

Markov chain Monte Carlo

Dr. Jarad Niemi

STAT 544 - Iowa State University

April 2, 2018

Markov chain construction

The techniques we have discussed thus far, e.g.

- Metropolis-Hastings
 - independent Metropolis-Hastings
 - random-walk Metropolis
 - Hamiltonian Monte Carlo
- Gibbs sampling
 - Slice sampling

form a set of techniques referred to as **Markov chain Monte Carlo** (MCMC).

Today we look at some practical questions involving the use of MCMC:

- What initial values should I use?
- How long do I need to run my chain?
- What can I do with the samples I obtain?

Markov chain Monte Carlo

An MCMC algorithm with transition kernel $K(\theta^{(t-1)}, \theta^{(t)})$ constructed to sample from $p(\theta|y)$ is the following:

1. Sample $\theta^{(0)} \sim \pi^{(0)}$.
2. For $t = 1, \dots, T$, perform the kernel $K(\theta^{(t-1)}, \theta^{(t)})$ to obtain a sequence $\theta^{(1)}, \dots, \theta^{(t)}$.

The questions can then be rephrased as

- What should I use for $\pi^{(0)}$?
- What should T be?
- What can I do with $\theta^{(1)}, \dots, \theta^{(t)}$?

Initial values

For ergodic Markov chains with stationary distribution $p(\theta|y)$, theory states that

$$\theta^{(t)} \xrightarrow{d} \theta \text{ where } \theta \sim p(\theta|y)$$

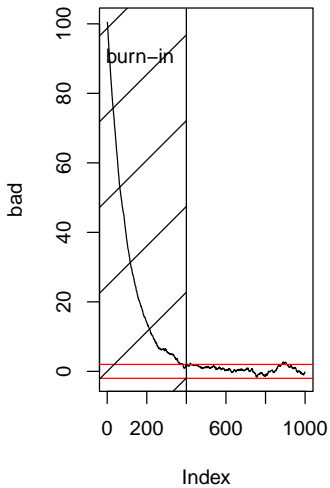
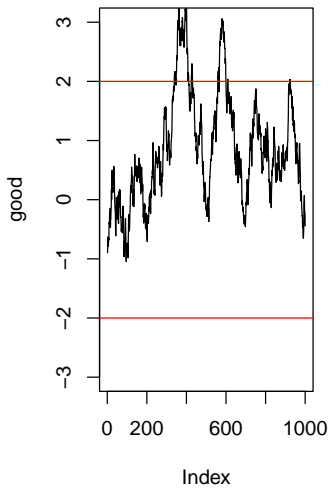
for **almost all** $\theta^{(0)}$ (all with Harris recurrence).

If $p(\theta^{(0)}|y) \ll p(\theta_{MAP}|y)$, then this can take a long time. For example, let

$$\theta^{(t)} = 0.99\theta^{(t-1)} + \epsilon_t \quad \epsilon_t \stackrel{iid}{\sim} N(0, 1 - .99^2)$$

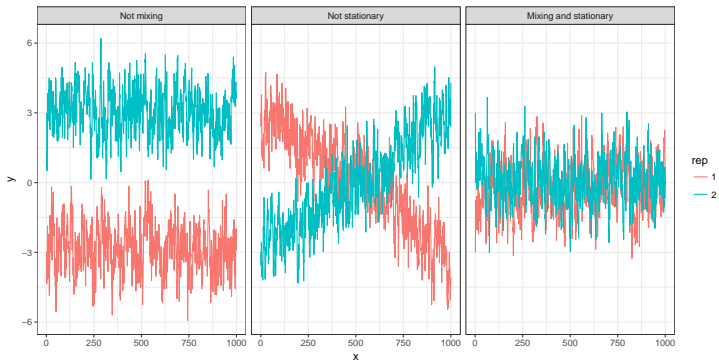
which has stationary distribution $p(\theta|y) \stackrel{d}{=} N(0, 1)$. If

- $\theta^{(0)} \sim p(\theta|y)$ then $\theta^{(t)} \dot{\sim} p(\theta|y)$ for all t , but if
- $\theta^{(0)}$ is very far from $p(\theta|y)$ then $\theta^{(t)} \dot{\sim} p(\theta|y)$ only for t very large.



How many iterations do I need for burn-in?

Imagine two different chains



Gelman-Rubin potential scale reduction factor

1. Start multiple chains with initial values that are **well dispersed values relative to $p(\theta|y)$** .
2. For each scalar estimand ψ of interest,
 - Calculate the between B and within W chain variances
 - Estimate the the marginal posterior variance of the estimand, i.e. $Var(\psi|y)$:

$$\widehat{Var}^+(\psi|y) = \frac{t-1}{t}W + \frac{1}{t}B$$

where t is the number of iterations.

- Calculate the potential scale reduction factor

$$\hat{R}_\psi = \sqrt{\frac{\widehat{Var}^+(\psi|y)}{W}}$$

3. If the \hat{R}_ψ are approximately 1, e.g. <1.1 , then **there is no evidence of non-convergence**.

Example potential scale reduction factors

```
[1] "Not mixing"  
Potential scale reduction factors:
```

```
      Point est. Upper C.I.  
[1,]      7.35      16.2
```

```
[1] "Not stationary"  
Potential scale reduction factors:
```

```
      Point est. Upper C.I.  
[1,]      2.62      5.31
```

```
[1] "Mixing and stationary"  
Potential scale reduction factors:
```

```
      Point est. Upper C.I.  
[1,]      1.01      1.04
```


Methods for finding good initial values

From <http://users.stat.umn.edu/~geyer/mcmc/burn.html>:

Any point you don't mind having in a sample is a good starting point.

Methods for finding good initial values:

- burn-in: throw away the first X iterations
- Start at the MLE, i.e. $\operatorname{argmax}_{\theta} p(y|\theta)$
- Start at the MAP (maximum aposterior), i.e. $\operatorname{argmax}_{\theta} p(\theta|y)$

How many iterations should I run (post 'convergence')?

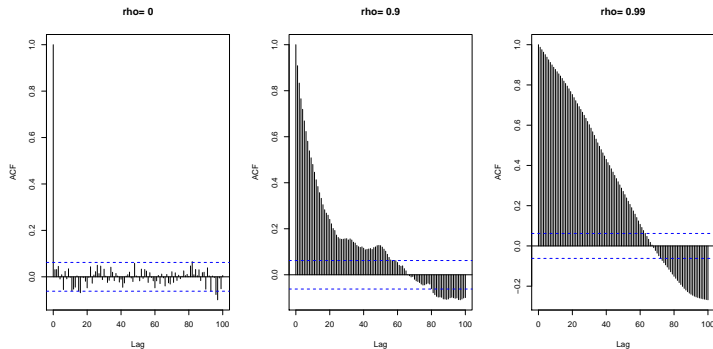
Compute the effective sample size, i.e. how many independent samples would we need to get the equivalent precision of our estimates?

```
d = ddply(data.frame(rho=c(0,.9,.99)), .(rho), function(x) data.frame(x=rwalk(1000,0,x$rho)))
ddply(d, .(rho), summarize,
      effective_size = round(coda::effectiveSize(x)))
```

	rho	effective_size
1	0.00	1000
2	0.90	35
3	0.99	6

BDA3 a total of 100-2000 effective samples. But this really depends on what you want to estimate. If you are interested in estimating probabilities of rare events, i.e. tail probabilities, you may need many more samples.

Autocorrelation function



Monte Carlo integration

Consider approximating the integral via it's Markov chain Monte Carlo (MCMC) estimate, i.e.

$$E_{\theta|y}[h(\theta)|y] = \int_{\Theta} h(\theta)p(\theta|y)d\theta \quad \text{and} \quad \hat{h}_T = \frac{1}{T} \sum_{t=1}^{(t)} h\left(\theta^{(t)}\right).$$

where $\theta^{(t)}$ is the t^{th} iteration from the MCMC. Under regularity conditions,

- SLLN: \hat{h}_T converges almost surely to $E[h(\theta)|y]$.
- CLT: **under stronger regularity conditions**,

$$\hat{h}_T \xrightarrow{d} N\left(E[h(\theta)|y], \sigma^2/T\right)$$

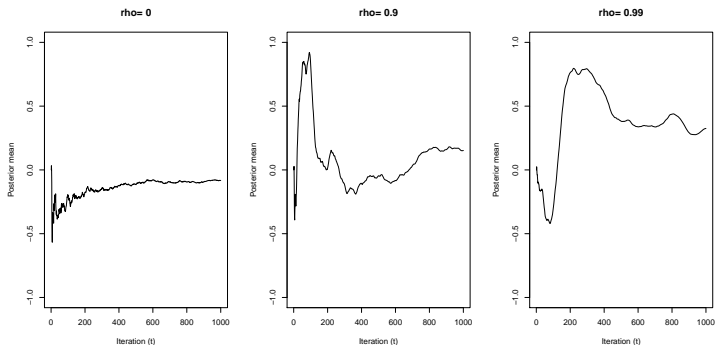
where

$$\sigma^2 = Var[h(\theta)|y] \left(1 + 2 \sum_{k=1}^{\infty} \rho_k\right)$$

where ρ_k is the k^{th} autocorrelation of the $h(\theta)$ values.

Sequential estimates

```
opar = par(mfrow=c(1,3))
d_ply(d, .(rho), function(x)
  plot(cumsum(x$x)/1:length(x$x), type="l", ylim=c(-1,1),
    ylab="Posterior mean", xlab="Iteration (t)", main=paste("rho=", x$rho[1]))
)
```



```
par(opar)
```

Treat the MCMC samples as samples from the posterior

Use `mcmcse::mcse` to estimate the MCMC variance

```
# Mean
ddply(d, .(rho), function(x) round(as.data.frame(mcmcse::mcse(x$x)),2))
```

	rho	est	se
1	0.00	-0.08	0.03
2	0.90	0.15	0.14
3	0.99	0.33	0.15

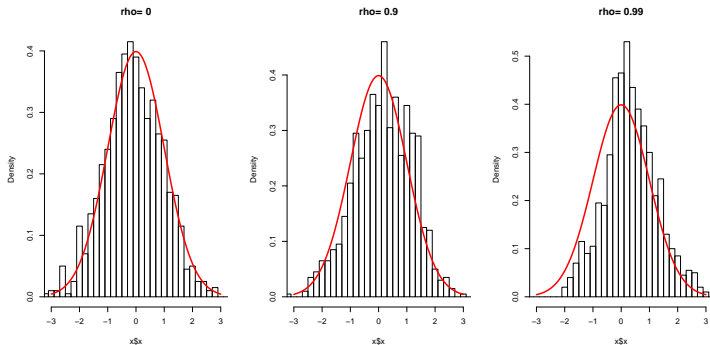
```
# Quantiles
ddply(d, .(rho), function(x) round(as.data.frame(mcmcse::mcse.q(x$x, .025)),2))
```

	rho	est	se
1	0.00	-2.06	0.08
2	0.90	-2.01	0.26
3	0.99	-1.53	0.24

```
ddply(d, .(rho), function(x) round(as.data.frame(mcmcse::mcse.q(x$x, .975)),2))
```

	rho	est	se
1	0.00	1.92	0.10
2	0.90	1.94	0.17
3	0.99	2.41	0.34

Treat the MCMC samples as samples from the posterior



A wasteful approach

The Gelman approach in practice is the following

1. Run an initial chain or, in some other way, approximate the posterior.
2. (Randomly) choose initial values for multiple chains well dispersed relative to this approximation to the posterior.
3. Run the chain until all estimands of interest have potential scale reduction factors less than 1.1.
4. Continuing running until you have a total of around 4,000 effective draws.
5. Discard the first half of all the chains.

Assuming this approach correctly diagnosis convergence or lack thereof, it seems computationally wasteful since

- You had to run an initial chain, but then threw it away.
- You threw away half of your later iterations.

One really long chain

From <http://users.stat.umn.edu/~geyer/mcmc/one.html>

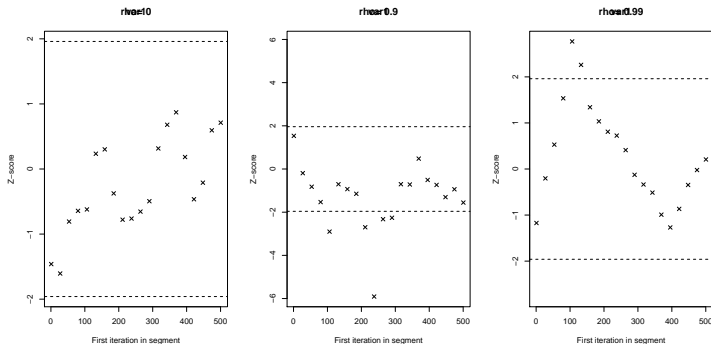
If you can't get a good answer with one long run, then you can't get a good answer with many short runs either.

1. Start a chain at a reasonable starting value.
2. Run it for many iterations (and keep running it).

If you really want a convergence diagnostic, you can try Geweke's which tests for equality of means in the first and last parts of the chain.

Geweke diagnostic

```
# Z-score for test of equality of means
par(mfrow=c(1,3))
d_ply(d, .(rho), function(x) geweke.plot(mcmc(x$x), auto=F, main=paste("rho=", x$rho[1])))
```



```
par(opar)
```

Thinning

You will hear of people **thinning** their Markov chain by only recording every n^{th} observation.

This has the benefit of reducing the autocorrelation in the retained samples.

But should only be used if memory or hard drive space is a limiting factor.

Thinning

```
sq = seq(10,1000,by=10)
ddply(d, .(rho), summarize, full=effectiveSize(x), thinned=effectiveSize(x[sq]))
```

	rho	full	thinned
1	0.00	1000.000000	103.29644
2	0.90	35.405683	39.37303
3	0.99	6.435595	16.21098

```
# Calculate standard error
ddply(d, .(rho), function(x) {
  rbind(data.frame(s="full", mcmcse::mcse(x$x)), data.frame(s="thinned", mcmcse::mcse(x$x[sq])))
})
```

	rho	s	est	se
1	0.00	full	-0.08223773	0.02789422
2	0.00	thinned	-0.16062926	0.08828254
3	0.90	full	0.15262785	0.13563890
4	0.90	thinned	0.19911506	0.16895797
5	0.99	full	0.32514041	0.15349268
6	0.99	thinned	0.32068486	0.26609394

Alternative use for burn-in

For MCMC algorithms that have tuning parameters, use burn-in (warm-up) to tune tuning parameters.

Suppose the target distribution is $N(0, 1)$ and we are performing a random-walk Metropolis with a normal proposal. The variance of this proposal is a tuning parameter and we can tune it during burn-in:

- if a proposal is accepted, then likely our variance is too small and therefore we should increase it
- if a proposal is rejected, then likely our variance is too big and therefore we should decrease it

Alternative use for burn-in

```
rw = function(n, theta0, tune=1, autotune=TRUE) {
  theta = rep(theta0, n)
  for (i in 2:n) {
    theta_prop = rnorm(1, theta[i-1], tune)
    logr = dnorm(theta_prop, log=TRUE) - dnorm(theta[i-1], log=TRUE)

    # This tuning tunes to an acceptance rate of 50%
    if (log(runif(1))<logr) {
      theta[i] = theta_prop
      if (autotune) tune = tune*1.1
    } else {
      theta[i] = theta[i-1]
      if (autotune) tune = tune/1.1
    }
  }
  return(list(theta=theta, tune=tune))
}

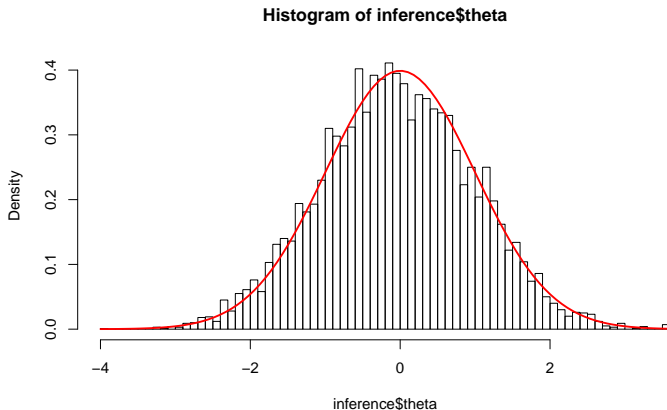
# Tune during burn-in
burnin = rw(1000, 0)
burnin$tune

[1] 1.61051

# Turn off tuning after burn-in for theory to hold
inference = rw(10000, burnin$theta[1000], burnin$tune, autotune=FALSE)
```

Alternative use for burn-in

```
hist(inference$theta, 100, prob=T)  
curve(dnorm, col="red", add=TRUE, lwd=2)
```



Summary

Since computing time/power is not very limited these days, my suggestion is

1. Run one long chain and continue running it
2. Run multiple chains according to suggestions in BDA
 - a. Start multiple chains with initial values relative to the posterior learned by the long chain
 - b. Monitor the potential scale reduction factor until < 1.1 for all quantities of interest
 - c. Monitor traceplots and cumulative mean plots
 - d. Discard burn-in (first half is probably overkill)
 - e. Run until effective sample size is around 2000
3. Use all samples for posterior inference

If things are not going well,

1. Check for identifiability of the parameters in your model.
2. Construct a better sampler.

A simple model

Let

$$Y_i \stackrel{\text{ind}}{\sim} N(\mu, \sigma^2) \quad \text{and} \quad p(\mu, \sigma) \propto Ca^+(\sigma; 0, 1)$$

In RStan,

```
model = "
data {
  int<lower=1> n;
  real y[n];
}
parameters{
  real mu;
  real<lower=0> sigma;
}
model {
  sigma ~ cauchy(0,1);
  y ~ normal(mu,sigma);
}
"
```

RStan

```
y = rnorm(10)
m = stan_model(model_code = model)
r = sampling(m, list(n=length(y), y=y))
```

Warning: There were 1 divergent transitions after warmup. Increasing adapt_delta above 0.8 may help. See <http://mc-stan.org/misc/warnings.html#divergent-transitions-after-warmup>

Warning: Examine the pairs() plot to diagnose sampling problems

```
r
```

Inference for Stan model: 6c86a547f723283854dd490525d54ee4.
 4 chains, each with iter=2000; warmup=1000; thin=1;
 post-warmup draws per chain=1000, total post-warmup draws=4000.

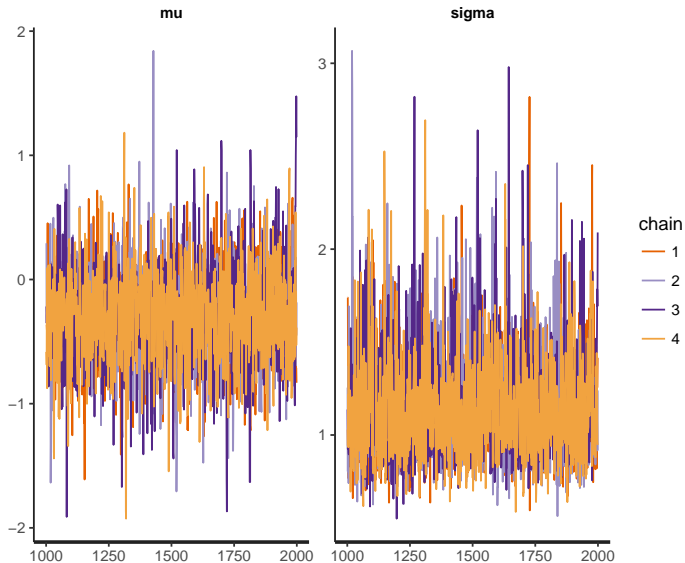
	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
mu	-0.30	0.01	0.38	-1.06	-0.55	-0.30	-0.05	0.45	2415	1
sigma	1.15	0.01	0.29	0.74	0.95	1.10	1.29	1.85	1701	1
lp__	-6.89	0.03	1.12	-9.93	-7.29	-6.56	-6.13	-5.83	1528	1

Samples were drawn using NUTS(diag_e) at Tue Apr 11 09:51:46 2017.
 For each parameter, n_eff is a crude measure of effective sample size,
 and Rhat is the potential scale reduction factor on split chains (at
 convergence, Rhat=1).

```
laply(extract(r, c("mu", "sigma")), function(x) length(unique(x))/length(x)) # Acceptance rate
```

```
[1] 0.85025 0.85025
```

RStan plot



Hierarchical binomial model

Recall the game-wise Andre Dawkins 3-point percentage data set with a hierarchical binomial model:

$$\begin{aligned} Y_i &\overset{ind}{\sim} \text{Bin}(n_i, \theta_i) \\ \theta_i &\overset{ind}{\sim} \text{Be}(\alpha, \beta) \\ p(\alpha, \beta) &\propto (\alpha + \beta)^{-5/2} \end{aligned}$$

In RStan,

```
hierarchical_binomial_model = "
data {
  int<lower=1> N; int<lower=0> n[N]; int<lower=0> y[N];
}
parameters {
  real<lower=0,upper=1> theta[N]; real<lower=0> alpha; real<lower=0> beta;
}
transformed parameters {
  real<lower=0,upper=1> mu; real<lower=0> eta;
  eta = alpha+beta; mu = alpha/eta;
}
model {
  target += -5*log(alpha+beta)/2;
  theta ~ beta(alpha,beta);
  y ~ binomial(n,theta);
}
"
```

Initial run

```
m = stan_model(model_code = hierarchical_binomial_model)

r = sampling(m,
  data = list(N = nrow(dawkins),
              n = dawkins$attempts,
              y = dawkins$made),
  pars = c("alpha", "beta", "mu", "eta"),
  chains = 1)
```

Warning: There were 29 divergent transitions after warmup. Increasing adapt_delta above 0.8 may help. See

<http://mc-stan.org/misc/warnings.html#divergent-transitions-after-warmup>

Warning: There were 1 chains where the estimated Bayesian Fraction of Missing Information was low. See <http://mc-stan.org/misc/warnings.html#bfmi-low>

Warning: Examine the pairs() plot to diagnose sampling problems

RStan

r

Inference for Stan model: dc6c37e2deab377b892b94e3b0b4a542.

1 chains, each with iter=2000; warmup=1000; thin=1;

post-warmup draws per chain=1000, total post-warmup draws=1000.

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
alpha	11.80	1.69	10.79	1.48	4.12	8.22	15.53	44.55	41	1.00
beta	13.35	1.92	12.23	1.75	4.85	9.18	16.88	49.51	41	1.01
mu	0.47	0.00	0.05	0.37	0.44	0.47	0.50	0.58	365	1.00
eta	25.16	3.61	22.86	3.38	9.11	17.34	32.99	94.55	40	1.01
lp__	-98.61	1.28	9.67	-117.22	-105.78	-98.53	-91.26	-81.23	57	1.00

Samples were drawn using NUTS(diag_e) at Mon Apr 2 16:01:54 2018.

For each parameter, n_eff is a crude measure of effective sample size,
and Rhat is the potential scale reduction factor on split chains (at
convergence, Rhat=1).

```
laply(rstan::extract(r, c("mu","eta")), function(x) length(unique(x))/length(x)) # Acceptance rate
```

```
[1] 0.983 0.983
```

RStan plot

