

Bayesian Statistical Methods for Astronomy

Part II: Markov Chain Monte Carlo

David A. van Dyk

Department of Statistics, University of California, Irvine
Statistics Section, Imperial College London

INAF - Osservatorio Astrofisico di Arcetri, September 2014

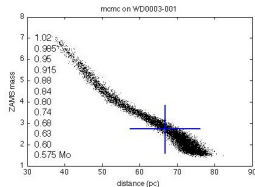
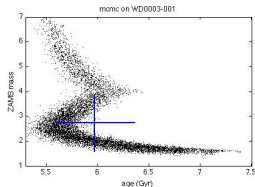
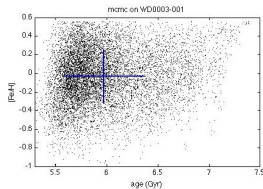
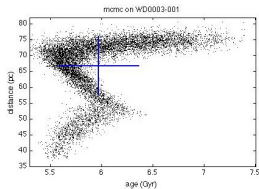
Outline

- 1 **Background**
 - Complex Posterior Distributions
 - Monte Carlo Integration
 - Markov Chains
- 2 **Basic MCMC Jumping Rules**
 - Metropolis Sampler
 - Metropolis Hastings Sampler
 - Basic Theory
- 3 **Practical Challenges and Advice**
 - Diagnosing Convergence
 - Choosing a Jumping Rule
 - Transformations and Multiple Modes
- 4 **The Gibbs Sampler and Data Augmentation**
 - The Gibbs Sampler
 - Data Augmentation

Outline

- 1 **Background**
 - Complex Posterior Distributions
 - Monte Carlo Integration
 - Markov Chains
- 2 **Basic MCMC Jumping Rules**
 - Metropolis Sampler
 - Metropolis Hastings Sampler
 - Basic Theory
- 3 **Practical Challenges and Advice**
 - Diagnosing Convergence
 - Choosing a Jumping Rule
 - Transformations and Multiple Modes
- 4 **The Gibbs Sampler and Data Augmentation**
 - The Gibbs Sampler
 - Data Augmentation

Complex Posterior Distributions

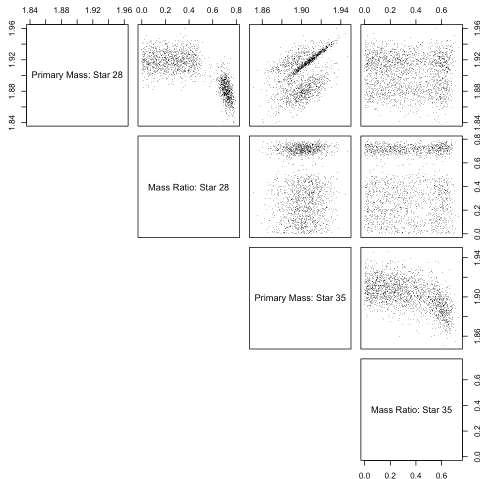


Highly non-linear relationship among stellar parameters.

Complex Posterior Distributions

Highly non-linear relationships among stellar parameters.

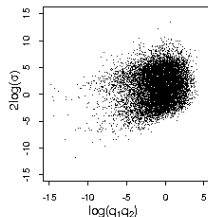
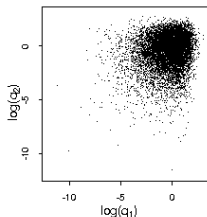
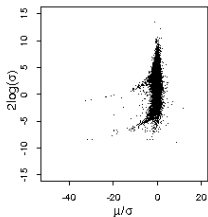
Complex Posterior Distributions



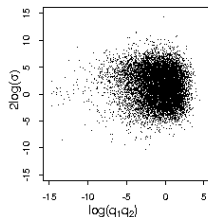
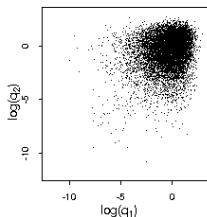
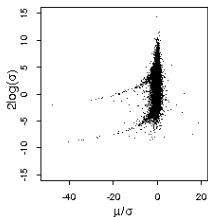
The classification of certain stars as field or cluster stars can cause multiple modes in the distributions of other parameters.

Complex Posterior Distributions

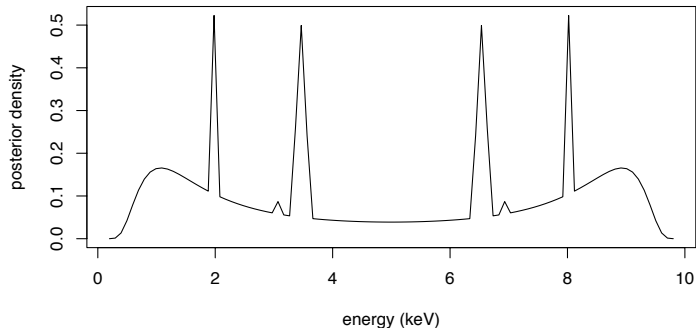
Standard Algorithm
one degree of freedom



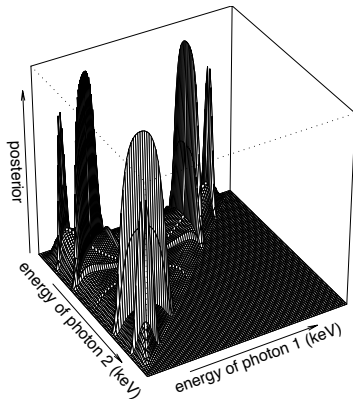
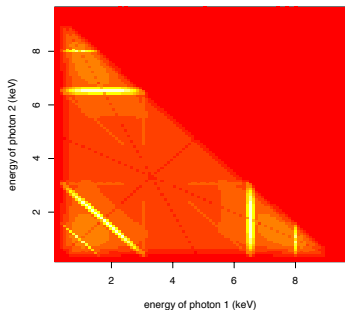
Marginal Augmentation
one degree of freedom



Complex Posterior Distributions

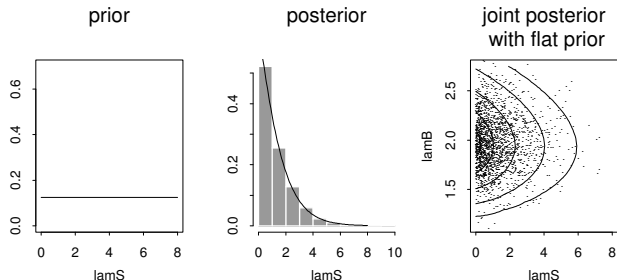


Complex Posterior Distributions



Simulating from the Posterior

- We can *simulate* or *sample* from a distribution to learn about its contours.
- With the sample alone, we can learn about the posterior.
- Here, $Y \sim \text{Poisson}(\lambda_S + \lambda_B)$ and $Y_B \sim \text{Poisson}(c\lambda_B)$.



Using Simulation to Evaluate Integrals

Suppose we want to compute

$$I = \int g(\theta)f(\theta)d\theta,$$

where $f(\theta)$ is a probability density function.

If we have a sample

$$\theta^{(1)}, \dots, \theta^{(n)} \sim f(\theta),$$

we can estimate I with

$$\hat{I}_n = \frac{1}{n} \sum_{i=1}^n g(\theta^{(i)}).$$

In this way we can compute means, variances, and the probabilities of intervals.

We Need to Obtain a Sample

Our primary goal:

Develop methods to obtain a sample from a distribution

- The sample may be independent or dependent.
- Markov chains can be used to obtain a dependent sample.
- In a Bayesian context, we typically aim to sample the *posterior* distribution.

*We first discuss independent methods:
Rejection Sampling & The Grid Method*

Rejection Sampling

Suppose we cannot sample $f(\theta)$ directly, but can find $g(\theta)$ with

$$f(\theta) \leq Mg(\theta)$$

for some M .

- 1 Sample $\tilde{\theta} \sim g(\theta)$.
- 2 Sample $u \sim \text{Unif}(0, 1)$.
- 3 If

$$u \leq \frac{f(\tilde{\theta})}{Mg(\tilde{\theta})}, \text{ i.e., if } uMg(\tilde{\theta}) \leq f(\tilde{\theta})$$

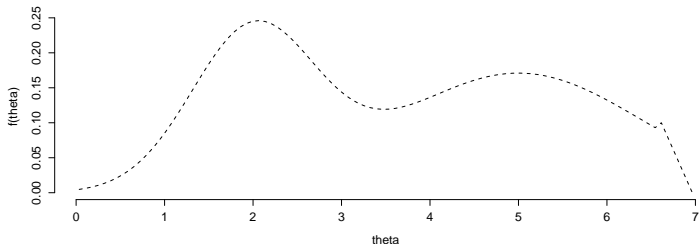
accept $\tilde{\theta}$: $\theta^{(t)} = \tilde{\theta}$.

Otherwise reject $\tilde{\theta}$ and return to step 1.

How do we compute M ?

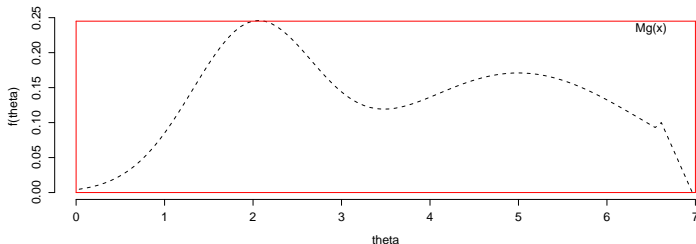
Rejection Sampling

Consider the distribution:



We must bound $f(\theta)$ with some unnormalized density, $Mg(\theta)$.

Rejection Sampling



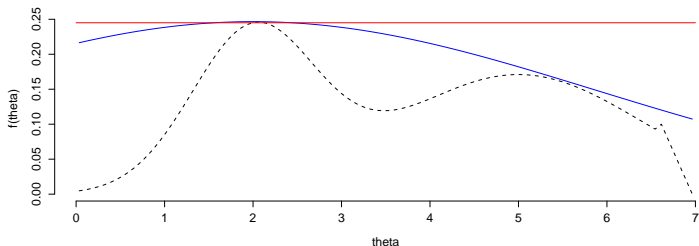
- Imagine that we sample uniformly in the red rectangle:

$$\theta \sim g(\theta) \text{ and } y = uMg(\theta)$$

- Accept samples that fall below the dashed density function.

How can we reduce the wait for acceptance??

Rejection Sampling

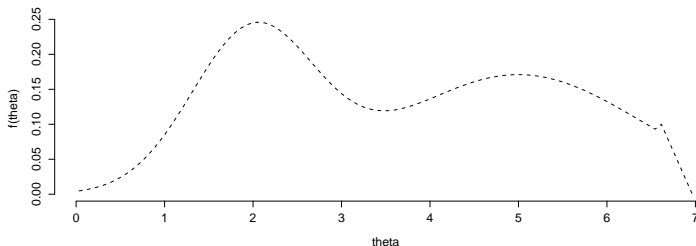


How can we reduce the wait for acceptance??

Improve $g(\theta)$ as an approximation to $f(\theta)$!!

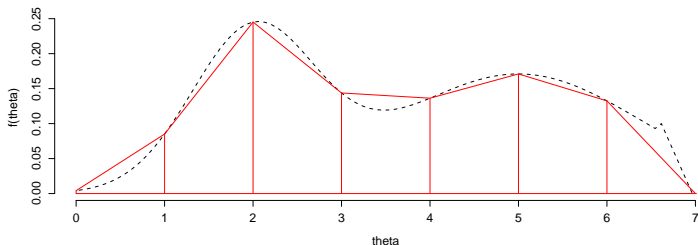
The Grid Method

The Grid method is a brute force / last resort method to sample from a density:



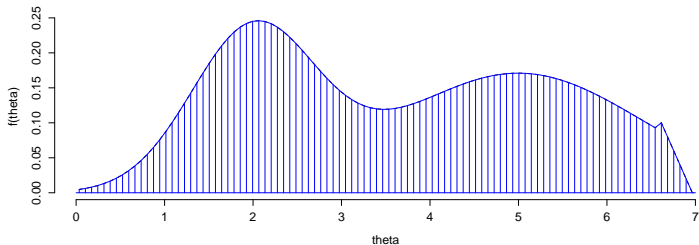
The Grid Method

- 1 Evaluate the density on a grid.
- 2 Compute the areas of the resulting trapezoids.
- 3 Sample from a multinomial distribution with probabilities proportional to the areas.



How can we improve the approximation??

The Grid Method



How can we improve the approximation??

Use a finer grid!!

Limitations?

What is a Markov Chain

Definition

A Markov chain is a sequence of random variables,

$$\theta^{(0)}, \theta^{(1)}, \theta^{(2)}, \dots$$

such that

$$p(\theta^{(t)} | \theta^{(t-1)}, \theta^{(t-2)}, \dots, \theta^{(0)}) = p(\theta^{(t)} | \theta^{(t-1)}).$$

A Markov chain is generally constructed via

$$\theta^{(t)} = \varphi(\theta^{(t-1)}, U^{(t-1)})$$

with $U^{(1)}, U^{(2)}, \dots$ independent.

What is a Stationary Distribution?

Definition

A stationary distribution is any distribution $f(x)$ such that

$$f(\theta^{(t)}) = \int p(\theta^{(t)} | \theta^{(t-1)}) f(\theta^{(t-1)}) d\theta^{(t-1)}$$

If we

- 1 have a sample from the stationary dist'n and
- 2 update the Markov chain,

then the next iterate also follows the stationary dist'n.

In practice we cannot obtain even one sample for the stationary dist'n.

What does a Markov Chain at Stationarity Deliver?

Under regularity conditions, the density at iteration t ,

$$f^{(t)}(\theta|\theta^{(0)}) \rightarrow f(\theta) \quad \text{and} \quad \frac{1}{n} \sum_{t=1}^n h(\theta^{(t)}) \rightarrow E_f[h(\theta)]$$

- The Markov chain converges to its stationary distribution.
- After sufficient burn-in, we treat $\{\theta^{(t)}, t = N_0, \dots, N\}$ as a *correlated* sample from the stationary distribution.
- This is an *approximation*: Use MCMC samples with care!
- Convergence diagnostics are critical.

We aim to find a Markov Chain with Stationary Dist'n equal to the Target Dist'n.

Outline

- 1 Background
 - Complex Posterior Distributions
 - Monte Carlo Integration
 - Markov Chains
- 2 **Basic MCMC Jumping Rules**
 - Metropolis Sampler
 - Metropolis Hastings Sampler
 - Basic Theory
- 3 Practical Challenges and Advice
 - Diagnosing Convergence
 - Choosing a Jumping Rule
 - Transformations and Multiple Modes
- 4 The Gibbs Sampler and Data Augmentation
 - The Gibbs Sampler
 - Data Augmentation

The Metropolis Sampler

Draw $\theta^{(0)}$ from some starting distribution.

For $t = 1, 2, 3, \dots$

Sample: θ^* from $J_t(\theta^*|\theta^{(t-1)})$

Compute: $r = \frac{p(\theta^*|y)}{p(\theta^{(t-1)}|y)}$

Set: $\theta^{(t)} = \begin{cases} \theta^* & \text{with probability } \min(r, 1) \\ \theta^{(t-1)} & \text{otherwise} \end{cases}$

Note

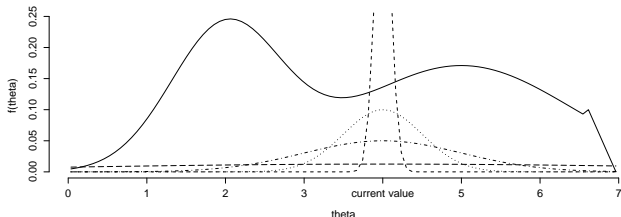
- J_t must be symmetric: $J_t(\theta^*|\theta^{(t-1)}) = J_t(\theta^{(t-1)}|\theta^*)$.
- If $p(\theta^*|y) > p(\theta^{(t-1)}|y)$, *jump!*

The Random Walk Jumping Rule

Typical choices of $J_t(\theta^*|\theta^{(t-1)})$ include

- Unif $(\theta^{(t-1)} - k, \theta^{(t-1)} + k)$
- Normal $(\theta^{(t-1)}, kl)$
- $t_{df}(\theta^{(t-1)}, kl)$

J_t may change, but may not depend on the history of the chain.



How should we choose k ? Replace l with M ? How?

An Example

A simplified model for high-energy spectral analysis.

- Model:

Consider a perfect detector:

- 1 1000 energy bins, equally spaced from 0.3keV to 7.0keV,
- 2 $Y_i \sim \text{Poisson}(\alpha E_i^{-\beta})$, with $\theta = (\alpha, \beta)$,
- 3 E_i is the energy, and
- 4 $(\alpha, \beta) \stackrel{\text{indep.}}{\sim} \text{Unif}(0, 100)$.

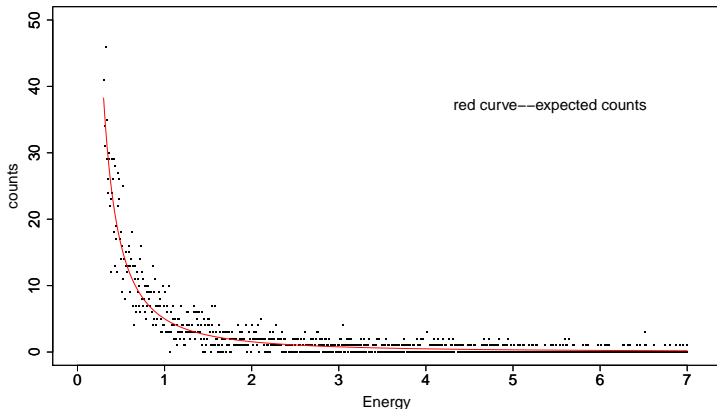
- The Sampler:

We use a Gaussian Jumping Rule,

- centered at the current sample, $\theta^{(t)}$
- with standard deviations equal 0.08 and correlation zero.

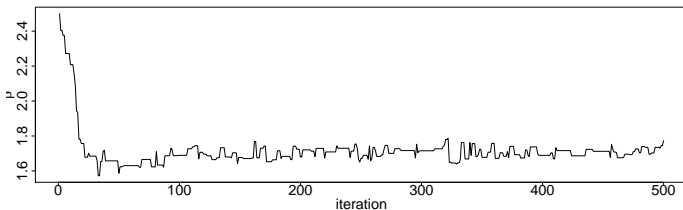
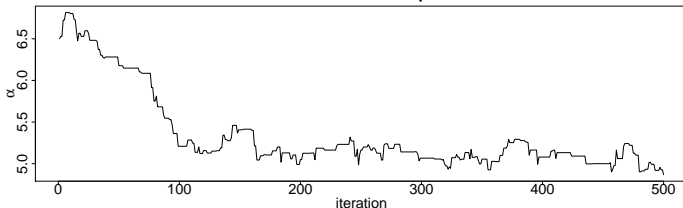
Simulated Data

2288 counts were simulated with $\alpha = 5.0$ and $\beta = 1.69$.



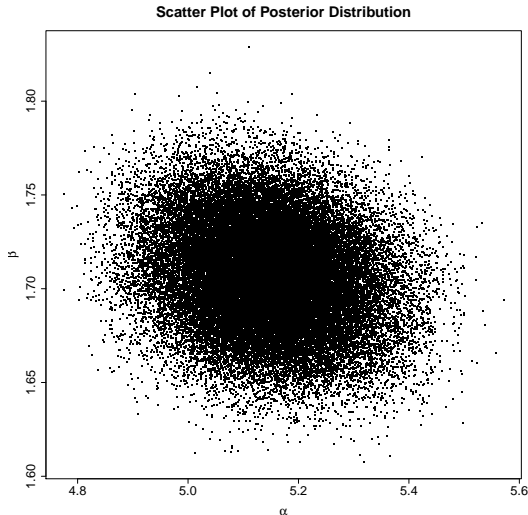
Markov Chain Trace Plots

Time Series Plot for Metropolis Draws



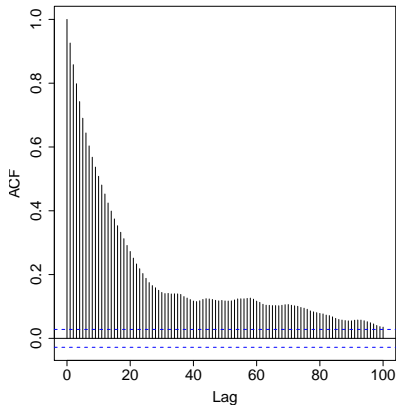
Chains “stick” at a particular draw when proposals are rejected.

The Joint Posterior Distribution

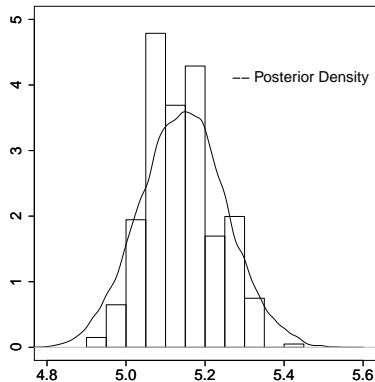


Marginal Posterior Dist'n of the Normalization

Autocorrelation for alpha



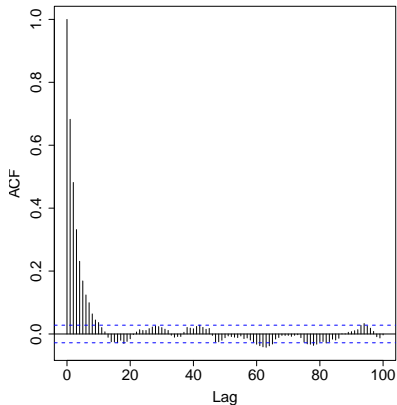
Hist of 500 Draws excluding Burn-in



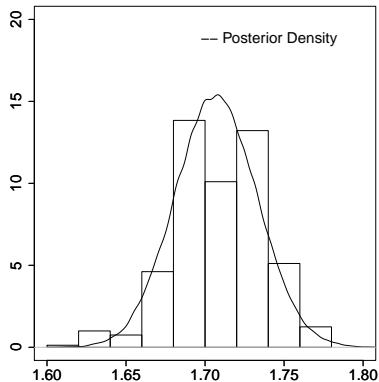
$E(\alpha|Y) \approx 5.13$, $SD(\alpha|Y) \approx 0.11$, and a 95% CI is (4.92, 5.41)

Marginal Posterior Dist'n of Power Law Param

Autocorrelation for beta



Hist of 500 Draws excluding Burn-in



$E(\beta|Y) \approx 1.71$, $SD(\beta|Y) \approx 0.03$, and a 95% CI is (1.65, 1.76)

The Metropolis-Hastings Sampler

A more general Jumping rule:

Draw $\theta^{(0)}$ from some starting distribution.

For $t = 1, 2, 3, \dots$

Sample: θ^* from $J_t(\theta^* | \theta^{(t-1)})$

Compute: $r = \frac{p(\theta^* | y) / J_t(\theta^* | \theta^{(t-1)})}{p(\theta^{(t-1)} | y) / J_t(\theta^{(t-1)} | \theta^*)}$

Set: $\theta^{(t)} = \begin{cases} \theta^* & \text{with probability } \min(r, 1) \\ \theta^{(t-1)} & \text{otherwise} \end{cases}$

Note

- J_t may be any jumping rule, it needn't be symmetric.
- The updated r corrects for bias in the jumping rule.

The Independence Sampler

Use an approximation to the posterior as the jumping rule:

$J_t = \text{Normal}_d(\text{MAP estimate, Curvature-based Variance Matrix}).$

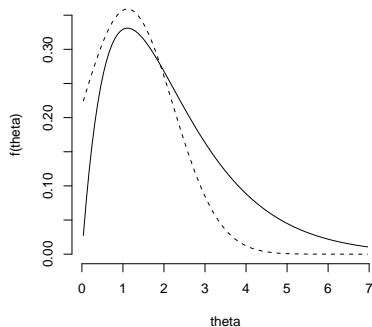
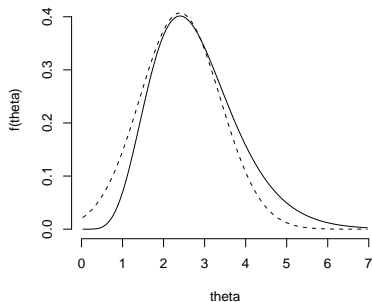
MAP estimate = $\text{argmax}_{\theta} p(\theta|y)$

$$\text{Variance} \approx \left[-\frac{\partial^2}{\partial \theta \cdot \partial \theta} \log p(\theta|Y) \right]^{-1}$$

Note: $J_t(\theta^* | \theta^{(t-1)})$ does not depend on $\theta^{(t-1)}$.

The Independence Sampler

The Normal Approximation may not be adequate.



- We can inflate the variance.
- We can use a heavy tailed distribution, e.g., lorentzian or t .

Example of Independence Sampler

A simplified model for high-energy spectral analysis.

- We use the same model and simulated data.
- This is a simple *loglinear model*, a special case of a *Generalized Linear Model*:

$$Y_i \sim \text{Poisson}(\lambda_i) \quad \text{with} \quad \log(\lambda_i) = \log(\alpha) - \beta \log(E_i).$$

- The model can be fit with the `glm` function in R:

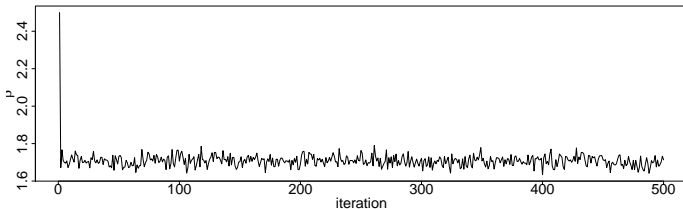
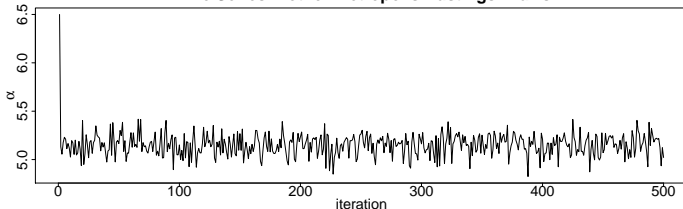

```
> glm.fit = glm( Y~I(-log(E)), family=poisson(link="log") )
> glm.fit$coef      #### best fit of (log(alpha), beta)
> vcov( glm.fit )  #### variance-covariance matrix
```
- Returns MLE of $(\log(\alpha), \beta)$ and variance-covariance matrix.

Example of Independence Sampler

- Alternatively, we can fit (α, β) directly with a general (but less stable) mode finder.
- Requires coding likelihood, specifying starting values, etc.
- Choose parameterization to improve Gaussian approx.
 - MLE is invariant to transformations.
 - Variance matrix of transform is computed via *delta method*.
- We use the general mode finder:
 $J_t = \text{Normal}_2(\text{MAP est}, \text{Curvature-based Variance Matrix}).$

Markov Chain Trace Plots

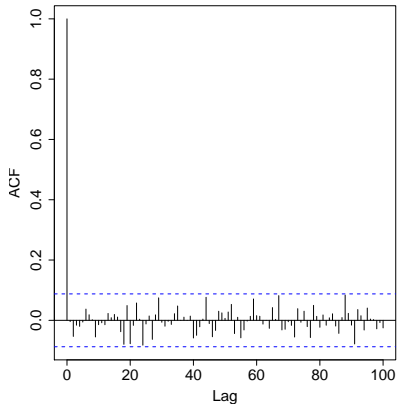
Time Series Plot for Metropolis Hastings Draws



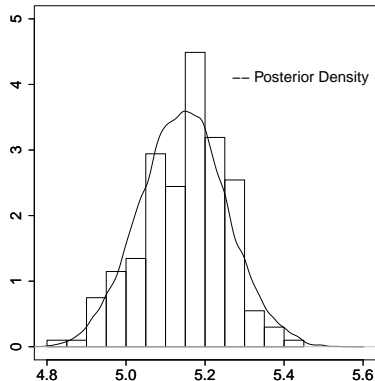
Very little “sticking” here: acceptance rate is 98.8%.

Marginal Posterior Dist'n of the Normalization

Autocorrelation for alpha



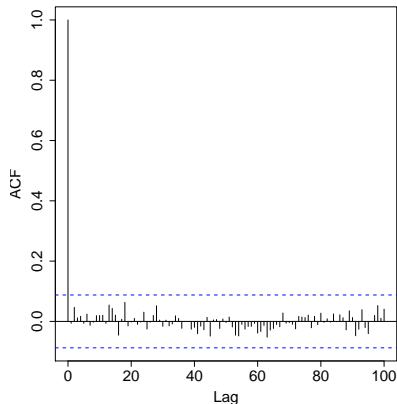
Hist of 500 Draws excluding Burn-in



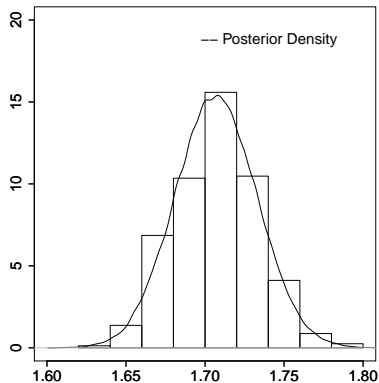
Autocorrelation is essentially zero: nearly independent sample!!

Marginal Posterior Dist'n of Power Law Param

Autocorrelation for beta



Hist of 500 Draws excluding Burn-in



This result depends critically on access to a very good approximation to the posterior distribution.

Convergence to Stationarity

Consider a finite state space \mathcal{S} with arbitrary elements i and j .

- Let $p_{ij}(t) = \Pr(\theta^{(t)} = j | \theta^{(0)} = i)$.
- Ergodic Theorem: If a Markov chain is *positive recurrent* and *aperiodic* then its stationary distribution is the unique distribution $\pi()$ such that

$$\sum_i p_{ij}(t)\pi(i) = \pi(j) \text{ for all } j \text{ and } t \geq 0.$$

We say the Markov chain is ergodic and the following hold:

- 1 $p_{ij}(t) \rightarrow \pi(j)$ as $t \rightarrow \infty$ for all i and j .
- 2

$$\Pr \left[\frac{1}{n} \sum_{t=1}^n h(\theta^{(t)}) \rightarrow E_{\pi}(h(\theta)) \right] = 1$$

Convergence to Stationarity

Definitions:

- 1 Chain is *irreducible* if for all i, j there is t with $p_{ij}(t) > 0$.

Let τ_{ij} be the time of first return, $\min\{t > 0 : \theta^{(t)} = i | \theta^{(0)} = i\}$.

- 2 Chain is *recurrent* if $\Pr[\tau_{ij} < \infty] = 1$ for all i .
- 3 Chain is *positive recurrent* if $E[\tau_{ij}] < \infty$ for all i .

Fact: Irreducible chain with a stationary dist'n is pos recurrent.

So we need our chain to

- 1 be irreducible,
- 2 be aperiodic, and
- 3 have the posterior distribution as a stationary distribution.

Outline

- 1 Background
 - Complex Posterior Distributions
 - Monte Carlo Integration
 - Markov Chains
- 2 Basic MCMC Jumping Rules
 - Metropolis Sampler
 - Metropolis Hastings Sampler
 - Basic Theory
- 3 Practical Challenges and Advice
 - Diagnosing Convergence
 - Choosing a Jumping Rule
 - Transformations and Multiple Modes
- 4 The Gibbs Sampler and Data Augmentation
 - The Gibbs Sampler
 - Data Augmentation

Has this Chain Converged?

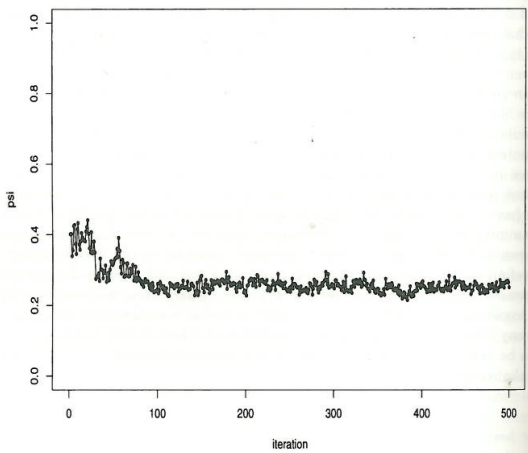


Image credit: Gelman (1995) In "MCMC in Practice" (Editors: Gilks, Richardson, and Spiegelhalter).

Has this Chain Converged?

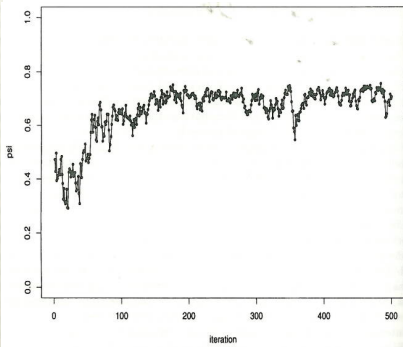
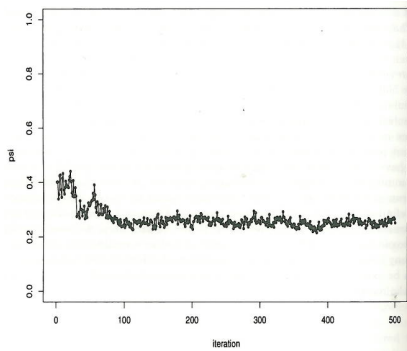
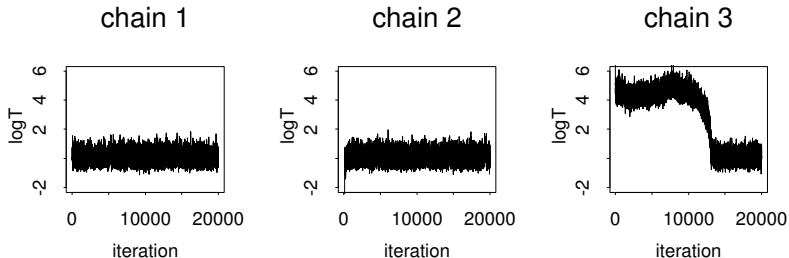


Image credit: Gelman (1995) In "MCMC in Practice" (Editors: Gilks, Richardson, and Spiegelhalter).

Comparing multiple chains can be informative!

Using Multiple Chains



- Compare results of multiple chains to check convergence.
- Start the chains from distant points in parameter space.
- Run until they appear to give similar results
 - ... or they find different solutions (multiple modes).

The Gelman and Rubin “R hat” Statistic

Consider M chains of length N : $\{\psi_{nm}, n = 1, \dots, N\}$.

$$B = \frac{N}{M-1} \sum_{m=1}^M (\bar{\psi}_{\cdot m} - \bar{\psi}_{\cdot \cdot})^2$$

$$W = \frac{1}{M} \sum_{m=1}^M s_m^2 \quad \text{where} \quad s_m^2 = \frac{1}{N-1} \sum_{n=1}^N (\psi_{nm} - \bar{\psi}_{\cdot m})^2$$

Two estimates of $\text{Var}(\psi | Y)$:

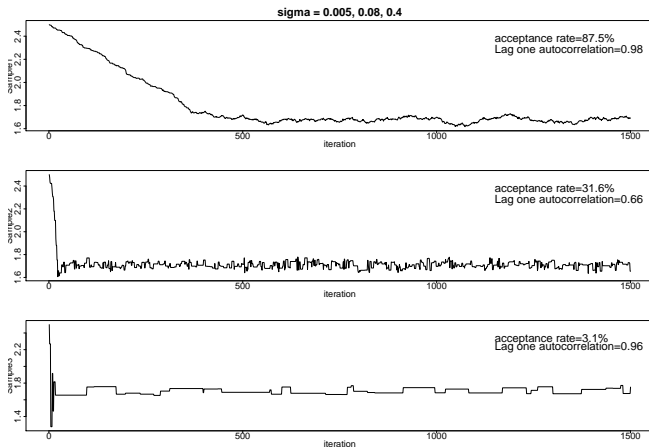
- 1 W : under estimate of $\text{Var}(\psi | Y)$ for any finite N .
- 2 $\widehat{\text{var}}^+(\psi | Y) = \frac{N-1}{N} W + \frac{1}{N} B$: over estimate of $\text{Var}(\psi | Y)$.

$$\hat{R} = \sqrt{\frac{\widehat{\text{var}}^+(\psi | Y)}{W}} \downarrow 1 \quad \text{as the chains converge.}$$

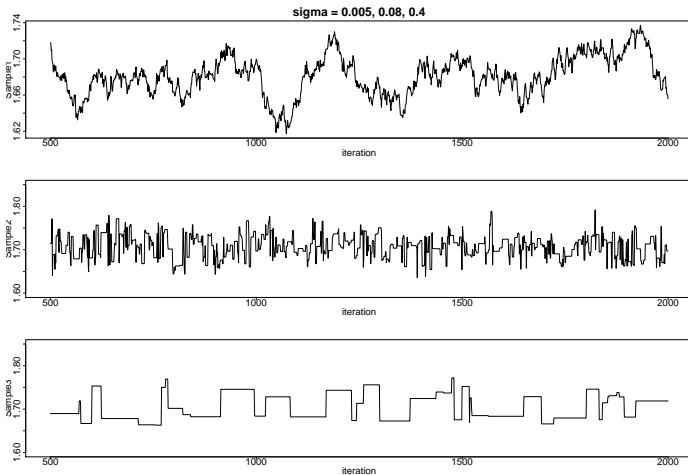
Compute with `coda` package in R: <http://cran.r-project.org/web/packages/coda/index.html>

Choice of Jumping Rule with Random Walk Metropolis

Spectral Analysis: effect on burn in of power law parameter



Higher Acceptance Rate is not Always Better!



Aim for 20% (vectors) - 40% (scalars) acceptance rate

Statistical Inference and Effective Sample Size

• Point Estimate: $\bar{h}_n = \frac{1}{n} \sum h(\theta^{(t)})$ (estimate of $E(h(\theta)|x)$!!)

• Variance Estimate: $\text{Var}(\bar{h}_n) \approx \frac{\sigma^2}{n} \frac{1+\rho}{1-\rho}$ with (not $\text{var}(\theta)$!!)

$\sigma^2 = \text{Var}(h(\theta))$ estimated by $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{t=1}^n [h(\theta^{(t)}) - \bar{h}_n]^2$,

$\rho = \text{corr}[h(\theta^{(t)}), h(\theta^{(t-1)})]$ estimated by

$$\hat{\rho} = \frac{1}{n-1} \frac{\sum_{t=2}^n [h(\theta^{(t)}) - \bar{h}_n][h(\theta^{(t-1)}) - \bar{h}_n]}{\sqrt{\sum_{t=1}^{n-1} [h(\theta^{(t)}) - \bar{h}_n]^2 \sum_{t=2}^n [h(\theta^{(t)}) - \bar{h}_n]^2}}$$

• Interval Estimate: $\bar{h}_n \pm t_d \sqrt{\text{Var}(\bar{h}_n)}$ with $d = n \frac{1-\rho}{1+\rho} - 1$

The *effective sample size* is $n \frac{1-\rho}{1+\rho} \dots$

...all computed with `coda` in R.

Illustration of the Effective Sample Size

Sample from $N(0, 1)$

with random walk Metropolis with $J_t = N(\theta^{(t)}, \sigma)$.

What is the Effective Sample Size here? and σ ?

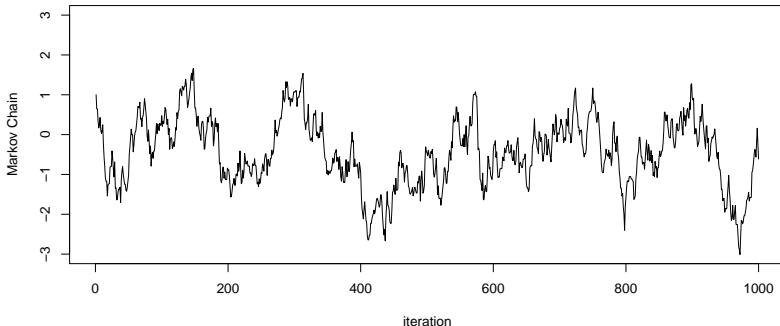


Illustration of the Effective Sample Size

What is the Effective Sample Size here? and σ ?

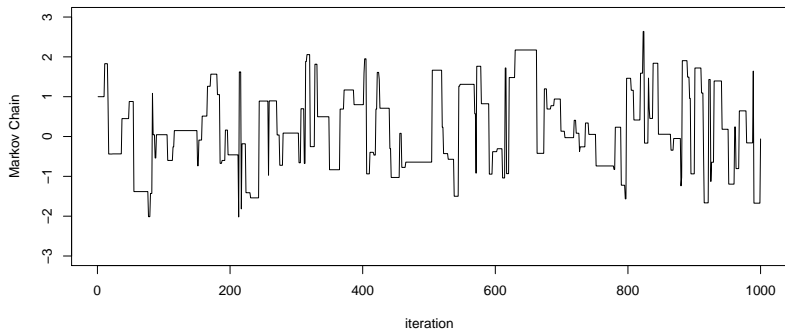


Illustration of the Effective Sample Size

What is the Effective Sample Size here? and σ ?

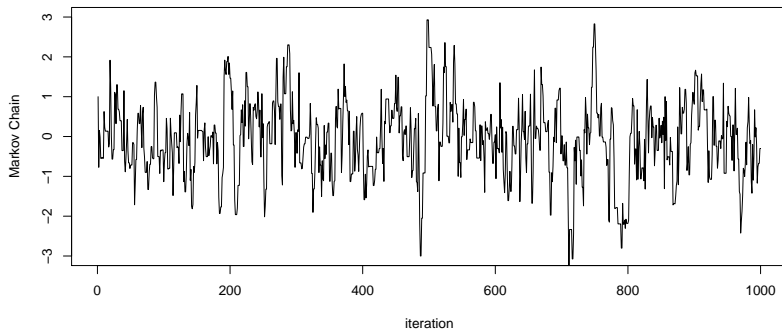


Illustration of the Effective Sample Size

What is the Effective Sample Size here? and σ ?

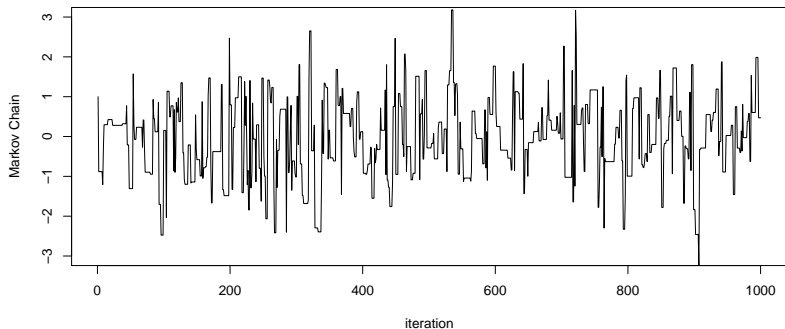


Illustration of the Effective Sample Size

What is the Effective Sample Size here? and σ ?

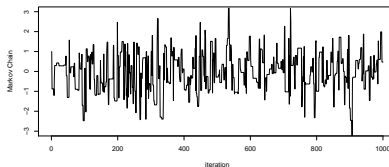
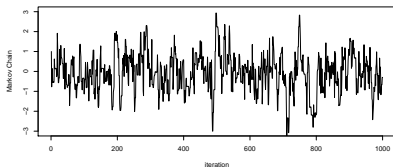
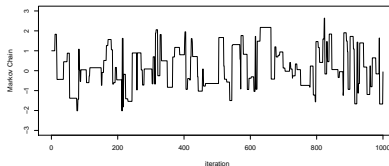
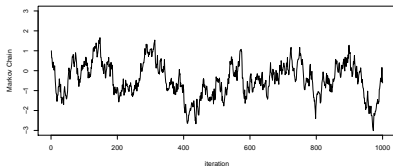
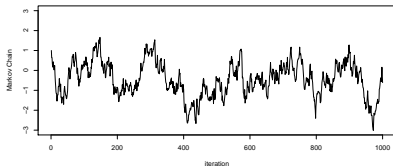
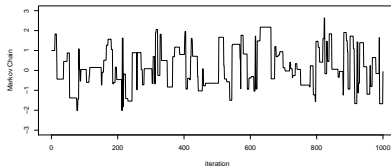


Illustration of the Effective Sample Size

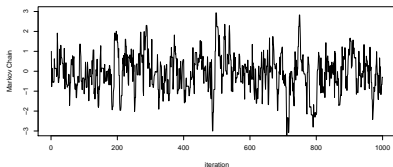
Effective Sample = 20; $\sigma = 0.25$.



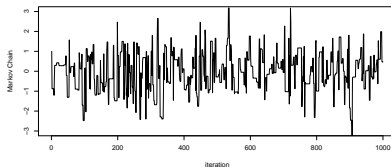
Effective Sample = 75; $\sigma = 0.10$.



Effective Sample = 100; $\sigma = 1$.

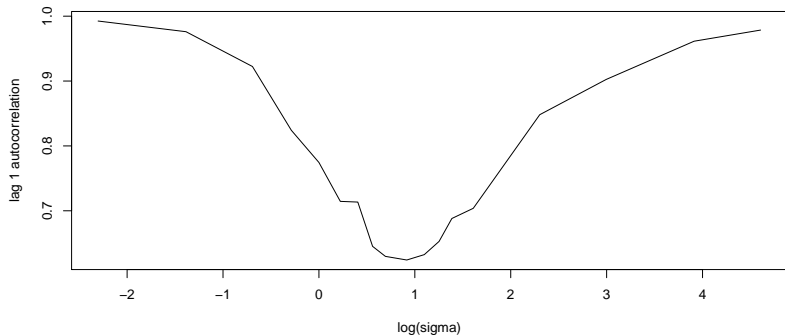


Effective Sample = 216; $\sigma = 3.5$.



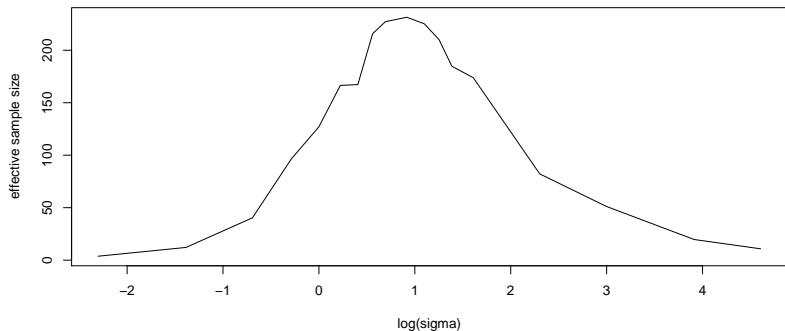
Lag One Autocorrelation

Small Jumps versus Low Acceptance Rates



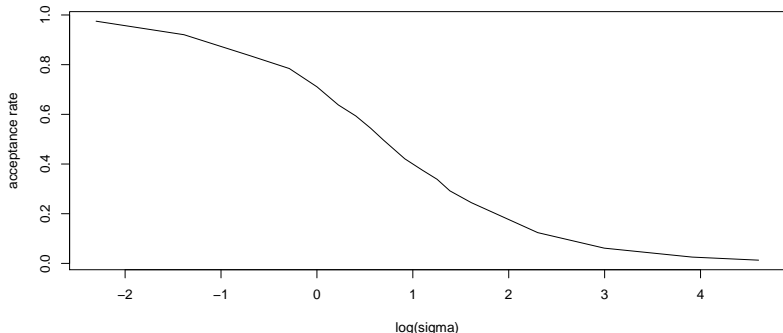
Effective Sample Size

Balancing the Trade-Off



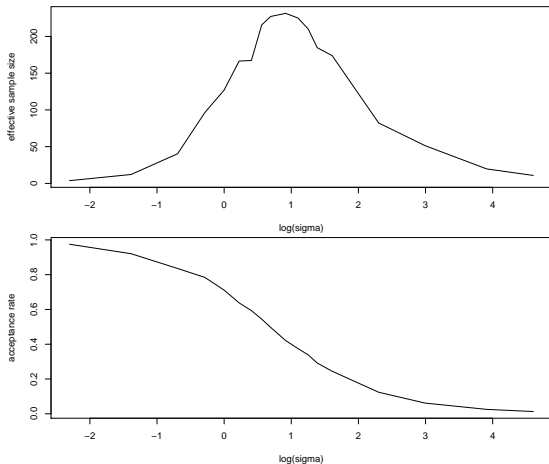
Acceptance Rate

Bigger is not always Better!!



High acceptance rates only come with small steps!!

Finding the Optimal Acceptance Rate

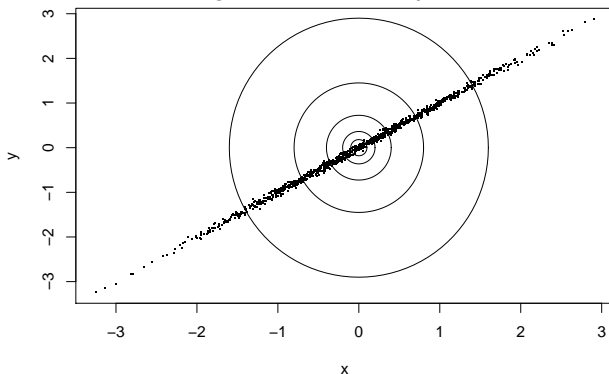


Random Walk Metropolis with High Correlation

A whole new set of issues arise in higher dimensions...

Tradeoff between high autocorrelation and high rejection rate:

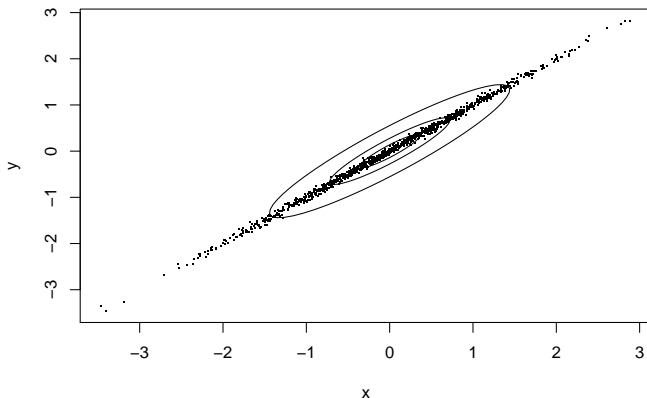
- more acute with high posterior correlations
- more acute with high dimensional parameter



Random Walk Metropolis with High Correlation

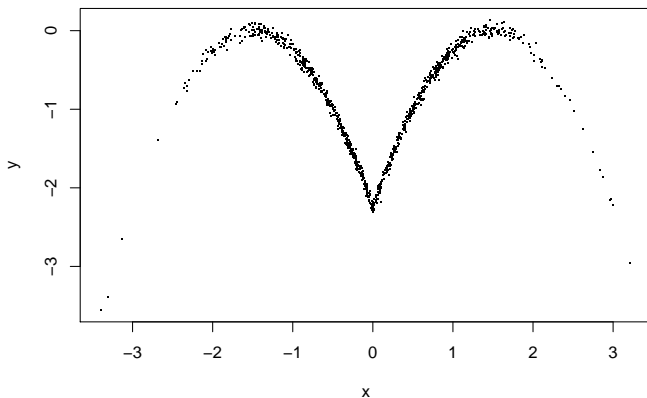
In principle we can use a correlated jumping rule, but

- the desired correlation may vary, and
- is often difficult to compute in advance.



Random Walk Metropolis with High Correlation

What random walk jumping rule would you use here?

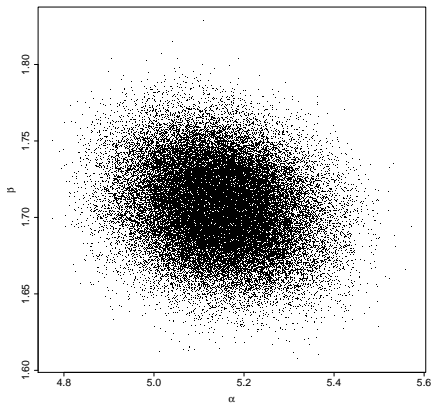


Remember: you don't get to see the distribution in advance!

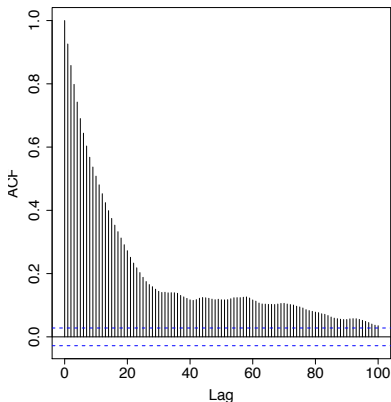
Parameters on Different Scales

Random Walk Metropolis for Spectral Analysis:

Scatter Plot of Posterior Distribution



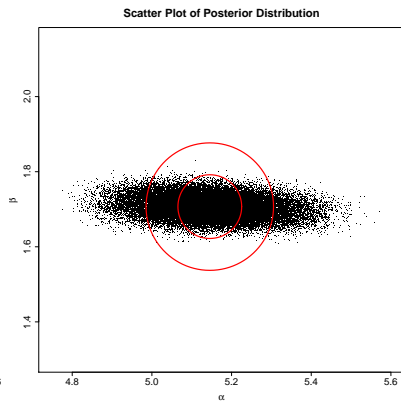
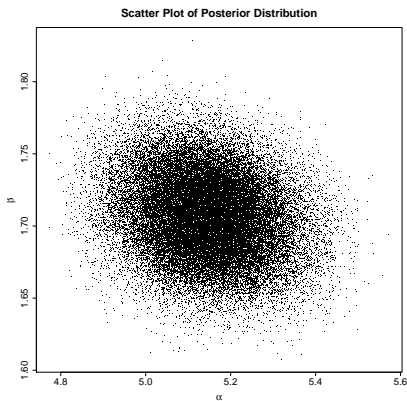
Autocorrelation for alpha



Why is the Mixing SO Poor?!??

Parameters on Different Scales

Consider the Scales of α and β :

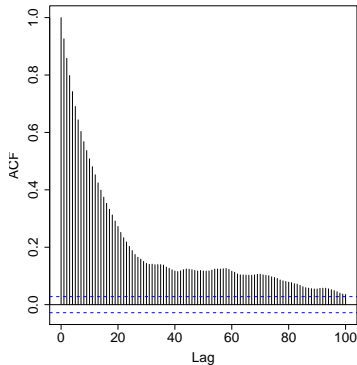


A new jumping rule: std dev for $\alpha = 0.110$, for $\beta = 0.026$, and $\text{corr} = -0.216$.

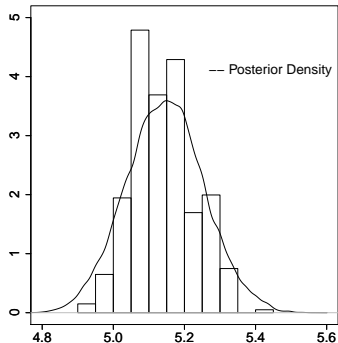
Improved Convergence

Original Jumping Rule:

Autocorrelation for alpha

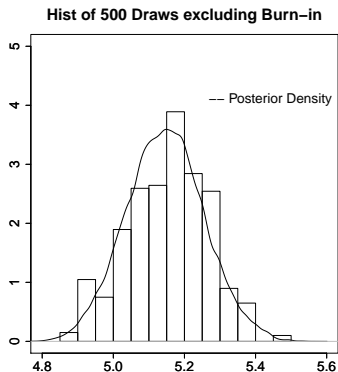
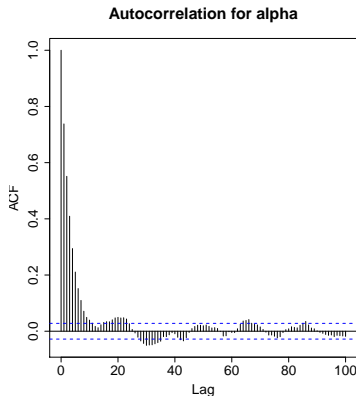


Hist of 500 Draws excluding Burn-in



Improved Convergence

Improved Jumping Rule:



Original Eff Sample Size = 19, Improved Eff Sample Size = 75, with $n = 500$.

Parameters on Different Scales

With Jumping Rule: $\text{NORM}(\theta^{(t-1)}, kM)$, or better $t_{\text{df}}(\theta^{(t-1)}, kM)$.

Try:

- 1 Using the variance-covariance matrix from a standard fitted model for M
... at least when standard mode-based model-fitting software is available.
- 2 New adaptive methods that allow the jumping rule to evolve on the fly.¹

Always: Aim for acceptance rate of

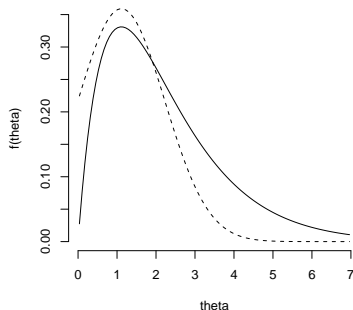
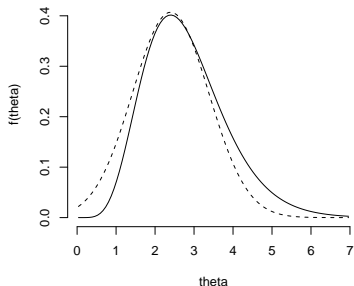
~20% (multivariate update) or ~40% (univariate update).

¹ E.g., "Optimal proposal distributions and adaptive MCMC" by JS Rosenthal in Handbook of Markov Chain Monte Carlo (CRC Press, 2011).

Transforming to Normality

Parameter transformations can greatly improve MCMC.

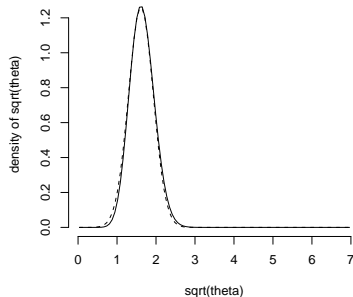
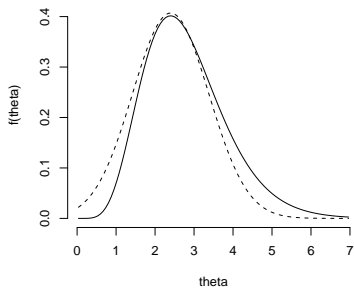
Recall the Independence Sampler:



The normal approximation is not as good as we might hope...

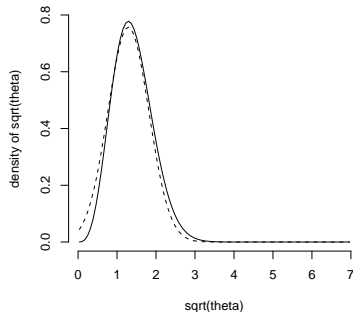
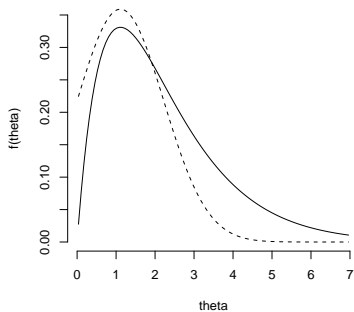
Transforming to Normality

But if we use the square root of θ :



Transforming to Normality

And...



The normal approximation is much improved!

Transforming to Normality

Working with with Gaussian or symmetric distributions leads to more efficient Metropolis and Metropolis Hastings Samplers.

General Strategy:

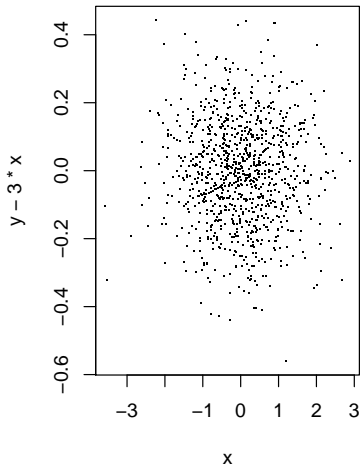
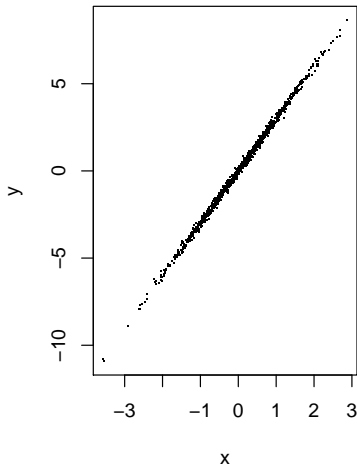
- Transform to the Real Line.
- Take the log of positive parameters.
- If the log is “too strong”, try square root.
- Probabilities can be transformed via the logit transform:

$$\log(p/(1 - p)).$$

- More complex transformations for other quantities.
- *Try out various transformations using an initial MCMC run.*
- Statistical advantages to using normalizing transforms.

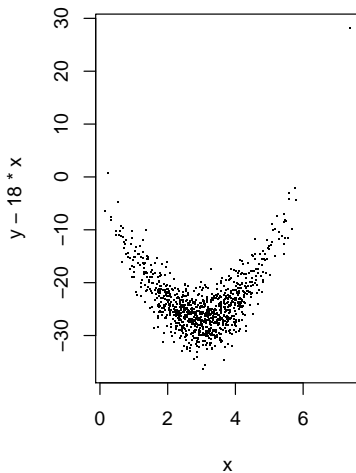
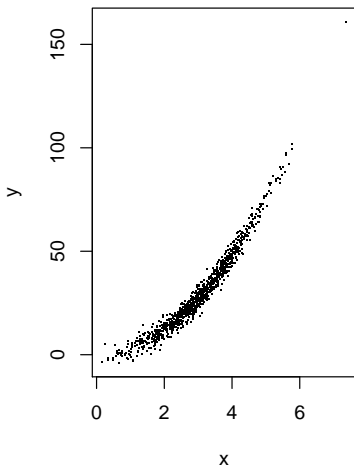
Removing Linear Correlations

Linear transformations can remove linear correlations



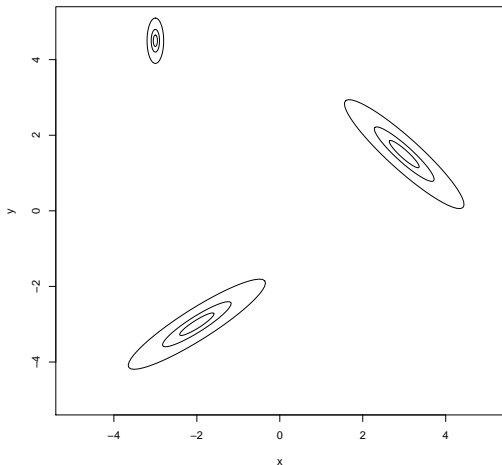
Removing Linear Correlations

... and can help with non-linear correlations.



Multiple Modes

- Scientific meaning of multiple modes.
- Do not focus only on the major mode!
- “Important” modes.
- Challenging for Bayesian and Frequentist methods.
- Consider Metropolis & Metropolis Hastings.
- Value of excess dispersion and multiple starting values.



Multiple Modes

- 1 Use a mode finder to “map out” the posterior distribution.
 - 1 Design a jumping rule that accounts for all of the modes.
 - 2 Run separate chains for each mode.
- 2 Use one of several sophisticated methods tailored for multiple modes.
 - 1 Adaptive Metropolis Hastings. Jumping rule adapts when new modes are found (van Dyk & Park, MCMC Hdbk 2011).
 - 2 Parallel Tempering.
 - 3 Nested Sampling (Skilling, 2006, *Bayesian Analysis*)
 - 4 Many other specialized methods.

Outline

- 1 Background
 - Complex Posterior Distributions
 - Monte Carlo Integration
 - Markov Chains
- 2 Basic MCMC Jumping Rules
 - Metropolis Sampler
 - Metropolis Hastings Sampler
 - Basic Theory
- 3 Practical Challenges and Advice
 - Diagnosing Convergence
 - Choosing a Jumping Rule
 - Transformations and Multiple Modes
- 4 The Gibbs Sampler and Data Augmentation
 - The Gibbs Sampler
 - Data Augmentation

Breaking a Complex Problem into Simpler Pieces

- Ideally we sample directly from $p(\theta|Y)$ without Metropolis.
- This may not work in complex problems.
- **BUT** in some cases we can split $\theta = (\theta_1, \theta_2)$ so that

$$p(\theta_1|\theta_2, Y) \text{ and } p(\theta_2|\theta_1, Y)$$

are both easy to sample although $p(\theta|Y)$ is not.

- The *Two-Step Gibbs Sampler*, starting with some $\theta^{(0)}$,

For $t = 1, 2, 3, \dots$

Draw: $\theta_1^{(t)} \sim p(\theta_1|\theta_2^{(t-1)}, Y)$

Draw: $\theta_2^{(t)} \sim p(\theta_2|\theta_1^{(t)}, Y)$

An Example

Recall Simple Spectral Model: $Y_i \sim \text{Poisson}(\alpha E_i^{-\beta})$.

Using $p(\alpha, \beta) \propto 1$,

$$\begin{aligned} p(\theta|Y) &\propto \prod_{i=1}^n e^{-[\alpha E_i^{-\beta}]} [\alpha E_i^{-\beta}]^{Y_i} \\ &= e^{-\alpha \sum_{i=1}^n E_i^{-\beta}} \alpha^{\sum_{i=1}^n Y_i} \prod_{i=1}^n E_i^{-\beta Y_i} \end{aligned}$$

So that

$$\begin{aligned} p(\alpha|\beta, Y) &\propto e^{-\alpha \sum_{i=1}^n E_i^{-\beta}} \alpha^{\sum_{i=1}^n Y_i} \\ &= \text{Gamma} \left(\sum_{i=1}^n Y_i + 1, \sum_{i=1}^n E_i^{-\beta} \right) \end{aligned}$$

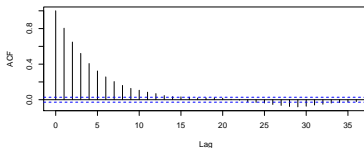
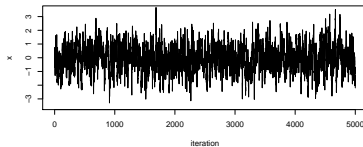
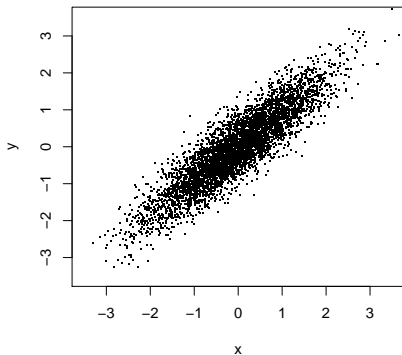
Embedding Other Samplers within Gibbs

In this case $p(\beta|\alpha, Y)$ is not a standard distribution:

$$p(\beta|\alpha, Y) \propto e^{-\alpha \sum_{i=1}^n E_i^{-\beta}} \prod_{i=1}^n E_i^{-\beta Y_i}$$

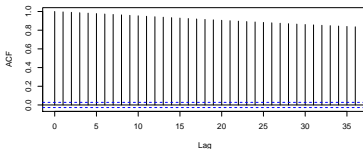
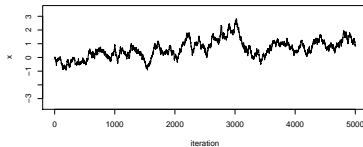
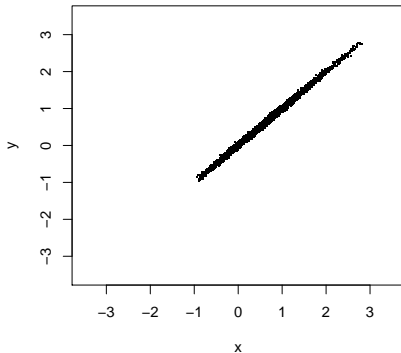
- We can use a Metropolis or Metropolis-Hastings step to update β within the Gibbs sampler.
- The result is known as Metropolis within Gibbs Sampler.
- **Advantage:** Metropolis tends to perform poorly in high dimensions. Gibbs reduces the dimension.
- **Disadvantage:** Case-by-case probabilistic calculations.
(But always need case-by-case algorithmic development and tuning.)

When Will Gibbs Sampling Work Well?



autocorrelation = 0.81, effective sample size = 525

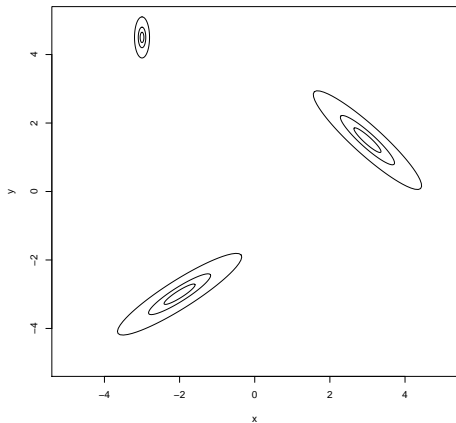
When Will Gibbs Sampling Work Poorly?



autocorrelation = 0.998, effective sample size = 5

High Posterior Correlations are Always Problematic.

Multiple Modes



How will the Gibbs Sampler Handle Multiple modes?

The General Gibbs Sampler

- 1 In general we break θ into P subvectors $\theta = (\theta_1, \dots, \theta_P)$.
- 2 The Complete Conditional Distributions are given by

$$p(\theta_p | \theta_1, \dots, \theta_{p-1}, \theta_{p+1}, \dots, \theta_P, Y), \text{ for } p = 1, \dots, P$$

- 3 The *Gibbs Sampler*, starting with some $\theta^{(0)}$,

For $t = 1, 2, 3, \dots$

Draw 1: $\theta_1^{(t)} \sim p(\theta_1 | \theta_2^{(t-1)}, \dots, \theta_P^{(t-1)}, Y)$

\vdots

Draw p : $\theta_p^{(t)} \sim p(\theta_p | \theta_1^{(t)}, \dots, \theta_{p-1}^{(t)}, \theta_{p+1}^{(t-1)}, \dots, \theta_P^{(t-1)}, Y)$

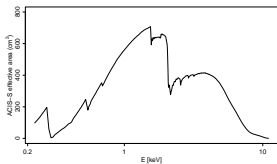
\vdots

Draw P : $\theta_P^{(t)} \sim p(\theta_P | \theta_1^{(t)}, \dots, \theta_{P-1}^{(t)}, Y)$

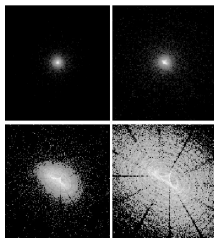
- 4 Determining the partition of θ is a matter of skill and art.

Example: Calibration Uncertainty in High Energy Astrophysics

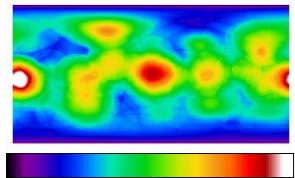
- Analysis is highly dependent on *Calibration Products*:
 - Effective area records sensitivity as a function of energy
 - Energy redistribution matrix can vary with energy/location
 - Point Spread Functions can vary with energy and location
 - Exposure Map shows how effective area varies in an image



A CHANDRA effective area.



Sample Chandra psf's
(Karovska et al., ADASS X)

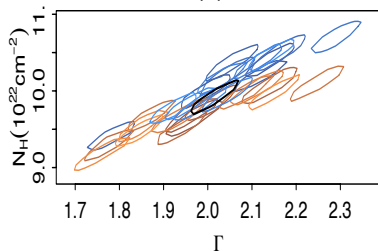
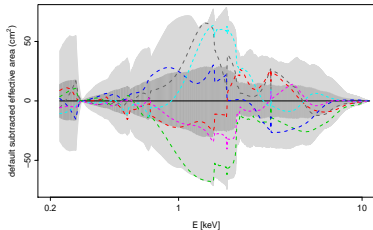


EGERT exposure map
(area \times time)

Example: Calibration Uncertainty

Derivation of Calibration Products

- Prelaunch ground-based and post-launch space-based empirical assessments.
- Aim to capture deterioration of detectors over time.
- Complex computer models of subassembly components.
- Calibration scientists provide a sample representing uncertainty



Example: Calibration Uncertainty

We wish to incorporate uncertainty represented in Calibration sample into a Fully Bayesian Analysis.

- **PyBLoCXS (Python Bayesian Low Count X-ray Spectral)**: provides a MCMC output for spectral analysis with *known* calibration products.
- Can we leverage PyBLoCXS for calibration uncertainty?
- Gibbs Sampler:
 - Draw 1: Update A (effective area) given θ (parameter).
 - Draw 2: Update θ given A with PyBLoCXS.

Power of Gibbs Sampling: breaks a problem into easier parts.

How do we draw A ?

We have only a calibration sample, not a formal model.

We use Principal Component Analysis to represent uncertainty:

$$A \sim A_0 + \bar{\delta} + \sum_{j=1}^m e_j r_j \mathbf{v}_j,$$

A_0 : default effective area,

$\bar{\delta}$: mean deviation from A_0 ,

r_j and \mathbf{v}_j : first m principle component eigenvalues & vectors,

e_j : independent standard normal deviations.

Capture 95% of variability with $m = 6 - 9$.

A Prototype Fully Bayesian Sampler

An MH within Gibbs Sampler:

STEP 1: $e \sim \mathcal{K}(e|e', \theta')$ via MH with limiting dist'n $p(e|\theta, Y)$

STEP 2: $\theta \sim \mathcal{K}(\theta|e', \theta')$ via MH with limiting dist'n $p(\theta|e, Y)$

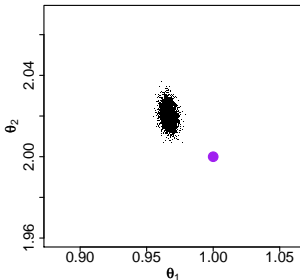
- STEP 1: Gaussian Metropolis jumping rule centered at e' .
- STEP 2: Simplified pyBLoCXS (no rmf or background).

A Simulation.

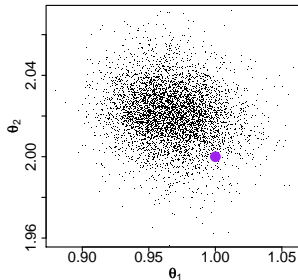
- Sampled 10^5 counts from a power law spectrum: e^{-2E} .
- A_{true} is 1.5σ from the center of the calibration sample.

Sampling From the Full Posterior

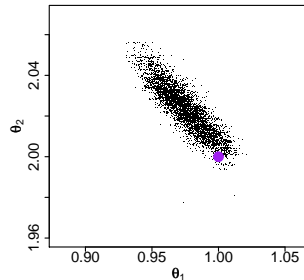
Default Effective Area



Pragmatic Bayes



Fully Bayes



θ_1 = normalization, θ_2 = power law param, purple bullet = truth

Citations:

- 1 Lee, Kashyap, van Dyk, Connors, Drake, Izem, Meng, Min, et al. (2011). Accounting for Calibration Uncertainties in X-ray Analysis: Effective Areas in Spectral Fitting. *The Astrophysical Journal*, **731**, 126–144.
- 2 Xu, van Dyk, Kashyap, Siemiginowska, Connors, Drake, Meng, et al. (2014). A Fully Bayesian for Jointly Fitting Instrumental Calibration and X-ray Spectral Models. *The Astrophysical Journal*, to appear.

Example: Transformations are Key

Fitting Computer Models for Stellar Evolution

- A complex computer model predicts observed *photometric magnitudes* of a stellar cluster as a function of
 - M_j : stellar masses, and
 - Θ : cluster composition, age, distance, and absorption:

$$\mathbf{G}(M_j, \Theta)$$

- We assume indep Gaussian errors with known variances:

$$L_0(\mathbf{M}, \Theta | \mathbf{X}) = \prod_{i=1}^N \left(\prod_{j=1}^n \left[\frac{1}{\sqrt{2\pi\sigma_{ij}^2}} \exp \left(-\frac{(x_{ij} - G_j(M_{i1}, \Theta))^2}{2\sigma_{ij}^2} \right) \right] \right).$$

Example: Stellar Evolution

Model Extensions:

- Binary stars: The luminosities of component stars sum.
- Field stars: Contaminate the data and magnitudes don't follow the pattern of the cluster.
- Initial Final Mass Relation is fit to combine stellar evolution models for the main sequence and for white dwarfs.
- A combination of informative and non-informative priors.

Citations:

- 1 van Dyk, D. A., DeGennaro, S., Stein, N., Jeffreys, W. H., von Hippel, T. Statistical Analysis of Stellar Evolution. *The Annals of Applied Statistics* **3**, 117-143, 2009.
- 2 DeGennaro, S., von Hippel, T., Jefferys, W., Stein, N., van Dyk, D., and Jeffery, E. Inverting Color-Magnitude Diagrams to Access Precise Cluster Parameters: A New White Dwarf Age for the Hyades. *The Astrophysical Journal*, **696**, 12–23, 2009.
- 3 Jeffery, E., von Hippel, T., DeGennaro, S., van Dyk, D., Stein, N., and Jeffreys, W. H., The White Dwarf Age of NGD 2477. *The Astrophysical Journal*, **730**, 35–44, 2011.

Stellar Evolution: MCMC Strategy

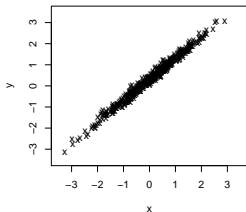
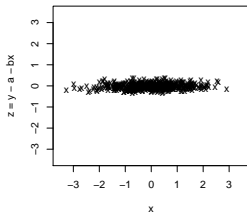
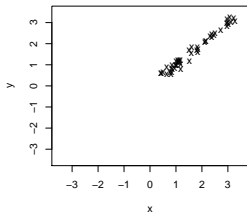
Metropolis within Gibbs Sampling

- $3N + 5$ parameters, none with closed form update.
- Strong posterior correlations among the parameters.

Strong Linear and Non-Linear Correlations Among Parameters

- Static and/or dynamic (power) transformations remove non-linear relationships.
- A series of preliminary runs is used to evaluate and remove linear correlations.
- We tune a linear transformation to the correlations of the posterior distribution on the fly.
- Results in a dramatic improvement in mixing.

Dynamic transformations

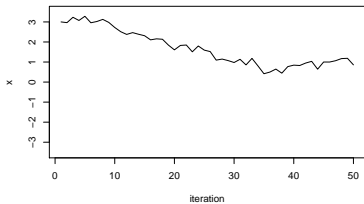


A toy example:

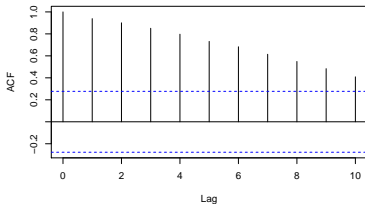
- 1 Initial Gibbs run shows high autocorrelation, panel 1.
- 2 Fit $y = \alpha + \beta x$ and transform $Z = Y - \hat{\alpha} - \hat{\beta}X$.
- 3 Rerun Gibbs, but sampling $p(X|Z)$ and $p(Z|X)$, panel 2.
- 4 Transform back to X, Y , panel 3.

Results for Toy Example

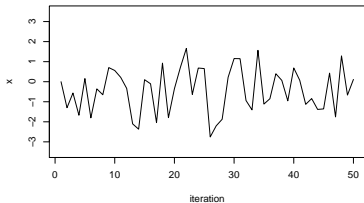
trace plot for initial run



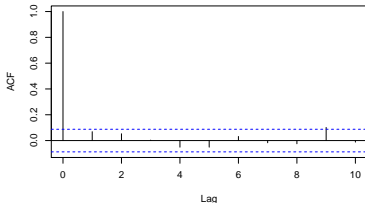
acf for initial run



trace plot for final run

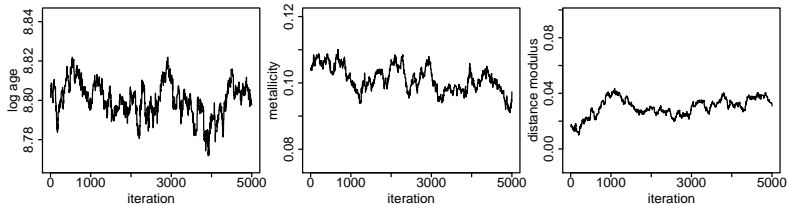


acf for final run

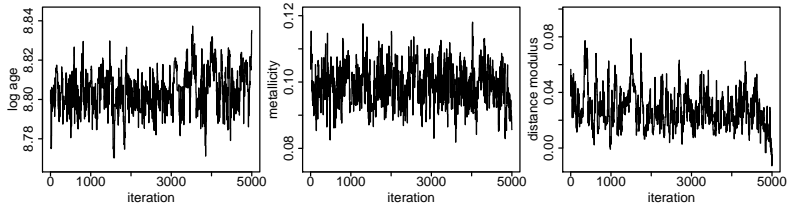


Results for Stellar Evolution Model

Initial Burn-in Period



After Dynamic Transformation



Data Augmentation

- We can sometimes simplify computation by including other unknown quantities in the model.
- Canonical Examples: *Missing Data* in Sample Surveys.
- Component photon energies of piled events (spectral analysis).
- If we had *Complete Data* analysis would be easier.
- More generally: there may quantities that we never *expected to observe*, but had we observed them, data analysis would be easier.

We call such quantities *Augmented Data* and their use in statistical computation *The Method of Data Augmentation*.

Handling Background with DA

Simple Example: Backgd contamination in single bin detector.

- Contaminated source counts: $Y = Y_S + Y_B$
- Background counts: X
- Background exposure is 24 times the source exposure.
- We observe Y and X .

A Poisson Multi-Level Model:

LEVEL 1: $Y|Y_B, \lambda_S \sim \text{Poisson}(\lambda_S) + Y_B$.

LEVEL 2: $Y_B|\lambda_B \sim \text{Pois}(\lambda_B)$ and $X|\lambda_B \sim \text{Pois}(24\lambda_B)$.

LEVEL 3: Specify a prior distribution on λ_B and λ_S .

Handling Background with DA

A Poisson Multi-Level Model:

LEVEL 1: $Y|Y_B, \lambda_S \sim \text{Poisson}(\lambda_S) + Y_B$.

LEVEL 2: $Y_B|\lambda_B \sim \text{Pois}(\lambda_B)$ and $X|\lambda_B \sim \text{Pois}(24\lambda_B)$.

LEVEL 3: Specify a prior distribution on λ_B and λ_S .

Data Augmentation

- Formulate model in terms of “missing data”.
- If Y_B were known.
- If λ_B and λ_S were known.

With Y_B we simplify the relationships among the quantities.

The Data Augmentation Sampler

A Two-Step Gibbs Sampler:

STEP 1: Sample Y_B given (λ_S, λ_B) , X , and Y .

$$Y_B \sim \text{Binomial} \left(Y, \frac{\lambda_B}{\lambda_S + \lambda_B} \right)$$

STEP 2: Sample (λ_S, λ_B) given X , Y_B , and Y_S .

$$\lambda_B \sim \text{Gamma}(X + Y_B + 1, 24 + 1)$$

$$\lambda_S \sim \text{Gamma}(Y_S + 1, 1)$$

The power of data augmentation is that it separates a complex problem into a series of simpler parts... just like Gibbs Sampler.

Details of STEP 1

$$\begin{aligned}
 p(Y_B, | \lambda_B, \lambda_S, Y) &\propto p(Y_B, Y | \lambda_B, \lambda_S) \\
 &= p(Y | \lambda_B, \lambda_S, Y_B) \times p(Y_B | \lambda_B, \lambda_S) \\
 &= \frac{e^{-\lambda_S} \lambda_S^{Y-Y_B}}{(Y-Y_B)!} \times \frac{e^{-\lambda_B} \lambda_B^{Y_B}}{Y_B!} \\
 &\propto \frac{1}{(Y-Y_B)! Y_B!} \lambda_S^{Y-Y_B} \lambda_B^{Y_B} \\
 &\propto \frac{Y!}{(Y-Y_B)! Y_B!} \left(\frac{\lambda_S}{\lambda_S + \lambda_B} \right)^{Y-Y_B} \left(\frac{\lambda_B}{\lambda_S + \lambda_B} \right)^{Y_B} \\
 &= \text{Binomial} \left(Y, \frac{\lambda_B}{\lambda_S + \lambda_B} \right)
 \end{aligned}$$

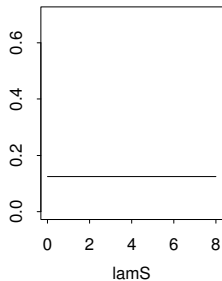
Requires case-by-case probability calculations.

Details of STEP 2

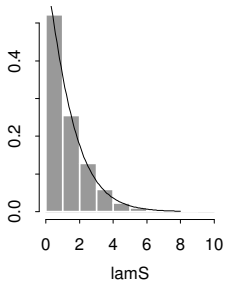
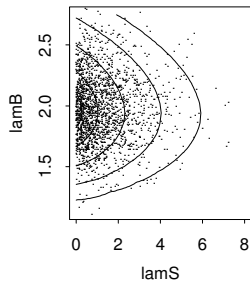
$$\begin{aligned}
 p(\lambda_S, \lambda_B, | Y_B, Y, X) &= p(\lambda_S, \lambda_B, | Y_S, Y_B, X) \\
 &\propto p(Y_S, Y_B, X | \lambda_B, \lambda_S) \\
 &= p(Y_S | \lambda_S) p(Y_B | \lambda_B) p(X | \lambda_B) \\
 &= \frac{e^{-\lambda_S} \lambda_S^{Y_S}}{Y_S!} \frac{e^{-\lambda_B} \lambda_B^{Y_B}}{Y_B!} \frac{e^{-24\lambda_B} (24\lambda_B)^X}{X!} \\
 &\propto \left(e^{-\lambda_S} \lambda_S^{Y_S} \right) \times \left(e^{-(24+1)\lambda_B} \lambda_B^{Y_B+X} \right) \\
 &\propto \gamma(Y_S + 1, 1) \times \gamma(X + Y_B + 1, 24 + 1)
 \end{aligned}$$

Results

prior



posterior

joint posterior
with flat prior

Here $Y = 1$ and $X = 48$.

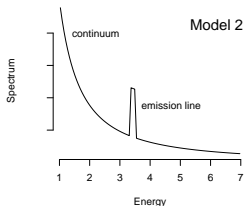
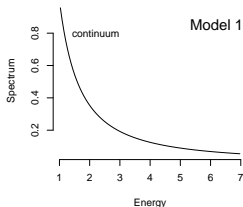
Handling a Spectral Emission Line

Recall the Power Law Spectral Model:

- $Y_i \sim \text{Poisson}(\alpha E_i^{-\beta}).$

Add a Spectral Emission Line:

- $Y_i \sim \text{Poisson}(\alpha E_i^{-\beta} + \gamma I\{i \in \mathcal{L}(\delta)\}).$
- $I\{i \in \mathcal{L}(\delta)\}$ is one if $i \in \mathcal{L}(\delta)$, otherwise it is zero.
- $\mathcal{L}(\delta) = \{\delta - 1, \delta, \delta + 1\}$
- $\theta_2 = (\alpha, \beta, \gamma, \delta)$



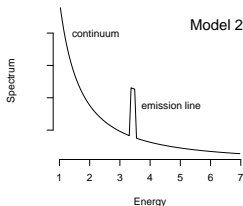
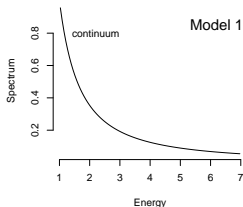
Handling a Spectral Emission Line

Continuum + Emission Line Model:

- 1 $Y_i \sim \text{Poisson} \left(\alpha E_i^{-\beta} + \gamma I\{i \in \mathcal{L}(\delta)\} \right)$.
- 2 An example of a *finite mixture model*.
- 3 Let Z_i be count in bin i due to line.
- 4 $Z_i | (Y_i, \theta_2) \sim$

$$\text{Binomial} \left(Y_i, \frac{\gamma I\{i \in \mathcal{L}(\delta)\}}{\gamma I\{i \in \mathcal{L}(\delta)\} + \alpha E_i^{-\beta}} \right)$$

- 5 Update α, β, γ , and δ given Z_i and continuum count = $X_i = Y_i - Z_i$?



A Metropolis within Gibbs Sampler

A Two-Step Sampler:

STEP 1: Sample Z_i given (θ_2, Y_i) , for $i = 1, \dots, n$.

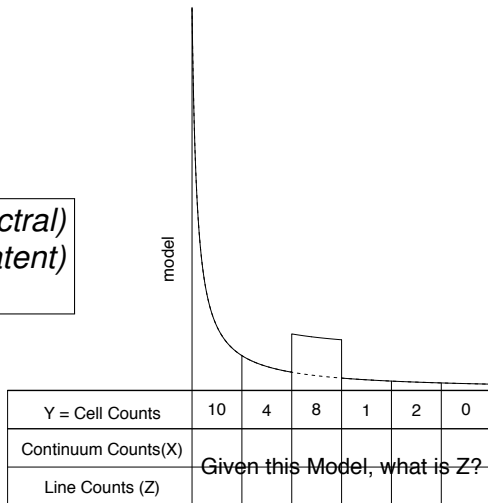
$$Z_i | (Y_i, \theta_2) \sim \text{Binomial} \left(Y_i, \frac{\gamma I\{i \in \mathcal{L}(\delta)\}}{\gamma I\{i \in \mathcal{L}(\delta)\} + \alpha E_i^{-\beta}} \right)$$

STEP 2: $p(\alpha, \beta, \gamma, \delta | X, Z) = p(\alpha, \beta | X) p(\gamma, \delta | Z)$
 $= p(\alpha, \beta | X) p(\gamma | \delta, Z) p(\delta | Z)$

- 1 Sample $p(\alpha, \beta | X)$ using Metropolis or MH.
- 2 $\gamma | (\delta, Z) \sim \text{gamma}(\sum Z_i, 3)$
- 3 Updating δ given Z is tricky.

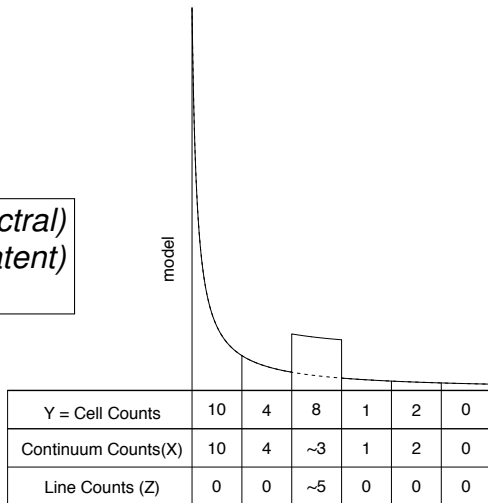
When Data Augmentation Fails

Consider a simple (spectral) model with the given (latent) cell counts.



When Data Augmentation Fails

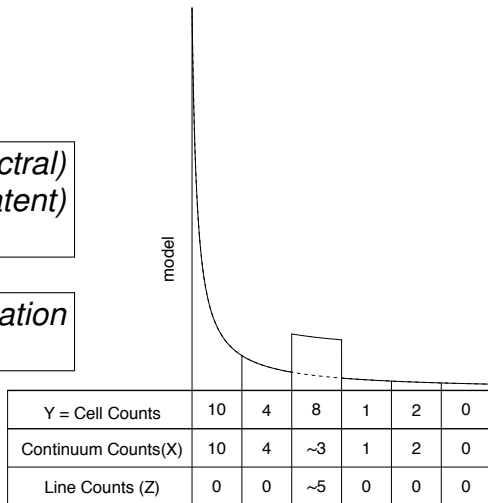
Consider a simple (spectral) model with the given (latent) cell counts.



When Data Augmentation Fails

Consider a simple (spectral) model with given (latent) cell counts.

Given Z what is the location of the emission line??

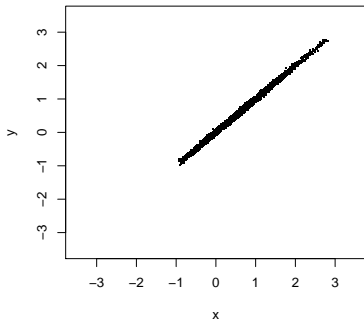


Handling a Spectral Emission Line

What Went Wrong?

*High Posterior Correlations
Are Always Problematic*

- Here Z and δ are highly correlated. In fact $\text{Var}(\delta|Z) = 0$.
- Given Z , δ will not change from iteration to iteration.



SOLUTION: Sample Z and δ in the same step.

An Improved Metropolis within Gibbs Sampler

A Two-Step Sampler:

STEP 1: Sample $p(Z, \delta | \alpha, \beta, \gamma, Y) = p(\delta | \alpha, \beta, \gamma, Y)p(Z | \theta_2, Y)$:

- 1 Sample δ given Y, α, β, γ using grid method:

$$p(\delta | \alpha, \beta, \gamma, Y) \propto p(Y | \theta_2).$$

- 2 For $i = 1, \dots, n$,

$$Z_i | (Y_i, \theta_2) \sim \text{Binomial} \left(Y_i, \frac{\gamma I\{i \in \mathcal{L}(\delta)\}}{\gamma I\{i \in \mathcal{L}(\delta)\} + \alpha E_i^{-\beta}} \right)$$

STEP 2: Sample $p(\alpha, \beta, \gamma | \delta, X, Z) = p(\alpha, \beta | X)p(\gamma | \delta, Z)$:

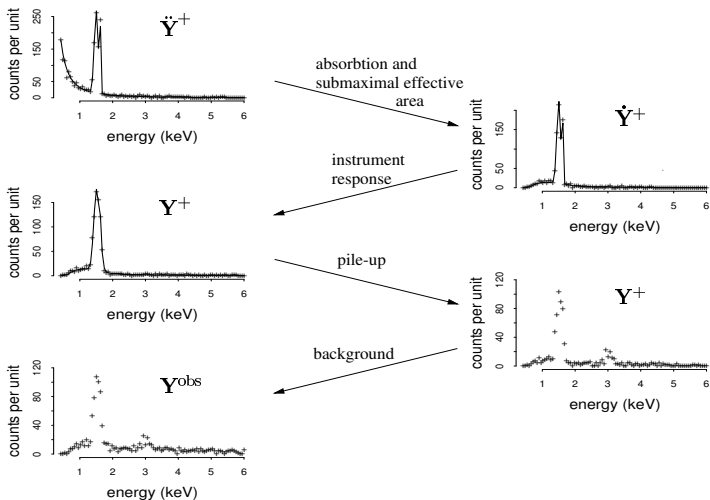
- 1 Sample $p(\alpha, \beta | X)$ using Metropolis or MH.
- 2 $\gamma | (\delta, X) \sim \text{gamma}(\sum Z_i, 3)$

Strategies for Implementing Gibbs Samplers

How we set up the complete conditional distributions can have a big impact on the performance of a Gibbs Sampler.

- 1 We have seen the potential effect of the choice of subsets:
 - $p(\vartheta|\varphi, \varsigma)$ and $p(\varphi, \varsigma|\vartheta)$ versus
 - $p(\vartheta, \varphi|\varsigma)$ and $p(\varsigma|\vartheta, \varphi)$
- 2 Combining steps into a single joint step is called *blocking*. This generally improves convergence:
 - $p(\vartheta|\varphi, \varsigma)$, $p(\varphi|\vartheta, \varsigma)$, and $p(\varsigma|\vartheta, \varphi)$ versus
 - $p(\vartheta, \varphi|\varsigma)$ and $p(\varsigma|\vartheta, \varphi)$
- 3 Removing a variable from the chain is called *collapsing*. This is also generally helpful:
 - $p(\vartheta, \varphi|\varsigma)$ and $p(\varsigma|\vartheta, \varphi)$ versus
 - $p(\vartheta|\varsigma)$ and $p(\varsigma|\vartheta)$
- 4 *Partial Collapsing* encompasses blocking and collapsing.

Example: Using DA for Spectral Analysis



Overview of Recommended Strategy

(Adopted from *Bayesian Data Analysis*, Section 11.10, Gelman et al. (2005), Second Edition)

- 1 Start with a crude approximation to the posterior distribution, perhaps using a mode finder.
- 2 Simulate directly, avoiding MCMC, if possible.
- 3 If necessary use MCMC with one parameter at a time updating or updating parameters in batches.
- 4 Use Gibbs draws for closed form complete conditionals.
- 5 Use metropolis jumps if complete conditional is not in closed form. Tune variance of jumping distribution so that acceptance rates are near 20% (for vector updates) or 40% (for single parameter updates).

Overview of Recommended Strategy- Con't

- 6 To improve convergence, use transformations so that parameters are approximately independent and/or approximately Gaussian.
- 7 Check for convergence using multiple chains.
- 8 Compare inference based on crude approximation and MCMC. If they are not similar, check for errors before believing the results of the MCMC.