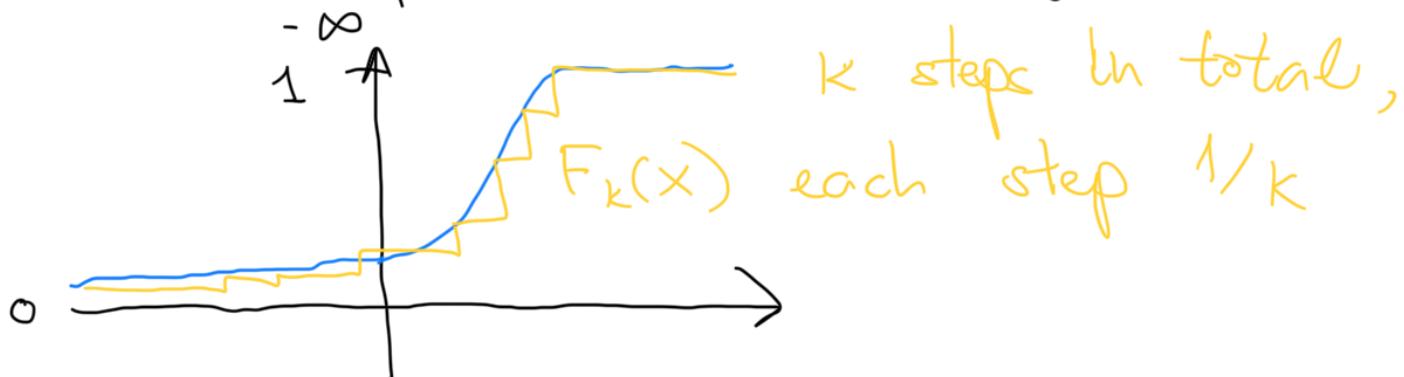


# Sampling algorithms . Monte Carlo

$x \in \mathbb{R}$ ,  $p(x)$  - a distribution of  $x$

$$F(t) = \int p(x) dx - \text{CDF of } x$$



## 1. Empirical distribution function

## Sum of delta functions

$$F_k(x) = \frac{1}{k} \sum_{i=1}^k I[x > x_i]$$

2. Dense steps for higher  $\frac{dF(t)}{dt} = p(t)$

3. Kolmogorov - Smirnov test:

$$D_k = \sup_{x \in \mathbb{R}} |F_k(x) - f(x)|$$

$\sqrt{k} D_k \rightarrow$  Kolmogorov distribution  
convergence in distribution

$$\underline{\text{Ex.1}} \quad p(y | \vec{x}) = \int p(y | \vec{x}, \vec{\theta}) p(\vec{\theta} | D) d\vec{\theta} \approx$$

↓      ↓      ↗  
 output    input      posterior

$$\approx \frac{1}{K} \sum_{i=1}^K p(y_i | \vec{x}_i, \vec{\theta}_i), \vec{\theta}_i \sim p(\vec{\theta}|D)$$

- unbiased estimate

- variance decreases  $O\left(\frac{1}{k}\right)$

$$\mathbb{E} f(x) = \int f(x) p(x) dx$$

$$\mathbb{E}_k f(x) = \int f(x) p_k(x) dx = \frac{1}{k} \sum_{i=1}^k f(x_i)$$

$$p_k(x) = \frac{1}{k} \sum_{i=1}^k \delta(x - x_i)$$

$$\begin{aligned} \mathbb{E} (\mathbb{E} f(x) - \mathbb{E}_k f(x))^2 &= \mathbb{E} (\mathbb{E} f(x) - \frac{1}{k} \sum_{i=1}^k f(x_i))^2 = \\ &= \mathbb{E} \left[ \frac{1}{k} \sum_{i=1}^k (\mathbb{E} f(x) - f(x_i)) \right]^2 = \frac{1}{k^2} k \mathbb{E} (\mathbb{E} f(x) - f(x_i))^2 + \\ &+ c \sum_{i \neq j} \mathbb{E} (\mathbb{E} f(x) - f(x_i)) (\mathbb{E} f(x) - f(x_j)) \xrightarrow{0} = \\ &= \frac{1}{k} \text{Var}_k (f(x)) \end{aligned}$$

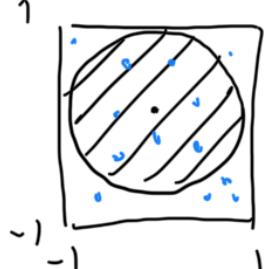
- We also need estimates for Variational inference

Ex.2 Find  $\pi$ , toy problem

$$\Theta_1, \Theta_2 \sim U([-1, 1]), \text{ ind.}$$

$$V_i = \left[ \left\| \begin{matrix} \Theta^{(i)} \\ m \end{matrix} \right\|_2 \leq 1 \right]$$

$$S_m = \frac{1}{m} \sum_{i=1}^m V_i \approx \frac{\pi}{4}$$



How to sample from  $p(\vec{\theta}|D)$ ?

We assume that we can sample from

Ex.0  $\sim U([0, 1])$  - uniform distribution with finite support

But typically  $p(\vec{\theta}|D)$  is more complex

Ex.1 Exponential distribution |  $p(x) = \lambda \exp(-\lambda x)$

- inverse sampling

$$x = -\frac{\ln u}{\lambda}$$

$$/ 1 - e$$

quantile

$$x \in \mathbb{R}$$

$F_x$  - CDF

$$F_x(x) = 1 - e^{-\lambda x}$$

$$\left| \begin{array}{l} u \in U(0,1) \\ F_x^{-1}(u) \sim p(x) \end{array} \right.$$

Ex. Gaussian distribution,

Box-Muller transformation

$$u_1 \sim U(0,1), u_2 \sim U(0,1)$$

$$\text{polar coordinates : } z = \sqrt{-2 \ln(u_1)} \\ \varphi = 2\pi u_2$$

$$\Rightarrow z \in (0, +\infty), \varphi \in (0, 2\pi)$$

$$\text{Cartesian coordinates : } x = z \sin \varphi \\ y = z \cos \varphi$$

$$\varphi \sim U(0, 2\pi)$$

$z = g(u_1)$  - monotonic one to one transf.

$$\ln(u_1) = -\frac{z^2}{2}, \quad u_1 = \exp(-\frac{z^2}{2}) = g^{-1}(z)$$

$$P(z_1 < z < z_2) = \int_{g^{-1}(z_1)}^{g^{-1}(z_2)} f_{u_1}(s) ds = \\ = - \int_{g^{-1}(z_1)}^{g^{-1}(z_2)} ds = \int_{z_1}^{z_2} -t \exp(-t^2/2) dt, \quad t = \sqrt{-2 \ln s}$$

$$\Rightarrow f_R(z) = z \exp(-z^2/2)$$

$$x = g_x(z, \varphi), \quad y = g_y(z, \varphi) \quad \boxed{A} \quad \begin{matrix} y_2 \\ g_y \\ x_1 \\ x_2 \end{matrix}$$

$$P(x_1 < x < x_2, y_1 < y < y_2) =$$

$$= \iint_A f_{x,y}(u, v) du dv = \iint_A f_z(u) f_\varphi(v) du dv =$$

$$= \iint_A u \exp(-u^2/2) \frac{1}{2\pi} du dv \quad (1)$$

$$p = u \cos(v), \quad q = u \sin(v)$$

$$u = \sqrt{p^2 + q^2}, \quad v = \arctan(q/p)$$

$$dudv = \left| \det \begin{bmatrix} \frac{\partial(u,v)}{\partial(p,q)} \end{bmatrix} \right| dp dq = \frac{1}{\sqrt{p^2 + q^2}} dp dq$$

$$P(x_1 < x < x_2, y_1 < y < y_2) = \iint_{x_1, y_1}^{x_2, y_2} \frac{\sqrt{p^2 + q^2}}{2\pi} dp dq \quad \text{from (1)}$$

$$\cdot \underbrace{\exp(- (p^2 + q^2)/2)}_{x_2} \frac{1}{\sqrt{p^2 + q^2}} dp dq = \underbrace{\frac{1}{\sqrt{2\pi}}}_{y_2} \exp(- p^2/2) dp \underbrace{\int_{y_1}^{y_2} \frac{1}{\sqrt{2\pi}} \exp(- q^2/2) dq}_{\text{change of var.}} =$$

$\Rightarrow x, y$  are independent and

$$f_x(x) = \frac{1}{\sqrt{2\pi}} \exp(-x^2/2)$$

$$f_y(y) = \frac{1}{\sqrt{2\pi}} \exp(-y^2/2)$$

For other distributions we can get

something similar, if we know their analytical form, but

in general for  $p(\vec{\theta} | D)$  we know only numerator and it's hard to design something similar

Importance Sampling

$$\mathbb{E}_p f(\vec{\theta}) \approx \sum_{k=1}^K p(\vec{\theta}_k) f(\vec{\theta}_k), \quad \vec{\theta}_k \sim \text{uniform}$$

$$\mathbb{E}_p f(\vec{\theta}) = \int f(\vec{\theta}) p(\vec{\theta}) d\vec{\theta} = \int f(\vec{\theta}) \frac{p(\vec{\theta})}{q(\vec{\theta})} q(\vec{\theta}) d\vec{\theta} =$$

$$= \mathbb{E}_q f(\vec{\theta}) \frac{p(\vec{\theta})}{q(\vec{\theta})} \approx \frac{1}{K} \sum_{i=1}^K f(\vec{\theta}_i) \frac{p(\vec{\theta}_i)}{q(\vec{\theta}_i)},$$

We use the proposal distribution  $q(\vec{\theta})$

$$\vec{\theta}_i \sim q(\vec{\theta}), \quad z_i := \frac{p(\vec{\theta}_i)}{q(\vec{\theta}_i)} \text{ - importance weights}$$

But typically we know  $p(\vec{\theta})$  up to normalization constant (e.g.  $p(\vec{\theta}|D)$ )

$$p(\vec{\theta}) = \frac{\tilde{p}(\vec{\theta})}{Z_p} \leftarrow \text{normalization constant}$$

$$\mathbb{E}_p f(\vec{\theta}) = \int f(\vec{\theta}) p(\vec{\theta}) d\vec{\theta} =$$

$$= \frac{Z_q}{Z_p} \frac{1}{K} \sum_{i=1}^K \tilde{z}_i f(\vec{\theta}_i), \quad \vec{\theta}_i \sim q(\vec{\theta}) = \frac{\tilde{q}(\vec{\theta})}{Z_q},$$

$$\tilde{z}_i = \frac{\tilde{p}(\vec{\theta}_i)}{\tilde{q}(\vec{\theta}_i)}$$

$$\frac{Z_p}{Z_q} = \frac{1}{Z_q} \int \tilde{p}(\vec{\theta}) d\vec{\theta} = \int \frac{\tilde{p}(\vec{\theta})}{\tilde{q}(\vec{\theta})} q(\vec{\theta}) d\vec{\theta}$$

$$\approx \frac{1}{K} \sum_{i=1}^K \tilde{z}_i$$

$\Rightarrow$

$$\mathbb{E}_p f(\vec{\theta}) \approx \sum_{i=1}^K \omega_i f(\vec{\theta}_i), \quad \omega_i = \frac{\tilde{z}_i}{\sum_{j=1}^K \tilde{z}_j}$$

- big variance

- no diagnostics

-  $q(\vec{\theta})$  should be close to  $p(\vec{\theta})$

# Mazkov Chain Monte Carlo

we don't know

Metropolis algorithm : sampling from  
 $p(\vec{\theta}) = \frac{\tilde{p}(\vec{\theta})}{Z_p}$

1. Initialize  $\vec{\theta}^{(0)}$

2. Repeat  $K$  times

2.1. Sample  $\vec{\theta}^*$  from  
the proposal distribution  $q(\vec{\theta}^* | \vec{\theta}^{(t)})$

2.2. Accept  $\vec{\theta}^*$  with probability

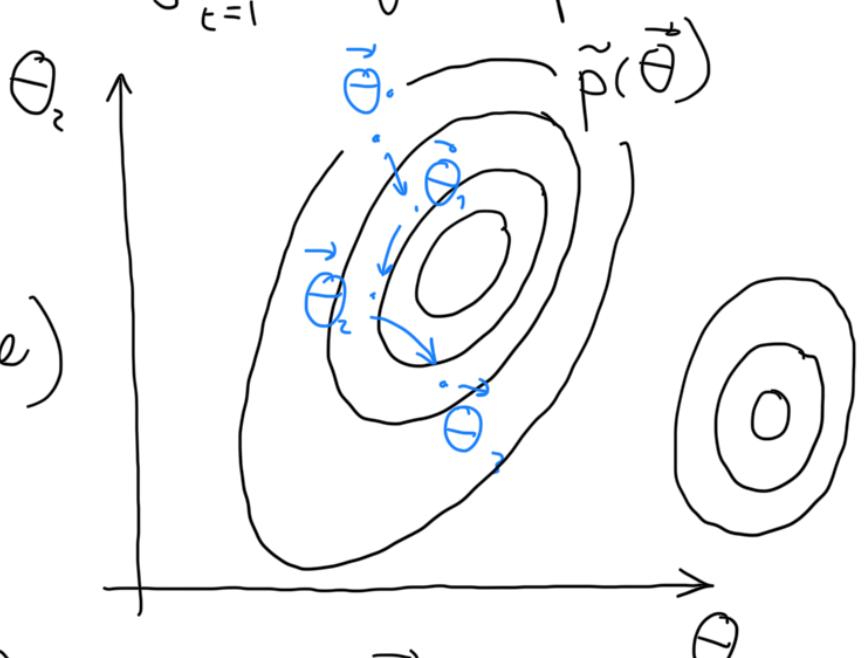
$$a = A(\vec{\theta}^*, \vec{\theta}^{(t)}) = \min(1, \frac{\tilde{p}(\vec{\theta}^*)}{\tilde{p}(\vec{\theta}^{(t)})})$$

$$\vec{\theta}^{(t+1)} = \begin{cases} \vec{\theta}^{(t)}, & u < a \\ \vec{\theta}^*, & u > a \end{cases}, \quad u \in U[0, 1]$$

Now we have  $\{\vec{\theta}_t\}_{t=1}^K$  of dependent

samples (we can just take

every  $m$ -th sample)



$$1. \frac{\tilde{p}(\vec{\theta})}{\tilde{p}(\vec{\theta}')} = \frac{p(\vec{\theta}) \cdot Z_p}{p(\vec{\theta}') \cdot Z_p} = \frac{p(\vec{\theta})}{p(\vec{\theta}')}$$

$$2. q(\vec{\theta} | \vec{\theta}') = q(\vec{\theta}' | \vec{\theta}) - \text{symmetric}$$

If  $q(\vec{\theta}' | \vec{\theta})$  is positive for all  $\vec{\theta}, \vec{\theta}'$ ,

then the distribution of  $\vec{\theta}_t$  goes

to  $p(\vec{\theta})$  for  $t \rightarrow \infty$  (sufficient cond.)

$$\text{Ex. } q(\vec{\theta}' | \vec{\theta}) = \mathcal{N}(\vec{\theta}', \beta^2 I_d) =$$

$$= \frac{1}{(\sqrt{2\pi} \beta)^d} \exp \left( -\frac{1}{2\beta^2} \|\vec{\theta} - \vec{\theta}'\|_2^2 \right)$$

We have a single parameter  $\beta^2$

$\beta^2$  is small - high correlation between  
 $\vec{\theta}^{(t)}, \vec{\theta}^{(t+1)}$

$\beta^2$  is large - low acceptance prob.

### Metropolis algorithm II (more general)

$\tilde{p}(\vec{\theta})$  - unnormalized density  
we can calculate

$q(\vec{\theta}' | \vec{\theta})$  - (unsymmetric) proposal distribution

1. Initial  $\vec{\theta}^{(0)}$  from some dist.

2. K iterations

2.1.  $\vec{\theta}^* \sim q(\vec{\theta}' | \vec{\theta}^{(t)})$

2.2. Accept with probability  
 $A(\vec{\theta}^*, \vec{\theta}^{(t)}) = \min(1, \frac{\tilde{p}(\vec{\theta}^*) q(\vec{\theta}^{(t)} | \vec{\theta}^*)}{\tilde{p}(\vec{\theta}^{(t)}) q(\vec{\theta}^* | \vec{\theta}^{(t)})})$

Let's prove, that this

approach is correct

*cancels if  $q(\vec{\theta}' | \vec{\theta})$  is symmetric*

### Markov chains

$\vec{\theta}^{(1)}, \dots, \vec{\theta}^{(t)}, \vec{\theta}^{(t+1)}$  - a series of random variables

Def. A first order Markov chain is a series of random variables s.t.

$$D(\vec{A}^{(t+1)} | \vec{\theta}^{(t)}) = D(\vec{\theta}^{(t+1)} | \vec{\theta}^{(t)})$$

the conditional density depends only on the previous state.

Ex.

A. independent  $\vec{\theta}^{(t)}$

$$p(\vec{\theta}^{(t+1)} | \vec{\theta}^{(1)}, \dots, \vec{\theta}^{(t)}) = p(\vec{\theta}^{(t+1)})$$

B. Gaussian Random Walk

$$p(\vec{\theta}^{(t+1)} | \vec{\theta}^{(1)}, \dots, \vec{\theta}^{(t)}) = N(\vec{\theta} | \vec{\theta}^{(t)}, \sigma^2 I)$$

How to fully specify Markov Chain

1. Initial distribution  $p(\vec{\theta}^{(0)})$

2. Transition distribution

$$p(\vec{\theta}^{(m+1)} | \vec{\theta}^{(m)}) = T_m(\vec{\theta}^{(m)}, \vec{\theta}^{(m+1)})$$

If there's no dependence on  $m \Rightarrow$

homogeneous Markov Chains.

We'll consider only homogeneous M.C.

Marginal distribution for  $\vec{\theta}^{(m+1)}$ :

$$p(\vec{\theta}^{(m+1)}) = \sum_{\vec{\theta}^{(m)}} p(\vec{\theta}^{(m+1)} | \vec{\theta}^{(m)}) p(\vec{\theta}^{(m)})$$

Def. A distribution is stationary (invariant), if each step in the chain doesn't change it

$$p^*(\vec{\theta}) = \sum_{\vec{\theta}'} T(\vec{\theta}', \vec{\theta}) p^*(\vec{\theta}')$$

Ex. Identity transitions  $\Rightarrow$  all are invariant

A sufficient condition is called the detailed balance property for  $p^*(\vec{\theta})$

$$p^*(\vec{\theta}) p(\vec{\theta}' | \vec{\theta}) = p^*(\vec{\theta}') p(\vec{\theta} | \vec{\theta}')$$

Stat. The invariance property follows from

the detailed balance property

$$\sum_{\vec{\theta}'} p(\vec{\theta} | \vec{\theta}') p^*(\vec{\theta}') = \sum_{\vec{\theta}} p(\vec{\theta}' | \vec{\theta}) p^*(\vec{\theta}) = \\ = p^*(\vec{\theta}) \sum_{\vec{\theta}'} p(\vec{\theta}' | \vec{\theta}) = p^*(\vec{\theta})$$

A Markov chain that respects the detailed balance is said to be reversible.

We design a Markov chain for sampling :

1. The desired distribution is invariant

2. Irrespective of the choice of the initial distribution

$$p(\vec{\theta}^{(t)}) \rightarrow p^*(\vec{\theta}), \text{ as } t \rightarrow \infty$$

Ergodicity holds.

M.C. will be ergodic under weak restrictions on the invariant dist. and transition probabilities.

M.-H. invariant distribution

accept. probability  $A_k(\vec{\theta}^*, \vec{\theta}^{(t)}) = \min \left( 1, \frac{\tilde{p}(\vec{\theta}^*) q_k(\vec{\theta}^{(t)} | \vec{\theta}^*)}{\tilde{p}(\vec{\theta}^{(t)}) q_k(\vec{\theta}^* | \vec{\theta}^{(t)})} \right)$  k possible options

Stat.  $p(\vec{\theta})$  is an invariant distribution of M.C.

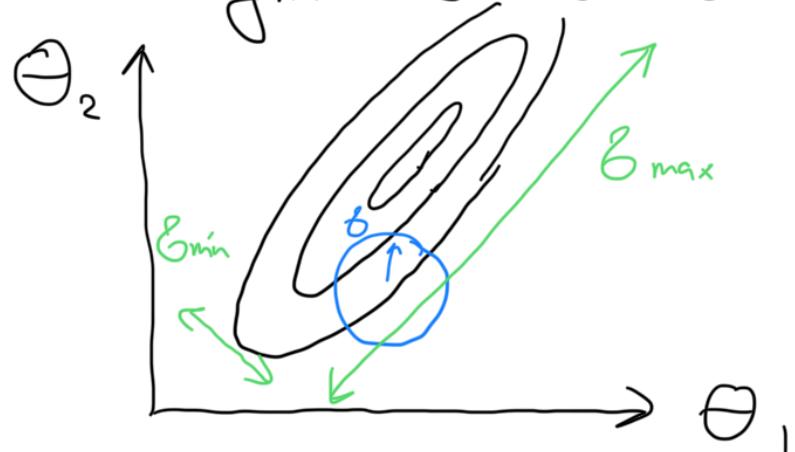
We'll check the detailed balance

$$p(\vec{\theta}) q_k(\vec{\theta} | \vec{\theta}') A_k(\vec{\theta}', \vec{\theta}) =$$

$$\begin{aligned}
 &= \min(p(\vec{\theta}) q_k(\vec{\theta} | \vec{\theta}'), p(\vec{\theta}') q_k(\vec{\theta} | \vec{\theta}')) = \\
 &= \min(p(\vec{\theta}') q_k(\vec{\theta} | \vec{\theta}'), p(\vec{\theta}) q_k(\vec{\theta}' | \vec{\theta})) = \\
 &= p(\vec{\theta}') q_k(\vec{\theta} | \vec{\theta}') A_k(\vec{\theta}, \vec{\theta}') \quad \blacktriangleleft
 \end{aligned}$$

Ergodicity - also holds

Ex. Non-symmetric Gaussian



number of steps  
to have "independent"  
samples:  
 $(\epsilon_{\max}, \epsilon_{\min})^2$

Gibbs sampling

$$p(\vec{\theta}) = p(\theta_1, \dots, \theta_d)$$

$$1. \text{ initialize } \vec{\theta}^{(0)} = (\theta_1^{(0)}, \dots, \theta_d^{(0)})$$

$$2. \tau = 1, \dots, T$$

- sample  $\theta_1^{(\tau+1)} \sim p(\theta_1 | \theta_2^{(\tau)}, \dots, \theta_d^{(\tau)})$
- sample  $\theta_2^{(\tau+1)} \sim p(\theta_2 | \theta_1^{(\tau+1)}, \theta_3^{(\tau)}, \dots, \theta_d^{(\tau)})$
- ⋮
- sample  $\theta_d^{(\tau+1)} \sim p(\theta_d | \theta_1^{(\tau+1)}, \theta_2^{(\tau+1)}, \dots, \theta_{d-1}^{(\tau+1)})$
- Markov chain type sampling

Stat. Is it invariant

for  $p(\vec{\theta})$

☒  $p(\vec{\theta}')$  - we change 1st component

We'll prove, that marginal and conditional  
are the same

$$\sim \vec{\theta}'_1 \sim \vec{\theta}'_1 \vec{\theta}'_2 \sim \vec{\theta}'_2 \sim \dots \sim \vec{\theta}'_d \sim \dots \sim \vec{\theta}'_1 \vec{\theta}'_2 \dots$$

$$p(\nu) = p(\nu_1 | \sigma_1) p(\sigma_1) = p(\nu_1 | \sigma_1) \cdot p(\vec{\theta}_1) = p(\vec{\theta})$$

Stat. It is a Metz.-Hast. algorithm

We consider

$$q_k(\vec{\theta}^* | \vec{\theta}) = p(\theta_{1k}^* | \vec{\theta}_{1k})$$

- conditional

$$A_k(\vec{\theta}^*, \vec{\theta}) = \frac{p(\vec{\theta}^*) q_k(\vec{\theta} | \vec{\theta}^*)}{p(\vec{\theta}) q_k(\vec{\theta}^* | \vec{\theta})} = \frac{p(\theta_{1k}^* | \vec{\theta}_{1k}) p(\vec{\theta}_{1k})}{p(\theta_{1k} | \vec{\theta}_{1k}) p(\vec{\theta}_{1k})} = 1$$

## Design choices MCMC

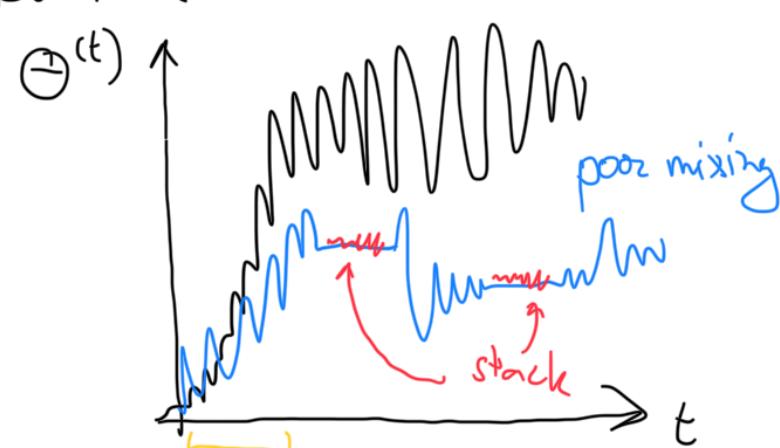
thinning: • skip steps, more independence

discarding: • burn-in period: don't depend on the initial distribution

adaptive proposals

multi-start

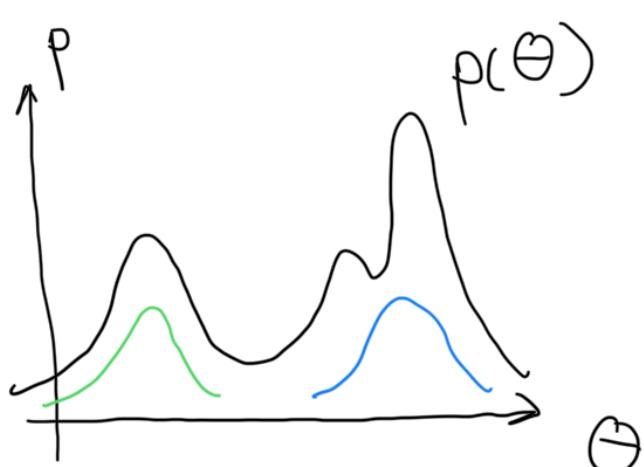
- global, local maxima
- proposal distribution



## Diagnosis of MCMC

- eye test
- formal tests :

- Geweke
- Heidelberg-Welch
- Raftery-Lewis



Adaptive proposal :



$t$

Next time:

- Hamiltonian Monte - Carlo
- NUT