Lec 19    Math 341   9/20/18

Previously... we needed to sample from $k(\theta|x)$ so we used grid sampling. This involved choices: $\theta_{min}$, $\theta_{max}$, $\Delta\theta$, to make g grid pts. We then

calculate $\{k(\theta_{min}|x), k(\theta_{min}+\Delta\theta|x), \ldots, k(\theta_{min}+(g-1)\Delta\theta|x), k(\theta_{max}|\theta)\}$

then we approximal $\int k(\theta|x) \approx \sum k(\theta_g|x)\Delta\theta = \frac{1}{c}$

and now we have $\{p(\theta_{min}|x), \ldots, p(\theta_{max}|x)\}$

$$\|$$

$$c\,k(\theta_{min}|x)$$

And we can sample from $\theta|x$ using the discrete r.v. sampling technique

# Grid Sampling Disadvantages

① Numerically unstable.

Computers have minimum values of #'s / max. values of #'s.

② How to pick $\theta_{min}$, $\theta_{max}$, $\Delta\theta$ ?

Bad decision for $\theta_{min}$, $\theta_{max}$ ⟹ You miss a part
of the support of the
parameter!

Bad decision for $\Delta\theta$ ⟹ Bad resolution ⟹ non-realistic samples

③ Let's say $\theta_{min} = 0$, $\theta_{max} = 1$, $\Delta\theta = 0.001$, $|G| = 10,000 = 10^5$

What if $\theta$ had 10 dimensions? ⟹ $|G| = 10^{5^{10}} = 10^{50}$ ⟹ IMPOSSIBLE
for a
computer!

⟹ Grid sampling only good in low dimensions
if you know the effective support of $\theta$ (i.e. where
most of the support lies) and if you know the
shape so you can pick a reasonable $\Delta\theta$.

It would be nice to fix these problems with a new method....

Recall

$$X_1, \ldots, X_n \mid \theta, \sigma^2 \overset{iid}{\sim} N(\theta, \sigma^2)$$

$$\theta \sim N(\mu_0, \tau_0^2)$$

$$\sigma^2 \sim \text{InvGamma}\left(\frac{n_0}{2}, \frac{n_0 \sigma_0^2}{2}\right)$$

$$P(\theta, \sigma^2 \mid X) \propto h(\theta, \sigma^2 \mid X) \quad \text{non-conjugate}$$

but

$$P(\theta \mid X, \sigma^2) = N(\theta_p, \sigma_\theta^2)$$

$$P(\sigma^2 \mid X, \theta) = \text{InvG}\left(\frac{n_0 + n}{2}, \frac{n_0 \sigma_0^2 + n \hat{\sigma}^2}{2}\right)$$

Can you use $P(\theta \mid X, \sigma^2)$ & $P(\sigma^2 \mid X, \theta)$ to solve for $P(\theta, \sigma^2 \mid X)$?

$$P(A, B) = P(A \mid B) P(B) = P(B \mid A) P(A)$$
$$P(\theta, \sigma^2 \mid X) = P(\theta \mid \sigma^2, X) P(\sigma^2 \mid X) = P(\sigma^2 \mid \theta, X) P(\theta \mid X)$$

not possible unless either $P(\theta \mid X)$ or $P(\sigma^2 \mid X)$

and those are not possible... so no!

However... what if you use an iterative algorithm?

① begin at $\theta_0$

② Draw $\sigma_0^2$ from $P(\sigma^2 \mid X, \theta = \theta_0)$

③ Draw $\theta_1$ from $P(\theta \mid X, \sigma^2 = \sigma_0^2)$

④ Draw $\sigma_1^2$ from $P(\sigma^2 \mid X, \theta = \theta_1)$

will "converge"

AKA "Gibbs sampling" or the "Gibbs sampler"

This is different than the N-R and E-M alg's. Hey?

Newton resolve

   solve for $f(x) = 0$   one value

E-M

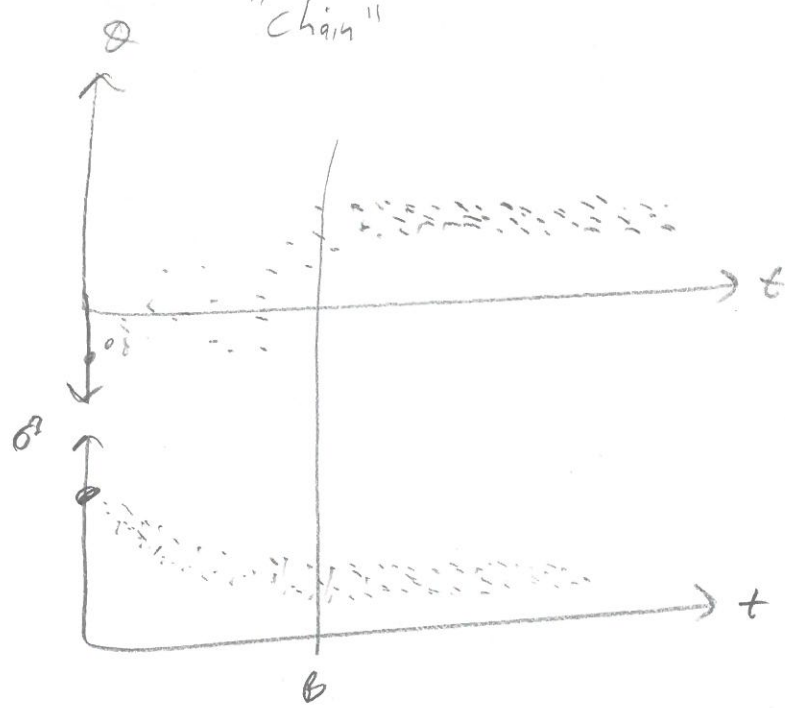   solve for $\hat{\Theta}_{MAP}$   is one value (or vector)

Here:

$$P(\Theta, \sigma^2 | X) \quad \dots \quad \text{entire posterior!!}$$

I dunno look like:

$$\left\langle \begin{pmatrix} \Theta_0 \\ \sigma_0 \end{pmatrix}, \begin{pmatrix} \Theta_1 \\ \sigma_1 \end{pmatrix}, \begin{pmatrix} \Theta_2 \\ \sigma_2 \end{pmatrix}, \dots \begin{pmatrix} \Theta_t \\ \sigma_{t\dots} \end{pmatrix} \right\rangle \quad \text{where } t \text{ is iteration } \#$$
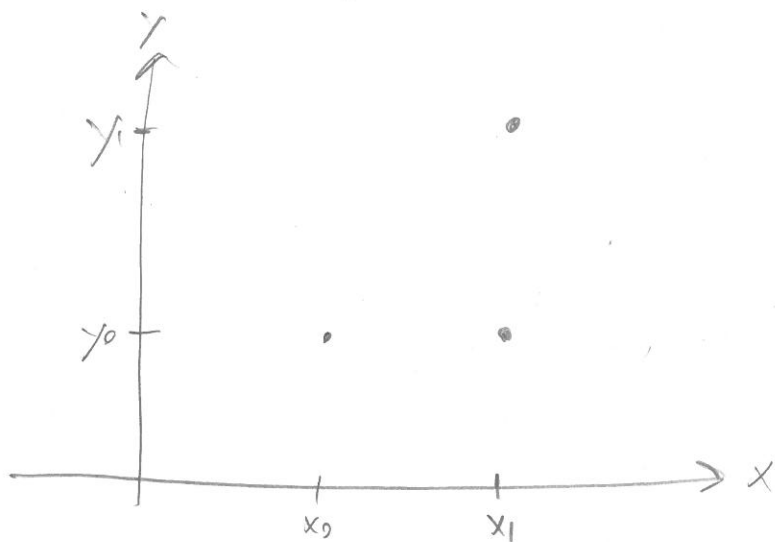
"chain"



let $b = \max_j b_j$
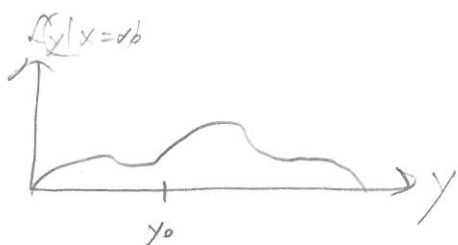
s.t. $b_j$ is the convergent

pt of $\Theta_j$

Where did algorithm converge?

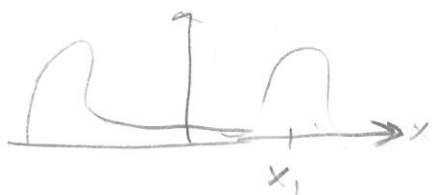we call $t = b$ the burn-in point. Kind of like E-M N-R or E-M

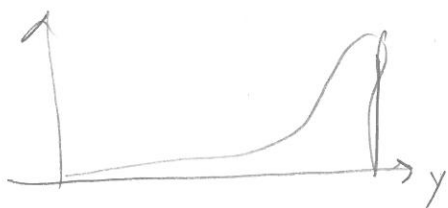Plus, consider way. You seek $f(x,y)$ but only know $f(x|y)$ & $f(y|x)$.
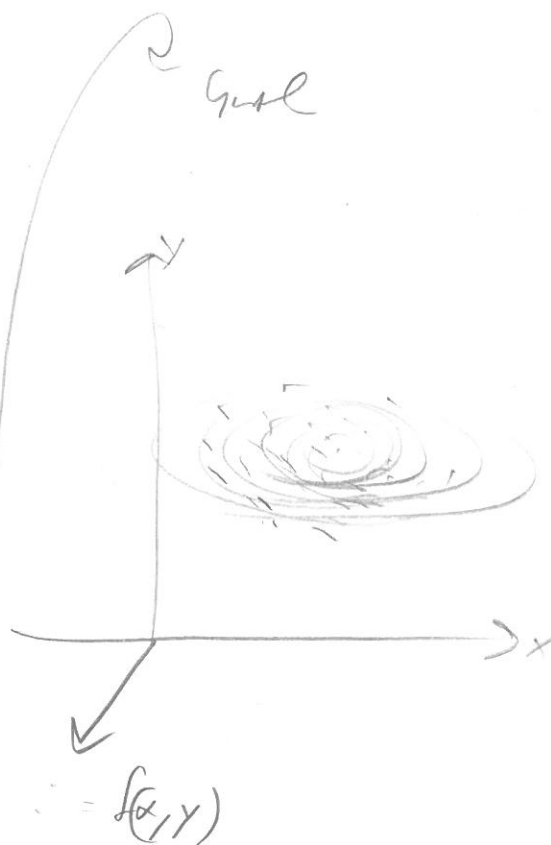


Begin with $x_0$. Draw $y_0$ from $f(y|x=x_0)$



then draw $x_1$ from $f(x|y=y_0)$



then draw $y_1$ from $f(y|x=x_1)$
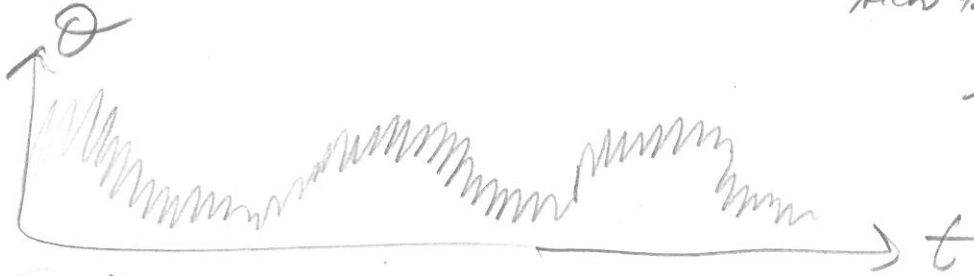


$\vdots$

Goal



$= f(x,y)$

If you only care about $f(x)$, you collapse all $y$'s by just deleting the second dimension



The Main problems with this type of sampling

① Bad mixing



$\vec{\theta}$ may be part of a sea of several minima's with multiple modes

Lacks ability to traverse $supp[\vec{\theta}]$ well



the problem pic area for $\theta_1$

The sampler will get stuck in any of these modes.
Sol: make many chains!
Start from all different starting pts, ...

problematic with big $dim(\theta)$ !

$\Rightarrow$ BTF problems...
because of its serial adequacy

A smaller (but flexible) problem is as follows

Init $\sigma_0^2$

draw $\theta_0$ for $p(\theta | Y, \sigma_0^2)$

draw $\sigma_1^2$ for $p(\sigma^2 | Y, \theta = \theta_0)$

draw $\theta_1$ for $p(\theta | Y, \sigma^2 = \sigma_1^2)$

Is $\theta_1$ related to $\theta_0$? Yes...

Is $\theta_{1000}$ relar to $\theta_{999}$? Yes... After Burn-in (B) still!!

the $\theta_{1000}$ and $\theta_{999}$ are not "independent samples." Since $\text{Corr}[\theta_{1000}, \theta_{999}] \neq 0$

recall $\text{Cou}[X, Y] = \dfrac{\text{Cov}(X,Y)}{SE(X)SE(Y)} = \dfrac{E\left[(X-m_X)(Y-m_Y)\right]}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$

est. by $r := \dfrac{S_{xy}}{S_x S_y} = \dfrac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i-\bar{x})^2}\sqrt{\sum(y_i-\bar{y})^2}}$

we can use an auto-correlation

$auto = self$

autocorrelation for lag 1 between $\text{corr}[\theta_t, \theta_{t+1}]$

$r_{q_1} := \dfrac{\sum\limits_{t=B}^{B+S-1} (\theta_t - \bar{\theta})(\theta_{t+1} - \bar{\theta})}{\sum\limits_{t=B}^{B+S} (\theta_t - \bar{\theta})^2}$  s.t. $\bar{\theta} = \dfrac{1}{S}\sum\limits_{t=B}^{B+S}\theta_t$