

### (New) 3 The Metropolis algorithm

As said, we want to make the acceptance ratio as large as possible to sample as many different states as possible. The Metropolis algorithm does this. \* Consider again

$$\frac{P(\mu \rightarrow \nu)}{P(\nu \rightarrow \mu)} = \frac{P_\nu}{P_\mu} = \frac{g(\mu \rightarrow \nu) A(\mu \rightarrow \nu)}{g(\nu \rightarrow \mu) A(\nu \rightarrow \mu)}$$

If you choose  $g(\mu \rightarrow \nu) = g(\nu \rightarrow \mu)$ , then

$$\frac{P_\nu}{P_\mu} = \frac{A(\mu \rightarrow \nu)}{A(\nu \rightarrow \mu)}$$

or equivalently

$$P_\nu A(\nu \rightarrow \mu) = P_\mu A(\mu \rightarrow \nu)$$

The choice made by Metropolis was

$$A(\mu \rightarrow \nu) = \begin{cases} 1 & ; P_\nu > P_\mu \\ P_\nu / P_\mu & ; \text{otherwise} \end{cases}$$

So let us for an example assume that  $P_\nu$  is larger than  $P_\mu$  then:  $A(\mu \rightarrow \nu) = 1$ ;  $A(\nu \rightarrow \mu) = \frac{P_\mu}{P_\nu}$   
and

$$P_\nu \cdot \frac{P_\mu}{P_\nu} = P_\mu \cdot 1 \Rightarrow P_\mu = P_\mu$$

Which means that Metropolis choice satisfied our requirement!

\* The point of Metropolis's algorithm is that we do not need to have a functional form for  $A$  that holds for all  $\mu$  and  $\nu$ 's.

E.g.

$$A(\mu \rightarrow \nu) = A_0 e^{-\frac{\beta}{2} (E_\nu - E_\mu)}$$

Would be such a functional choice. This would lead to very low acceptance ratios in cases where  $E_\nu$  is even slightly larger or not much smaller than  $E_\mu$  (depending on the algorithm of  $g$ ), because  $A$  needs to be  $\in [0, \dots, 1]$ , so it must be normalized with the largest decrease in  $E$  possible by some choice of  $g$ :  $A_0 = \exp \left[ \frac{\beta}{2} \Delta E_{\max} \right]$  ;  $\Delta E_{\max} < 0$   
 Anyway. What Metropolis observed was that our requirement

$$\frac{P(\mu \rightarrow \nu)}{P(\nu \rightarrow \mu)} = \frac{P_\nu}{P_\mu} = \frac{A(\mu \rightarrow \nu)}{A(\nu \rightarrow \mu)}$$

fixes  $A$  only for a given pair  $(\mu, \nu)$  and hence, there can be as many functions  $A$  as there are pairs  $(\mu, \nu)$ !

For a thermal system, the Metropolis algorithm is hence

$$A(\mu \rightarrow \nu) = \begin{cases} 1 & ; E_\nu < E_\mu \\ e^{-\beta(E_\nu - E_\mu)} & ; \text{otherwise} \end{cases}$$

In cosmology, we want to sample the posterior  $\text{prob}(\vec{\lambda} | \text{data})$ . Or equivalently the likelihood  $\text{prob}(\text{data} | \vec{\lambda})$  times the prior  $\text{prob}(\vec{\lambda} | \text{prior knowledge})$ . A cosmological MCMC therefore uses

$$A(\mu \rightarrow \nu) = \begin{cases} 1 & ; \text{Likelihood of } \nu > \text{Likelihood of } \mu \\ \frac{\text{prob}(\vec{\lambda}_\nu)}{\text{prob}(\vec{\lambda}_\mu)} & ; \text{otherwise} \end{cases}$$

In practice, a simple Metropolis search is performed by choosing a Gaussian proposal  $g(\mu \rightarrow \nu)$  like

$$\Omega_{\text{m}} h^2 (\frac{t}{t+1}) \text{ is drawn from } e^{-\frac{[\Omega_{\text{m}} h^2(t) - \Omega_{\text{m}} h^2(t+1)]^2}{2\sigma_{\text{m}}^2}}$$

etc (there are some refinements like choosing parameter combinations to fight degeneracies, but the idea remains the same) and accept according to  $A(\mu \rightarrow \nu)$  above, where the probabilities are computed by the likelihood routines available for the different experiments.

## (Sivia) 4 MODEL SELECTION

Suppose there is a debate which model <sup>describes</sup> fits the data best. Naively, we could look at how well each model fits the data. However, a model with many parameters will always be able to give better fits compared to a simple model. So how do we decide?

The story of A and B:

Mr. A has a theory; Mr. B also has one but with an adjustable parameter  $\lambda$ . Whose theory should we prefer given data  $\mathcal{D}$ ?

To settle the issue, we should look at the relative merit of the two theories. If the

$$\text{posterior ratio} = \frac{\text{prob}(A|\mathcal{D},I)}{\text{prob}(B|\mathcal{D},I)}$$

is  $> 1$  : prefer A

$< 1$  : prefer B

$\approx 1$  : no preference

Bayes theorem yields:

$$\frac{\text{prob}(A|\mathcal{D},I)}{\text{prob}(B|\mathcal{D},I)} = \frac{\text{prob}(\mathcal{D}|A,I)}{\text{prob}(\mathcal{D}|B,I)} \times \frac{\text{prob}(A|I)}{\text{prob}(B|I)}$$

Without  $\lambda$ , we can't quantify  $\text{prob}(\mathcal{D}|B,I)$ , so we proceed

← marginal.

$$\begin{aligned} \text{prob}(\mathcal{D} | \mathcal{B}, \mathcal{I}) &= \int \text{prob}(\mathcal{D}, \lambda | \mathcal{B}, \mathcal{I}) d\lambda \\ &= \int \underbrace{\text{prob}(\mathcal{D} | \lambda, \mathcal{B}, \mathcal{I})}_{\text{not an ordinary pdf.}} \cdot \text{prob}(\lambda | \mathcal{B}, \mathcal{I}) d\lambda \end{aligned}$$

↗ product rule

Assume that Mr B only knows that  $\lambda \in [\lambda_{\min}, \lambda_{\max}]$

so

$$\text{prob}(\lambda | \mathcal{B}, \mathcal{I}) = \frac{1}{\lambda_{\max} - \lambda_{\min}} \quad \text{for } \lambda_{\min} \leq \lambda \leq \lambda_{\max}$$

Assume further that  $\exists \lambda_0$  which yields best fit and well approximated by Gaussian

$$\text{prob}(\mathcal{D} | \lambda, \mathcal{B}, \mathcal{I}) = \text{prob}(\mathcal{D} | \lambda_0, \mathcal{B}, \mathcal{I}) \exp\left[-\frac{(\lambda - \lambda_0)^2}{2\delta\lambda^2}\right]$$

In addition, we can take the prior  $\text{prob}(\lambda | \mathcal{B}, \mathcal{I})$  out of the  $\int d\lambda$  integral (it does only depend on  $\lambda_{\max}, \lambda_{\min}$ )

So

$$\text{prob}(\mathcal{D} | \mathcal{B}, \mathcal{I}) = \frac{1}{\lambda_{\max} - \lambda_{\min}} \int_{\lambda_{\min}}^{\lambda_{\max}} \text{prob}(\mathcal{D} | \lambda, \mathcal{B}, \mathcal{I}) d\lambda$$

Furthermore, we assume that  $\lambda_{\max}, \lambda_{\min}$  does not truncate the  $\lambda$ -peak around  $\lambda_0$  so that we can perform the Gaussian integral yielding

$$\begin{aligned} \text{prob}(\mathcal{D} | \mathcal{B}, \mathcal{I}) &= \frac{1}{\lambda_{\max} - \lambda_{\min}} \int_{\lambda_{\min}}^{\lambda_{\max}} \text{prob}(\mathcal{D} | \lambda, \mathcal{B}, \mathcal{I}) d\lambda \\ &= \frac{\text{prob}(\mathcal{D} | \lambda_0, \mathcal{B}, \mathcal{I}) \times \delta\lambda \sqrt{2\pi}}{\lambda_{\max} - \lambda_{\min}} \end{aligned}$$

So

$$\frac{\text{prob}(A | D, I)}{\text{prob}(B | D, I)} = \underbrace{\frac{\text{prob}(A | I)}{\text{prob}(B | I)}}_{= 1 \text{ (b/fair)}} \times \underbrace{\frac{\text{prob}(D | A, I)}{\text{prob}(D | B, I)}}_{\text{ratio of goodness to fit}} \times \frac{\lambda_{max} - \lambda_{min}}{\delta \lambda} \uparrow \text{penalty}$$

penalty  $\gg 1$ , because  $\lambda_{max} - \lambda_{min}$  usually used larger than  $\delta \lambda$ . This is called a Ockham factor or Ockham's razor.

Problem: If no knowledge about additional parameter known then infinite penalty. No big problem in practice, though, because usually boundaries given and then goodness to fit decides. If goodness to fit undecided then Ockham effect kicks in.

If you take a look back at Bayes's theorem:

$$\text{prob}(A | D, B, I) = \frac{\text{prob}(D | A, B, I) \text{prob}(A | B, I)}{\text{prob}(D | B, I)}$$

You'll recognize that the denominator played a crucial role in ascertaining the relative merit of B over A. It is therefore called "evidence" or "marginal likelihood".

## Hypothesis testing

Suppose we have a hypothesis  $H_1$ . To quantify how much we believe that this is true given the data  $\mathcal{D}$  and information  $\mathcal{I}$ , we need to evaluate the posterior

$$\text{prob}(H_1 | \mathcal{D}, \mathcal{I}) = \frac{\text{prob}(\mathcal{D} | H_1, \mathcal{I}) \times \text{prob}(H_1 | \mathcal{I})}{\text{prob}(\mathcal{D} | \mathcal{I})}$$

Suppose we have another hypothesis  $H_2$ .

$$\frac{\text{prob}(H_1 | \mathcal{D}, \mathcal{I})}{\text{prob}(H_2 | \mathcal{D}, \mathcal{I})} = \frac{\text{prob}(\mathcal{D} | H_1, \mathcal{I})}{\text{prob}(\mathcal{D} | H_2, \mathcal{I})} \times \frac{\text{prob}(H_1 | \mathcal{I})}{\text{prob}(H_2 | \mathcal{I})}$$

let  $H_2 = \bar{H}_1$  then

$$\text{prob}(\mathcal{D} | \mathcal{I}) = \text{prob}(\mathcal{D} | H_1, \mathcal{I}) \text{prob}(H_1 | \mathcal{I}) + \text{prob}(\mathcal{D} | \bar{H}_1, \mathcal{I}) \text{prob}(\bar{H}_1 | \mathcal{I})$$

$$\text{where } \text{prob}(H_1 | \mathcal{I}) + \text{prob}(\bar{H}_1 | \mathcal{I}) = 1$$

The problem is that we can't compute  $\text{prob}(\mathcal{D} | \bar{H}_1, \mathcal{I})$ , because even if we know that it is not  $H_1$ , we still don't know what  $\bar{H}_1$  is. We would need a specific alternative to give quantitative answers.

Usual statistics would say that we should compute the  $\chi^2$  and compare to the expected deviation  $\sqrt{N}$  for  $N$  measurements.

We could say that the misfit statistics serves the purpose to think about alternatives.

### 4.3 Examples

Common problem: classification. Archeologist finds two skeletons differing by 2 million years. He wants to know if there is significant evolutionary change. Suppose the two sites yield  $N_1$  and  $N_2$  measurements represented by  $\vec{D}_1$  and  $\vec{D}_2$ .

Consider the following hypotheses:

A: no change over period of time. Both sets characterized by some (unknown) mean  $\mu$  and standard deviation  $\sigma$ .

B: There is change with unknown  $\mu_1, \mu_2$  and  $\sigma_1, \sigma_2$ .

We need to compute the evidence. So we need to compute  $\text{prob}(\vec{D}_1, \vec{D}_2 | A, I)$  and  $\text{prob}(\vec{D}_1, \vec{D}_2 | B, I)$ . Let's start.

$$\text{prob}(\vec{D}_1, \vec{D}_2 | A, I) = \iint \text{prob}(\vec{D}_1, \vec{D}_2 | \mu, \sigma, A, I) \text{prob}(\mu, \sigma | A, I) d\mu d\sigma$$

We use a uniform prior.

$$\text{prob}(\mu, \sigma | A, I) = \frac{1}{(\mu_{\max} - \mu_{\min}) \sigma_{\max}}$$

$N \equiv N_1 + N_2$  treating as independent  $N$  measurements:

$$\text{prob}(\vec{D}_1, \vec{D}_2 | \mu, \sigma, A, I) = (\sigma \sqrt{2\pi})^{-N} \exp\left[-\frac{1}{2\sigma^2} \sum_{k=1}^N (x_k - \mu)^2\right]$$



and the  $\int d\mu d\sigma$  integral yields

$$\int_{\mu_{\min}}^{\mu_{\max}} \int_0^{\sigma_{\max}} \exp\left[-\frac{1}{2} \left( \alpha(\mu - \mu_0)^2 + \beta(\sigma - \sigma_0)^2 \right)\right] d\mu d\sigma = \frac{2\pi}{\sqrt{\alpha\beta}}$$

$$\Rightarrow \text{prob}(\vec{D}_1, \vec{D}_2 | A, I) \approx \frac{(\sigma_0 \sqrt{2\pi})^{2-N} \exp(-N/2)}{(\mu_{\max} - \mu_{\min}) \sigma_{\max} N \sqrt{2}}$$

Having computed this, we also need

$$\text{prob}(\vec{D}_1, \vec{D}_2 | B, I) = \text{prob}(\vec{D}_1 | B, I) \times \text{prob}(\vec{D}_2 | B, I)$$

where we are allowed to factorize, because our hypothesis said that the two are independent.

Again, each of the two can be written as

$$\text{prob}(\vec{D}_i | B, I) = \iint \text{prob}(\vec{D}_i | \mu_i, \sigma_i, B, I) \text{prob}(\mu_i, \sigma_i | D, I) d\mu_i d\sigma_i$$

Both priors  $i=1,2$  can be set the same, because we have no evidence for the contrary.

So the result for each looks exactly like the one for A:

$$\text{prob}(\vec{D}_i | B, I) \approx \frac{(\sigma_0 \sqrt{2\pi})^{2-N_i} \exp(-N_i/2)}{(\mu_{\max} - \mu_{\min}) \sigma_{\max} N_i \sqrt{2}}$$

Finally, we divide the two evidences:

$$\frac{\text{prob}(\vec{D}_1, \vec{D}_2 | A, I)}{\text{prob}(\vec{D}_1, \vec{D}_2 | B, I)} = \frac{(\mu_{\max} - \mu_{\min}) \sigma_{\max}}{\pi \sqrt{2}} \frac{N_1 N_2 (\sigma_0)^{2-U}}{N (\sigma_{01})^{2-U_1} (\sigma_{02})^{2-U_2}}$$

### THE LIGHT BULB EXAMPLE

Consider 2 manufacturers of light bulbs. We know that their mean life-time and variance will be different. But which are better, given some samples we received?

What we need is the posterior probability of the hypothesis  $\mu_1 > \mu_2$ :

$$\text{prob}(\mu_1 > \mu_2 | \vec{D}_1, \vec{D}_2, I) = \int_0^{\infty} d\mu_1 \int_0^{\mu_1} d\mu_2 \text{prob}(\mu_1, \mu_2 | \vec{D}_1, \vec{D}_2, I)$$

Surely, the expected lifetime factories:

$$\text{prob}(\mu_1, \mu_2 | \vec{D}_1, \vec{D}_2, I) = \text{prob}(\mu_1 | \vec{D}_1, I) \text{prob}(\mu_2 | \vec{D}_2, I)$$

If the number of samples  $N_j$  is reasonably large, then

$$\text{prob}(\mu_j | \vec{D}_j, I) = \frac{\sqrt{N_j}}{\sqrt{2\pi} S_j} \exp \left[ -\frac{1}{2} \frac{N_j (\mu_j - \mu_{0j})^2}{S_j^2} \right]$$

Where  $\mu_{0j} = \frac{1}{N_j} \sum \text{lifetimes}$

$$S_j = \frac{1}{N_j - 1} \sum (\text{lifetimes} - \mu_{0j})^2$$

So

$$\text{prob}(\mu_1 > \mu_2 | \vec{D}_1, \vec{D}_2, \vec{I}) = \int_0^{\mu_1} \int_0^{\mu_2} d\mu_2 \frac{\sqrt{N_1 N_2}}{2\pi S_1 S_2} \exp\left[-\frac{1}{2} \frac{N_1 (\mu_1 - \mu_{01})^2}{S_1^2} - \frac{1}{2} \frac{N_2 (\mu_2 - \mu_{02})^2}{S_2^2}\right]$$

Make a change of variables  $Z = \mu_1 - \mu_2$ ,

then the integral becomes

$$\text{prob}(\mu_1 > \mu_2 | \vec{D}_1, \vec{D}_2) = \frac{1}{S_2 \sqrt{2\pi}} \int_0^{\infty} \exp\left[-\frac{(Z - Z_0)^2}{2S_2^2}\right] dZ$$

$$\text{where } Z_0 = \mu_{01} - \mu_{02} \quad S_2^2 = \frac{S_1^2}{N_1} + \frac{S_2^2}{N_2}$$

So if the difference  $Z_0 = \mu_{01} - \mu_{02}$  (which we inferred from our samples) equals  $S_2$ ,

then  $\text{prob}(\mu_1 > \mu_2) = 0.84$  and so on...