

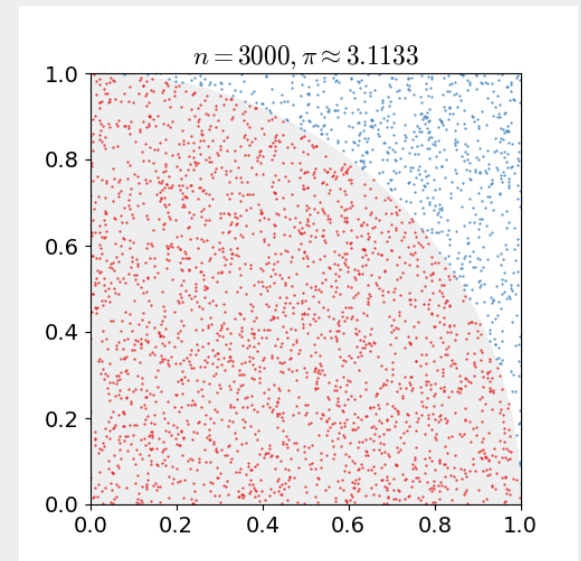
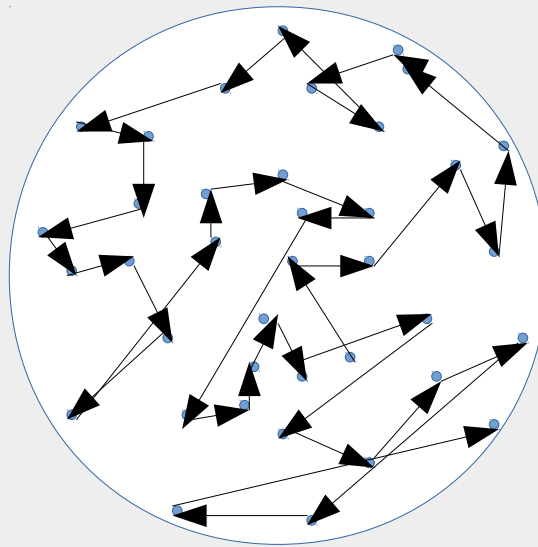
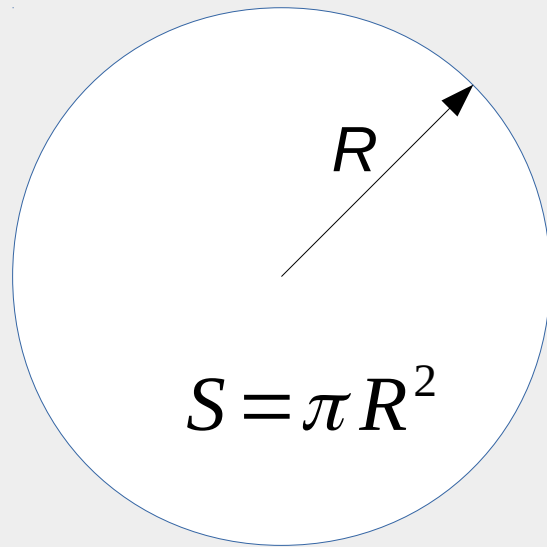
Stochastic methods in Mathematical Modelling

Lecture 17. Markov Chain Monte Carlo methods.



Random sampling instead of exact solution

Alternatively we could do a random walk and sample the space in a correlated way

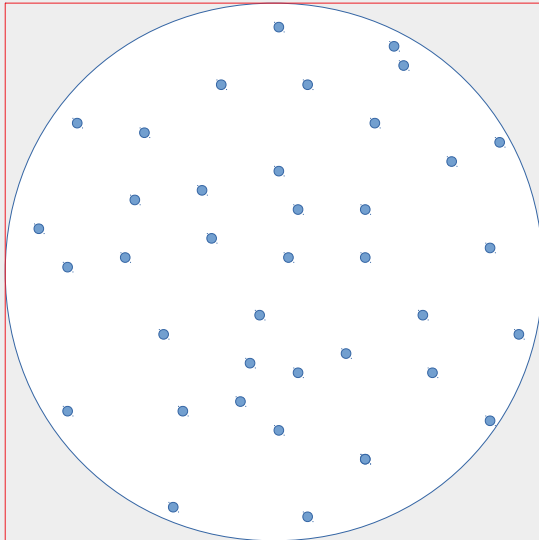


Monte Carlo. Two ways

Usually the goal is to sample some complicated distribution $f(x)$ or find $E[f(x)]$

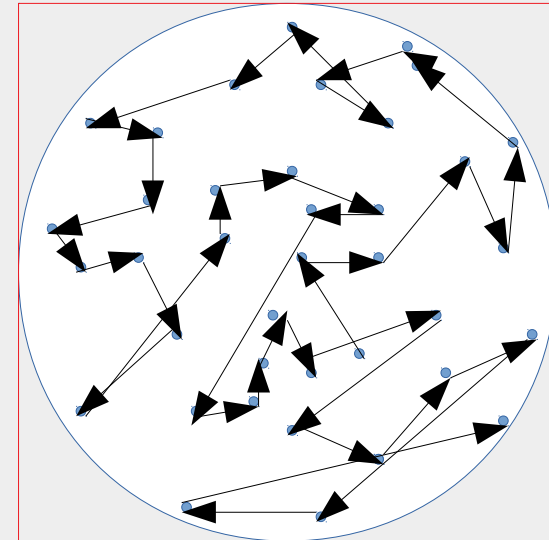
Direct sampling

Independent samples from a distribution



Markov Chain Monte Carlo (MCMC)

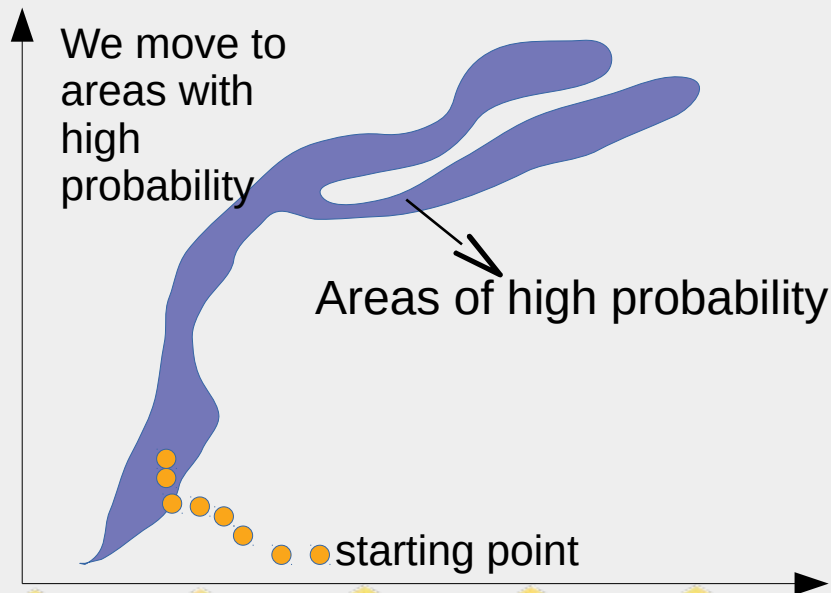
Draws are correlated according to a Markov chain



Usually the goal is to sample some complicated distribution $f(x)$ or find $E_f[h(X)]$ or an integral $\int f(x) dx$

In low dimensional problems rejection sampling (RS) and importance sampling (IS) are sufficiently good because it is possible to find the proposal density $P(x)$ similar to $f(x)$. However, in high dimensional spaces, for large and complex problems it is difficult to create a single density $Q(x)$ that has this property.

Intuitive idea of MCMC



(a) the region of high probability tends to be “connected”, i.e. one can get from one point to another without going through a low-probability region

(b) we tend to be interested in the expectations of functions that are relatively smooth and have lots of “symmetries”, hence need only a small number of representative points

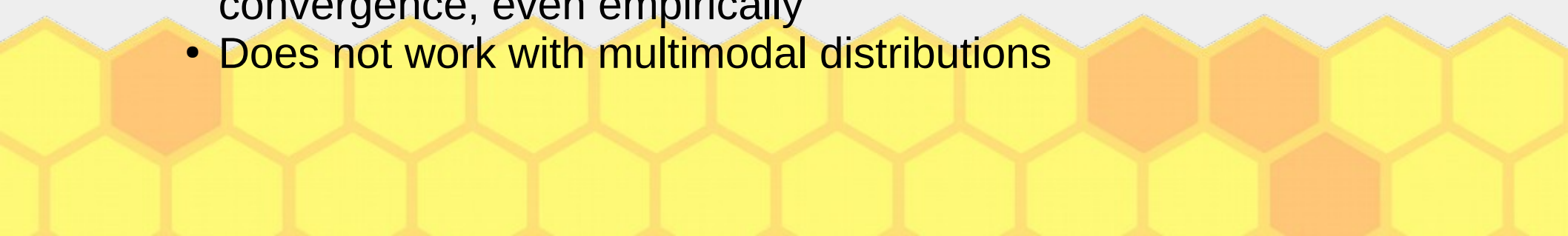
Skoltech Advantages/Disadvantages of MCMC:

Skolkovo Institute of Science and Technology

Advantages:

- applicable even when we can't directly draw samples
- works for complicated distributions in high-dimensional spaces, even when we don't know where the regions of high probability are
- relatively easy to implement
- fairly reliable

Disadvantages:

- the initial samples may follow a very different distribution, especially if the starting point is in a region of low density. A burn-in (thermalisation) period is typically necessary, where an initial number of samples (e.g. the first 1000 or so) are thrown away.
 - can be very difficult to assess accuracy and evaluate convergence, even empirically
 - Does not work with multimodal distributions
- 



MCMC methods

- Metropolis-Hastings method
- Gibbs Sampling
- Slice sampling
- Hamiltonian (or Hybrid) Monte Carlo (HMC)
- etc.



MC MC. Metropolis-Hastings method

we start $X(0)$

our goal is to approach f $X(t) \xrightarrow{t \rightarrow \infty} \underline{f}$

we want $X(t+s)$ to be independent of $X(t)$

$A = E[f(V(x))]$ → function of „interest“

$$\hat{A} = \frac{1}{L} \sum_{t=1}^L V(x(t))$$

If ergodicity holds for the Markov chain
 $\forall X(0) \quad \hat{A} \rightarrow A \quad \text{for } L \rightarrow \infty$



Let's consider ~~et~~of discrete states \mathcal{S} :

$$x_1, \dots, x_d$$

$$f = (f_1, \dots, f_d), \quad f_k \equiv f(x_k)$$

Transition probability: $P_{xy} = P(x \leftarrow y) = P(x_{t+1} = x | x_t = y)$
we would like to make f to be stationary

$$f(x) = \sum_{y \in \mathcal{S}} f(y) P_{xy}$$

distr. of Markov chain

Detailed balance condition

$$f(x) P_{yx} = f(y) P_{xy} \quad \forall x, y$$



Metropolis method

1) P_{xy} have to follow detailed balance (DB)

2) Let's build RW assuming $f(x)$ are given



$$a(x) = P(x \rightarrow x+1)$$

$$b(x) = P(x \rightarrow x)$$

$$c(x) = P(x \rightarrow x-1)$$

$$a(x) + b(x) + c(x) = 1$$

$$0 \leq a(x) \leq 1, 0 \leq b(x) \leq 1, 0 \leq c(x) \leq 1 \quad \forall x$$

Let's also assume

$$P(d \rightarrow d+1) = 0 = a(d) = 0$$

$$a(x), c(x) > 0$$

$$P(1 \rightarrow 0) = c(1) = 0$$

DB condition: $f(x) a(x) = f(x+1) c(x+1)$

$$f(x) a(x) = f(x+1) c(x+1)$$

Case (1) $f(x) \leq f(x+1)$ \downarrow $c(x+1) = a(x) \frac{f(x)}{f(x+1)}$

if $a(x) = 1$ $c(x+1) = \frac{f(x)}{f(x+1)} \leq 1$

(2) $f(x) \geq f(x+1)$ $c(x+1) = 1$ $a(x) = \frac{f(x+1)}{f(x)} \leq 1$

\Downarrow we could find a, b, c

BUT: $a + b + c = 1 \Rightarrow$ we could get $f(x) < 0$

Hence we need to adjust it:



Let's try to be less greedy:

$$a(x) \leq \frac{1}{2}, \quad c(x) \leq \frac{1}{2}$$

then $b(x) \geq 0$

$$(i) \quad f(x) \geq f(x+1) \quad c(x+1) = \frac{1}{2}$$

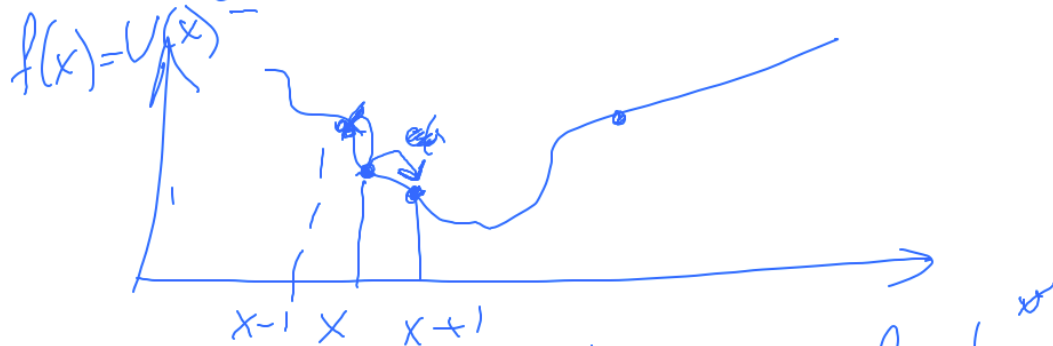
$$a(x) = \frac{f(x+1)}{f(x)} \cdot \frac{1}{2} \neq \frac{1}{2}$$

(ii) $f(x) \leq f(x+1)$ $a(x) = 1/2$

$$C(x+1) = \frac{f(x)}{f(x+1)} \cdot \frac{1}{2}$$

the pvd stay $a(d) = 0 = c(d)$

$$b(x) = 1 - a(x) - c(x)$$



1) proposal stage

2) acceptance/rejection stage



(i): we propose a move with $p = \frac{1}{2}$
if $x \rightarrow x_1$ we accept it
with $p = 1$

$x_i \rightarrow x_{i+1}$ we accept it
with $p = \frac{f(x_{i+1})}{f(x_i)} \leq 1$

In general : 1) $K_{xy} \geq 0$ the proposal

$$\sum_x K_{xy} = 1$$

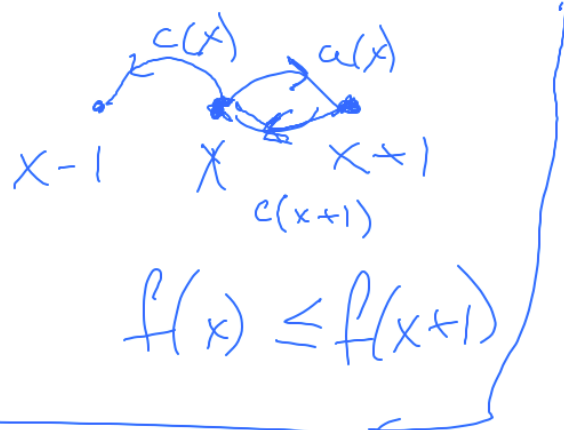
2) The acceptance rates R_{xy}
to fulfil DB \Rightarrow

$$f(x) K_{yx} R_{yx} = f(y) K_{xy} R_{xy}$$

$$0 \leq R_{yx}, K_{yx} \leq 1$$

we would like R_{xy} to be close to 1:

$$R_{xy} = \min \left(1, \frac{f(y) K_{yx}}{f(x) K_{xy}} \right)$$



Convergence speed
crude estimate

from the theory of MCMC:
after T steps the distance
(the number of visited states)

grows as $\sqrt{T} \varepsilon (\sim d)$
↑ length of the move

If we have d states
 $T \sim \left(\frac{d}{\varepsilon}\right)^2$

a rough estimate
of the number of states
to produce independent
samples.

1. Initialisation:

Pick an initial state x_0

2. Iteration:

(i) Generate a random candidate state x' according to $K(x',x)$.

(ii) Calculate the acceptance probability $R_{x',x} = \min\left(1, \frac{f(x) K_{x,x'}}{f(x') K_{x',x}}\right)$

3. Accept or reject:

(i) generate a uniform random number u from $U[0,1]$

if $u \leq R_{x',x}$, then accept the new state and set $x_{t+1} = x'$

if $u > R_{x',x}$, then reject the new state, and copy the old state forward $x_{t+1} = x_t$.

4. Increment: set $t=t+1$

Note: It is not the same as acceptance-rejection algorithm. In the latter when we *reject* a proposed value the number of samples does not increase. In MH method when the rejection happens the chain still takes a time step

3. Accept or reject:

(i) generate a uniform random number u from $U[0,1]$

if $u \leq R_{x',x}$, then accept the new state and set $x_{t+1} = x'$

if $u > R_{x',x}$, then reject the new state, and copy the old state forward $x_{t+1} = x_t$

Comment: It is not the same as acceptance-rejection algorithm. In the latter when we *reject* a proposed value the number of samples *does not increase*. In MH method when the rejection happens the chain *still* takes a time step! In that case two consecutive states are the same

in $\frac{1}{n} \sum_{k=1}^n f(X_k)$ some terms can be the same



An example of slow convergence

Target distribution

$$P(x) = 1/21, x \in \{0, 1, \dots, 20\}$$

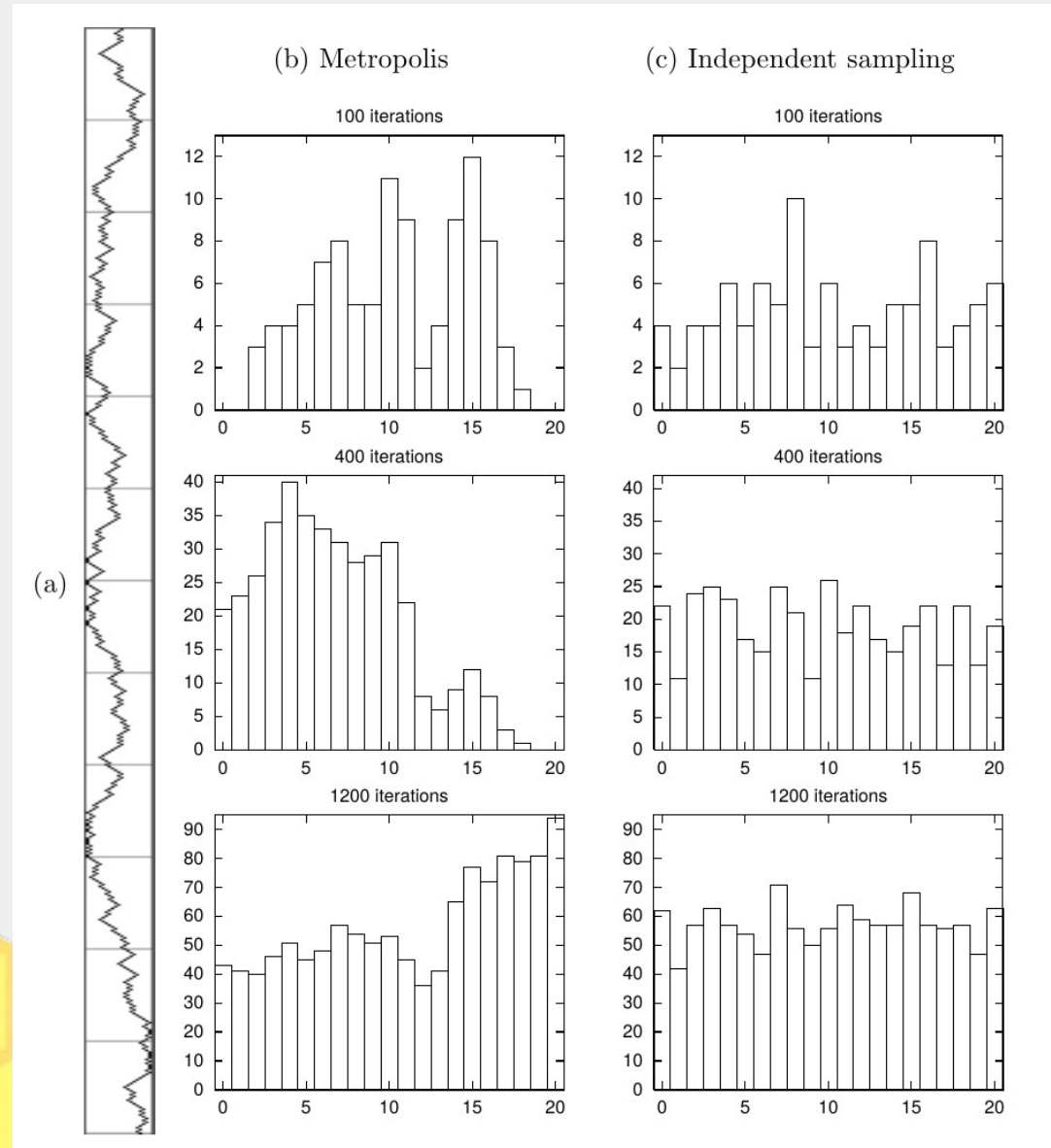
$$P(x) = 0, \text{otherwise}$$

Proposal distribution

$$K(x, x') = 1/2, x' = x \pm 1$$

$$K(x, x') = 0, \text{otherwise}$$

and rejections at the ends!



Metropolis for Ising model

$S_i = \pm 1$ $\uparrow \downarrow$; the set of spins

$$E = - \sum J_{ij} S_i S_j - \sum S_j h_j$$

Simple case: nearest-neighbour interactions & $J_{ij} = 1$

$$E = - \sum_{\langle ij \rangle} S_i S_j$$

Minimum is achieved when all spins ~~are~~ point in 1 direction

Boltzmann distr:

$$T \rightarrow 0 \ (\beta \rightarrow \infty) \quad P_{eq}(s) \neq 0 \quad \text{only for a few states close to minimum energy state (ground state)}$$
$$P_{eq}(s) = \frac{1}{Z} e^{-\beta E(s)}, \quad Z = \sum_s e^{-\beta E(s)}$$
$$\beta = \frac{1}{T}$$



$T \rightarrow \infty$ ($\beta \rightarrow 0$) all states are equiprobable

$\uparrow\uparrow\downarrow\uparrow\uparrow\downarrow\uparrow$ > the same energy but diff ^{micro} states

$\uparrow\downarrow\uparrow\uparrow\downarrow\uparrow\uparrow$

Macrostates: W is a number of microstates with the same energy

$$\frac{W}{Z} e^{-\beta E} = \frac{1}{Z} e^{-\beta E + \ln W} =$$

$\ln W = S \rightarrow$ entropy

$$= \frac{1}{Z} e^{-\beta \underbrace{(E - TS)}_F}$$

the most likely state corresponds to minimum of F is free energy



N spins $\uparrow\uparrow\downarrow\uparrow$

2^N combinations or states

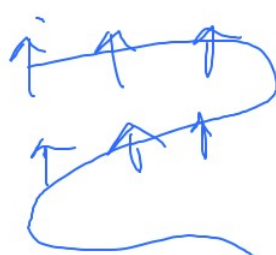
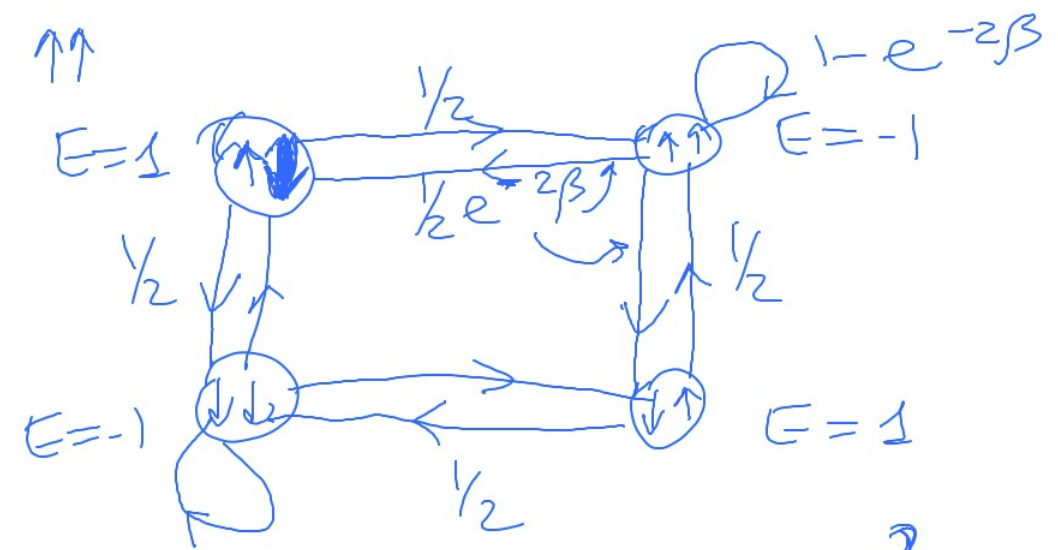
1) Brute force approach requires 2^N
inefficient

2) Metropolis-Hastings for Ising
Let's construct the Markov chain on a hypercube
with 2^N vertices

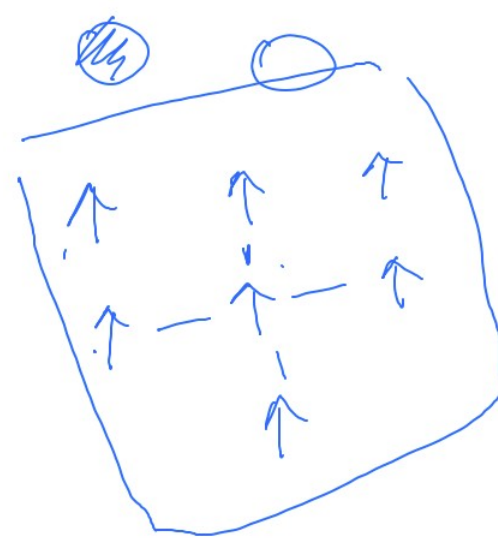
a) Choose a state at random

b) Compute the prob. to flip the spin as

$$P = \begin{cases} 1, & \Delta E \leq 0 \\ e^{-\beta \Delta E}, & \Delta E \geq 0 \end{cases}$$



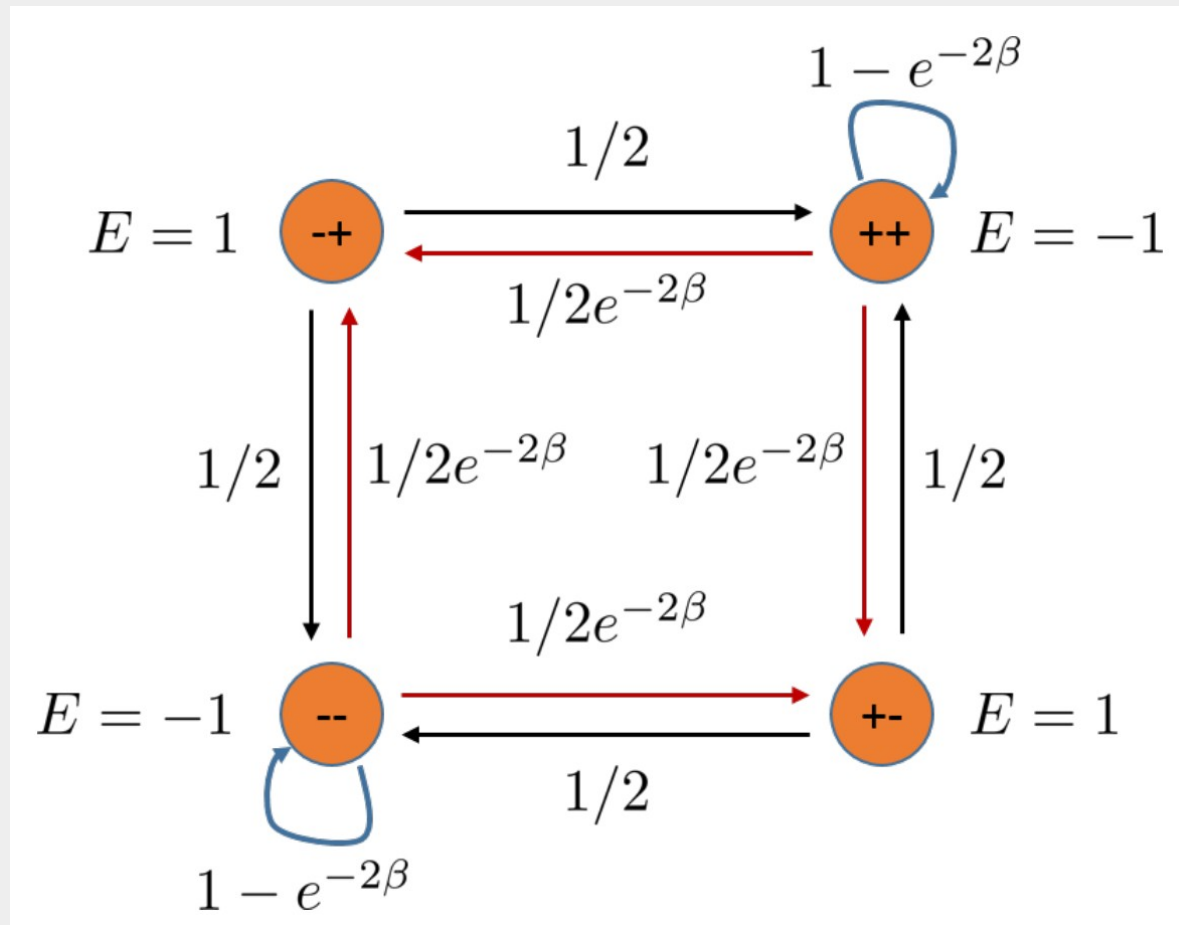
N^2 vs 2^N



2^N



Ising model. Metropolis-Hastings Markov chain example for two spins



An alternative

Gibbs (Glauber) sampling

other way to fulfil the DB

$$\begin{cases} P_{\downarrow} + P_{\uparrow} = 1 & P_{\uparrow} \\ \frac{P_{\uparrow}}{P_{\downarrow}} = e^{-\beta \Delta E} & P_{\downarrow} \end{cases} \quad N^2$$

Probability flux \uparrow to \downarrow :

$$J_{\uparrow\downarrow} = \frac{1}{Z} e^{-\beta E(\uparrow)} P_{\downarrow}$$

$$J_{\downarrow\uparrow} = \frac{1}{Z} e^{-\beta E(\downarrow)} P_{\uparrow}$$

$$\boxed{J_{\uparrow\downarrow} = J_{\downarrow\uparrow}}$$



Gibbs sampling method

(also called Glauber dynamics if applied for Ising model)

Gibbs sampling can be viewed as a Metropolis method in which a sequence of proposal distributions Q are defined in terms of the conditional distributions of the joint distribution $P(x)$. Used when the conditional distributions are easier to sample than $P(x)$

$$\begin{aligned}x_1^{(t+1)} &\sim P(x_1 \mid x_2^{(t)}, x_3^{(t)}, \dots, x_K^{(t)}) \\x_2^{(t+1)} &\sim P(x_2 \mid x_1^{(t+1)}, x_3^{(t)}, \dots, x_K^{(t)}) \\x_3^{(t+1)} &\sim P(x_3 \mid x_1^{(t+1)}, x_2^{(t+1)}, \dots, x_K^{(t)}), \text{ etc.}\end{aligned}$$

For big problems it may be more efficient to sample groups of variables jointly, that is to use several proposal distributions

$$\begin{aligned}x_1^{(t+1)}, \dots, x_a^{(t+1)} &\sim P(x_1, \dots, x_a \mid x_{a+1}^{(t)}, \dots, x_K^{(t)}) \\x_{a+1}^{(t+1)}, \dots, x_b^{(t+1)} &\sim P(x_{a+1}, \dots, x_b \mid x_1^{(t+1)}, \dots, x_a^{(t+1)}, x_{b+1}^{(t)}, \dots, x_K^{(t)})\end{aligned}$$

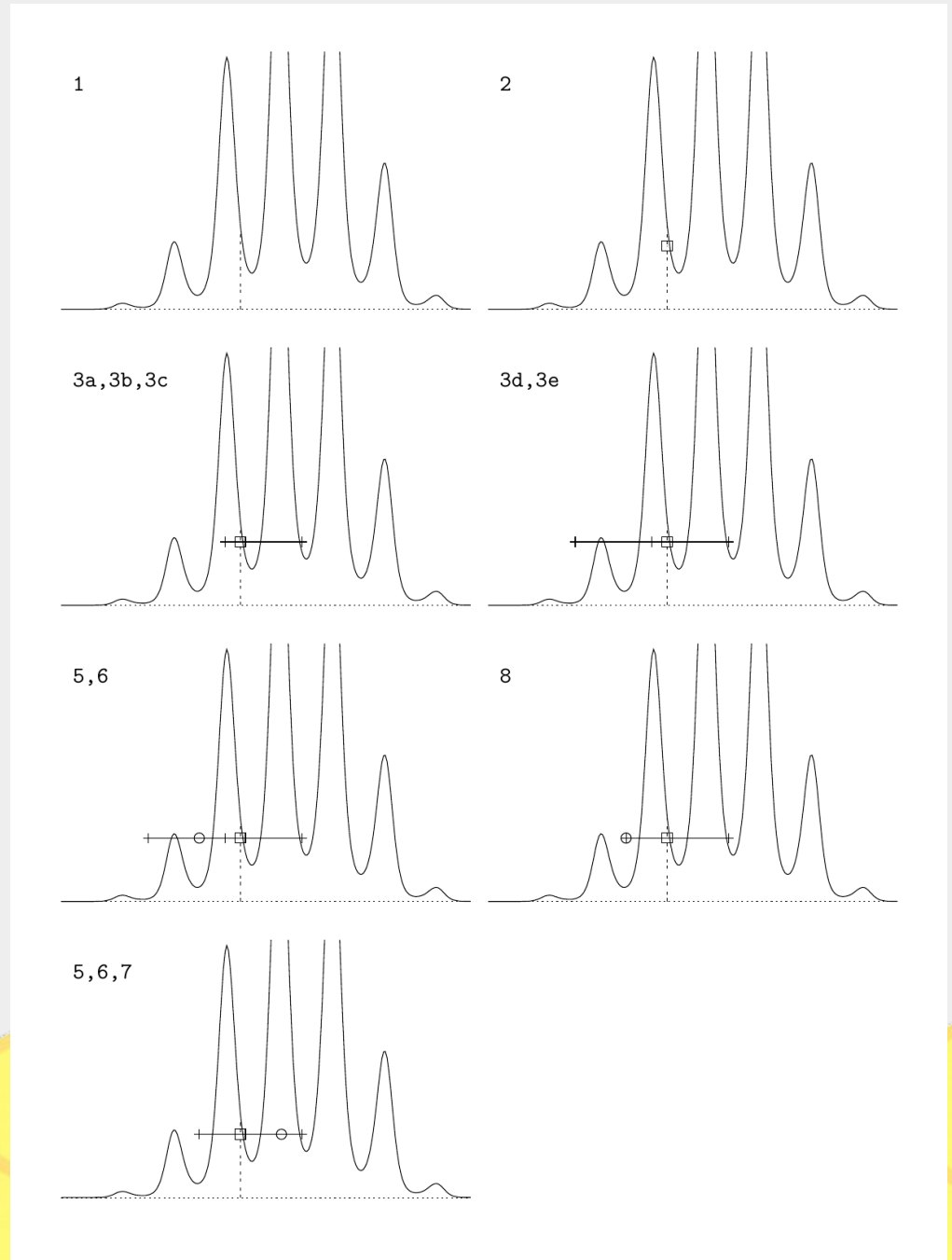
Slice sampling method

Can be applied to any distribution
when one can evaluate $P(x)$ at any point x

- 1: evaluate $P^*(x)$
- 2: draw a vertical coordinate $u' \sim \text{Uniform}(0, P^*(x))$
- 3: create a horizontal interval (x_l, x_r) enclosing x
- 4: loop {
- 5: draw $x' \sim \text{Uniform}(x_l, x_r)$
- 6: evaluate $P^*(x')$
- 7: if $P^*(x') > u'$ break out of loop 4-9
- 8: else modify the interval (x_l, x_r)
- 9: }

step size is self-tuning

D. J. C. Mackay, Information theory,
Inference, and learning algorithms. 2003



Hamiltonian MC method

The Hamiltonian Monte Carlo method is a Metropolis method, applicable to continuous state spaces, that makes use of gradient information to reduce sampling inefficiency of a simple random walk.

$$P(\mathbf{x}) = \frac{e^{-E(\mathbf{x})}}{Z}$$

We can also compute the gradient of $E(\mathbf{x})$. It indicates the direction with states with higher probability. Add momentum variables \mathbf{p}

$$H(\mathbf{x}, \mathbf{p}) = E(\mathbf{x}) + K(\mathbf{p})$$

$$P_H(\mathbf{x}, \mathbf{p}) = \frac{1}{Z_H} \exp[-H(\mathbf{x}, \mathbf{p})] = \frac{1}{Z_H} \exp[-E(\mathbf{x})] \exp[-K(\mathbf{p})]$$

$$p(x, v) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{\det \Sigma}} \exp(-E(x) + \frac{1}{2} v^T \Sigma^{-1} v)$$

Hamiltonian MC method

1. A proposal for \mathbf{p} is obtained from a marginal distribution for momenta. This proposal is always accepted

2.

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{p} \\ \dot{\mathbf{p}} &= -\frac{\partial E(\mathbf{x})}{\partial \mathbf{x}}.\end{aligned}$$

Because of the persistent motion of \mathbf{x} in the direction of the momentum \mathbf{p} during each dynamical proposal, the state of the system tends to move a distance that goes *linearly* with the computer time, rather than as the square root (ballistic vs diffusive motion).

The second proposal is accepted in accordance with the Metropolis rule.



Hamiltonian MC method

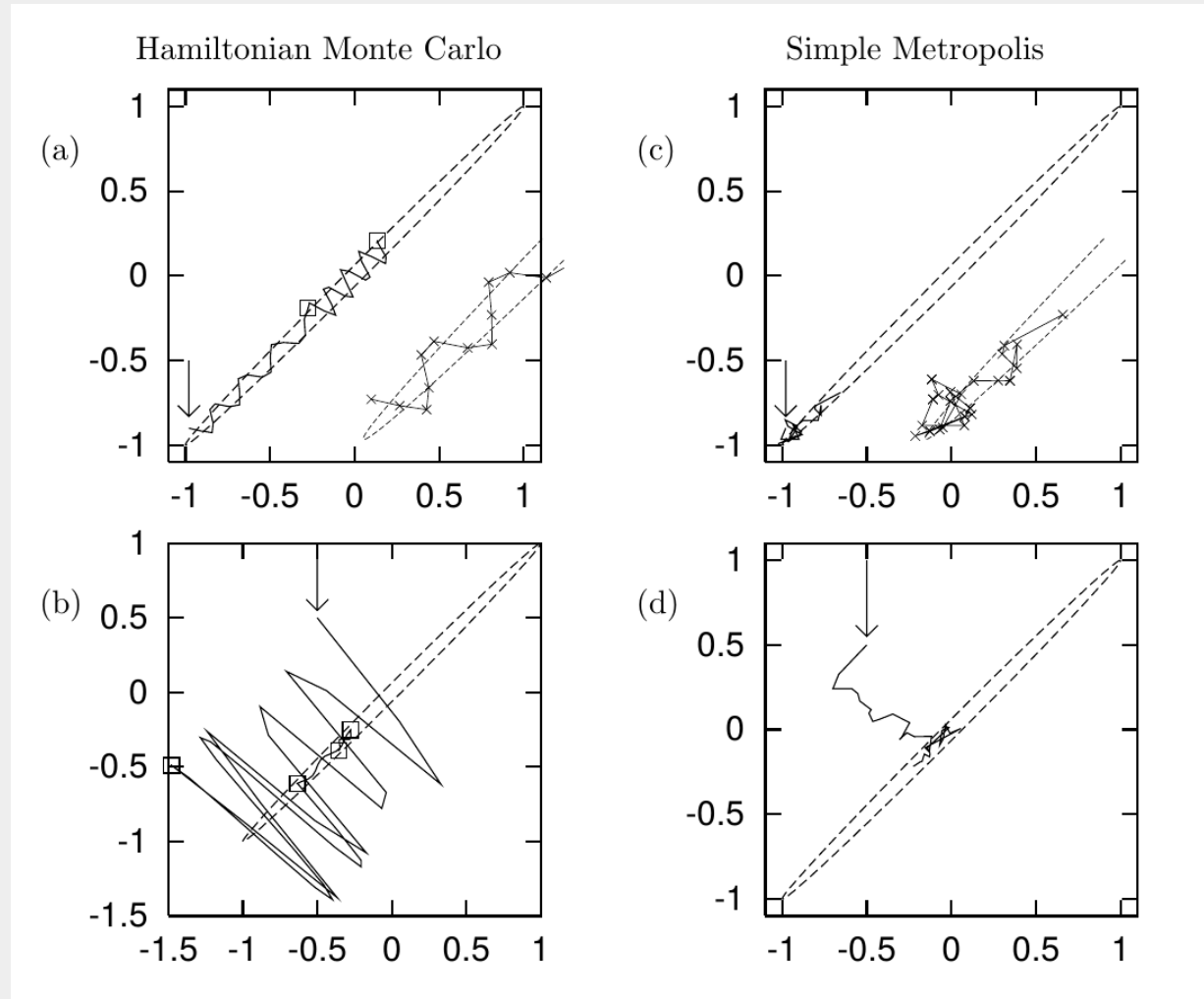


Figure 30.2. (a,b) Hamiltonian Monte Carlo used to generate samples from a bivariate Gaussian with correlation $\rho = 0.998$. (c,d) For comparison, a simple random-walk Metropolis method, given equal computer time.

D. J. C. Mackay, Information theory, Inference, and learning algorithms. 2003

Hamiltonian MC method

- Euler Integrator:

$$\begin{cases} v(t + \epsilon) = p(t) - \epsilon \frac{\partial H}{\partial q(t)} \\ q(t + \epsilon) = q(t) - \epsilon \frac{\partial H}{\partial v(t+\epsilon)} \end{cases}$$

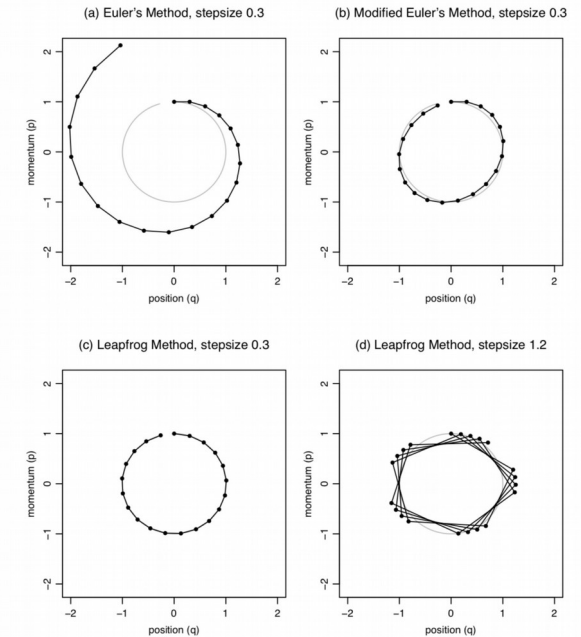
- Leapfrog :

$$\begin{cases} v(t + \frac{\epsilon}{2}) = v(t) - \frac{\epsilon}{2} \frac{\partial H}{\partial q(t)} \\ q(t + \epsilon) = q(t) - \epsilon \frac{\partial H}{\partial v(t+\frac{\epsilon}{2})} \\ v(t + \epsilon) = v(t + \frac{\epsilon}{2}) - \frac{\epsilon}{2} \frac{\partial H}{\partial q(t+\epsilon)} \end{cases}$$

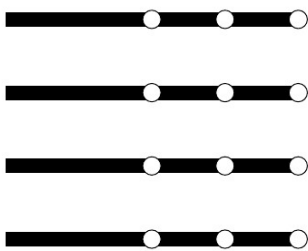
- Problems:

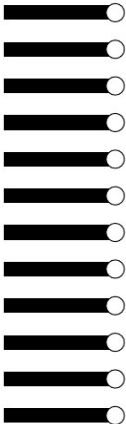
- Approximate Hamiltonian
- Discretization error
- one mode exploration
- slow mixing
- slow burn-in
- not mixing across levels
- geometric of space

Comparison of Integrators



(1) 

(2) 

(3) 

Error sources of MCMC:

1) Burn-in period

2) The variance for correlated samples

$$\text{cov}(X, Y) \equiv E[XY] - E[X]E[Y] = E[(X - \mu_X)(Y - \mu_Y)]$$

$$\text{var}\left(\sum_{i=1}^N Y_i\right) = \frac{1}{N^2} \sum_{i,j=1}^N \text{cov}(X_i, X_j)$$

we get into stationary state & cov depends $|i-j|$

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N f(x_n)$$

$$\text{var}(\hat{\mu}) = \frac{1}{N^2} \sum_{k,n=1}^N \text{cov}(f(x_k), f(x_n))$$

If N is large $\text{var}(f(x_n))$ should be close to $\frac{\sigma_f^2}{N}$

$$\text{cov}(f(x_k), f(x_n)) = [\text{var}(f(x_k)) \text{var}(f(x_n))]^{\frac{1}{2}} \text{corr}(f(x_k), f(x_n)) \\ \approx \sigma^2(f) \text{corr}(f(x_k), f(x_n))$$

$$\text{Var}(\hat{\mu}) = \frac{\sigma^2(f)}{N^2} \sum_{k=1}^N \sum_{n=1}^N \text{corr}(f(x_k), f(x_n))$$

$$\sum_{n=1}^{\infty} \text{corr}(f(x_k), f(x_n)) \approx \sum_{n=-\infty}^{n=+\infty} \text{corr}(f(x_k), f(x_{k+n})) =$$

$$= 1 + 2 \sum_{n=1}^{\infty} \text{corr}(f(x_k), f(x_{k+n})) = \tau(f)$$

autocor time

$$\text{Var}(\hat{\mu}) \approx \frac{\sigma^2(f)}{N^2} \sum_{k=1}^N \tau(f) = \frac{\sigma^2(f) \tau(f)}{N}$$

size of a batch

$\frac{N}{\tau(f)}$

Convergence

Independence test

$$\lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{n=0}^N \sum_{n=1}^N f(x_n) g(x_{n-1}) \rightarrow \mathbb{E}[f(x)] \mathbb{E}[g(x)]$$

Asymptotic convergence

$$\lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{n=0}^N \sum_{n=1}^N f(x_n) \rightarrow \mathbb{E}[f(x)]$$

Parametric convergence

The target estimate is reached only in a special limit with respect to a special parameter

$$\lim_{s \rightarrow s_*} \lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{n=0}^N \sum_{n=1}^N f_s(x_n) \rightarrow \mathbb{E}[f_{s_*}(x)]$$

Literature

1. D. J. C. Mackay, Information theory, inference, and learning algorithms. Cambridge: Cambridge University Press, 2003.
<http://www.inference.org.uk/itprnn/book.html>
2. <https://math.nyu.edu/faculty/goodman/teaching/>
3. <https://people.eecs.berkeley.edu/~jordan/courses/260-spring10/>
4. <https://www.math.arizona.edu/~tgk/mc/book.pdf>

