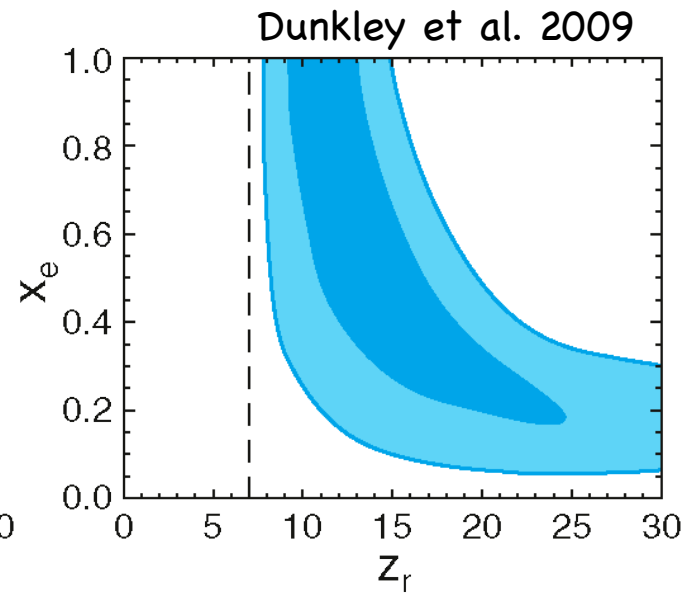
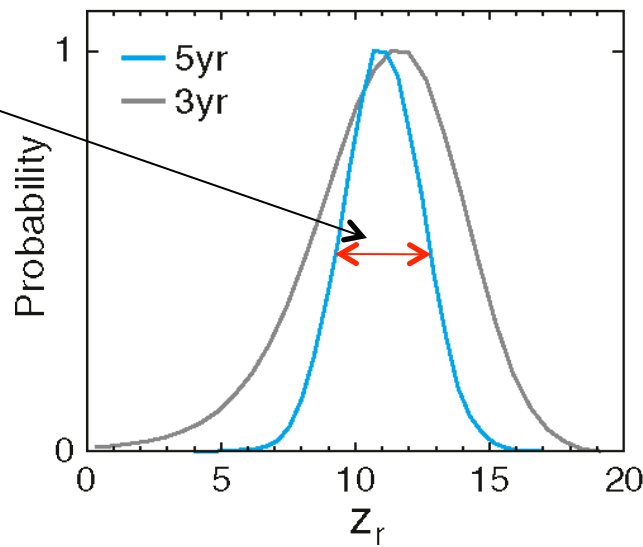


Marginalisation

Marginal probability: posterior probability of a given parameter regardless of the value of the others. It is obtained by integrating the posterior over the parameters that are not of interest.

$$p(\vartheta_2 | x) = \int p(\theta | x) d\theta_1 d\theta_3 \dots d\theta_n$$

Marginal errors characterise the width of the marginal posterior distributions.



Pitfalls of Bayesian inference

- There is no “correct” way to choose a prior. Bayesian inference requires skills to translate subjective prior beliefs into a mathematically formulated prior.
- It can produce posterior distributions which are heavily influenced by priors.
- It often comes with high computational costs.

Pitfalls of frequentist inference

The likelihood principle

- The likelihood principle is a conjecture that has been formulated based on the work of Fisher (in the 1920s) as well as Barnard, Savage and Birnbaum (in the 1960s).
- It states that, in the inference for the model parameters θ_i after the data \mathbf{x} have been observed, all the relevant experimental information is contained in the likelihood function.
- A second statement is that two likelihood functions contain the same information about the model parameters if they are proportional to each other.

Testing fairness of a coin I

- Experiment 1: flipping a coin one gets the sequence HTHHTHHHHHHT
Can you conclude that there is evidence against an unbiased coin?

Frequentist analysis: (how rare something as extreme should happen if the experiment is repeated many times and the coin is fair?)

If the experiment consists of flipping a coin 12 times, then

binomial statistical model:

$$P(x | \theta) = \binom{12}{x} \theta^x (1 - \theta)^{12-x}$$

Null hypothesis H_0 : $\theta = 0.5$ (θ = prob. to get head)

Alternative hypothesis: H_1 : $\theta > 0.5$

Let's compute the probability to find equal or more extreme results by repeating the experiment (the so-called p-value)

$$P(x \geq 9 | \theta) = \sum_{x=9}^{12} \binom{12}{x} \theta^x (1 - \theta)^{12-x} = 0.073$$

In 7.3 per cent of the cases we will find 9 heads or more. This is above the conventional 5% C.L. threshold. **We cannot reject H_0 at the 5% C.L.**

Testing fairness of a coin II

- Experiment 2: flipping a coin one gets the sequence HTHHTHHHHHHT
Can you conclude that there is evidence against an unbiased coin?

Frequentist analysis: (how rare something as extreme should happen if the experiment is repeated many times and the coin is fair?)

If the experiment consists of flipping a coin until one gets 3 tails, then negative binomial statistical model:

$$P(x | \theta) = \binom{3+x-1}{3-1} \theta^x (1-\theta)^3$$

Null hypothesis H_0 : $\theta = 0.5$ (θ = prob. to get head)

Alternative hypothesis: H_1 : $\theta > 0.5$

Let's compute the probability to find equal or more extreme results by repeating the experiment

$$P(x \geq 9 | \theta) = \sum_{x=9}^{12} P(x | \theta) = 0.0327$$

In 3.27 per cent of the cases we will find 9 heads or more. This is below the conventional 5% C.L. threshold. **We can reject H_0 at the 5% C.L.**

- Frequentist inference depends on the question asked although the data are exactly the same.
- P-values change according to the decision procedure (i.e. the procedure with which the experiment is repeated).
- This is a violation of the likelihood principle and a source of criticism towards frequentist statistics (from the Bayesians)
- In Bayesian statistics one would build the likelihood function and the posterior probability for θ . Do this as an exercise in the two cases!

How can we do this in practice?

Markov Chain Monte Carlo

**Andrey Andreyevic Markov
(1856-1922)**

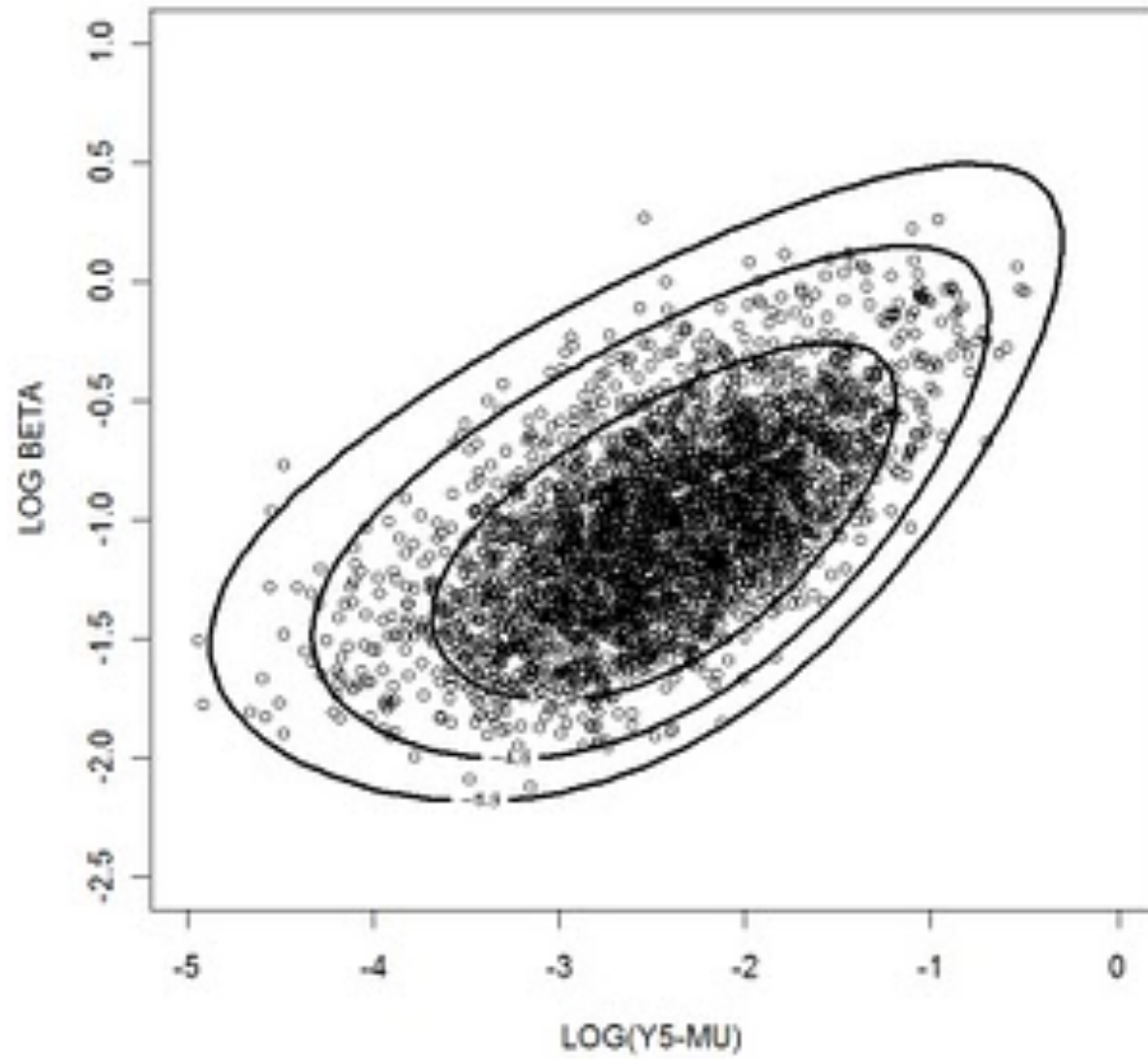


**Monte Carlo Casino
(1863-now)**



Markov Chain Monte Carlo

- **WHAT?** A numerical simulation method
- **AIM:** Sampling a given distribution function (known as the target density)
i.e. generate a finite set of points in some parameter space that are drawn from a given distribution function.
- **HOW?** By building a Markov chain that has the desired distribution as its equilibrium distribution



Markov chains

- A Markov chain is a sequence of random variables (or vectors) X_i (where i is an integer index: $i=0, \dots, N$) with the property that the transition probability

$$P(x_{N+1} \mid x_0, \dots, x_N) = P(x_{N+1} \mid x_N)$$

This means that the future of the chain does not depend on the entire past but only on the present state of the process.

Monte Carlo

- The term Monte Carlo method refers, in a very general meaning, to any numerical simulation which uses a computer algorithm explicitly dependent on a series of (pseudo) random numbers
- The idea of Monte Carlo integration was first developed by Enrico Fermi in the 1930s and by Stanislaw Ulam in 1947

$$\int f(x)p(x)dx \approx \frac{1}{N} \sum_{i=1}^N f(x_i) \quad [\text{where the } x_i \text{ are samples from } p(x)]$$

- Ulam and von Neumann used it for classified work at Los Alamos and as a “code name” for the project chose “Monte Carlo” as a reference to the famous Casino in Monaco.

MCMC and Bayesian statistics

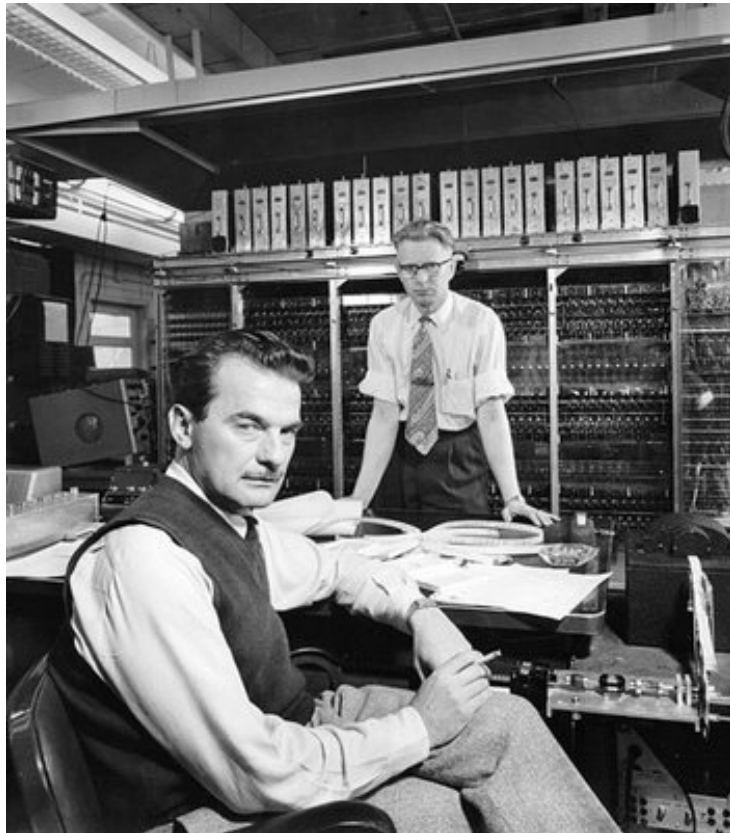
- The MCMC method has been very successful in modern Bayesian computing.
- In general (with very few exceptions) posterior densities are too complex to work with analytically.
- With the MCMC method, it is possible to generate samples from an arbitrary posterior density and to use these samples to approximate expectations of quantities of interest.
- Most importantly, the MCMC is guaranteed to converge to the target distribution under rather broad conditions, regardless of where the chain was initialized.
- Furthermore, if the chain is run for very long time (often required) you can recover the posterior density to any precision.
- The method is easily applicable to models with a large number of parameters (although the “curse of dimensionality” often causes problems in practice).

MCMC algorithm

- Choose a random initial starting point in parameter space, and compute the target density
- Repeat:
 - ✓ Generate a step in parameter space from a proposal distribution, generating a new trial point for the chain.
 - ✓ Compute the target density at the new point, and accept it or not with the Metropolis-Hastings algorithm (see next slide).
 - ✓ If the point is not accepted, the previous point is repeated in the chain.
- End Repeat

The Metropolis algorithm

Nicholas Constantine Metropolis
(1915-1999)



“Equation of state calculation by fast
computing machines”

Metropolis et al. (1953)

- After generating a new MCMC sample using the proposal distribution, calculate

$$r = \text{probability of acceptance} = \min\left(\frac{f(\theta_{new})}{f(\theta_{old})}, 1\right)$$

- Then sample u from the uniform distribution $U(0,1)$
- Set $\theta_{t+1} = \theta_{new}$ if $u < r$; otherwise set $\theta_{t+1} = \theta_t$
- Note that the number of iterations keeps increasing regardless of whether a proposed sample is accepted.

The Metropolis algorithm

- It can be demonstrated that the Metropolis algorithm works.
- The proof is beyond the scope of this course but, if you are curious, you can check standard statistics textbooks including Roberts (1996) and Liu (2001).
- You are not limited to a symmetric random-walk proposal distribution in establishing a valid sampling algorithm. A more general form, now known as the Metropolis-Hastings algorithm, was proposed by Hastings (1970). In this case:

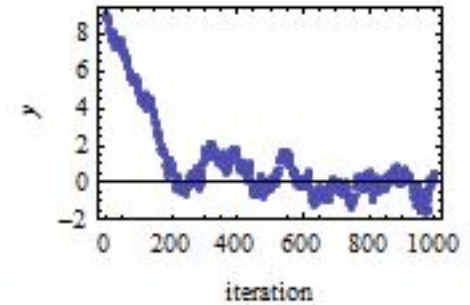
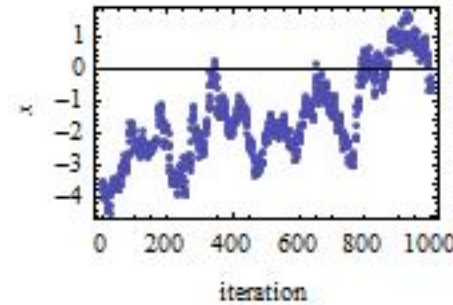
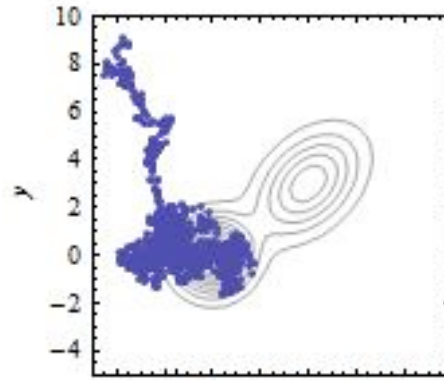
$$r = \text{probability of acceptance} = \min\left(\frac{f(\theta_{new})q(\theta_t | \theta_{new})}{f(\theta_{old})q(\theta_{new} | \theta_t)}, 1\right)$$

The proposal distribution

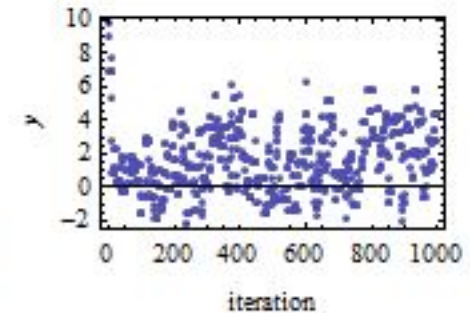
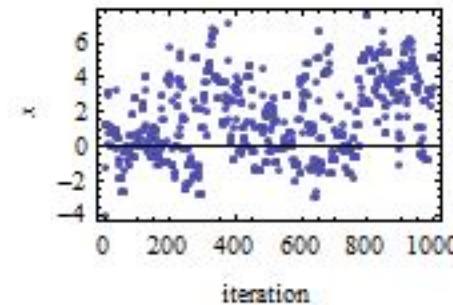
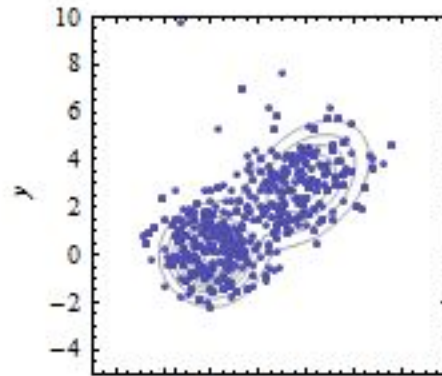
- If one takes too small steps, it takes long time to explore the target and the different entries of the chain are very correlated
- If one takes too large steps, almost all trials are rejected and the different entries of the chain are very correlated
- There is an optimal proposal distribution (easy to identify if we knew already the target density)

Effect of the sampling distribution

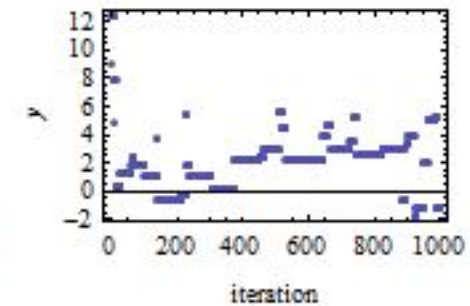
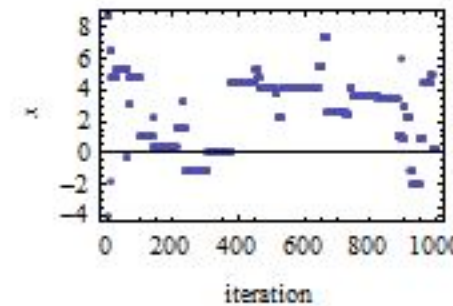
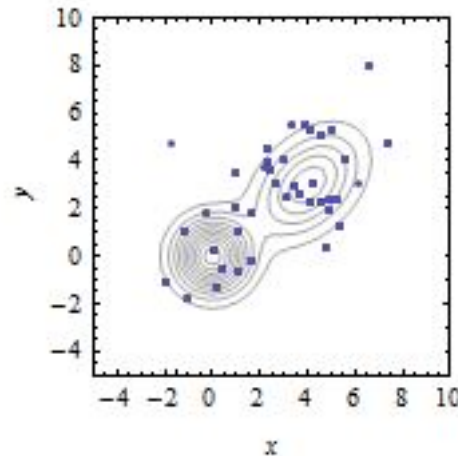
Gaussian proposal distribution with $\sigma = 0.2$, acceptance rate = 85.1%



Gaussian proposal distribution with $\sigma = 2.2$, acceptance rate = 37.9%



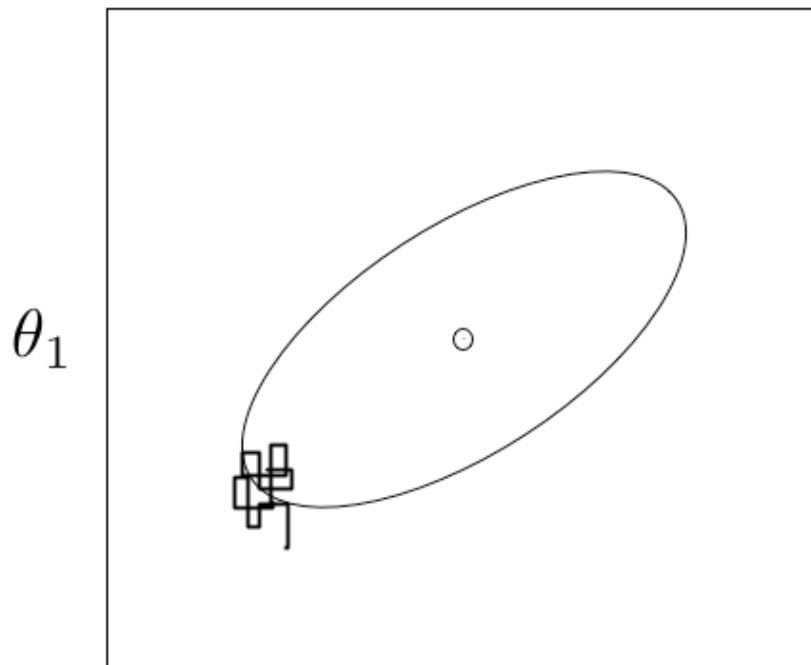
Gaussian proposal distribution with $\sigma = 10.2$, acceptance rate = 4.1%



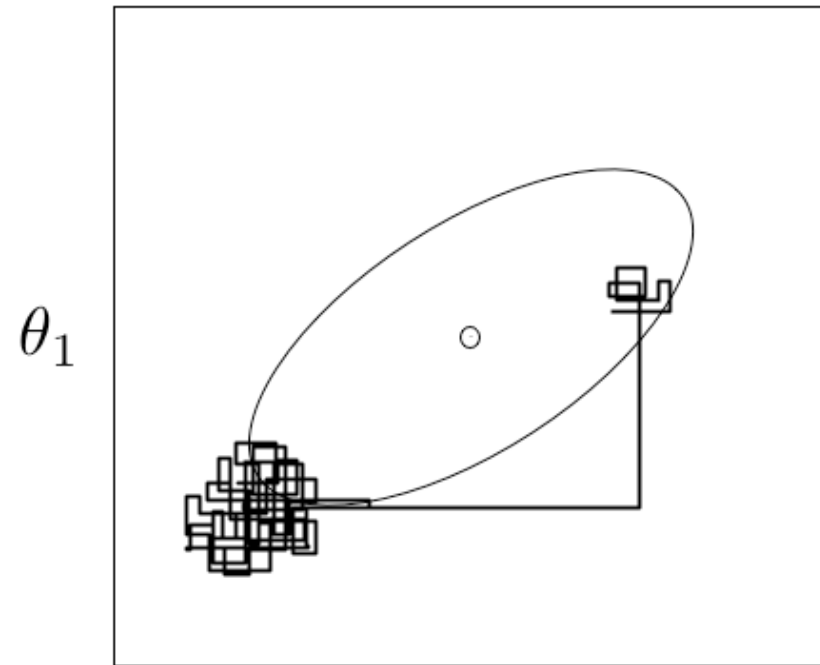
Mixing

Mixing refers to the degree to which the Markov chain explores the support of the posterior distribution. Poor mixing may stem from inappropriate proposals (if one is using the Metropolis–Hastings sampler) or from attempting to estimate models with highly correlated variables.

Bad mixing

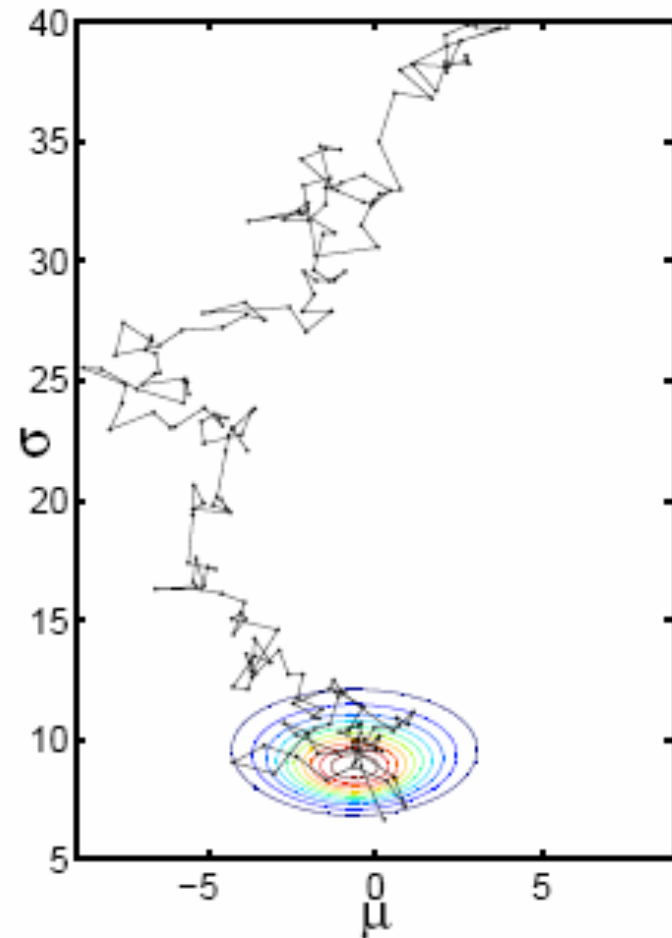


Metastability



Burn-in

- Mathematical theorems guarantee that the Metropolis algorithm will asymptotically converge to the target distribution independently of its starting point.
- However, there will be an initial transient of unknown length during which the chain reaches its stationary state.
- In practice, you have to assume that after N_b iterations, the chain converged and started sampling from its target distribution.
- The value of N_b is called the burn-in number.



Issues with MCMC

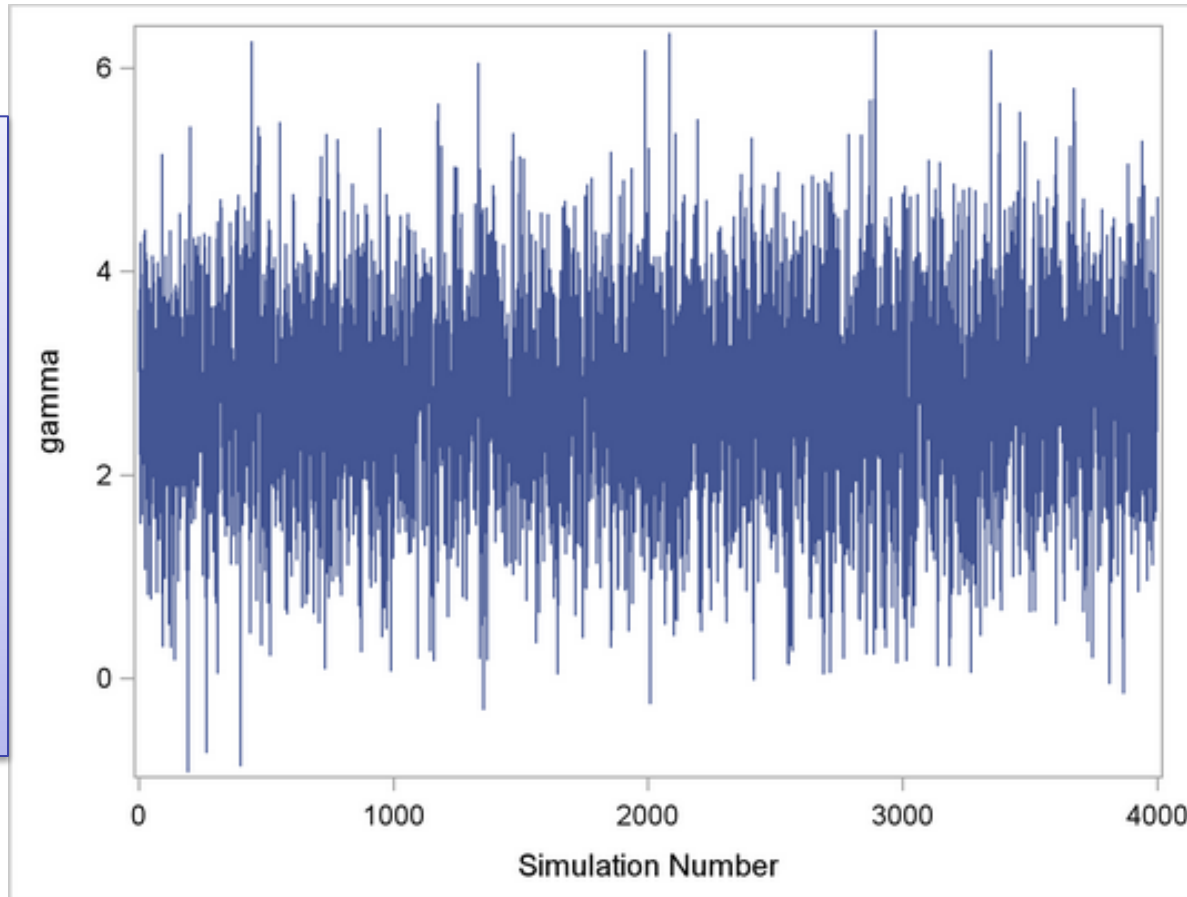
- You have to decide whether the Markov Chain has reached its stationary distribution
- You have to decide the number of iterations to keep after the Markov Chain has reached stationarity
- Convergence diagnostics help to resolve these issues. Note, however, that most diagnostics are designed to verify a necessary but NOT sufficient condition for convergence.

Visual analysis via Trace Plots

- The simplest diagnostic is obtained by plotting the value of one model parameter versus the simulation index (i.e. the first point in the Markov chain has index 1, the second 2, and so on).
- This is called a Trace Plot.
- As we will see, a trace tells you if a longer burn-in period is needed, if a chain is mixing well, and gives you an idea about the stationary state of the chain.
- Trace plots must be produced for all the parameters, not only for those of interest! If some of parameters have bad mixing you cannot get accurate posterior inference for parameters that appear to have good mixing.

Example I

If the chain has reached stationarity the mean and the variance of the trace plot should keep relatively constant.

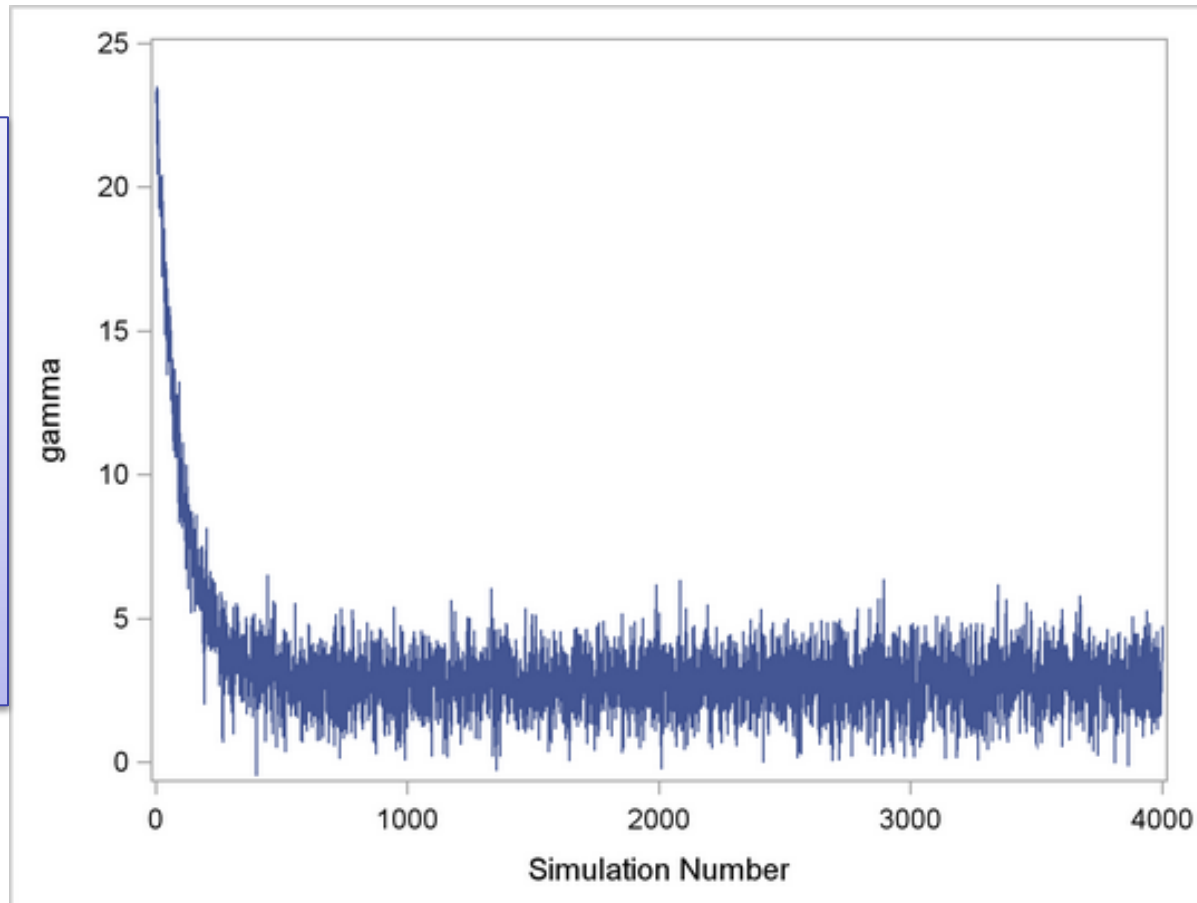


A chain that mixes well traverses the posterior space rapidly, and it can jump from a remote region of the posterior to another in relatively few steps.

The figure displays a "perfect" trace plot, not easy to achieve in high-dimensions

Example II

This chain starts at a very remote location and makes its way to the targeting distribution .

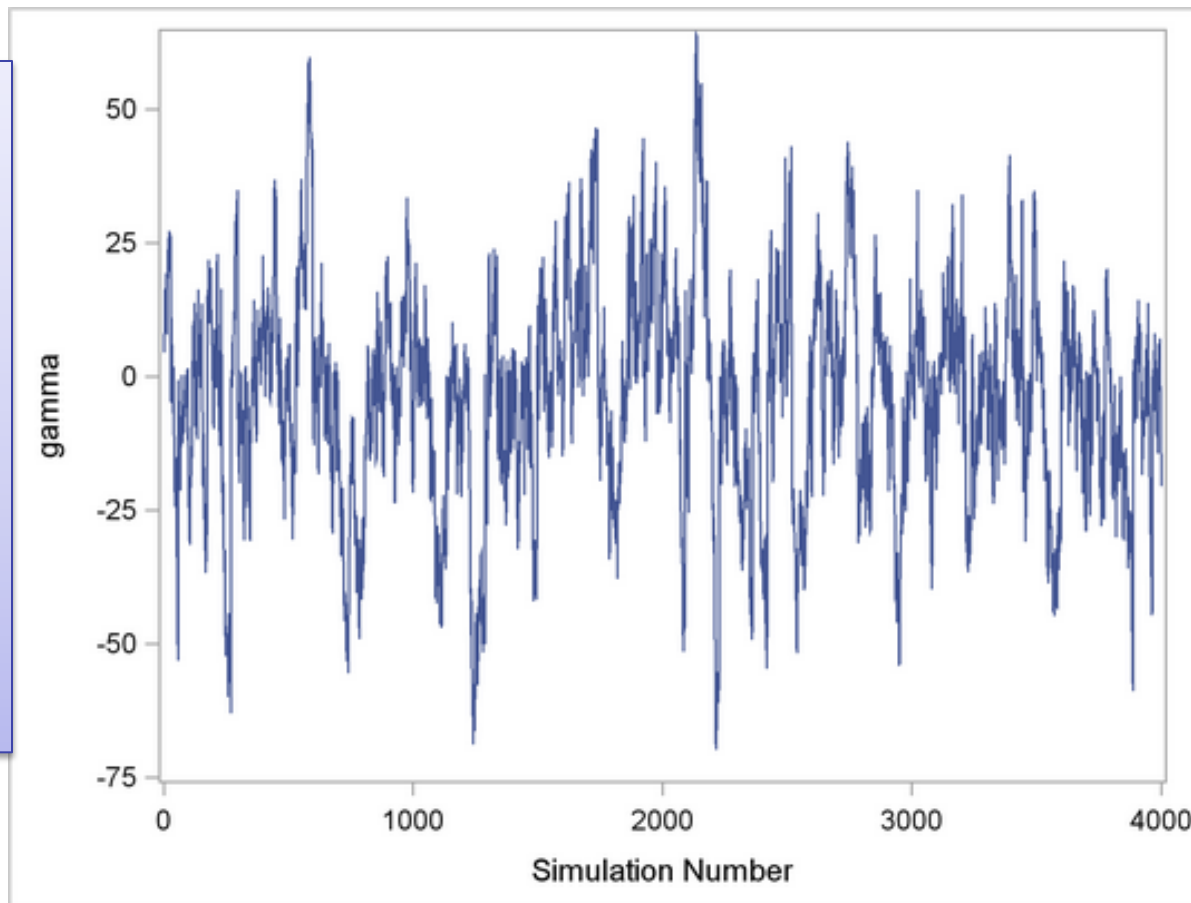


This chain mixes well locally and travels relatively quickly to the target distribution, reaching it in a few hundred iterations.

If you have a chain like this, increase the burn-in sample size.

Example III

This trace plot shows marginal mixing. The chain is taking small steps and does not traverse its distribution quickly.

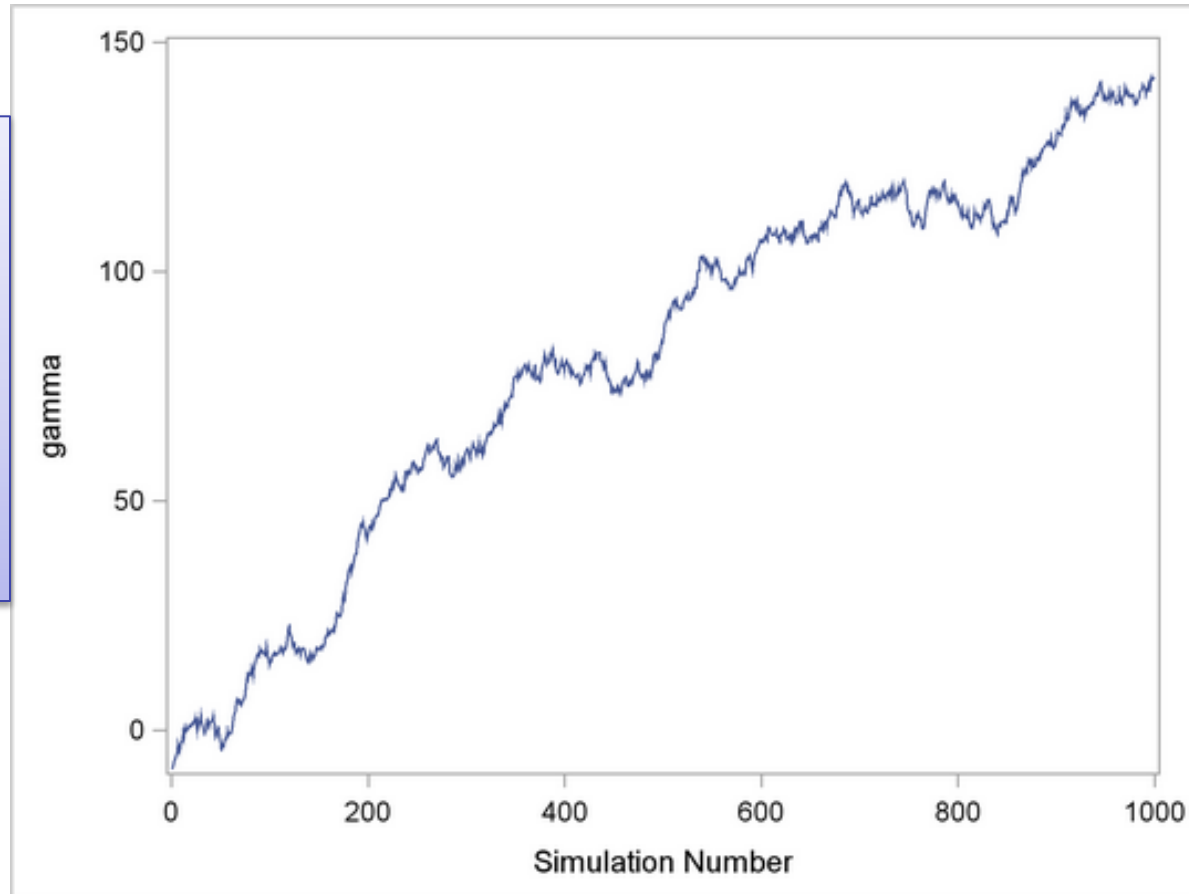


This type of trace plot is typically associated with high correlation among the samples. The chain takes too long to forget where it was before.

In order to obtain a given number of independent samples you need to run the chain for much longer.

Example IV

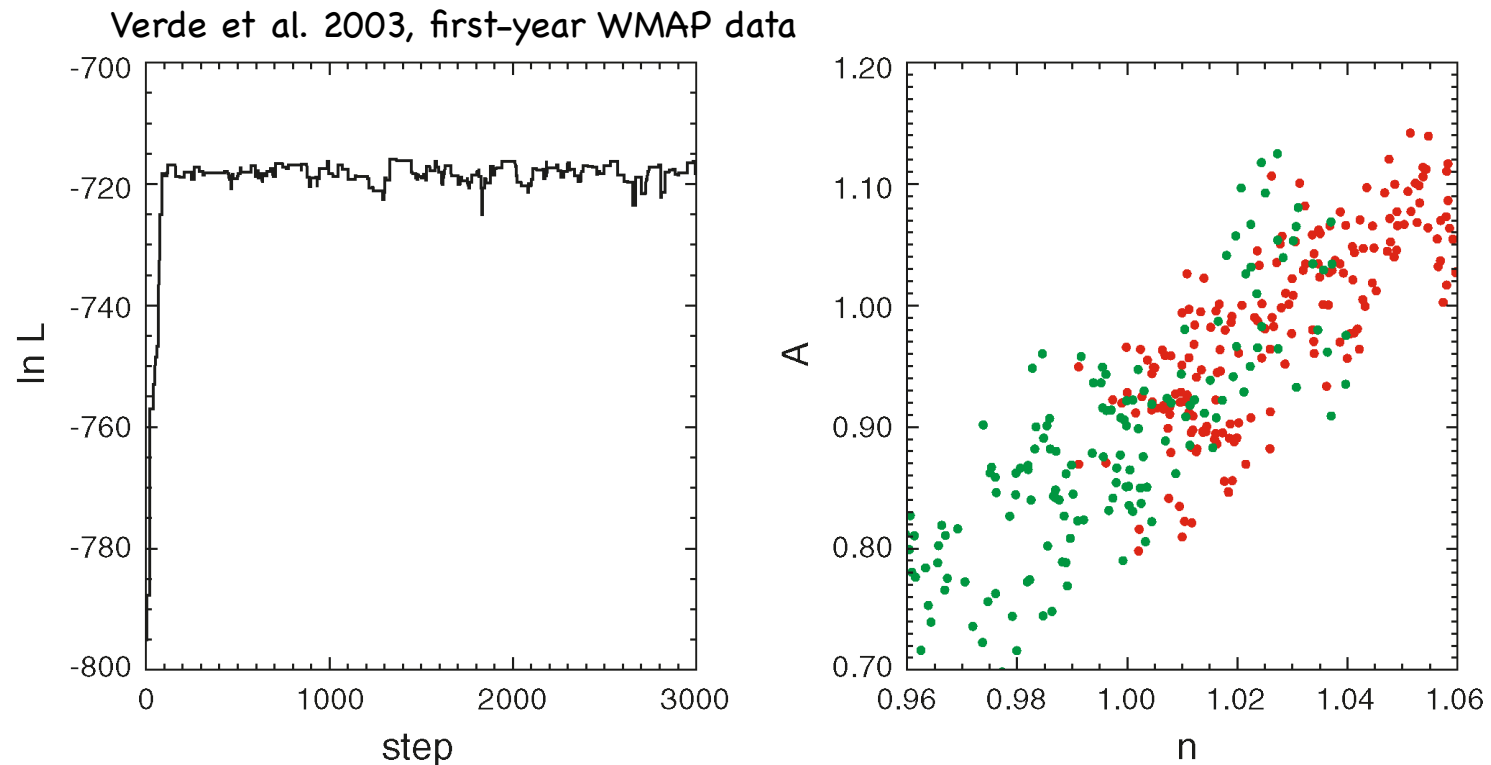
This chain has serious problems. It mixes very slowly, and it does not give any evidence of convergence.



You would want to try to improve the mixing of this chain. For example, you might consider changing the proposal distribution or reparameterizing your model.

This type of chain is entirely unsuitable for making parameter inferences!

Convergence



Although the trace plot on the left may appear to indicate that the chain has converged after a burn-in of a few hundred steps, in reality it has not fully explored the posterior surface.

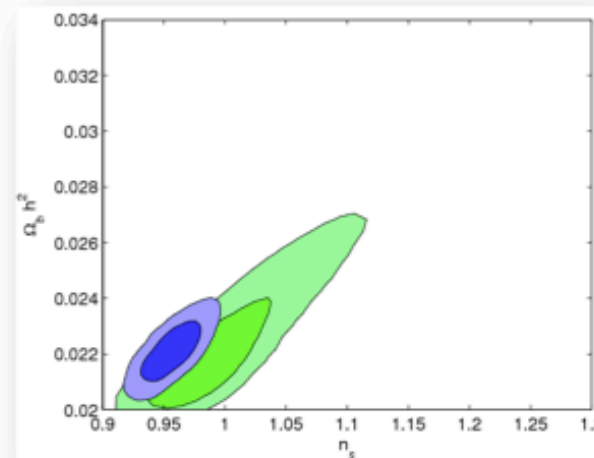
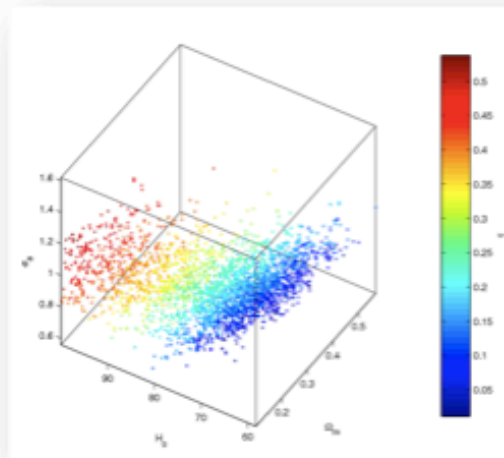
This is shown on the right where two chains of the same length are plotted. Using either of these two chains at this stage will give incorrect results for the best-fit cosmological parameters and their errors.

Statistical diagnostics

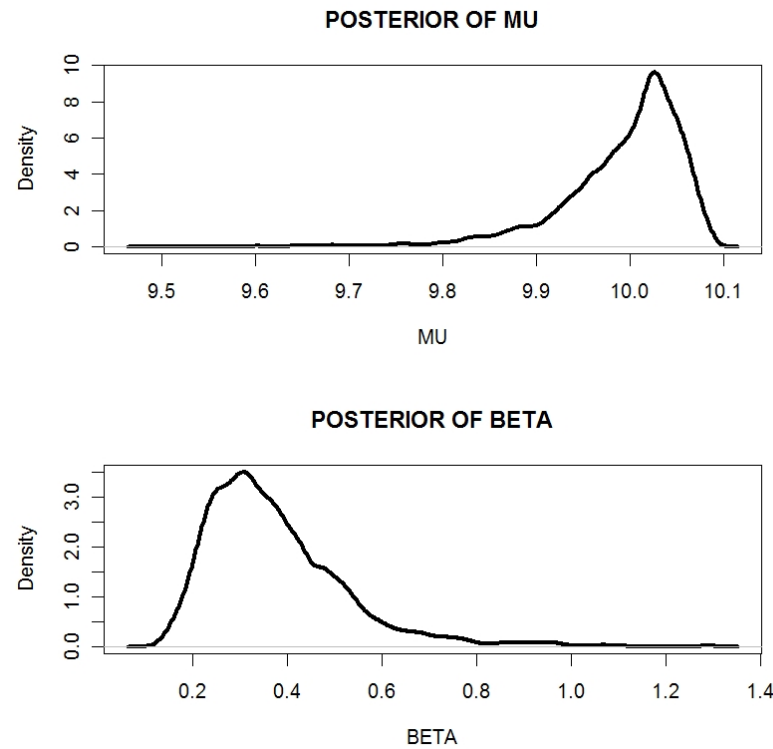
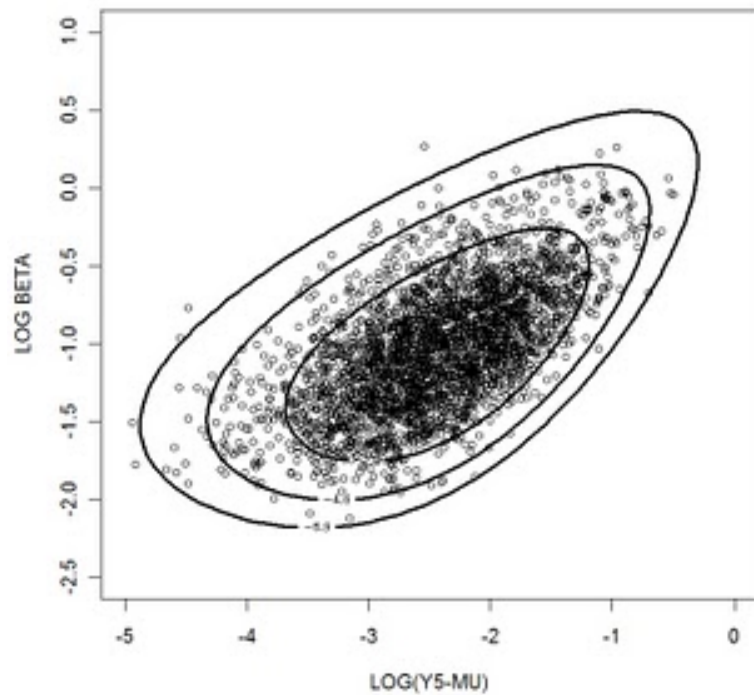
- **Gelman-Rubin:** uses parallel chains with dispersed initial values to test whether they all converge to the same target distribution.
- **Geweke:** tests whether the mean estimates of the parameters have converged by comparing means from the early and latter part of the Markov chain.
- **Raftery-Lewis:** Evaluates the accuracy of the estimated percentiles by reporting the number of samples needed to reach the desired accuracy.
- And many, many, more...

Marginalisation

- Marginalisation is trivial
 - Each point in the chain is labelled by all the parameters
 - To marginalise, just ignore the labels you don't want

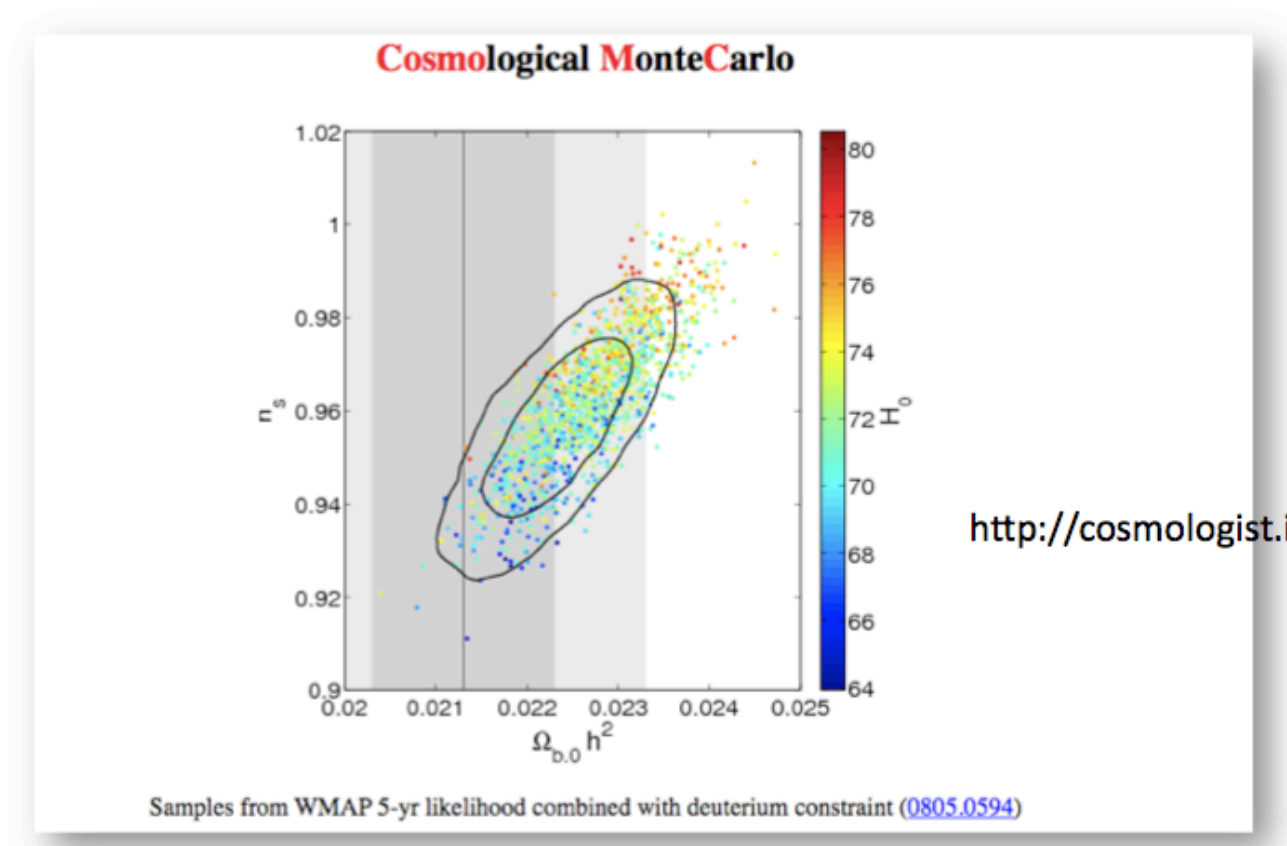


How to plot the results



Public codes for MCMC in cosmology

CosmoMC



Public codes for MCMC in cosmology

<http://baudren.github.io/montepython.html>

Monte Python

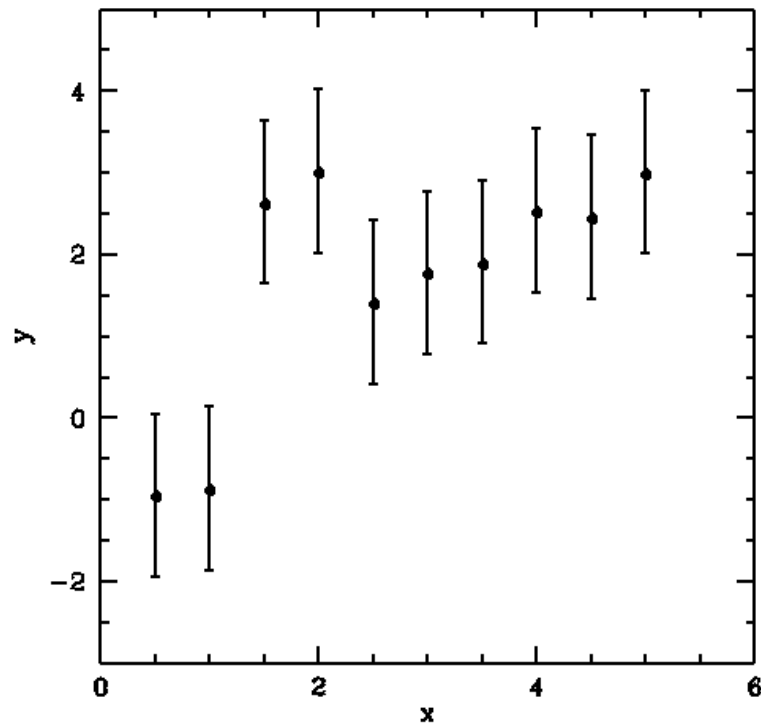
The Monte Carlo code for **CLASS** in Python

[Download the latest version](#)

Bayesian model comparison

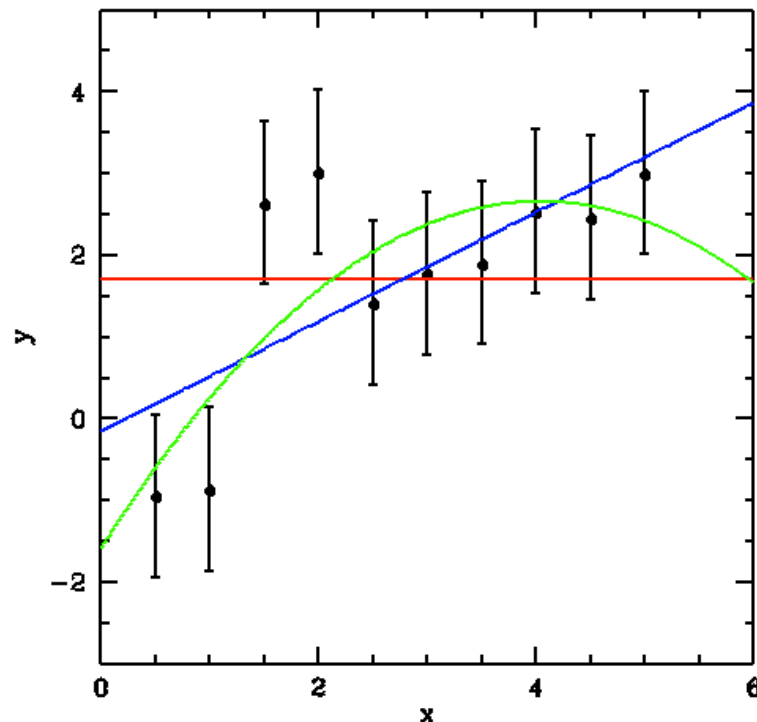
- Suppose you have some data
- You want to fit them with some model
- More than a model is available with a different number of free parameters (e.g. vanilla Λ CDM, Λ CDM + massive neutrinos, Λ CDM + curvature, CDM + dynamic dark energy)
- How do you choose the “best model”?
- Central problem in learning: how to balance “goodness of fit” criteria against the complexity of models
- We want to avoid “overfitting”

An example



- What function would you use to fit these data?
- Constant? $y=c$ (1 parameter)
- Linear? $y=bx+c$ (2 parameters)
- Quadratic? $y=ax^2+bx+c$ (3 parameters)

Maximum likelihood solution



- **Constant** fit: $\chi^2_{\min}=19.33$
- **Linear** fit: $\chi^2_{\min}=10.02$
- **Quadratic** fit: $\chi^2_{\min}=7.79$
- Which one should be preferred?
- I could get $\chi^2_{\min}=0$ by using a polynomial of degree 10

Bayes factor

Let's write Bayes theorem for the models in odds form:

$$K = \frac{P(M_1 | \vec{x})}{P(M_2 | \vec{x})} = BF \times \frac{\pi(M_1)}{\pi(M_2)}$$

$$BF = \frac{P(\vec{x} | M_1)}{P(\vec{x} | M_2)} = \frac{\int L(\vec{x} | \vec{\vartheta}_1, M_1) \pi(\vec{\vartheta}_1 | M_1) d\vec{\vartheta}_1}{\int L(\vec{x} | \vec{\vartheta}_2, M_2) \pi(\vec{\vartheta}_2 | M_2) d\vec{\vartheta}_2}$$

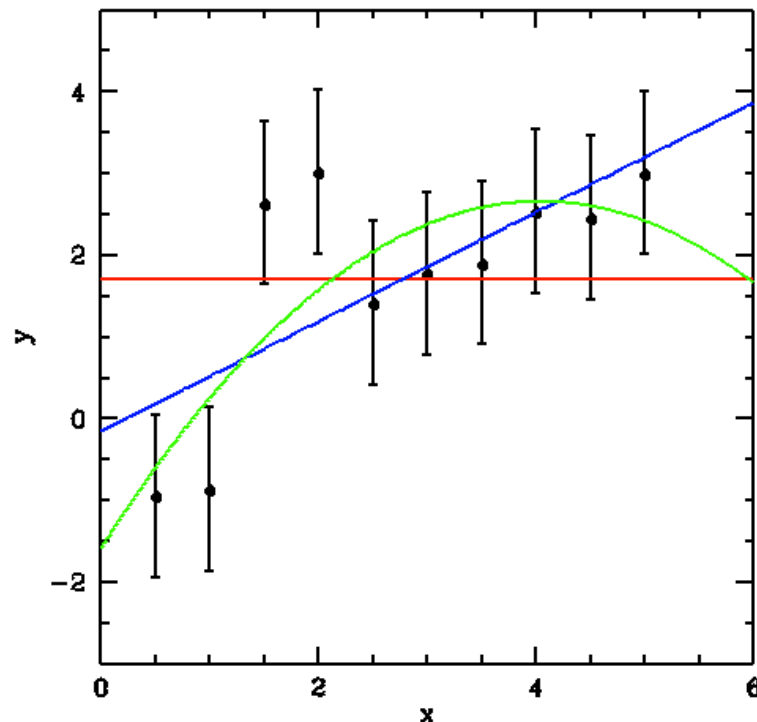
In words: posterior odds = Bayes factor x prior odds
where the Bayes factor (BF) is given by the evidence ratio for the two models.

The prior ratio is often taken as unity. The evidence ratio penalizes for unnecessary complexity in the models: models are penalized if a small part of their prior parameter range matches the data (Occam factor).

Jeffrey's scale

Odds ratio	Strength of evidence
$1 < K < 3$	Barely worth mentioning
$3 < K < 10$	Substantial
$10 < K < 30$	Strong
$30 < K < 100$	Very strong
$K > 100$	Decisive

Evidence and Bayes factors



Evidence:

- **Constant** fit: 5.04×10^{-5}
- **Linear** fit: 2.92×10^{-3}
- **Quadratic** fit: 1.93×10^{-3}
- The linear fit is slightly preferred to the quadratic even though it has a worse χ^2_{\min} (**BF**=1.5 while **BF**=57.9)
- I generated the data adding Gaussian noise with unit variance to the relation $y=x/2$ which was indeed linear

