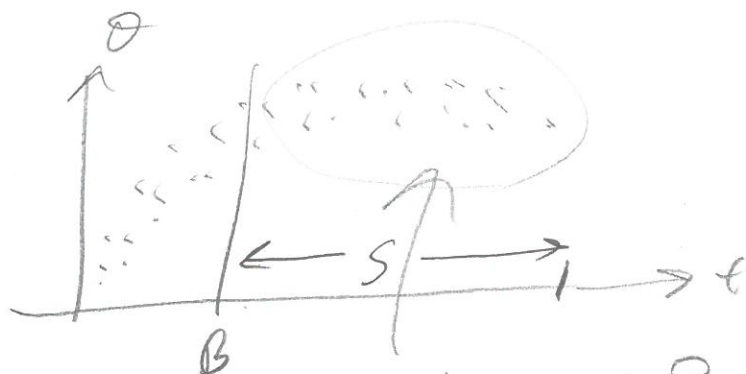


Lee 20 Math 341 5/2/18

Gibbs Sampling



are these iid? No!

They are dependent on the previous realizations.

How do we measure dependence?

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \in [-1, 1] \quad \text{if } r \approx 0 \Rightarrow X, Y \text{ mostly independent}$$

Use this to see how one iteration depends on the previous...

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

$$\text{where } \bar{\theta} := \frac{1}{S} \sum_{t=B+1}^{B+S} \theta_t$$

this is called "auto-correlation" with lag = 1
self!

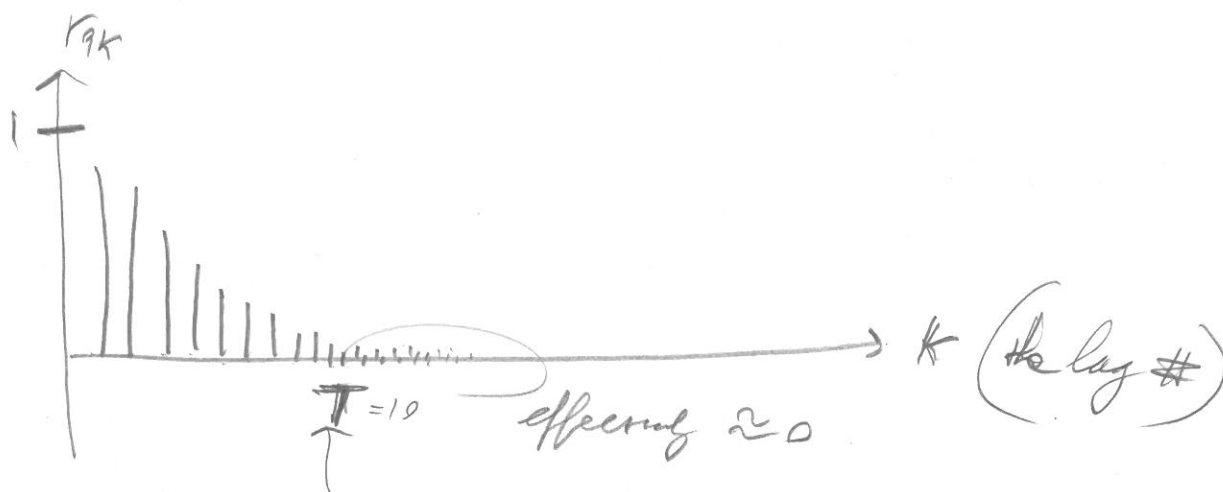
It measures how θ_t varies with θ_{t-1} - What do you think $r_{q,1}$ will be? Large & Positive!

Autocorrelation for lag 2 is:

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

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

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



at around T intervals ago, the data effectively are independent

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

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

Let $l=1 \dots L$ be the index on the brand not desired drug. This is almost as good as having $p(\theta|x)$ directly:

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

$\hat{\theta}_{max} = \text{re}(\theta|x) = \text{order for smallest to largest } \theta_{(1)}, \dots, \theta_{(L)}$

$CR_{\theta, 1-\alpha} = \text{order for smallest to largest}$

$\left[\theta_{(\frac{\alpha}{2}L)}, \theta_{(1-\frac{\alpha}{2}L)} \right]$
rank ... rank

old term $\theta_{(\frac{L}{2})}$
(rank)

$p_{nd} = P(H_d|x) = P(\theta \in \Theta_d|x) \approx \frac{1}{L} \sum_{l=1}^L \mathbb{1}_{\theta_l \in \Theta_d}$

e.g. $P(\theta > 3|x) \approx \frac{1}{L} \sum_{l=1}^L \mathbb{1}_{\theta_l > 3}$ 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 $l \in \{1, \dots, L\}$ to get θ_l

② Draw x^* from $P(x^*|\theta = \theta_l)$

③ Repeat steps 1,2 over and over...

Def: Symmetric Sleep Gibbs Sampler

Assume posterior $P(\theta_1, \dots, \theta_p|x)$ 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 θ 's except θ_j

Step 1: Initialize $\vec{\theta}_0 = \langle \theta_{1,0}, \theta_{2,0}, \dots, \theta_{p,0} \rangle$

Step 2: Sample $\theta_{1,1}$ from $P(\theta_1 | \theta_2 = \theta_{2,0}, \dots, \theta_p = \theta_{p,0}, x)$
 Sample $\theta_{2,1}$ from $P(\theta_2 | \theta_1 = \theta_{1,1}, \theta_3 = \theta_{3,0}, \dots, \theta_p = \theta_{p,0}, x)$
 \vdots
 Sample $\theta_{p,1}$ from $P(\theta_p | \theta_1 = \theta_{1,1}, \dots, \theta_{p-1} = \theta_{p-1,1}, x)$

Step 3: record $(\theta_{1,1}, \dots, \theta_{p,1})$ as a "sample"
 repeat step 2 for many times...

The claim is that the samples, given $t \gg 0$ come from $P(\theta_1, \dots, \theta_p | x)$. Proof...

Def. Consider X_0, X_1, X_2, \dots a sequence of r.v.'s (scalar or vector) with support \mathcal{X} .

If $P(X_t^{EA} | X_{t-1}, X_{t-2}, X_{t-3}, \dots, X_{t-s}) = P(X_t^{EA} | X_{t-1}) \quad \forall t, s$

then X_1, X_2, \dots is called a "Markov Chain".

Is the Gibbs sampler a Markov Chain?? YES!!! "MCMC"

trans. kernel if distr. / terms kind if cont. $\forall A \subset \mathcal{X}$

Def. A Markov chain's "invariant distribution" is defined as:

$$P(X_{t+1}) = \int P(X_{t+1}, x_t) dx_t = \int P(X_{t+1} | x_t) P(x_t) dx_t$$

\mathcal{X}
margin out the effect of the previous value

$$P(x_t) = P(x_t)$$

known

$$= \int P(x_{t+1} | x_t) P(x_t | x_{t-1}) P(x_{t-1}) dx_t$$

$$= \int \left(\prod_{i=0}^{t-1} P(x_{i+1} | x_i) \right) P(x_0) dx$$

cool equation huh?

Thm. for any starting distr. $P(X_0)$,

$$P(X_t) = \lim_{t \rightarrow \infty} \int \prod_{i=0}^{t-1} P(X_i | X_{i-1}) P(X_0 = x) dx$$

\Rightarrow doesn't matter where you start, given enough "hops" or time, you wind up in the same steady state distr. AKA "long-term" / "limiting" / "stationary" distr.

Def: A jdf $P(X_1, \dots, X_p)$ has the positivity cond.

$$\forall j, P(X_j) > 0 \quad \forall x_j \in \text{supp}(X_j)$$

Thm: Consider $P(X_1, \dots, X_p)$ which has the positivity cond. Then, $\forall \vec{a} \in \text{supp}(\vec{X})$,

$$P(X_1, \dots, X_p) \propto \prod_{j=1}^p \frac{P(X_j | X_1, \dots, X_{j-1}, X_{j+1} = a_{j+1}, \dots, X_p = a_p)}{P(X_j = a_j | X_1, \dots, X_{j-1}, X_{j+1} = a_{j+1}, \dots, X_p = a_p)}$$

Corr: If $P(X_1, \dots, X_p)$ has the positivity cond.

$$\Rightarrow P(X_j | X_{-j}) > 0 \quad \forall x_j \in \mathcal{X} \quad \text{i.e. all cond. densities are nonzero.}$$

We need this for the proof...

What is the transition kernel for the Gibbs sampler?

$$P(\vec{\theta}_{t+1} | \vec{\theta}_t, X) = P(\theta_{t+1,1}, \dots, \theta_{t+1,p} | \theta_{t,1}, \dots, \theta_{t,p}, X)$$

I will not write X anymore...
I'll just
pretend
it's
there