

Algorithm: Systematic Sweep

Gibbs Sampler for $P(\theta_1, \dots, \theta_p | X)$, the unknown posterior w/ p parameters.

Here all conditionals, $P(\theta_i | \theta_{-i})$ where $\theta_{-i} = \{\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_p\}$ are known and can be "easily" sampled from.

step 1: Initialize $\theta = \langle \theta_{0,1}, \theta_{0,2}, \dots, \theta_{0,p} \rangle$ (conditional on data X)

step 2: Sample $\theta_{1,1}$ from $P(\theta_{1,1} | \theta_2 = \theta_{0,2}, \dots, \theta_p = \theta_{0,p})$

Sample $\theta_{1,2}$ from $P(\theta_{1,2} | \theta_1 = \theta_{1,1}, \theta_3 = \theta_{0,3}, \dots, \theta_p = \theta_{0,p})$

\vdots

Sample $\theta_{1,p}$ from $P(\theta_{1,p} | \theta_1 = \theta_{1,1}, \dots, \theta_{p-1} = \theta_{1,p-1})$

step 3: Repeat 2 until "convergence".

(Not Required) Proof: Consider X_0, X_1, X_2, \dots a sample of random variables. Each have Sample X . If $P(X_t \in A | X_{t-1}, X_{t-2}, \dots, X_0) = P(X_t \in A | X_{t-1}) \forall t, \forall A \subset \mathcal{X}$ then the sample sequence is called a ~~discrete-time~~ Markov Chain. The Gibbs sampler is a Markov Chain.

This is why the Gibbs sampler is ~~the~~ a form of "Markov Chain Monte Carlo" ^{transition kernel} _{mcmc}.

$$P(X_{t+1}) = \int P(X_{t+1}, X_t) dx = \int P(X_{t+1} | X_t) P(X_t) dx$$

If $P(X_{t+1}) \stackrel{x}{=} P(X_t)$ then this distribution is deemed the "invariant" "equilibrium" "stationary" "long-term".

$$P(X_{t+1}) = P(X_t | X_{t+1}) P(X_{t+1} | X_{t-2}) \dots P(X_1 | X_0) P(X_0)$$

$$P(X) = \lim_{t \rightarrow \infty} \int \overset{\text{integrating whole thing}}{dx_0} \quad \text{you'll get invariant distribution.}$$

$$= P(\theta_{t+1,1} | \theta_{t+2}, \dots, \theta_{t,p}) \cdot P(\theta_{t+2} | \theta_{t+1,1}, \theta_{t+3}, \dots, \theta_{t,p}) \cdot \dots \cdot P(\theta_{t+1,p-1} | \theta_{t+1,1}, \dots, \theta_{t,p-1}) \cdot P(\theta_{t+1,p} | \theta_{t+1,1}, \dots, \theta_{t+1,p-1}, \theta_{t+1,p-1})$$

(vector notation) $P(\vec{\theta}_{t+1}) = \int P(\vec{\theta}_{t+1} | \vec{\theta}_t) \cdot P(\vec{\theta}_t) d\vec{\theta}_t$

(scalar notation) $P(\theta_{t+1,1}, \dots, \theta_{t+1,p}) = \int \int \dots \int_{\Theta_1} \int_{\Theta_2} \dots \int_{\Theta_p} P(\theta_{t+1,1}, \dots, \theta_{t+1,p}) \dots d\theta_1 \dots d\theta_p$

$$= \int \dots \int_{\Theta_2} \int_{\Theta_p} \text{kernel} \int_{\Theta_1} P(\theta_{t+1}, \theta_{t+2}, \dots, \theta_p | \theta_1) d\theta_1$$

"1"

$$P(\theta_t, \dots, \theta_p) d\theta_1 \dots d\theta_p$$

$$= \int_{\Theta_3} \int_{\Theta_p} \text{rest of kernel} \left[\int_{\Theta_2} P(\theta_{t+1,1} | \theta_{t,2}, \dots, \theta_{t,p}) P(\theta_{t,2}, \dots, \theta_{t,p}) d\theta_t \right] d\theta_3 \dots d\theta_p$$

$$\underbrace{P(\theta_{t+1,1}, \theta_{t,2}, \theta_{t,p})}_{P(\theta_{t,2}, \theta_{t,3}, \dots, \theta_{t,p})}$$

$$= \int_{\Theta_4} \dots \int_{\Theta_1} \text{rest of kernel} \left[\int_{\Theta_3} P(\theta_{t+1,2} | \theta_{t,3,1}, \theta_{t,3}, \dots, \theta_{t,p}) P(\theta_{t,3,1}, \theta_{t,3}, \dots, \theta_{t,p}) d\theta_4 \dots d\theta_p \right]$$

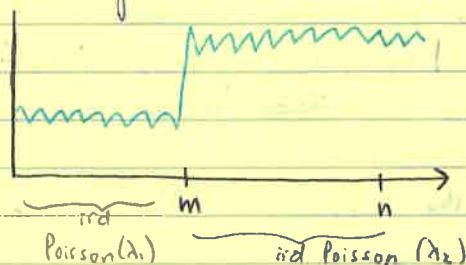
$$\underbrace{P(\theta_{t+1,1}, \theta_{t+1,2}, \theta_{t,3}, \dots, \theta_{t,p})}_{P(\theta_{t+1,1}, \theta_{t+1,2}, \theta_{t+1,3}, \dots, \theta_{t,p})}$$

and the process repeats...

proof of eventual convergence

$$= P(\theta_{t+1,1}, \dots, \theta_{t+1,p})$$

Change Point Model



Parameters: λ_1 - mean of "first process"
 λ_2 - mean of "second process"
 m - "change point."

Priors: $P(\lambda_1) = \text{Gamma}(\alpha, \beta)$

$P(\lambda_2) = \text{Gamma}(\alpha, \beta)$

$P(m) = \text{Uniform} \{0, \dots, n\} = \frac{1}{n} \forall m$

Posterior: $P(\lambda_1, \lambda_2, m | X_1, \dots, X_n) \propto P(X_1, \dots, X_n | \lambda_1, \lambda_2, m) \cdot P(\lambda_1, \lambda_2, m)$

$$\propto \left(\prod_{i=1}^m \frac{e^{-\lambda_1} \lambda_1^{x_i}}{x_i!} \right) \left(\prod_{i=m+1}^n \frac{e^{-\lambda_2} \lambda_2^{x_i}}{x_i!} \right) (\lambda_1^{\alpha-1} e^{-\beta \lambda_1}) (\lambda_2^{\alpha-1} e^{-\beta \lambda_2})$$

$$\propto e^{-m \lambda_1} \lambda_1^{\sum_{i=1}^m x_i} e^{-m \lambda_2} \lambda_2^{\sum_{i=m+1}^n x_i} \lambda_1^{\alpha-1} e^{-\beta \lambda_1} \lambda_2^{\alpha-1} e^{-\beta \lambda_2}$$

$$= e^{-\frac{(m+\beta) \lambda_1}{\lambda_1^{\alpha+1}}} e^{-\frac{(n-m+\beta) \lambda_2}{\lambda_2^{b+\alpha-1}}}$$

"Unknown Distribution" that's the best we can do...

Need Conditionals:

$$P(\lambda_1 | X_1, \dots, X_n, \lambda_2, m) \propto e^{-(m+\beta)\lambda_1} \lambda_1^{a+\alpha-1} \propto \text{Gamma}(a+\alpha, m+\beta)$$

$$P(\lambda_2 | X_1, \dots, X_n, \lambda_1, m) \propto e^{-(n-m+\beta)\lambda_2} \lambda_2^{b+\alpha-1} \propto \text{Gamma}(b+\alpha, n-m+\beta)$$

$$P(m | X_1, \dots, X_n, \lambda_1, \lambda_2) \propto e^{-(m+\beta)\lambda_1} e^{-(n-m+\beta)\lambda_2} \propto e^{-m\lambda_1} e^{\lambda_2 m}$$

$$\propto e^{-m(\lambda_1 - \lambda_2)} \lambda_1^a \lambda_2^b \leftarrow \text{dependent on } m.$$

$$\propto P \underbrace{\quad}_{\text{call } h(m)}$$

$$\propto \frac{h(m)}{\sum_{k=0}^{\infty} h(k)}$$

then. Pick λ_1 and a starting point. Plug in to get next round.
and repeat. ✓

$$\left\langle \begin{bmatrix} \lambda_{0,1} \\ \lambda_{0,2} \\ m_0 \end{bmatrix}, \begin{bmatrix} \lambda_{1,1} \\ \lambda_{1,2} \\ m_1 \end{bmatrix}, \dots \right\rangle$$

Drawing vertical line through the 3
constitutes 1 data point.

All have same burn-in.

converges quickly.

Discard the data points before burn-in.

Problem: When do you discard data points.
 $\approx 4\%5$.

(thing)

When does it dip below significance level.
Then Posterior \rightarrow get (.Regions).

Color on Github.

\rightarrow Dark line = mean. ($\hat{\theta}_{mmse}$)

the real mean $\neq \hat{\theta}_{mmse}$ but still
in the Credible Region.

Machine Learning Club ✓