

## Data analysis (Binary systems)

For binary systems we can model shape of the signal.

If the signal weak we can search this shape in the noise  $\rightarrow$  (correlation  $\rightarrow$  matched filter) matched filter. The shape (signal) which we see is called template (waveform).

- Weak signal  $\Rightarrow$  we have two hypothesis
  - (1) data = noise (~~noise~~)
  - (2) data = noise + signal  $H_1$

$\rightarrow$  model testing (hypothesis testing)

- Frequentist approach : signal is deterministic & noise corrupts the signal.

Newman - Pearson criterion ; maximization of detection statistic for a chosen significance of the test.

$\rightarrow$  Likelihood ratio (proportional to  $P(H_1)/P(H_0)$ ) - most powerful detection statistic in Gaussian noise

- we fix acceptable false alarm probability

$$P(Y \geq Y_{\text{threshold}} | H_0) = \alpha \quad - \text{probability of detection stat}$$

in  $H_0$  (data = noise) to be above threshold  $Y_{\text{threshold}} = \alpha$ .

$\rightarrow$  We choose  $\alpha \rightarrow$  determine  $Y_{\text{thr}}$ .  $\rightarrow$

$$y = \frac{P(d | H_1, \theta_i)}{P(d | H_0, \theta_0)} \leftarrow \text{likelihood} \quad \begin{aligned} &\text{detection probability } P(Y \geq Y_{\text{thr}} | H_1) \\ &\text{assigned to observed } y \end{aligned}$$

$y < Y_{\text{thr}} \rightarrow$  not a detection ;  $y \geq Y_{\text{thr}} \rightarrow$  we have significant

of detection  $P(Y \geq Y_{\text{thr}} | H_1)$  given false alarm probability

$\alpha$  (say 1% or smaller).  $\rightarrow$  maximum likelihood estimation of parameters

• Bayesian approach : We have to assign probability to our hypothesis before looking at the data  $\rightarrow$  prior probability  $P_0$ ,

We ~~then~~ treat parameters of the signal as random variables & try to estimate the probability density function  $b$  for parameter based on observation

$$P(H_i|D) = P(D|H_i) \frac{P(H_i)}{P(D)}$$

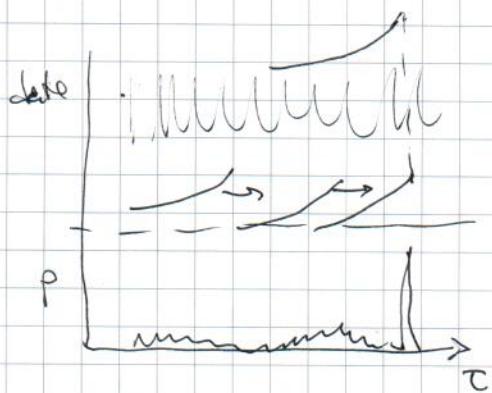
Both approaches are used in DA. Let us build likelihood f-n

~~form~~ \* Matched filter : search for signal of known form.

- assume we know the form of the signal exactly but we do not know the arrival time : correlation of the signal with the data

$$g(f) = 4 \operatorname{Re} \int \frac{\tilde{J}(f) \tilde{h}^*(f)}{S_n(f)} dt$$

$$\tilde{J}(f) = \int d(t) e^{j2\pi ft} dt$$



$\tilde{J}(f) \rightarrow$  data in Fourier domain  
 $\tilde{h}^*(f) \rightarrow$  template in Fourier domain  
 c.c.

$t \rightarrow$  time shift

But! We ~~are~~ not equally sensitive at all frequencies if signal has very low or

very high frequencies we will <sup>4</sup>see them with our detector unless signal has very strong amplitude

→ we use  $'S_n(f) \rightarrow'$  noise power Spectral density as weights at ~~different frequencies~~

$$\text{SNR}^2(f) = 4 \operatorname{Re} \int_{f_{\max}}^f \frac{\tilde{J}(f) \tilde{h}^*(f)}{S_n(f)} e^{-j2\pi ft} dt$$

← determined by detector

Allows search over time (shift)

"2"  $S_n(f) \rightarrow$  one-sided PSD (for positive freq. only)

$$(2) \rightarrow \int \frac{\tilde{J}(f) \tilde{h}^*(f) + \tilde{J}^*(f) \tilde{h}(f)}{S_n(f)} dt$$

(to have real  $e^{-j2\pi ft}$  is absorbed in  $\tilde{h}(f)$ )

In practice the signal depends on parameters  $\{\theta_i\}$  which we do not know a priori. Noise might also depend on parameters!

Assume that we know the noise & it has Gaussian distribution & noise has zero mean!  $\langle \text{noise} \rangle = \langle n \rangle = 0$

$\langle \cdot \rangle$  over noise realizations or assuming ergodicity  $\langle \cdot \rangle_t$  over time  
 $d(t) = n(t) + s(t, \theta_t)$    
 noise signal parameters

we compute  $h(\vec{\theta}_1), h(\vec{\theta}_2), h(\vec{\theta}_3), \dots$

$$d(t) - h(\vec{\theta}_i) \rightarrow \text{residuals} \quad \text{if } h(\vec{\theta}_t) = s(t, \vec{\theta}_t) \rightarrow d(t) - h(t, \vec{\theta}_t) = n(t)$$

$$\text{otherwise } d(t) - h(t, \vec{\theta}_i) = n + [s(t, \theta_t) - h(t, \vec{\theta}_i + \vec{\theta}_t)] \neq n(t) \quad \text{"deformed signal"}$$

Likelihood  $\propto n \rightarrow$  maximum likelihood  $\rightarrow$  the residuals represent the Gaussian noise

$$P(n) = \prod_{i=1}^N P(n(i) | f_{\theta_0}(\text{white, Gaussian})) = \frac{1}{(6\sqrt{2\pi})^N} e^{-\frac{1}{2\sigma^2} \sum_i n(i)^2}$$

Gaussian distribution with zero mean &  $\sigma$ -variance  
 Non white  $\rightarrow$  different variance @ different frequencies

$$P(d = n | f_{\theta_0}) \sim e^{-\frac{1}{2} \langle n | n \rangle} = e^{-\frac{1}{2} \langle d | d \rangle} \quad \text{generalization}$$

$$(a | b) = \int_{-\infty}^{+\infty} \hat{a}(f) \hat{b}^*(f) + \hat{a}^*(f) \hat{b}(f) S_n(f) \leftarrow \text{non-white}$$

$$\text{likelihood } P(d | f_{\theta_1}, \theta_1) \sim e^{-\frac{1}{2} (d - h(\theta_1))^2} \quad \text{if } h(\theta_1) = s \text{ (signal)}$$

$f$  = likelihood ratio  
 (frequentist approach)

$$\frac{P(d | f_{\theta_1}, \theta_1)}{P(d | f_{\theta_0})} = e^{+\frac{1}{2} \ell \cdot \frac{1}{2} (d - h(\theta_1))^2 - \frac{1}{2} (h(\theta_1))^2}$$

if  $d = n$  ( $\text{H}_0$ )

$$\mathbb{E} \log L = \langle (d | h(\theta_i)) - \frac{1}{2} \langle (h(\theta_i) | h(\theta_i)) \rangle \stackrel{\text{H}_0}{=} \langle n | h(\theta_i) \rangle - \frac{1}{2} \langle h | h \rangle > 0$$

$$= -\frac{1}{2} \langle (h | h) \rangle < 0$$

if  $d = n+s$

$$\langle \log L \rangle = \langle s | h(\theta_i) \rangle - \frac{1}{2} \langle h(\theta_i) | h(\theta_i) \rangle$$

$$\langle \log L \rangle_{\max \theta_i} = \left[ \text{when } \theta_i = \theta_{\text{tr}} \text{ and } s = h(\theta_{\text{tr}}) \right] = \frac{1}{2} \langle s | s \rangle = \frac{1}{2} \langle h(\theta_{\text{tr}}) | h(\theta_{\text{tr}}) \rangle > 0$$

$$2 \langle \log L \rangle_{\max \theta_i} = \text{SNR}^2$$

- maximization of the likelihood ratio (or likelihood)  $\rightarrow \hat{\{\theta\}} \rightarrow$  maximum likelihood estimators of parameters  $\rightarrow$  unbiased if averaged over noise realizations  $\langle \hat{\{\theta\}} \rangle = \{\theta_{\text{tr}}\}$
- for a given noise realization  $\{\hat{\theta}_i\} + \{\theta_{\text{tr}}\}$  - estimator  $\rightarrow$  how close  $\hat{\theta}_i$  to  $\theta_{\text{tr}}$   $\rightarrow$  depends on the strength of the signal
- the stronger signal  $\rightarrow$  less influence of the noise  $\rightarrow \hat{\theta}_i$  closer to  $\theta_{\text{tr}}$
- If  $s \neq h(\theta_{\text{tr}})$   $\rightarrow$  systematic error in modelling G-W statistic  
 $s \approx h(\theta_{\text{tr}}) \rightarrow \min_{\theta_i} |(s - h(\theta_i))| \leq |(s - h(\theta_{\text{tr}}))| \rightarrow \hat{\theta}_i \approx \theta_{\text{tr}}$   
 $|\theta_{\text{tr}} - \hat{\theta}_i| \rightarrow$  bias in parameter estimation  $\&$  we still might detect the signal  $\$$  but with bias in parameters (trying to correct mis-modelling by adjusting parameters) effectiveness

faithfulness :  $\frac{(s | h(\theta_{\text{tr}}))}{\sqrt{(s | s)(h(\theta_{\text{tr}}) | h(\theta_{\text{tr}}))}}$  - overlap  $[-1, 1]$

1 - perfect match

0 - very bad

-1  $\rightarrow \hat{\theta}_i = -s$

- maximisation of likelihood - optimization problem.

PSO, GA, grid-based methods

Consider only grid based: We want to cover parameter space  $\{\theta\}$  by grid points at equal distance between them. Grid should be not very coarse (loose signal) or bad parameter estimation or should be not very fine (computationally expensive).

Distance  $\rightarrow$  that we can see  $\theta_i$  - as coordinates on the parameter space  $\rightarrow$  introduce interval & metric

$$\begin{aligned} ds^2 &= \|\hat{h}(\theta_i + d\theta_i) - \hat{h}(\theta_i)\| \approx \text{assume } d\theta_i \text{ is infinitesimal} \approx \\ &= (\hat{h}(\theta_i + d\theta_i) - \hat{h}(\theta_i))^\top (\hat{h}(\theta_i + d\theta_i) - \hat{h}(\theta_i)) \approx \left( \frac{\partial \hat{h}}{\partial \theta_i} \mid \frac{\partial \hat{h}}{\partial \theta_i} \right) d\theta_i^2 \\ &= g_{ij} dx^i dx^j \end{aligned}$$

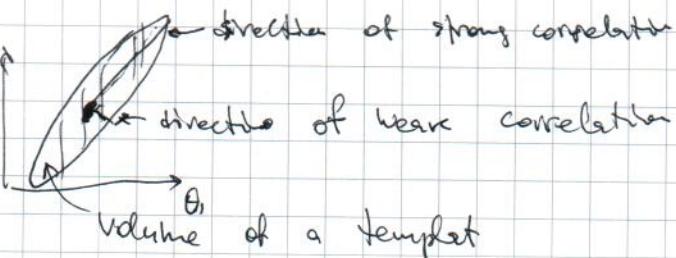
$\downarrow$  metric on the parameter space!

$$g_{ij} = \left( \frac{\partial \hat{h}}{\partial \theta_i} \mid \frac{\partial \hat{h}}{\partial \theta_j} \right)$$

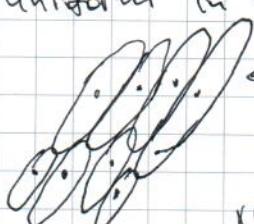
$$\hat{h}: (\hat{h}^\top \hat{h}) = 1 \text{ (normalized templates)}$$

$ds^2 \rightarrow$  tells us how similar (correlated) two nearby templates  $\downarrow$  in parameter space

$$\begin{aligned} ds^2(\theta) &: \\ &\geq 0.9 \\ &= 0.01 \end{aligned}$$



at correlation between two nearby templates  $\rightarrow$  we want grid to be uniform in this distance



$\leftarrow$  ellipses are fns of central point or in geometrical language  $g_{ij} = g_{ij}(\theta_i)$

Non-trivial problem in  $N_{\text{dim}} > 2$

Need coordinates (parametrization = combination of physical parameters such that coordinates are close to Cartesian):

$m_1 \& m_2 \rightarrow$  really bad, the  $\frac{m_1}{m_2} \rightarrow$  better  $\rightarrow \dots$

Assume we have solved placement problem  $\Rightarrow$  set of grid points in parameter space  $\{\theta_i\}$  → generate GW signal @ each point.

$$(1 - ds) = 0.03 \rightarrow (\hat{h}(\theta_i) | \hat{h}(\theta_i + \delta\theta_i)) = 0.97 \rightarrow < 10\% \text{ loss in event rate}$$

→ Template bank → compute likelihood  $\sim (d | h(\theta_i)) - \frac{1}{2} (h(\theta_i) | h(\theta_i + \delta\theta_i))$   
at each point → find maximum over  $\{\theta_i\}$

→ low latency strategy for searching GW signal in LIGO/Virgo & estimation of significance of a candidate.

- Bayesian analysis: expensive computationally → often used when we know that  $H_1$  is true. The power of this method → allows testing several models (different signal models, non-GW features)

Consider several models  $M_i$  each parametrized by own set of parameters  $\vec{\theta}_i$  (vector of parameters for  $i$ -th model)

$$P(M_i | d) = \frac{P(d | M_i) \pi(M_i)}{P(d) \pi(M_i)}$$

& for a given model the posterior pdf:

$$\text{posterior } P(\vec{\theta}_i | d, M_i) = \frac{P(d | \vec{\theta}_i, M_i) \pi(\vec{\theta}_i)}{P(d | M_i)} \leftarrow \begin{array}{l} \text{likelihood} \\ \leftarrow \text{prior on parameters} \\ \leftarrow \text{evidence of } M_i \end{array}$$

$$P(M_i | d) = \int d\vec{\theta}_i P(d | \vec{\theta}_i, M_i) \pi(\vec{\theta}_i) \quad \left| \begin{array}{l} \text{normalization} \\ \text{of pdf } \int p(\vec{\theta}_i) d\vec{\theta}_i \end{array} \right.$$

$$\text{then } P(M_i | d) = \left[ \int d\vec{\theta}_i P(d | \vec{\theta}_i, M_i) \pi(\vec{\theta}_i) \right] \frac{\pi(M_i)}{P(d)} \leftarrow \begin{array}{l} \text{probability} \\ \text{of model } M_i \end{array}$$

This problem is to evaluate  $P(d)$  in rare cases when we have limited number of models that are mutually exclusive  $P(d) = \sum_i P(d | M_i) \pi(M_i) \rightarrow$  hardly ever possible

⇒ usually consider posterior odd ratio: →

$$O_{ij}(d) = \frac{P(M_i|d)}{P(M_j|d)} = \underbrace{\left[ \frac{\int d\vec{\theta}_i P(d|\vec{\theta}_i, M_i) \pi(\vec{\theta}_i)}{\int d\vec{\theta}_j P(d|\vec{\theta}_j, M_j) \pi(\vec{\theta}_j)} \right]}_{\text{Bayes factor}} \underbrace{\frac{\pi(M_i)}{\pi(M_j)}}_{\text{prior odds}}$$

If we do not have reference of one model over other → we assign equal  $\pi$ -priors to models → Bayes factor tells us if the data prefers one or another model.

The problem is to set the threshold on the Bayes factor besides common sense: > 10 - 20 → likely  
 > 100 strong evidence > 1000 definite, few - mat conclusion  
 In reality (weak signal, small deviations from GR) - we do not expect very high Bayes factor: often used frequentist approach to set the threshold.

→ Problem to evaluate  $p(d|M_i) = \int p(d|\vec{\theta}, M_i) \pi(\vec{\theta}|M_i) d\vec{\theta}$

→ multidimensional integral

& <sup>another</sup> ~~permitted~~ problem to get the posterior:

$$p(\vec{\theta} | d, M_i) = \frac{p(d|\vec{\theta}, M_i) \pi(\vec{\theta}|M_i)}{p(d|M_i)}$$



PDF for all parameters → can marginalize for each parameter & say about mean, median, confidence interval, can compare with prior → how informative were observations (did they update our prior knowledge)

How to ~~say~~ get  $p(\vec{\theta}|d, M_i)$  &  $p(d|M_i)$

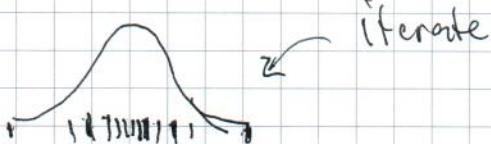
$p(d\vec{m})$  - evaluation of multidimensional integrals

- Vegas algorithm  $\rightarrow$  monte-carlo integration

$$I = \int f(\vec{x}) d\vec{x} \approx \frac{V}{N} \sum_{i=1}^N f(\vec{x}_i)$$

$$\sigma_I^2 = \frac{1}{N-1} \sum_{i=1}^N (f(\vec{x}_i) - \bar{f})^2$$

Divide parameter space into equal hypercubes, choose random points in each hypercube (box), evaluate integral, resize cubes according to density  $g = \frac{|f|}{I(f)}$ : make boxes smaller (denser) in the parts of parameter space where we have largest contribution to the integral. In less dense (larger) other parts



- Markov chain Monte Carlo

We construct Markov chain: stochastic process where the next point  $\vec{\theta}_{k+1}$  in the chain depends only on previous point  $\vec{\theta}_k$ .

(i) We want chain to move towards the region of parameter space with high likelihood

(ii) We ~~want~~ to introduce transition probability

$P(\vec{\theta}_{k+1} | \vec{\theta}_k) \rightarrow$  way of moving chain from  $\vec{\theta}_k$  to  $\vec{\theta}_{k+1}$ . (iii) If transition probability satisfies the balance equation:

$$\Lambda(\vec{\theta}_k) P(\vec{\theta}_{k+1} | \vec{\theta}_k) = \Lambda(\vec{\theta}_{k+1}) P(\vec{\theta}_k | \vec{\theta}_{k+1}) \text{ then}$$

after some "burn-in" length chain samples from the distribution  $\Lambda(\vec{\theta}_k)$

Consider a particular implementation (Metropolis & Hastings)

→ particular way of building the transitional probability which satisfies balance equation. (i) introduce a proposal distribution  $q(\vec{\theta}_{k+1} | \vec{\theta}_k)$  (could be arbitrary) & (ii) introduce the acceptance probability:

$$\alpha(\vec{\theta}_{k+1} | \vec{\theta}_k) = \min \left\{ 1, \frac{\underbrace{p(d | \vec{\theta}_{k+1})}_{\text{likelihood ratio}} \underbrace{q(\vec{\theta}_k | \vec{\theta}_{k+1})}_{\text{prior of }} \underbrace{\pi(\vec{\theta}_{k+1})}_{\text{prior}}}{\underbrace{p(d | \vec{\theta}_k)}_{\text{likelihood ratio}} \underbrace{q(\vec{\theta}_{k+1} | \vec{\theta}_k)}_{\text{prior of }} \underbrace{\pi(\vec{\theta}_k)}_{\text{prior}}} \right\}$$

\* Easier to understand if we assume uniform priors & use symmetric proposal distribution ( $q(\vec{\theta}_k | \vec{\theta}_{k+1}) = q(\vec{\theta}_{k+1} | \vec{\theta}_k)$ ) - Metropolis ratio →

$$\tilde{\alpha}(\vec{\theta}_{k+1} | \vec{\theta}_k) = \min \left\{ 1, \frac{p(d | \vec{\theta}_{k+1})}{p(d | \vec{\theta}_k)} \right\}$$

$\tilde{\alpha} \rightarrow$  probability of acceptance  $\vec{\theta}_{k+1}$ : if  $p(d | \vec{\theta}_{k+1}) > p(d | \vec{\theta}_k)$  point is always accepted, if  $p(d | \vec{\theta}_{k+1}) < p(d | \vec{\theta}_k)$  point is accepted with probability  $\tilde{\alpha}$  (likelihood ratio) → (draw random number  $x \in U(0, 1)$  accept if  $x \geq \tilde{\alpha}$ ) → chain moves towards the maximum likelihood, but (!) it can also go down the slope.



However, the theorem tells us that the chain will sample from posterior pdf (after some burn-in length) however of the proposal distribution  $q(\vec{\theta}_{k+1} | \vec{\theta}_k)$  however

- ① the efficiency of the sampling very strongly depends on the proposal: \* ideally proposal is similar / close to the posterior, if we know multimodality → take it <sup>into</sup> proposal proposal can be as complex as you want

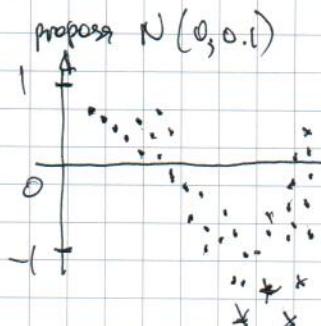
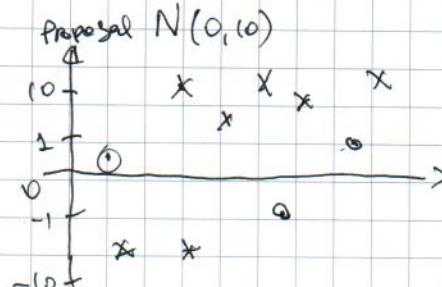
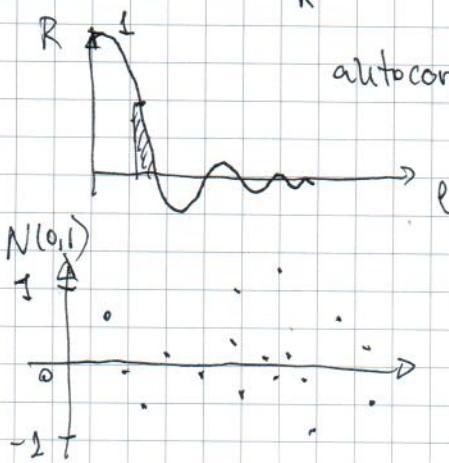
# of independent samples (uncorrelated) : defined by auto correlation length of the chain

$$R_{xx}(l) = \sum_n x(n)x^*(n-l)$$

$$\int_{-\infty}^{\infty} f(u) f^*(u+\tau) du = \int_{-\infty}^{\infty} f(u) f^*(u-\tau) du$$

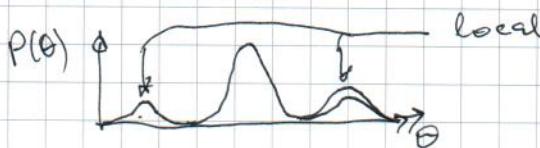
autocorrelation length  $\rightarrow$  not unique definition

$$R = \frac{1}{2}; R = 0; R = \frac{1}{e}$$



high acceptance  
but highly correlated chain  
(large autocorrelation length)  
low computational efficiency

## ② Multimodal posterior distribution



maxima  $\leftarrow$  chain can get stuck there, it will move  $\rightarrow$  find the global / all maxima if we wait > 100 years

### - Simulated annealing

$$\text{likelihood} \rightarrow \int L \propto e^{-\frac{1}{2} \langle h | h \rangle}$$

$\leftarrow$  similar to Maxwell distribution

$$f(E) \sim \exp\left[-\frac{E}{kT}\right]$$

$$f(v) \sim e^{-\frac{mv^2}{2kT}}$$

introduce "temperature"

$$e^{\frac{-\langle d| h(\theta) - \frac{1}{2} \langle h | h \rangle}{kT}} \sim e^{-\frac{E''}{kT}}$$

$T' = 1 \rightarrow$  likelihood we introduce  $T > 1 \rightarrow$  "heating the likelihood"  $\rightarrow$  makes likelihood surface smoother by increasing the noise level

$f(\theta)$  $T=1$  $f(\theta)$  $T \approx 100$ 

← easy to find near maximum

Simulated annealing: we start with high temperature  $\rightarrow$  can explore large part of parameter space (random walk) until we find the largest maximum in the likelihood  $\rightarrow$  temperature is gradually decreases to 0.

### - parallel tempering

We run multiple chains (always good) with in parallel "with different temperature  $T_i$  & make chains "talk" to each other. We can explore parameter space simultaneously @ different scales (high temper.  $\rightarrow$  large scale exploration, mid temper.  $\rightarrow$  mid scale, low temper.  $\rightarrow$  vicinity of maximum).

We perform cross-talk between chains by exchanging the points  $\vec{\theta}_k^i \leftrightarrow \vec{\theta}_k^j$  (or equivalently swap the temperature) with probability

$$p = \min \left\{ 1, e^{(E_i - E_j) \left( \frac{1}{T_i} - \frac{1}{T_j} \right)} \right\}$$

$$= \frac{\exp \left( -\frac{E_j}{T_i} - \frac{E_i}{T_j} \right)}{\exp \left( -\frac{E_i}{T_i} - \frac{E_j}{T_j} \right)}$$

& we define the frequency of cross-talk (how often we swap attempt to swap)

In addition parallel tempering allows us to evaluate the evidence  $\rightarrow$  introduce  $\beta = \frac{1}{T}$ , The evidence corresponding to each chain with the temperature  $T$  is

$$Z(\beta) = \int d\vec{\theta} L(\vec{\theta}, T) \pi(\vec{\theta}) \quad (\text{T assumed to be continuous variable})$$

We are interested in chain with  $T=1$  (true likelihood!)

$\rightarrow Z(1) \leftarrow$  evidence we are interested in. Take derivative!

$$\frac{d \log Z(\beta)}{d\beta} = \frac{1}{Z(\beta)} \int \frac{\partial}{\partial \beta} L(\theta, \beta) \pi(\theta) d\theta \rightarrow \log \langle L(\theta, T=1) \rangle - \bar{L}(\theta, \beta)$$

$$\langle f(x) \rangle = \frac{\int f(x) p(x) dx}{\int p(x) dx} \leftarrow \text{average over pdf } p(x) \text{ - normalization}$$

$$\frac{d \log Z(\beta)}{d\beta} = \langle \log \bar{L}(\theta, T=1) \rangle_\beta \leftarrow \text{averaging over posterior at temperature } T=\frac{1}{\beta}$$

Integrating  $\int_0^1 \frac{d \log Z(\beta)}{d\beta} d\beta = \log Z(1) \Big|_0^1 ; Z(0) = \int \pi(\theta) d\theta = 1$

$$\rightarrow \log Z(1) = \int_0^1 \langle \log \bar{L}(\theta, T=1) \rangle_\beta d\beta \text{ & can be evaluated numerically if we have}$$

properly chosen Temperature ladder.

(Thermodynamic integration)

- Nested models  $\vec{\theta}_0$  is subset of  $\vec{\theta}$  (descend from old non-splitting BMs)

$$B \approx \frac{P(\vec{\theta} = \vec{\theta}_0 | d)}{\pi(\vec{\theta}_0)}$$

