

Recall

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

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

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

$$P(\theta, \sigma^2 \mid x) \propto K(\theta, \sigma^2 \mid x) \text{ non-conjugate}$$

but

$$P(\theta \mid x, \sigma^2) = N(\theta_p, \sigma_p^2)$$

$$P(\sigma^2 \mid x, \theta) = \text{InvGamma}\left(\frac{n_0 + n}{2}, \frac{n_0 \sigma_0^2 + n \bar{\sigma}^2}{2}\right)$$

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

$$P(A \cap B) = P(A \mid B) P(B) = P(B \mid A) P(A)$$

$$P(\theta, \sigma^2 \mid x) = P(\theta \mid \sigma^2, x) P(\sigma^2 \mid x) = P(\sigma^2 \mid \theta, x) P(\theta \mid x)$$

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  $P(\sigma^2 \mid x, \theta = \theta_1)$

until "convergence"

AKA "Gibbs sampling" or the "Gibbs sampler".

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

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

E-M

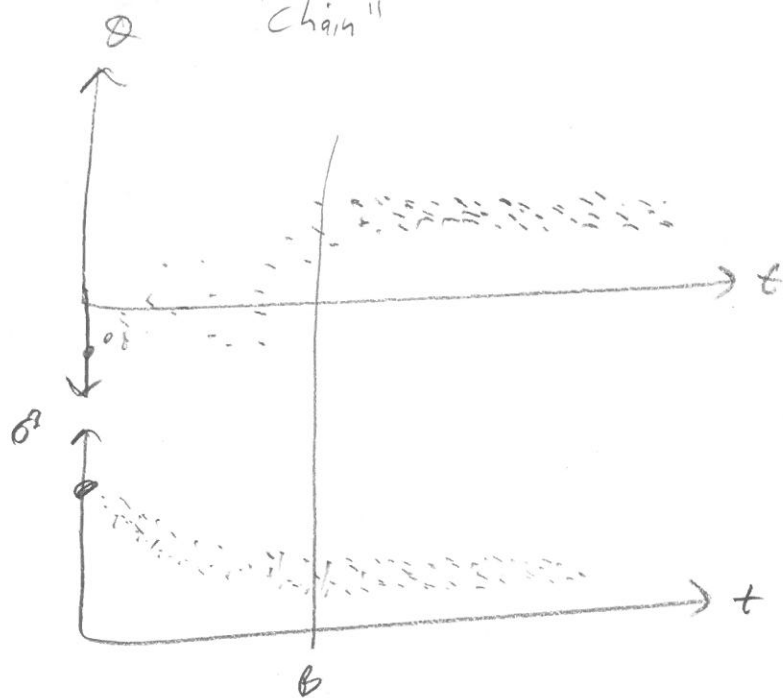
solves for  $\hat{\theta}_{MAP}$  is one value (or vector)

Here:

$p(\theta, \sigma^2 | x)$  ... entire posterior!!

Iterations look like:

$\langle \begin{bmatrix} \theta_0 \\ \sigma_0^2 \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \sigma_1^2 \end{bmatrix}, \begin{bmatrix} \theta_2 \\ \sigma_2^2 \end{bmatrix}, \dots, \begin{bmatrix} \theta_t \\ \sigma_t^2 \end{bmatrix} \rangle$  where  $t$  is iteration #



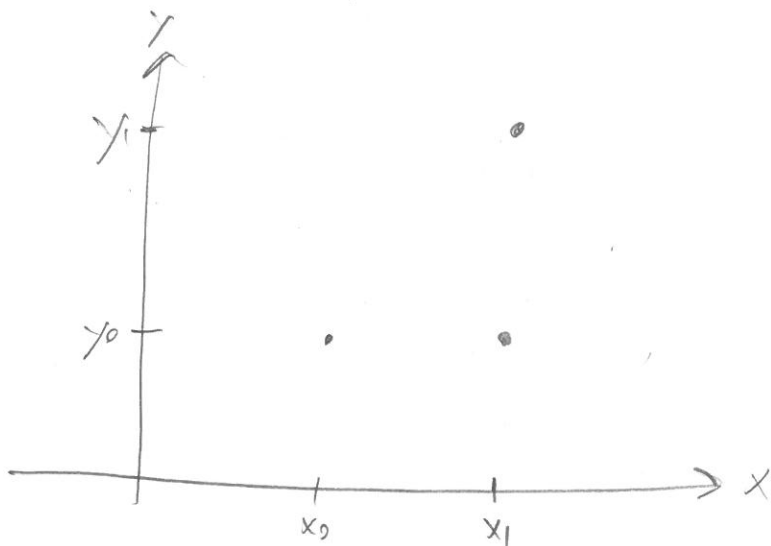
Let  $B = \max_j B_j$

s.t.  $B_j$  is the convergent  
 pt. of  $\theta_j$

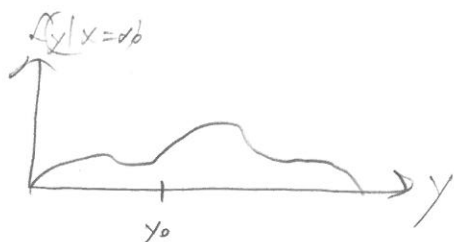
When did algorithm converge?

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

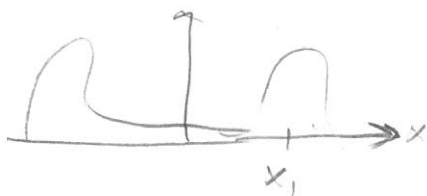
Pro-Quito: wir suchen  $f(x,y)$  nur aus  $f(x|y)$  &  $f(y|x)$ .



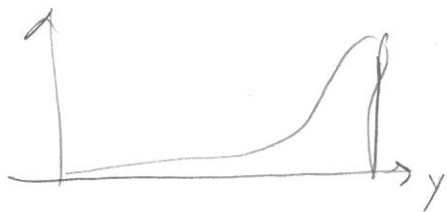
Beginn mit  $x_0$ . Dann  $y_0$  für  $f(y|x=x_0)$



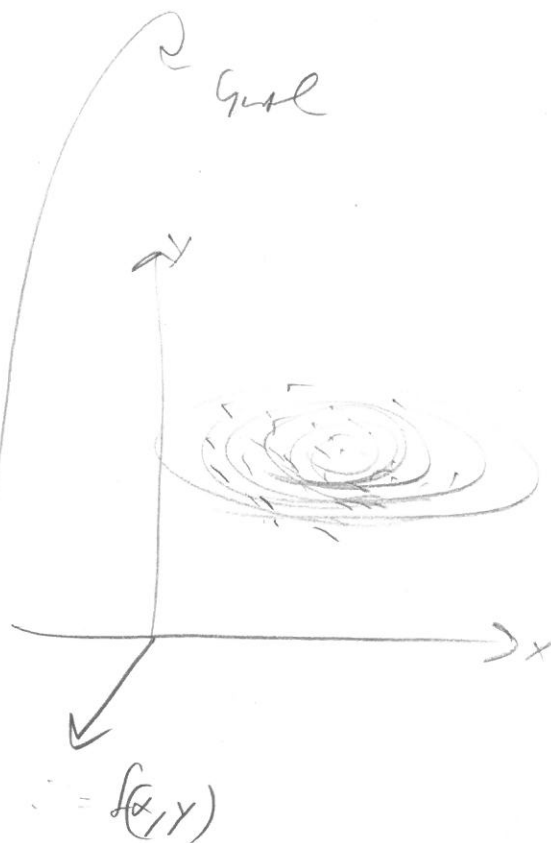
oder dann  $x_1$  für  $f(x|y=y_0)$



oder dann  $y_1$  für  $f(y|x=x_1)$



⋮

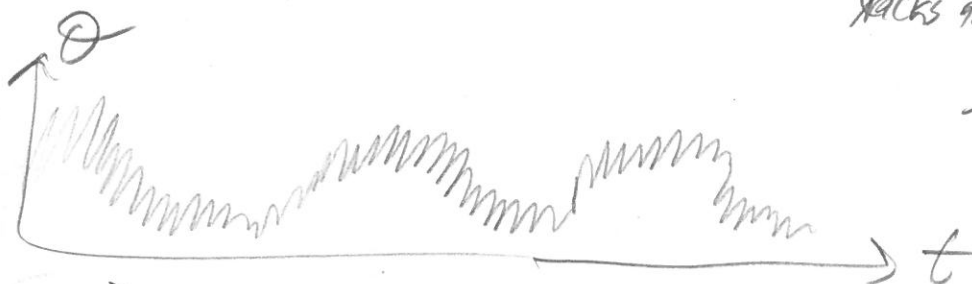


If you only care about  $f(x)$ , you collapse all  $y$ 's by just deleting the second dimension



The main problem with this type of system

① Bad model

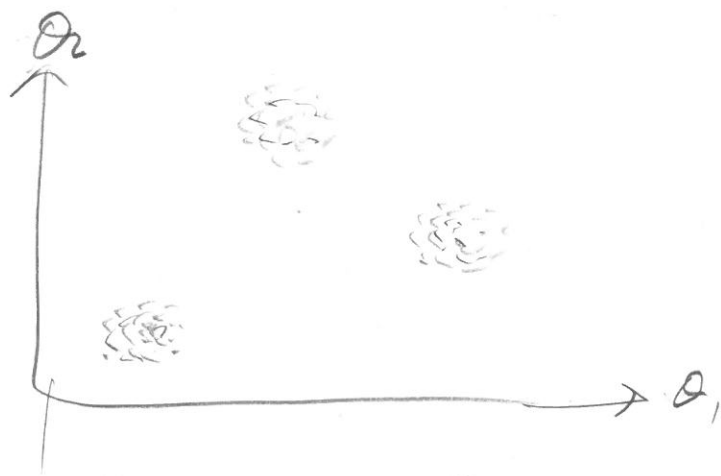


Lacks ability to generalize

$\text{Supp}[\vec{\theta}]$  well

$\vec{\theta}$  may be part of a set of

related distros with multiple modes



The system 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...

Worse if it's really adequately

the problem is not for  $\theta_1$

3

A smaller (but fixable) problem is as follows

- Int  $\sigma_0^2$
- draw  $\theta_0$  for  $P(\theta | x, \sigma_0^2)$
- draw  $\sigma_1^2$  for  $P(\sigma^2 | x, \theta = \theta_0)$
- draw  $\theta_1$  for  $P(\theta | x, \sigma_1^2)$

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

Is  $\theta_{1000}$  related 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  $\text{Corr}[X, Y] = \frac{\text{Cov}(X, Y)}{\text{SE}(X)\text{SE}(Y)} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}}$

est. by  $r := \frac{s_{xy}}{s_x s_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$

he can be this quasi-random  
quasi = self

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

$$r_1 := \frac{\sum_{t=B}^{B+S-1} (\theta_t - \bar{\theta})(\theta_{t+1} - \bar{\theta})}{\sum_{t=B}^{B+S} (\theta_t - \bar{\theta})^2}$$

est.  $\bar{\theta} = \frac{1}{S} \sum_{t=B}^{B+S} \theta_t$

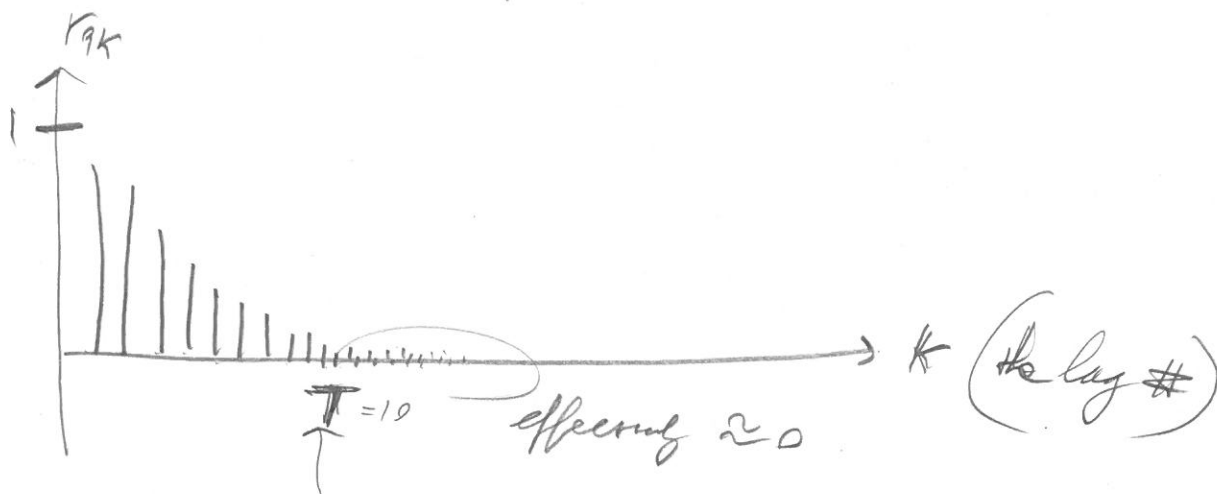
Autocorrelation for lag 2 is:

$$r_{q2} := \frac{\sum_{t=B}^{B+S-2} (\theta_t - \bar{\theta})(\theta_{t+2} - \bar{\theta})}{\sum_{t=B}^{B+S} (\theta_t - \bar{\theta})^2}$$

$$\vdots$$

$$r_{qK} := \frac{\sum_{t=B}^{B+S-K} (\theta_t - \bar{\theta})(\theta_{t+K} - \bar{\theta})}{\sum_{t=B}^{B+S} (\theta_t - \bar{\theta})^2}$$

At some  $K$ ,  $r_{qK} \approx 0$ . Why? Eventually the dependence is lost...  
How to see? Autocorrelation plot



at around  $T$  (brown age), the data effectively are independent

In order to make the chain represent all independent samples from process, we need to throw out all samples except those that are multiple of  $T$  after  $B$ . This is known as "thinning".

$\left\{ \begin{pmatrix} \theta_B \\ \sigma_B^2 \end{pmatrix}, \begin{pmatrix} \theta_{B+T} \\ \sigma_{B+T}^2 \end{pmatrix}, \begin{pmatrix} \theta_{B+2T} \\ \sigma_{B+2T}^2 \end{pmatrix}, \dots \right\} \rightarrow$  the burned and thinned chain.

Let  $l=1 \dots L$  be the index on the block and shared data. This is almost as good as having  $p(\theta|x)$  directly:

How to get  $\hat{\theta}_{MLE} = E(\theta|x) \approx \bar{\theta} = \frac{1}{L} \sum_{l=1}^L \theta_l$

$\hat{\theta}_{MLE} = \text{val}(\theta|x) = \text{order from smallest to largest } \theta_{(1)}, \dots, \theta_{(L)}$

$CR_{\theta, 1-\alpha} = \text{order from smallest to largest } \theta_{(1)}, \dots, \theta_{(L)}$   
 $\left[ \theta_{(\frac{\alpha}{2}L)}, \theta_{((1-\frac{\alpha}{2})L)} \right]$   
 rank ... rank

$p_{ul} = P(\theta_{ul}|x) = P(\theta \in \theta_{ul}|x) \approx \frac{1}{L} \sum_{l=1}^L \mathbb{I}_{\theta_l \in \theta_{ul}}$

e.g.  $P(\theta > 3|x) \approx \frac{1}{L} \sum_{l=1}^L \mathbb{I}_{\theta_l > 3}$  i.e. prop. of  $L$  s.t.  $\theta > 3$

$P(x^*|x) = \int P(x^*|\theta) P(\theta|x) d\theta$

to sample from this...

- ① Pick  $l \in \{1, \dots, L\}$
- ② Draw  $x^*$  from  $P(x^*|\theta = \theta_l)$
- ③ Repeat steps 1,2 over and over...

Def: Symmetric Sweep Gibbs Sampler

Assume parameter  $P(\theta_1, \dots, \theta_p|x)$  unknown but

$P(\theta_j | \theta_{-j}, x)$  s.t.  $\theta_{-j} = \{\theta_1, \dots, \theta_{j-1}, \theta_{j+1}, \dots, \theta_p\}$   
 is known  $\forall j$ .  
 i.e. all  $\theta$ 's except  $\theta_j$