

# LECTURE 18: SOME MCMC PRACTICALITIES

STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

---

Vinayak Rao

Purdue University

November 19, 2018

## SUMMARY SO FAR...

Independent samples from prob. distrib.  $p$  is often difficult.

MCMC addresses this by producing dependent samples.

- Begin with an arbitrary initialization  $X_0$ .
- Sequentially produce samples  $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_N$ .

If the chain is stationary w.r.t.  $p(x)$ , irreducible and aperiodic:

$$\frac{1}{S} \sum_{i=1}^S h(X_i) \rightarrow \mathbb{E}_p[h]$$

## SUMMARY SO FAR...

Independent samples from prob. distrib.  $p$  is often difficult.

MCMC addresses this by producing dependent samples.

- Begin with an arbitrary initialization  $X_0$ .
- Sequentially produce samples  $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_N$ .

If the chain is stationary w.r.t.  $p(x)$ , irreducible and aperiodic:

$$\frac{1}{S} \sum_{i=1}^S h(X_i) \rightarrow \mathbb{E}_p[h]$$

In practice,  $S$  is finite.

Assessing error is much harder

## HOW WELL DOES YOUR CHAIN MIX?

Are our MCMC samples representative of the overall posterior?

- Difficult with multimodal distributions.

Do we have enough samples to estimate expectations accurately?

- This is hard with Monte Carlo methods in general
- Trickier with MCMC because of correlation between samples.

Burn-in time: time to 'forget' the arbitrary initialization.

Typically deal with burn-in by discarding the first  $B$  samples  
(e.g.  $B = 1000$ )

Burn-in time: time to 'forget' the arbitrary initialization.

Typically deal with burn-in by discarding the first  $B$  samples (e.g.  $B = 1000$ )

Sometimes people deal with sample dependence by 'thinning' the Markov chain: E.g. Use every  $m$ th sample (e.g.  $m = 10$ )

Thinning is usually unnecessary and increases variance of estimates (unless you want to save memory/computation).

Burn-in time: time to 'forget' the arbitrary initialization.

Typically deal with burn-in by discarding the first  $B$  samples (e.g.  $B = 1000$ )

Sometimes people deal with sample dependence by 'thinning' the Markov chain: E.g. Use every  $m$ th sample (e.g.  $m = 10$ )

Thinning is usually unnecessary and increases variance of estimates (unless you want to save memory/computation).

However, it's worthwhile remembering that  $N$  MCMC samples correspond to a smaller number of independent samples.

## EFFECTIVE SAMPLE SIZE

A good diagnostic is the effective sample size (ESS):

$$N_{ESS} = \frac{N}{1 + 2 \sum_{k=1}^{\infty} \rho_k}$$

$\rho_k$  is the auto-correlation between  $X_i$  and  $X_{i+k}$ :

$$\rho_k = \frac{\mathbb{E}[(X_{i+k} - \mu)(X_i - \mu)]}{\sigma^2}$$

$(\mu, \sigma^2)$  are mean and var. of  $X_i$  under stationary distribution.



## EFFECTIVE SAMPLE SIZE

A good diagnostic is the effective sample size (ESS):

$$N_{ESS} = \frac{N}{1 + 2 \sum_{k=1}^{\infty} \rho_k}$$

$\rho_k$  is the auto-correlation between  $X_i$  and  $X_{i+k}$ :

$$\rho_k = \frac{\mathbb{E}[(X_{i+k} - \mu)(X_i - \mu)]}{\sigma^2}$$

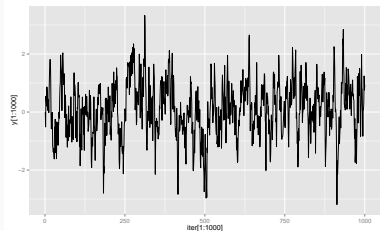
$(\mu, \sigma^2)$  are mean and var. of  $X_i$  under stationary distribution.

Comes from CLT for Markov chains:

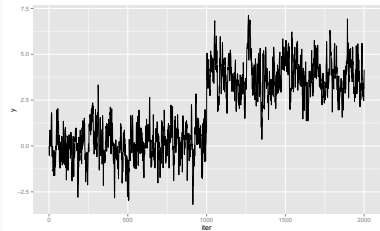
$$\left( \frac{1}{N} \sum_{i=1}^N f(X_i) - \mathbb{E}[f(X)] \right) \rightarrow \mathcal{N}(0, \sigma^2 / N_{ESS})$$

# EFFECTIVE SAMPLE SIZE

The coda package in R calculates this and other diagnostics.



ESS: 130.4

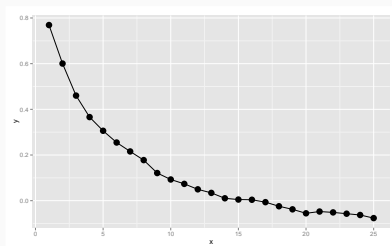


ESS: 9.21

```
> effectiveSize(data.frame(half=z[1:1000],full=z))  
      half      full  
260.997261  9.216991
```

Note: always useful to visualize traceplots.

## OTHER DIAGNOSTICS

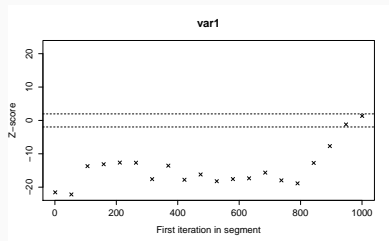
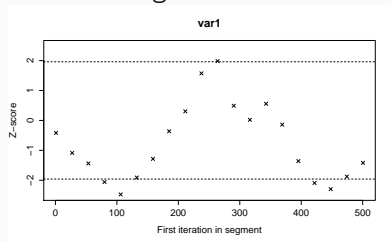


Correlation vs lag

```
> acf <- autocorr(mcmc(z[1:1000]),c(1:25))
```

# OTHER DIAGNOSTICS

Geweke diagnostic:



Compare 2 non-overlapping parts of the chain (in R `CODA` is the first 10% and last 50%, and test if their means come from the same distribution.

Can repeat, successively discarding initial parts.

```
> geweke.plot(mcmc(z[1:1000]))
```

```
> geweke.plot(mcmc(z))
```

**Gelman-Rubin diagnostic:** Run  $m \geq 2$  independent chains with overdispersed starting points (e.g. sampled from the prior)

- Calculate within-chain variance and between-chain variance.

**Gelman-Rubin diagnostic:** Run  $m \geq 2$  independent chains with overdispersed starting points (e.g. sampled from the prior)

- Calculate within-chain variance and between-chain variance.
- Former typically underestimates variance (bad mixing), and latter overestimates it (overdispersed initialization).
- If latter is much larger than former, run chain longer

> `gelman.diag`

## ONE LONG CHAIN VS MANY SHORTER CHAINS?

$M$  short chain of length  $N$  vs 1 chain of length  $MN$ :

# ONE LONG CHAIN VS MANY SHORTER CHAINS?

$M$  short chain of length  $N$  vs 1 chain of length  $MN$ :

Pros:

- Diverse initialization likely means better exploration of different modes.
- Allows easy parallelization



# ONE LONG CHAIN VS MANY SHORTER CHAINS?

$M$  short chain of length  $N$  vs 1 chain of length  $MN$ :

Pros:

- Diverse initialization likely means better exploration of different modes.
- Allows easy parallelization

Cons:

- Each chain still has a burn-in period  $B$ . Must discard  $MB$  samples vs  $B$  for a single chain.

Never mind mixing, how do we know our sampler is correct?!

After changing something, how do we know it's still correct?

Never mind mixing, how do we know our sampler is correct?!

After changing something, how do we know it's still correct?  
Can never be sure, but useful to run a few standard tests.

Never mind mixing, how do we know our sampler is correct?!

After changing something, how do we know it's still correct?

Can never be sure, but useful to run a few standard tests.

Do your results make sense for special cases?

Never mind mixing, how do we know our sampler is correct?!

After changing something, how do we know it's still correct?

Can never be sure, but useful to run a few standard tests.

Do your results make sense for special cases?

Compare different samplers: a Gibbs and MH sampler should give similar results, but unlikely to have same errors.

Never mind mixing, how do we know our sampler is correct?!

After changing something, how do we know it's still correct?

Can never be sure, but useful to run a few standard tests.

Do your results make sense for special cases?

Compare different samplers: a Gibbs and MH sampler should give similar results, but unlikely to have same errors.

On scaled down datasets, compare with simple Monte Carlo methods like rejection/importance sampling.

Never mind mixing, how do we know our sampler is correct?!

After changing something, how do we know it's still correct?  
Can never be sure, but useful to run a few standard tests.

Do your results make sense for special cases?

Compare different samplers: a Gibbs and MH sampler should give similar results, but unlikely to have same errors.

On scaled down datasets, compare with simple Monte Carlo methods like rejection/importance sampling.

Can you analytically calculate the posterior for 1 observation or 2 states or 2 time-periods?

## USING MCMC SAMPLES:

Consider a Markov chain on  $(x, y, z)$  with stationary distrib.  $P()$ .

We obtain samples  $(x_1, y_1, z_1), (x_2, y_2, z_2), (x_3, y_3, z_3), \dots$



## USING MCMC SAMPLES:

Consider a Markov chain on  $(x, y, z)$  with stationary distrib.  $P()$ .

We obtain samples  $(x_1, y_1, z_1), (x_2, y_2, z_2), (x_3, y_3, z_3), \dots$

What is  $\mathbb{E}[f(x, y, z)]$ ?

$$\mathbb{E}[f(x, y, z)] \approx \frac{1}{N} \sum_{i=1}^N f(x_i, y_i, z_i)$$

## USING MCMC SAMPLES:

Consider a Markov chain on  $(x, y, z)$  with stationary distrib.  $P()$ .

We obtain samples  $(x_1, y_1, z_1), (x_2, y_2, z_2), (x_3, y_3, z_3), \dots$

What is  $\mathbb{E}[f(x, y, z)]$ ?

$$\mathbb{E}[f(x, y, z)] \approx \frac{1}{N} \sum_{i=1}^N f(x_i, y_i, z_i)$$

What is  $P(x = 1)$ ?

## USING MCMC SAMPLES:

Consider a Markov chain on  $(x, y, z)$  with stationary distrib.  $P()$ .  
We obtain samples  $(x_1, y_1, z_1), (x_2, y_2, z_2), (x_3, y_3, z_3), \dots$

What is  $\mathbb{E}[f(x, y, z)]$ ?

$$\mathbb{E}[f(x, y, z)] \approx \frac{1}{N} \sum_{i=1}^N f(x_i, y_i, z_i)$$

What is  $P(x = 1)$ ?

$$P(x = 1) = \mathbb{E}[\delta(x = 1)] \approx \frac{1}{N} \sum_{i=1}^N \delta(x_i = 1)$$

## USING MCMC SAMPLES:

Consider a Markov chain on  $(x, y, z)$  with stationary distrib.  $P()$ .  
We obtain samples  $(x_1, y_1, z_1), (x_2, y_2, z_2), (x_3, y_3, z_3), \dots$

What is  $\mathbb{E}[f(x, y, z)]$ ?

$$\mathbb{E}[f(x, y, z)] \approx \frac{1}{N} \sum_{i=1}^N f(x_i, y_i, z_i)$$

What is  $P(x = 1)$ ?

$$P(x = 1) = \mathbb{E}[\delta(x = 1)] \approx \frac{1}{N} \sum_{i=1}^N \delta(x_i = 1)$$

Can we do better? E.g. what if  $x$  is continuous and we want the density  $p(x = 1)$ ?

Suppose we can calculate  $P(x|y, z)$ .

This is the case if our Markov chain is a Gibbs sampler.

Suppose we can calculate  $P(x|y, z)$ .

This is the case if our Markov chain is a Gibbs sampler.

Then:

$$P(x = 1) = \int \int P(x = 1, y, z) dy dz$$

Suppose we can calculate  $P(x|y, z)$ .

This is the case if our Markov chain is a Gibbs sampler.

Then:

$$\begin{aligned} P(x = 1) &= \int \int P(x = 1, y, z) dy dz \\ &= \int \int P(x = 1|y, z) p(y, z) dy dz \end{aligned}$$

Suppose we can calculate  $P(x|y, z)$ .

This is the case if our Markov chain is a Gibbs sampler.

Then:

$$\begin{aligned} P(x = 1) &= \int \int P(x = 1, y, z) dy dz \\ &= \int \int P(x = 1|y, z) p(y, z) dy dz \\ &\approx \frac{1}{N} \sum_{i=1}^N P(x = 1|y_i, z_i) \end{aligned}$$



Suppose we can calculate  $P(x|y, z)$ .

This is the case if our Markov chain is a Gibbs sampler.

Then:

$$\begin{aligned} P(x = 1) &= \int \int P(x = 1, y, z) dy dz \\ &= \int \int P(x = 1|y, z) p(y, z) dy dz \\ &\approx \frac{1}{N} \sum_{i=1}^N P(x = 1|y_i, z_i) \end{aligned}$$

Typically, this estimate will have lower variance.