Want $P(\theta, \sigma^2 | x)$ and we have

$$P(\theta | x, \sigma^2) = N\left(\frac{\frac{\bar{x}\eta}{\sigma^2} + \frac{\mu_0}{t^2}}{\frac{\eta}{\sigma^2} + \frac{1}{t^2}}, \frac{1}{\frac{\eta}{\sigma^2} + \frac{1}{t^2}}\right)$$ and

$$P(\sigma^2 | x, \theta) = InvGamm\left(\frac{\eta_0 + \eta}{2}, \frac{\eta_0 \sigma_0^2 + \eta \hat{\sigma}_{MLE}^2}{2}\right)$$

Can we use these two uni-dimensional non-marginal posteriors to sample from the full posterior? Note that:

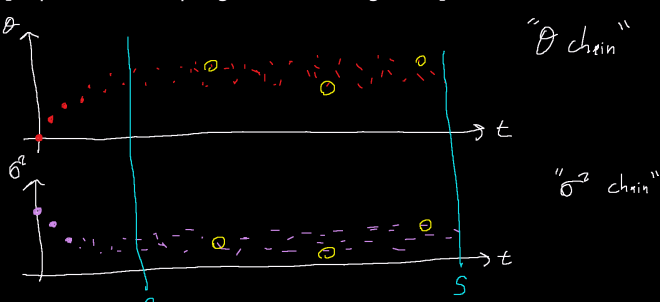$$P(\theta, \sigma^2 | x) = \frac{P(\theta | x, \sigma) \ P(\sigma^2 | x)}{P(\sigma^2 | x, \theta) \ P(\theta | x)}$$

Consider the following numerical sampling algorithm:

(1) Begin at $\theta_0$ = some reasonable value.
(2) Draw $\sigma^2\_1$ from $P(\sigma^2 | X, \theta = \theta\_0)$ via rinvgamma.
(3) Draw $\theta\_1$ from $P(\theta | X, \sigma^2 = \sigma^2\_1)$ via rnorm.
(4) Draw $\sigma^2\_2$ from $P(\sigma^2 | X, \theta = \theta\_1)$ via rinvgamma.
(5) Draw $\theta\_2$ from $P(\theta | X, \sigma^2 = \sigma^2\_2)$ via rnorm.
.....
[Repeat this sampling until "convergence"]



"$\theta$ chain"

"$\sigma^2$ chain"

"B" is the iteration number that refers to the "burn-in" i.e. the number of samples of this algorithm it takes for the algorithm to be converged. Once it's converged, it's considered [almost] ready to provide samples for your inference. There is another problem... these samples are dependent!! Why $\theta\_98$ is very related to $\theta\_97$ which is dependent on $\theta\_96$, etc... It seems like it's impossible to get iid samples from the posterior! That is technically true... but... we can remove ourselves by enough iterations that this dependence is negligible. How do we assess how many iterations? And negligibility?

Let's go back to basic stats for a minute. Consider two r.v.'s $X\_1, X\_2$. Then

$$\sigma_{12} := Cov[X, X_2] = E[(X_1 - \mu_1)(X_2 - \mu_2)]$$

$$Q := Corr[X, X_2] := \frac{\sigma_{12}}{\sigma_1 \sigma_2} \in [-1, +1] \quad \text{prob in 368.}$$

We have estimators/estimates for these parameters:

$$\hat{\sigma}_{12} \approx s_{12} := \frac{1}{n-1} \sum_{i=1}^{\eta} (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)$$

$$\hat{\rho} \approx r = \frac{s_{12}}{s_1 s_2} = \frac{\sum(x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)}{\sqrt{\sum(x_{1i} - \bar{x}_1)^2 \ \sum(x_{2i} - \bar{x}_2)^2}}$$

Now consider $X\_1, X\_2, ..., X\_S$ to be r.v.'s from an iterative process where $X\_1$ is the first iteration, $X\_2$ is the second iteration, etc where each iteration has a dependence on the previous iteration. We define "autocorrelation" which is correlation with a previous iteration. First, we have autocorrelation with the iteration directly before:

$$r_{t1} := \frac{\sum_{t=2}^{S}(x_t - \bar{x})(x_{t-1} - \bar{x})}{\sum_{t=1}^{S}(x_t - \bar{x})^2}, \quad \bar{x} = \frac{1}{S}\sum_{t=1}^{S} x_t$$

$$r_{t2} := \frac{\sum_{t=3}^{S}(x_t - \bar{x})(x_{t-2} - \bar{x})}{\sum_{t=1}^{S}(x_t - \bar{x})^2}$$

$$\vdots$$

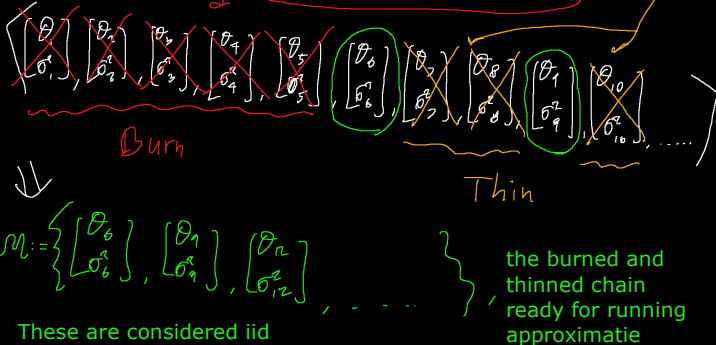In our sampling idea, we assess autocorrelations for both $\theta$ and $\sigma^2$:

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

$k = 1, 2, ...,$ the "lag"

How to assess the dependence? Look at the autocorrelation plot:
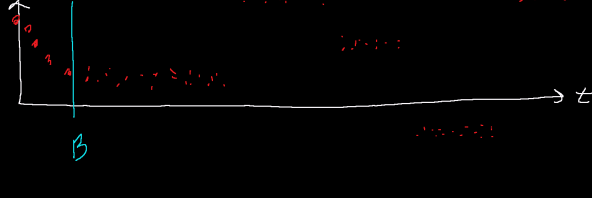


$T = 7$

Once the autocorrelation appears to be not significantly different from zero, you select that first $k$ and call it "T" for "thinning" distance. Then, you "thin" the chain of samples e.g. B = 5, T = 3:



Burn                                              Thin

$$M := \left\{ \begin{bmatrix} \theta_6 \\ \sigma_6^2 \end{bmatrix}, \begin{bmatrix} \theta_9 \\ \sigma_9^2 \end{bmatrix}, \begin{bmatrix} \theta_{12} \\ \sigma_{12}^2 \end{bmatrix}, \dots \right\}$$

the burned and thinned chain ready for running approximatie inference!

These are considered iid samples from the posterior $P(\theta, \sigma | X)$

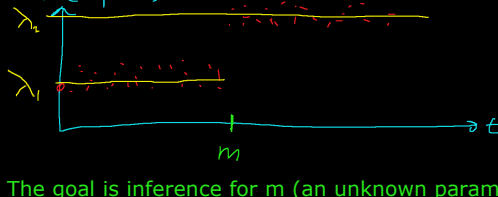There are still problems with this approach. For instance, the following happens frequently:



The problem here is the posterior has many pockets of density that may not be fully explored in this sampling scheme. Partial solution: (1) let the chain run a very long time and (2) begin the chain from different locations and then union the multiple chains together.

Another practical problem is that you either need to know all conditional distributions or be able to grid sample them accurately and without undue computational burden.

The algorithm we just discussed is quite famous and is called the "Systematic Sweep Gibbs Sampler". Here is the general algorithm below to sample from $P(\theta\_1, \theta\_2, ..., \theta\_p | X)$:

the $p$ conditional distributions

(1) Pick $\vec{\theta}_0 = \langle \theta_{0,1}, \theta_{0,2}, ..., \theta_{0,p} \rangle$

(2-1) Sample $\theta_{1,1}$ from $P(\theta_1 | \theta_2 = \theta_{0,2}, ..., \theta_p = \theta_{0,p})$

(2-2) Sample $\theta_{1,2}$ from $P(\theta_2 | \theta_1 = \theta_{1,1}, \theta_3 = \theta_{0,3}, ..., \theta_p = \theta_{0,p})$

...

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

(3) Record $\vec{\theta}_1 = \langle \theta_{1,1}, \theta_{1,2}, ..., \theta_{1,p} \rangle$ i.e. the result of step 2

(4) Repeat step 2 and 3 many times.

(5) Burn all the chains at the highest B value across all p chains.
(6) Thin the chains at the highest T value across all p chains.

─────────────────────────────────────

A real-world example finally! Change-point modeling. Assume there is a poisson number of phone calls with mean $\lambda\_1$ and then at some point in time, it changes to a poisson number of calls with mean $\lambda\_2$ which is diffferent than $\lambda\_1$.



X (# of calls)

The goal is inference for m (an unknown parameter that is the change time). There are two unknown nuisance parameters $\lambda\_1$ and $\lambda\_2$ which you don't need to infer. Thus p = 3 dimensions and the full posterior will be $P(m, \lambda\_1, \lambda\_2 | x)$. We will build a Gibbs sampler that will provide us inference next class.