# **LECTURE 17: SOME MCMC PRACTICALITIES** STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

Vinayak Rao Purdue University

November 17, 2016

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  $X_1 \rightarrow X_2 \rightarrow \ldots \rightarrow X_N$ .

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

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

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  $X_1 \rightarrow X_2 \rightarrow \ldots \rightarrow X_N$ .

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

$$\frac{1}{5}\sum_{i=1}^{5}h(X_i)\to \mathbb{E}_p[h]$$

In practice, S is finite. Assessing error is much harder

#### How much burn-in is required to forget initial state?

How much burn-in is required to forget initial state?

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? Tricky 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. 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 variance of the stationary distribution.

CLT for Markov chains:

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

$$N_{ESS} = \frac{N}{1 + 2\sum_{k=1}^{\infty} \rho_k},$$
$$\rho_k = \frac{\mathbb{E}[(f(X_{i+k}) - \mu)(f(X_i) - \mu)]}{\sigma^2}$$

 $(\mu, \sigma^2)$  are mean, variance of f(X) under stationary distribution.

CLT for Markov chains:

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

$$N_{ESS} = \frac{N}{1 + 2\sum_{k=1}^{\infty} \rho_k},$$
$$\rho_k = \frac{\mathbb{E}[(f(X_{i+k}) - \mu)(f(X_i) - \mu)]}{\sigma^2}$$

 $(\mu, \sigma^2)$  are mean, variance of f(X) under stationary distribution. Different variables/functions have different ESS. Often take the minimum of a few.

#### **EFFECTIVE SAMPLE SIZE**

The Coda package in Rcalculates this and other diagnostics.



ESS: 130.4

ESS: 9.21

Note: always useful to visualize traceplots.



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



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  $m \ge 2$  independent chains with overdispersed starting points (e.g. sampled from the prior)

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

**Gelman-Rubin diagnostic:** Run  $m \ge 2$  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).
- $\cdot\,$  If latter is much larger than former, run chain longer

**Gelman-Rubin diagnostic:** Run  $m \ge 2$  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(mcmc.list(mcmc(z[1:1000]),

```
mcmc(z[1001:2000])))
```

Potential scale reduction factors:

Point est. Upper C.I. [1,] 4.87 10.6

Potential scale reduction factor much larger than 1 is trouble

#### ONE LONG CHAIN VS MANY SHORTER CHAINS?

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

*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

*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.

# DEBUGGING MCMC

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

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

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

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?

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.

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.

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?

Generate a new dataset every MCMC iteration.

Generate a new dataset every MCMC iteration.

Every iteration, MCMC samples  $p(Y_{n+1}|X, Y_n)$ .

Generate a new dataset every MCMC iteration.

Every iteration, MCMC samples  $p(Y_{n+1}|X, Y_n)$ .

After this sample a new dataset  $p(X_{n+1}|Y_{n+1})$ .

Generate a new dataset every MCMC iteration.

Every iteration, MCMC samples  $p(Y_{n+1}|X, Y_n)$ . After this sample a new dataset  $p(X_{n+1}|Y_{n+1})$ . Overall, a Gibbs sampler on (X, Y)

Generate a new dataset every MCMC iteration.

Every iteration, MCMC samples  $p(Y_{n+1}|X, Y_n)$ .

After this sample a new dataset  $p(X_{n+1}|Y_{n+1})$ .

Overall, a Gibbs sampler on (X, Y)

What is the joint distribution?

What is the marginal distribution of Y?

Generate a new dataset every MCMC iteration.

Every iteration, MCMC samples  $p(Y_{n+1}|X, Y_n)$ .

After this sample a new dataset  $p(X_{n+1}|Y_{n+1})$ .

Overall, a Gibbs sampler on (X, Y)

What is the joint distribution?

What is the marginal distribution of Y?

James P. Hobert and Galin L. Jones, Honest Exploration of Intractable Probability Distributions via Markov Chain Monte Carlo Statist. Sci. Volume 16, Number 4 (2001), 312-334.

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

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

What is P(x = 1)?

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)?

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

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?

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:

$$P(x=1) = \int \int P(x=1,y,z) \mathrm{d}y \mathrm{d}z$$

Suppose we can calculate P(x|y, z). This is the case if our Markov chain is a Gibbs sampler. Then:

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

Suppose we can calculate P(x|y,z). This is the case if our Markov chain is a Gibbs sampler. Then:

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

Suppose we can calculate P(x|y,z). This is the case if our Markov chain is a Gibbs sampler. Then:

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

Typically, this estimate will have lower variance.