## Posterior Inference

① $E[\Theta|x] \approx \frac{1}{N_T} \sum_{i=1}^{N_T} \Theta_i$

② $\text{Quantile}[\Theta|x, p] = \{p \cdot N_T^{th} \text{ value after sorting}\}$

③ $95\% \; CR = \text{ditto}$

④ Hypothesis tests, etc ... same as before

Gibbs sampling offers inference for previously intractable models!

## Gibbs Sampling algorithm (systematic sweep)    Lecture START

$$P(\Theta_1,\ldots,\Theta_K|x) \propto P(x|\Theta_1,\ldots,\Theta_K)\, P(\Theta_1,\ldots,\Theta_K) \propto L(x|\Theta_1,\ldots,\Theta_K)$$

Assume

$\quad P(\Theta_i|\Theta_{-i}, x)$ can be sampled from $\forall i$

Then Gibbs Sampler is:

$\quad$ Step 1. Initialize $\vec{\Theta}_0 = [\Theta_{1,0}, \ldots, \Theta_{K,0}]$

$\quad$ Step 2: Draw $\Theta_{1,1}$ from $P(\Theta_1|\Theta_{2,0},\ldots,\Theta_{K,0}, x)$

$\qquad\qquad$ Draw $\Theta_{1,2}$ from $P(\Theta_2|\Theta_{1,1},\Theta_{3,0},\ldots\Theta_{K,0}, x)$

$\qquad\qquad\qquad \Theta_{1,K}$ from $P(\Theta_K|\Theta_1,\ldots\Theta_{1,K-1}, x)$

$\quad$ Step 3: Record $\vec{\Theta}_1$ as a sample    Step 4: Repeat 2-3 until "converged"

# Proof of Convergence

Def: Markov Chain on space $X$ is a seq. of r.v's $X_0, X_1, \ldots$ s.t.

$$P(X_{t+1} \in A \mid X_0, X_1, \ldots, X_t) = P(X_{t+1} \in A \mid X_t) \quad \forall A \subseteq X$$

once you're at a certain state, you lose memory of all previous states.

Further $\ldots = \int_A P(X_{t+1} \mid X_t) f(X_t) dx$ where $P(X_{t+1} \mid X_t)$ is called the "transition kernel"

If $X$ is ctbl, it's called the "transition matrix".

$f$ is the chain's

Def: invariant distribution if:

$$f(X_{t+1}) = \int_X f(X_t) P(X_{t+1} \mid X_t) dx$$

Thm: $\forall \underset{\substack{(\text{i.e. strong} \\ \text{positive})}}{g(X_0)}, \underset{\substack{(\text{invariant} \\ \text{distr})}}{f(a)} = \lim_{T \to \infty} \int_X \prod_{t=0}^{T} P(X_{t+1} \mid X_t) g(X_0) dx$ w/ reg. cond's.

$X \leftarrow$ avg. over everywhere $X_t$ could be!

$\Rightarrow$ No matter where the Markov chain begins, it converges to the same distr.

jdf

Def: $f(X_1, \ldots, X_k)$ satisfies the positivity condition if $f(X_i) > 0 \; \forall \; i = 1 \ldots k$

$\Rightarrow f(X_1, \ldots, X_k) > 0 \iff Supp[X_1, \ldots, X_k] = Supp[X_1] \times \ldots \times Supp[X_k]$

Thm: If $\forall x, f(X_1, \ldots, X_k)$ satisfies the pos. cond., then $\forall (q_1, \ldots, q_k) \in Supp[X_1, \ldots, X_k]$

$$\Rightarrow f(X_1, \ldots, X_k) \propto \prod_{i=1}^{K} \frac{f(X_j \mid X_1, \ldots, X_{j-1}, X_{j+1} = q_{j+1}, \ldots, X_k = q_k)}{f(X_j = q_j \mid \; '' \quad '' \quad '' \quad '' \quad '')} > 0$$

If jdf has pos,

Lem: All conditional densities are non-zero. If not, then $\perp$.

$\Rightarrow$ Before running Gibbs Sampler, ask... is the mass everywhere in the support???

**Thm:** the kernel of a Gibbs sampler can be expressed as:

Recall the Gibbs sampler in $K$ dimensions... the transition kernel would be:

$$P\left(\Theta_{t+1,1}, \ldots \Theta_{t+1,K} \mid \Theta_{t,1}, \ldots \Theta_{t,K}, X\right)$$

<u>future iteration</u>　　<u>previous iteration</u>　　data

$$= P\left(\Theta_{t+1,K} \mid \Theta_{t+1,1}, \ldots, \Theta_{t+1,K-1}, X\right) \cdot$$

$$P\left(\Theta_{t+1,K-1} \mid \Theta_{t+1,1}, \ldots \Theta_{t+1,K-2}, \Theta_{t,K}, X\right) \cdot$$

$$P\left(\Theta_{t+1,K-2} \mid \Theta_{t+1,1}, \ldots, \Theta_{t+1,K-3}, \Theta_{t,K-1}, \Theta_{t,K}, X\right) \cdot \; \cdots \;$$

$$P\left(\Theta_{t+1,2} \mid \Theta_{t+1,1}, \Theta_{t,2}, \ldots, \Theta_{t,K}, X\right) \cdot$$

$$P\left(\Theta_{t+1,1} \mid \Theta_{t,2}, \ldots \Theta_{t,K}, X\right)$$

$K$ steps

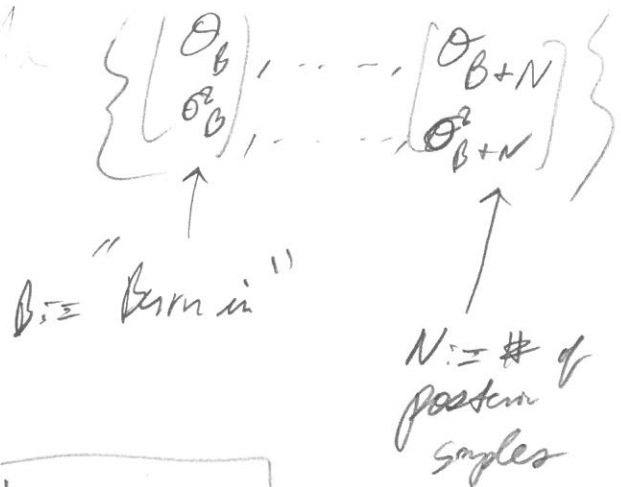| $\Theta_{1,0} \; \Theta_{2,0} \cdots \qquad \Theta_{K,0}$ | $t=0$ |
| $\Theta_{1} \; \Theta_{2,0} \cdots \qquad \Theta_{K,0}$ | $t=1.1$ |

$\Rightarrow$ a Gibbs sampler is a Markov chain

---

**Thm:** The Gibbs sampler converges to the posterior if $f\left(\Theta_1, \ldots \Theta_K \mid X\right)$ is the "invariant distn." **Proof:**

$$\overset{?}{f}\left(\Theta_1 = \Theta_{t+1,1}, \ldots \Theta_K = \Theta_{t+1,K} \mid X\right) \overset{?}{=} \int_{\Theta_1} \cdots \int_{\Theta_K} \int_{\Theta_{t,K}} f\left(\Theta_1 = \Theta_{t,1}, \ldots \mid X\right) (\text{kernel}) \, d\Theta_{t,1} \cdots d\Theta_{t,K}$$

$$\int_{\Theta_2} f\left(\Theta_1 = \Theta_{t+1,1}, \Theta_2 = \Theta_{t,2}, \ldots \Theta_K = \Theta_{t,K} \mid X\right) d\Theta_{t,2}$$

$\|$

$$P(A)\,P(B \mid A) \cdot P(B,A)$$

Recall since $f\left(\Theta_1, \ldots \Theta_K \mid X\right) > 0 \Rightarrow$ all cond's are $> 0$

$$= \int_{\Theta_2} \cdots \int_{\Theta_K} \int_{\Theta_1} f\left(\Theta_1 = \Theta_{t,1}, \ldots \Theta_K = \Theta_{t,K} \mid X\right) d\Theta_{t,1} \; \text{kernel} \; d\Theta_{t,2} \cdots d\Theta_{t,K} = \int_{\Theta_3} \cdots \int_{\Theta_K} \int_{\Theta_2} f\left(\Theta_2 = \Theta_{t,2}, \ldots \Theta_K = \Theta_{t,K} \mid X\right) P\left(\Theta_{t+1,1} \mid \Theta_{t,2}, \ldots \Theta_{t,K}, X\right) d\Theta_{t,2} \cdots d\Theta_{t,K}$$

last of kernel

$$= \int_{\Theta_4} \cdots \int_{\Theta_K} \int_{\Theta_3} f\left(\Theta_1 = \Theta_{t+1}, \Theta_3 = \Theta_{t,3}, \ldots, \Theta_K = \Theta_{t,K} \mid X\right) P\left(\Theta_{t+1,2} \mid \Theta_{t+1,1}, \Theta_{t,2}, \ldots \Theta_{t,K}, X\right) d\Theta_{t,3} \cdots d\Theta_{t,4}$$

last of kernel

$$\int_{\Theta_3} \int f\left(\Theta_1 = \Theta_{t+1}, \Theta_2 = \Theta_{t+1,2}, \Theta_3 = \Theta_{3,t}, \Theta_4 = \Theta_{4,t}, \ldots \Theta_K = \Theta_{t,K} \mid X\right) d\Theta_3$$

$$= \int_{\Theta_5} \cdots \int_{\Theta_K} \int_{\Theta_4} f\left(\Theta_1 = \Theta_{t+1}, \Theta_2 = \Theta_{t+1,2}, \Theta_4 = \Theta_{4,t}, \ldots \Theta_K = \Theta_{t,K} \mid X\right) \cdots d\Theta_{t,4} \cdots d\Theta_{t,K}$$

$\vdots$

$$= f\left(\Theta_1 = \Theta_{t+1,1}, \ldots \Theta_K = \Theta_{t+1,K} \mid X\right) \; ▨$$

So it converges to invariant distribution which is the jdf we seek.
Let $B$ be the iteration of convergence, thus...

$$\left\{ \begin{bmatrix} \theta_B \\ \sigma^2_B \end{bmatrix}, \cdots, \begin{bmatrix} \theta_{B+N} \\ \sigma^2_{B+N} \end{bmatrix} \right\}$$

$B := $ "Burn in"

$N := $ # of posterior samples

Constitute samples from $P(\theta, \sigma^2 | X)$, the density which previously defied simulation

---

**Gibbs problem #3** $\theta_{t+1}$ is dependent on $\theta_t$ which is dep on $\theta_{t-1}$, etc.

At what point are they independent?

Recall the $\text{Corr}[X,Y] := \dfrac{\text{Cov}(Y,Y)}{SE(X) \, SE(Y)} := \dfrac{E(X-\mu_x)(Y-\mu_Y)}{\sqrt{\text{Var}(X) \, \text{Var}(Y)}}$

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

Here we care about "autocorrelation"

auto-self (greek prefix)

automatic
↑ ↑
self move

or $\theta$, Autocorrelation for lag 1 is:

$$r_{a1} := \frac{\sum_{t=B}^{B+N-1} (\theta_t - \bar{\theta})(\theta_{t+1} - \bar{\theta})}{\sum_{t=B}^{B+N} (\theta_t - \bar{\theta})^2}$$
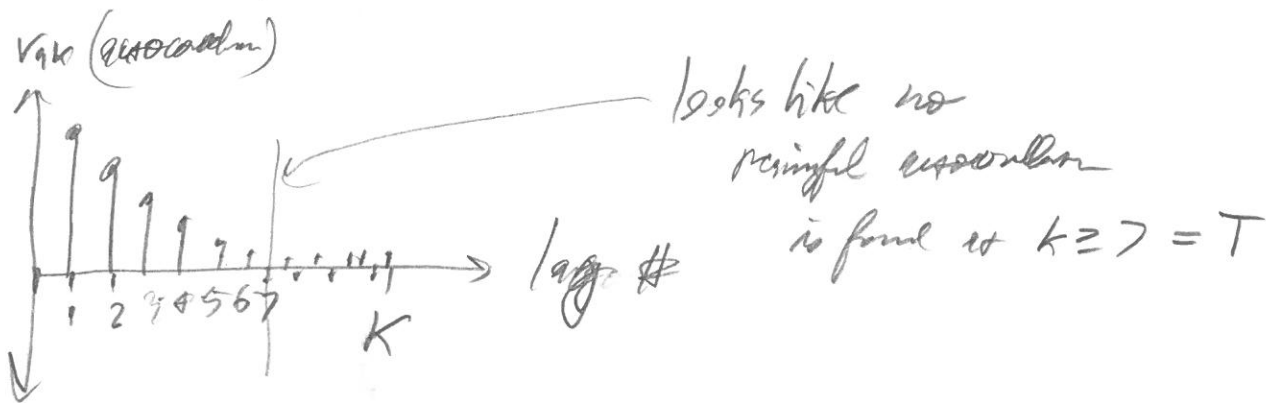
where $\bar{\theta} := \dfrac{1}{N} \sum_{t=B}^{B+N} \theta_t$

And autocorrelation for lag 2 is:

$$r_{a2} := \sum_{t=B}^{B+N-2} (\theta_t - \bar{\theta})(\theta_{t+2} - \bar{\theta}) \Big/ \sum_{t=B}^{B+N} (\theta_t - \bar{\theta})^2$$

And for the $k^{th}$ lag:

$$V_{ak} := \frac{\sum\limits_{t=b}^{b+N-k} \left(\theta_t - \overline{\theta}\right)\left(\theta_{t+k} - \overline{\theta}\right)}{\sum\limits_{t=b}^{b+N} \left(\theta_t - \overline{\theta}\right)^2}$$

Pick a max $k$, $K$ and look at



$V_{ak}$ (autocorrelation)

lag #

looks like no meaningful association is found at $k \geq 7 = T$

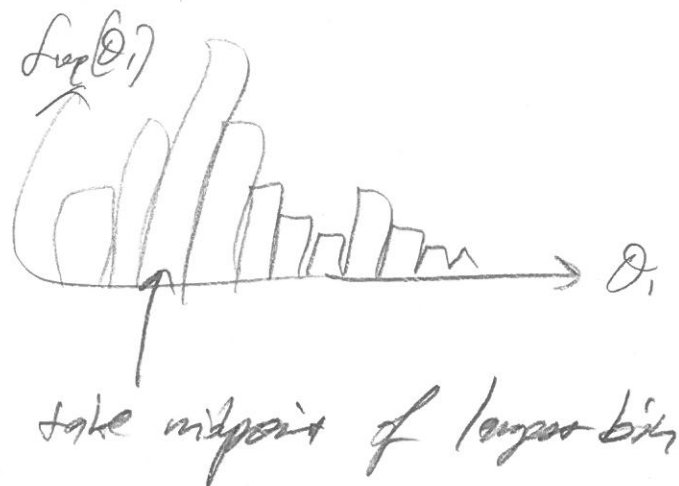$\Longrightarrow$ "Thin" the chains by throwing out all non-multiples of $7$:

$$\left\{ \begin{bmatrix} \theta_b \\ \sigma_b^2 \end{bmatrix}, \dots, \begin{bmatrix} \theta_{b+7} \\ \sigma_{b+7}^2 \end{bmatrix}, \dots \begin{bmatrix} \theta_{b+14} \\ \sigma_{b+14}^2 \end{bmatrix} \dots \right\}$$

$$\left\{ \begin{bmatrix} \theta_b \\ \sigma_b^2 \end{bmatrix}, \begin{bmatrix} \theta_{b+7} \\ \sigma_{b+7}^2 \end{bmatrix}, \begin{bmatrix} \theta_{b+14} \\ \sigma_{b+14}^2 \end{bmatrix}, \dots \right\}$$

Our "burned and thinned" chains

Which can be used as draws from joint or marginal densities.

$\hat{\theta}_{1,MAP}$ ?    Bin:



take midpoint of longest bin

$$P(X^* | X) ?$$

$$= \int P(x^* | \theta_1, \dots \theta_k) \, P(\theta_1, \dots, \theta_k | X) \, d\theta_1 \dots d\theta_k$$

## Procedure

① Run gibbs sampler, until convergence and get "burned-in" chain

② Sample $\vec{\theta}$ from chain $\vec{\theta}_1$

③ Sample $X^*$ from Lik. model, $X^*_1$ (assume often this is possible)

④ Repeat the above $\overset{\text{steps } 2-3}{N}$ times, $X^*_1, \dots, X^*_n$

what if you could draw...

$$P(x^* \geq c | x) \approx \frac{1}{N} \sum \mathbb{1}_{x^* \geq c}$$

First example was

$$X | \theta, \sigma^2 \sim N(\theta, \sigma^2)$$

$$\theta \sim N(\mu_0, \tau^2), \quad \sigma^2 \sim \text{InvGamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$$