P}(X^{n-1} = 1)$$

$$= \quad \mathbf{E_P}\,[\mathbf{E_P}\,[X^n{=}0 \mid X^{n-1}]\,]$$

$$= \quad \mathbf{E_P}\,[X^n{=}0\,]$$

$$= \quad \mathbf{P}(X^n = 0)\,.$$

The *Ergodic* Theorem:

*Definitions*:

- A *Markov chain*, $X_0, X_1, \ldots$ satisfies

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

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

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

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

- An irreducible Markov chain is *recurrent* if, for each set **A** having positive **P**-probability, with positive **P**-probability the chain visits **A** infinitely often.

- A Markov chain is *periodic* if for some integer $k > 1$, there is a partition into $k$ sets $\{\mathbf{A}_1, \ldots, \mathbf{A}_k\}$ such that

  $\mathbf{P}(X_{n+1} \in \mathbf{A}_{j+1} \mid x_n \in \mathbf{A}_j) = 1$ for all $j = 1, \ldots, k\text{-}1$ (mod k). That

  is, the chain cycles through the partition.

  Otherwise, the chain is *aperiodic*.

*Theorem*:  If the Markov chain $X_0, X_1, \ldots$  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

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

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

$$4.\ \ lim_{n\to\infty}\ \tfrac{1}{n}\sum_{i=i}^{n} h(X_i)\ =\ \int h(x)\, f(x)\ dx\ \ (= \mathbf{E_F}[h(x)]\ ),$$

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

A (now-familiar) puzzle.

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

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

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

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 $\theta$ be the statistical parameter for $X$ with $f(X|\theta)$ the exponential model. What "prior" density for $\theta$ yields the *posterior* $f(\theta \mid x) \propto x e^{-x\theta}$? 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 $\theta$ in the exponential distribution is the Gamma $\Gamma(\alpha, \beta)$.

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

Then the posterior for $\theta$ based on $x = (x_1, .., x_n)$, $n$ *iid* observations from the exponential distribution is

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

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

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

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

$f(\theta) \propto 1/\theta$. But then $E_F(\theta)$ 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:**

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

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

And $L(\theta ; z \mid x)$ is a *conditional likelihood* for $\theta$ with respect to $z$, given $x$;

which is based on $h(z \mid x, \theta)$: the conditional density for the data $z$, given $(x,\theta)$.

Then as $\qquad\qquad f(X \mid \theta) = f(X, Z \mid \theta) \ / \ h(Z \mid x, \theta)$

we have $\qquad\qquad log \ L(\theta ; x) = log \ L(\theta ; (x, z)) - log \ L(\theta ; z \mid x) \quad$ (*)

$$As \ below, \ we \ use \ the \ \mathbf{EM} \ algorithm \ to \ compute \ the \ \mathbf{mle}$$
$$\hat{\theta} \ = \ argmax_{\Theta} \ L(\theta ; x)$$

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

(**E**-step) $\qquad\qquad Q(\theta \mid \pmb{x}, \hat{\theta}_0) \;=\; \int_Z [\pmb{log}\, \mathbf{L}(\theta\,;\,\pmb{x},\,\pmb{z})]\, \pmb{h}(\pmb{z} \mid \pmb{x}, \hat{\theta}_0)\, dz$

$$\text{and}$$

$$H(\theta \mid \pmb{x},\, \hat{\theta}_0) \;=\; \int_Z [\pmb{log}\, \mathbf{L}(\theta\,;\,\pmb{z} \mid \pmb{x})]\, \pmb{h}(\pmb{z} \mid \pmb{x},\, \hat{\theta}_0)\, dz.$$

then $\qquad\qquad \pmb{log}\, \mathbf{L}(\theta\,;\,\pmb{x}) \;=\; Q(\theta \mid \pmb{x},\, \theta_0) - H(\theta \mid \pmb{x},\, \theta_0),$

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

The **EM algorithm** is an iteration of
  i.  the **E**-step: determine the integral $Q(\theta \mid \pmb{x},\, \hat{\theta}_j)$,
  ii.  the **M**-step: define $\hat{\theta}_{j+1}$ as $\pmb{argmax}_\Theta\, Q(\theta \mid \pmb{x},\, \hat{\theta}_j)$.
Continue until there is convergence of the $\hat{\theta}_j$.

Now, for a *Generalized EM* algorithm.

Let be **P(Z)** any distribution over the augmented data **Z**, with density **p(z)**
*Define* the function **F** by**:**

$$F(\theta, P(Z)) = \int_Z [log\ \mathbf{L}(\theta; x, z)]\ p(z)dz - \int_Z log\ p(z)\ p(z)dz$$
$$= \mathbf{E}_P[log\ \mathbf{L}(\theta; x, z)] - \mathbf{E}_P[\ log\ p(z)]$$

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

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

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

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

Consider the *GEM* approach and, instead of maximizing the choice of $\mathbf{P(Z)}$ 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) $\mathbf{M}$-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 $U \sim$ Uniform $(0,1)$.  Thus the *inverse* tranformation $F_Y^{-1}(U)$ gives us the desired distribution for $Y$.

Examples:

- If $Y \sim$ 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($\lambda$)

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

    **Note** only even integer values possible here, alas!

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

- $Z = \dfrac{\sum_{j=1}^{a} ln(U_j)}{\sum_{j=1}^{a+b} ln(U_j)} \sim$ 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:

$$P(V \leq y, \, U \leq (1/c)\,f_Y(V)\,) \;\; = \;\; \int_0^y \int_0^{f_Y(v)/c} du\,dv$$

$$= \;\; (1/c) \int_0^y f_Y(v)\,dv$$

$$= \;\; (1/c)\, P(Y \leq y).$$

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 algorthim 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 $M = sup_y [f_Y(y) / f_V(y)]$ exists, i.e., $M < \infty$.

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

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

(i)    Generate values $(u, v)$.

(ii)    If $u < (1/M) f_Y(v) / f_V(y)$ 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*

$$P(Y \le y) = P(V \le y \mid \text{stop})$$

$$= P(V \le y \mid U < (1/M) f_Y(v) / f_V(y))$$

$$= \frac{P(V \le y, U < (1/M) f_Y(V) / f_V(V))}{P(U < (1/M) f_Y(V) / f_V(V))}$$

$$= \frac{\int_{-\infty}^{y} \int_{0}^{(1/M) f_Y(v) / f_V(v)} du \, f_V(v) dv}{\int_{-\infty}^{\infty} \int_{0}^{(1/M) f_Y(v) / f_V(v)} du \, f_V(v) dv}$$

$$= \int_{-\infty}^{y} f_Y(v) dv.$$

*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 $Y \sim f_Y(y)$, $V \sim f_V(v)$, $U \sim$ Uniform(0,1), with $(U,V)$ independent.

Assume only that $Y$ and $V$ have a common support.

*Metropolis Algorithm*:

Step$_0$: Generate $v_0$ and set $z_0 = v_0$.        For $i = 1, \dots.,$

Step$_i$: Generate $(u_i, v_i)$

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

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. Staistican* **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.