1977 Dempster, Laird, Rubin

But 1974 Rolf Sundberg found it first

1977- generalized method

Wu, 1983: proved convergence for a wide variety of parametric models

___

Newton Raphson: useful for solving $f(x) = c$. We used it for finding $\hat{\theta}_{MAP} = \text{argmax} \{k(\theta | x)\}$ by solving $k'(\theta | x) = 0$ where it was not solvable in closed form.

E-M: useful for cases where MLE's break down due to not in closed form but if you knew some latent data it would be easy.

Can this help with our semi-conjugate model?

$X_1, \ldots, X_n \overset{iid}{\sim} N(\theta, \sigma^2)$, $\theta \sim N(m_0, \tau^2)$, $\sigma^2 \sim \text{Inv-Gamma}\left(\frac{h_0}{2}, \frac{h_0 \sigma_0^2}{2}\right)$

$$P(\theta, \sigma^2 | X) \propto P(X | \theta, \sigma^2) P(\theta) P(\sigma^2)$$

$$\propto \left(\prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \theta)^2}\right) \left(e^{-\frac{1}{2\tau^2}(\theta - m_0)^2}\right) \left((\sigma^2)^{\frac{h_0}{2} - 1} e^{-\frac{h_0 \sigma_0^2}{2\sigma^2}}\right)$$

$$\propto (\sigma^2)^{-\frac{n+h_0}{2} - 1} e^{-\frac{\Sigma x_i^2 + h_0 \sigma_0^2}{2\sigma^2}} e^{\frac{\theta n \bar{x}}{\sigma^2}} e^{-\frac{n\theta^2}{2\sigma^2}} e^{-\frac{\theta^2}{2\tau^2}} e^{\frac{\theta m_0}{\tau^2}}$$

$$e^{\left(\frac{\mu_0}{\tau^2} + \frac{h\bar{x}}{\sigma^2}\right)\theta} + \left(-\frac{h}{2\sigma^2} - \frac{1}{2\tau^2}\right)\theta^2$$

with braces labeled $a$ and $b$.

Need $\quad e^{-\frac{1}{2V}(\theta - c)^2}$

$$\Rightarrow -\frac{1}{2V}(\theta - c)^2 = -\frac{1}{2V}(\theta^2 - 2\theta c + c^2) = -\frac{\theta^2}{2V} + \frac{\theta c}{V} - \frac{c^2}{2V}$$

$$\Rightarrow -\frac{1}{2V} = b \Rightarrow V = -\frac{1}{2b} = -\frac{1}{2\left(-\frac{h}{2\sigma^2} - \frac{1}{2\tau^2}\right)} = \frac{1}{\frac{h}{\sigma^2} + \frac{1}{\tau^2}}$$

$$\frac{c}{V} = a \Rightarrow c = aV = \frac{\frac{\mu_0}{\tau^2} + \frac{h\bar{x}}{\sigma^2}}{\frac{h}{\sigma^2} + \frac{1}{\tau^2}}$$

$$-\frac{c^2}{2V} = bc^2 = \frac{1}{\frac{h}{\sigma^2} + \frac{1}{\tau^2}}\left(\frac{\frac{\mu_0}{\tau^2} + \frac{h\bar{x}}{\sigma^2}}{\frac{h}{\sigma^2} + \frac{1}{\tau^2}}\right)^2 = Q \quad \text{s.t.} \quad Q \not\ni \, f(\theta)$$

$$\sqrt{2\pi V} \frac{1}{\sqrt{2\pi V}} e^{a\theta} e^{b\theta^2} e^{Q} e^{-Q}$$

$$= N(\theta_p, \sigma_p^2)$$

$$\Rightarrow P(\theta, \sigma^2 | x) \propto N(\theta_p, \sigma_p^2) \underbrace{(\sigma^2)^{-\frac{\eta_0 + h}{2} - 1} e^{-\frac{\{\Sigma x_i^2 + \eta_0\sigma_0^2\}}{2\sigma^2}}}_{K(\sigma^2 | x)} \overbrace{e^{-Q}\sqrt{\frac{2\pi}{\frac{h}{\sigma^2} + \frac{1}{\tau^2}}}}^{\text{other stuff}}$$

with "kernel of inv gamma" labeled above.

Resort to grid sampling of $K(\sigma^2|x)$ to approximate $P(\sigma^2|x)$ and then sample $\sigma^2$, sample $N(\theta_p, \sigma_p^2)$

Can we do better? Is there an iterative algorithm?

$p(\theta, \sigma^2 | X)$ non-std distr.

But $p(\theta | X, \sigma^2) = N(\theta_p, \sigma^2_p)$

$\sum(x_i - \theta)^2 = 4\hat{\sigma}^2$

$p(\sigma^2 | X, \theta) \propto p(X | \theta, \sigma^2) \, p(\sigma^2)$

$$\propto (\sigma^2)^{-\frac{h_0 + h}{2} - 1} e^{- \frac{\sum x_i^2 - h_0 \sigma_0^2 + 4\theta^2 - 2\theta_4 \bar{x}}{2\sigma^2}}$$

$$\propto InvGamma\left(\frac{h_0 + h}{2}, \frac{h_0 \sigma_0^2 + h\hat{\sigma}^2}{2}\right)$$

But in $\theta | X, \sigma^2 \Rightarrow \sigma^2$ is unknown    and

in $\sigma^2 | X, \theta \Rightarrow \theta$ is unknown!

Algorithm: similar to E-M:

Step 1: Guess $\sigma_0^2$ ... maybe use $s^2$

Step 2: Sample $\theta | X, \sigma^2 = \sigma_0^2$ to get $\theta_0$

Step 3: Sample $\sigma^2 | X, \theta = \theta_0$ to get $\sigma_1^2$

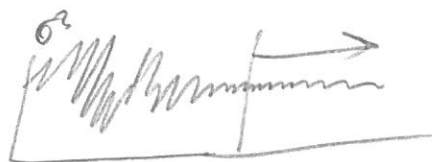Step 4: Sample $\theta | X, \sigma^2 = \sigma_1^2$ to get $\theta_1$

Step 5: Repeat Steps 3 & 4 to dozen...

$$\left\langle \begin{bmatrix} \theta_0 \\ \sigma_0^2 \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \sigma_1^2 \end{bmatrix}, \cdots, \begin{bmatrix} \theta_t \\ \sigma_t^2 \end{bmatrix} \right\rangle$$

where    $\langle \theta_0, \theta_1, \ldots \rangle$ and $\langle \sigma_0^2, \sigma_1^2, \ldots \rangle$ form

where t is large    two "chains"

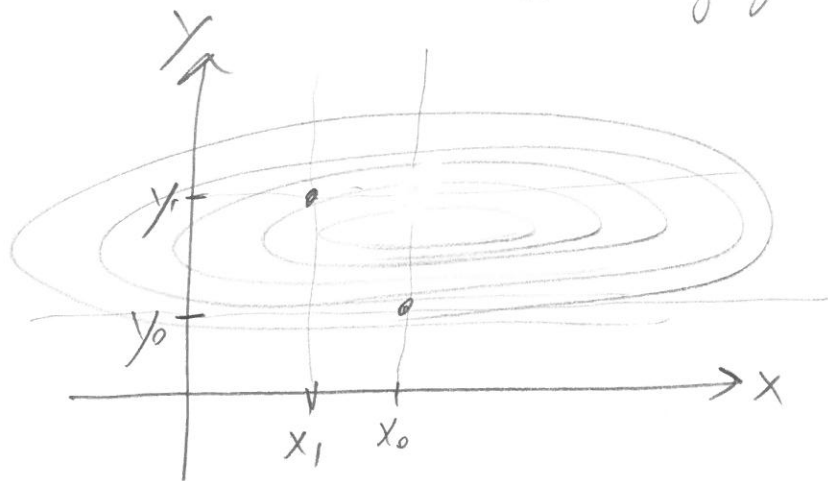Quit upon convergence of the chains max t s.t. all converge

Traceplots:



$R = \max_t \{ \text{both} \atop \text{(converge)} \}$
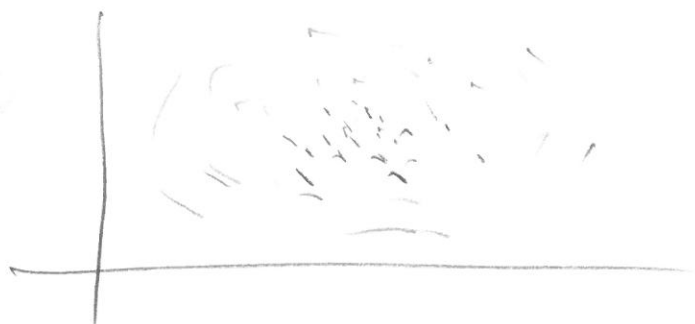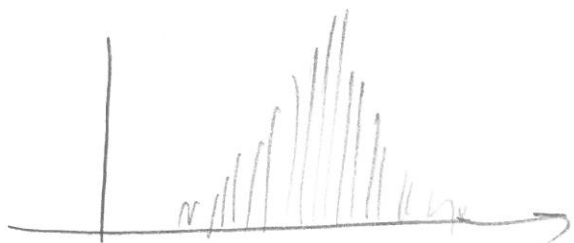
This Algorithm is called "Gibbs Sampling"

$f(x,y)$ is a density you wish to sample from but cannot. But...

You know $f(x|y)$ & $f(y|x)$

So begin w/ $y_0$: You sample $x_0 = f(x|y=y_0)$. Then you sample $y_1 = f(y|x=x_0)$ then $x_1 = f(x|y=y_1)$, etc.
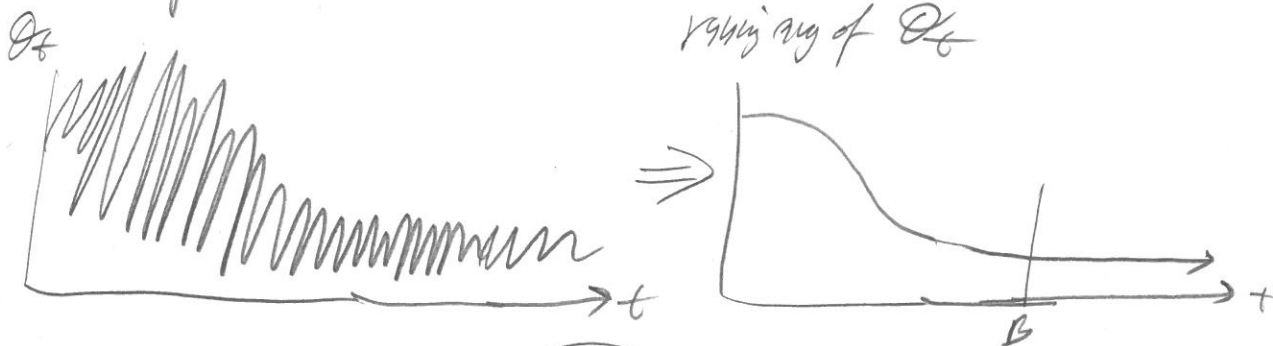
Eventually...



You sample the whole space. If you only care about $x$, then you ditch the $y$'s and have an estimate of $f(x)$:

## Problem #1    At what pt. did it converge?

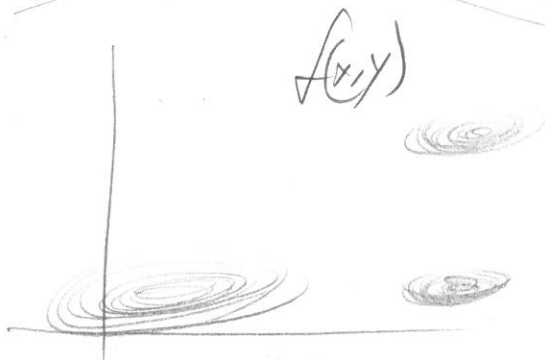Trace plot
of
$\theta_t$

running avg of $\theta_t$

$\Rightarrow$

It appears to converge (9) $t \geq B$.

B

"Good mixing"

But did it really converge?

If you let it run long enough you may see this,
corresponding to bad "mixing". Mixing is the ability of the chain to
effectively traverse the param space (i.e. the support).

$f(x,y)$

## Problem #2    If $f(x,y)$

has multiple modes, the Gibbs chain can get
stuck in a mode and never get out.
Diagnosis... run lots of chains from
multiple starting positions and look at
traceplots. As $dim(\theta)$ increases... this
becomes a bigger problem.

Hard density to Gibbs sample