

Probabilistic Programming Languages

Guillaume Baudart

MPRI 2023-202

MCMC Metropolis-Hastings

Probabilistic Programming Languages

Markov Chain Monte Carlo (MCMC)

Main idea

- Create a Markov chain that converge to the posterior distribution
- Iterate the process until convergence
- Generate samples to approximate the distribution

Pros

- Faster convergence
- Better results for high-dimensional models
- Advanced state-of-the-art optimizations (e.g., HMC, NUTS).

Cons

- Convergences?
- Traps: multimodal, funnel
- Samples correlation

Reminder: rejection sampling

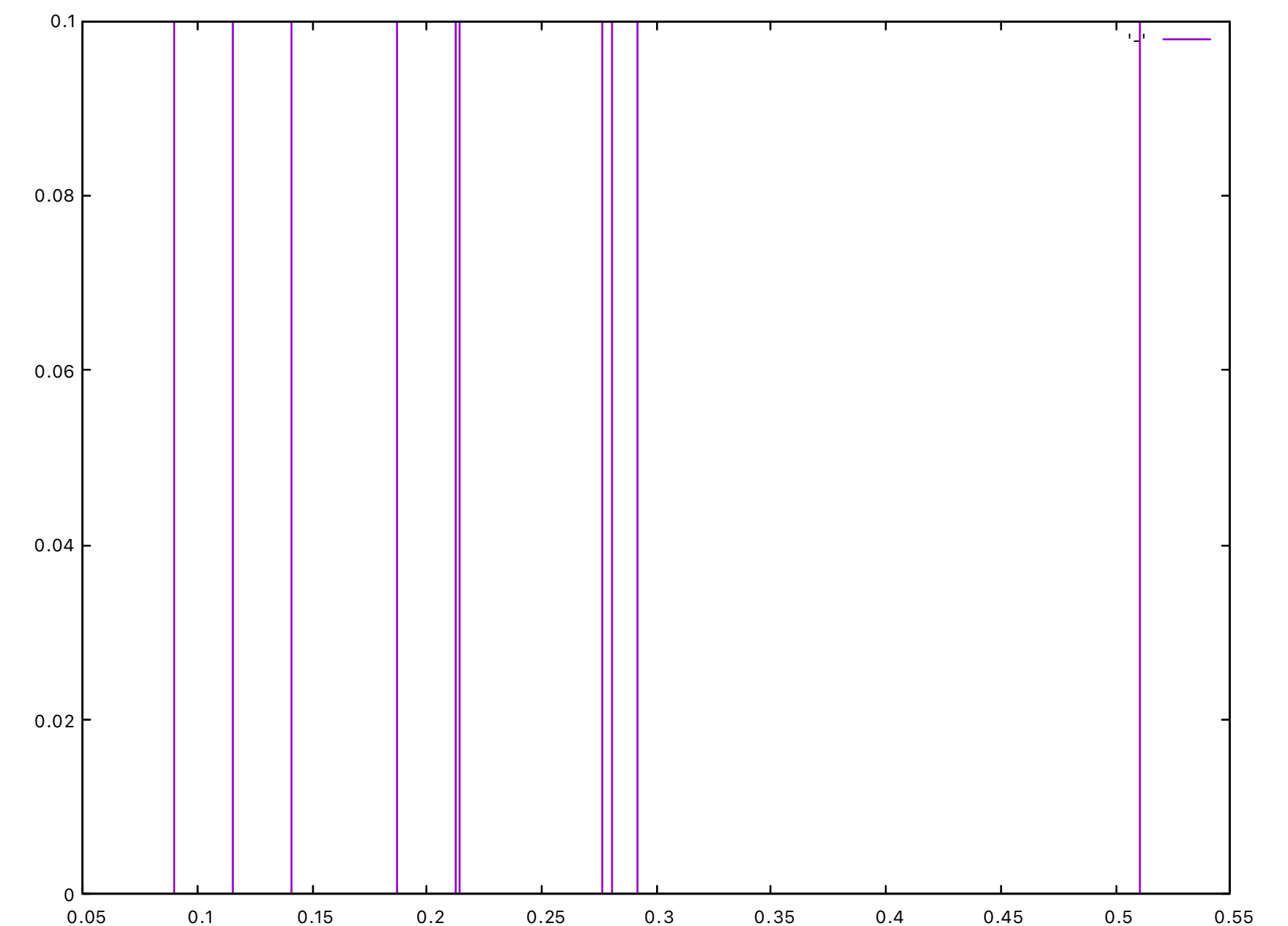
coin.ml

```
let coin prob data =  
  let z = sample prob (uniform ~a:0. ~b:1.) in  
  List.iter (observe prob (bernoulli ~p:z)) data;  
  z  
  
let _ =  
  let d = infer coin [ 1; 1; 0; 0; 0; 0; 0; 0; 0; 0 ] in  
  plot d
```

Executing the model generates one sample

- **sample**: draw from a distribution
- **assume/observe**: hard conditioning, reject invalid samples
- Terminates with *n* valid samples

10 particles



Reminder: rejection sampling

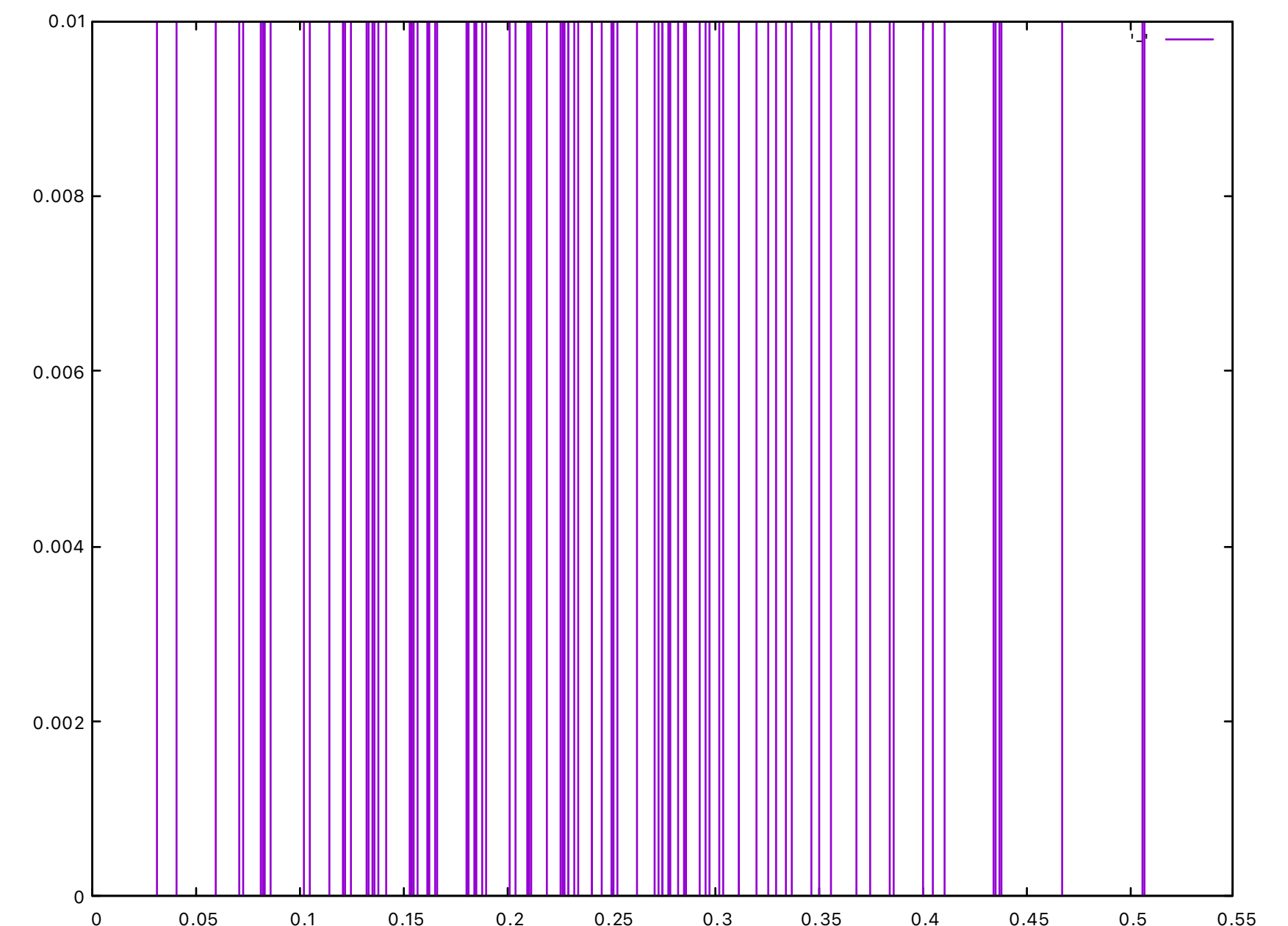
coin.ml

```
let coin prob data =  
  let z = sample prob (uniform ~a:0. ~b:1.) in  
  List.iter (observe prob (bernoulli ~p:z)) data;  
  z  
  
let _ =  
  let d = infer coin [ 1; 1; 0; 0; 0; 0; 0; 0; 0; 0 ] in  
  plot d
```

Executing the model generates one sample

- **sample**: draw from a distribution
- **assume/observe**: hard conditioning, reject invalid samples
- Terminates with *n* valid samples

100 particles



Reminder: rejection sampling

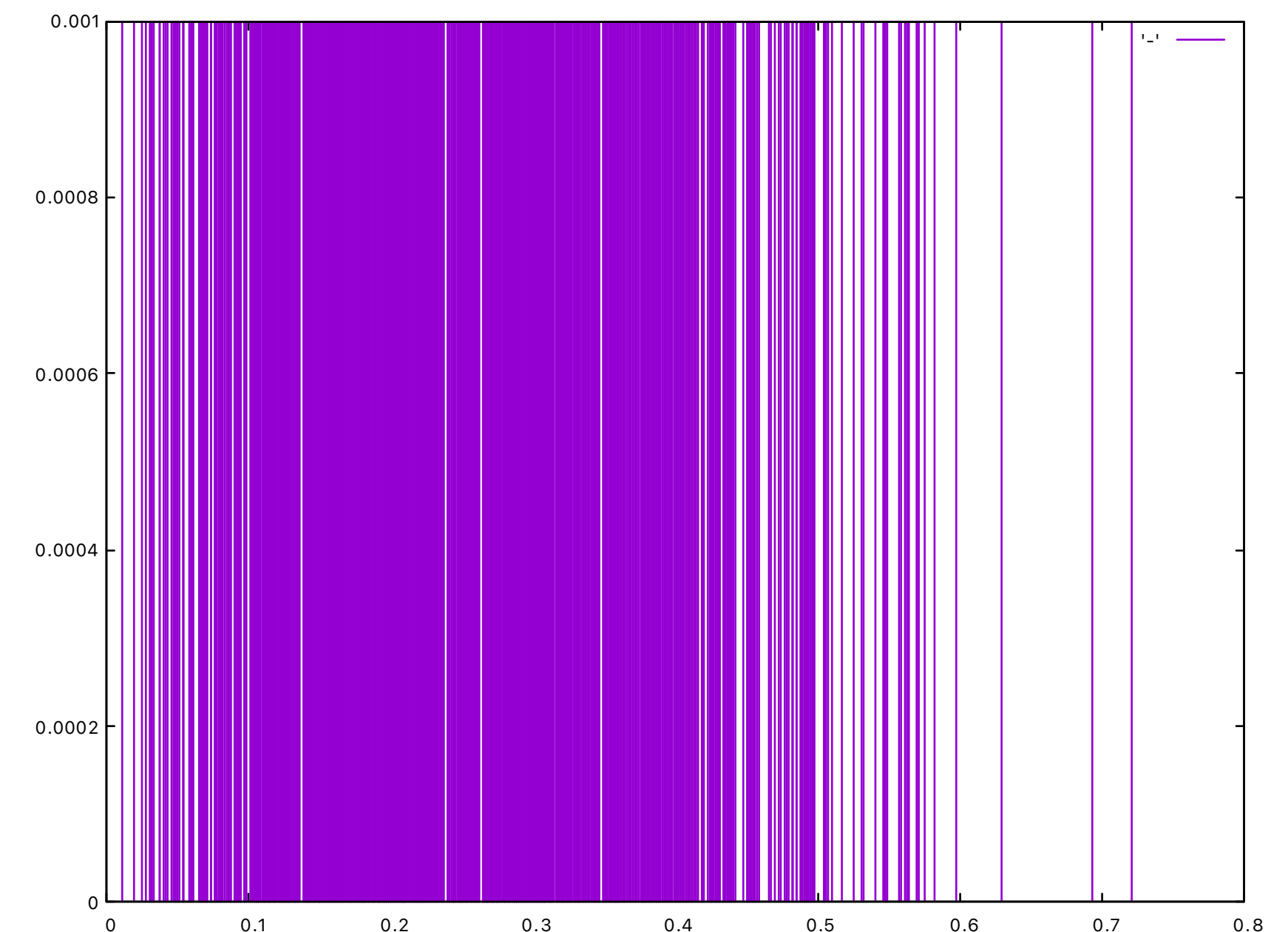
coin.ml

```
let coin prob data =  
  let z = sample prob (uniform ~a:0. ~b:1.) in  
  List.iter (observe prob (bernoulli ~p:z)) data;  
  z  
  
let _ =  
  let d = infer coin [ 1; 1; 0; 0; 0; 0; 0; 0; 0; 0 ] in  
  plot d
```

Executing the model generates one sample

- **sample**: draw from a distribution
- **assume/observe**: hard conditioning, reject invalid samples
- Terminates with *n* valid samples

1000 particles



Reminder: rejection sampling

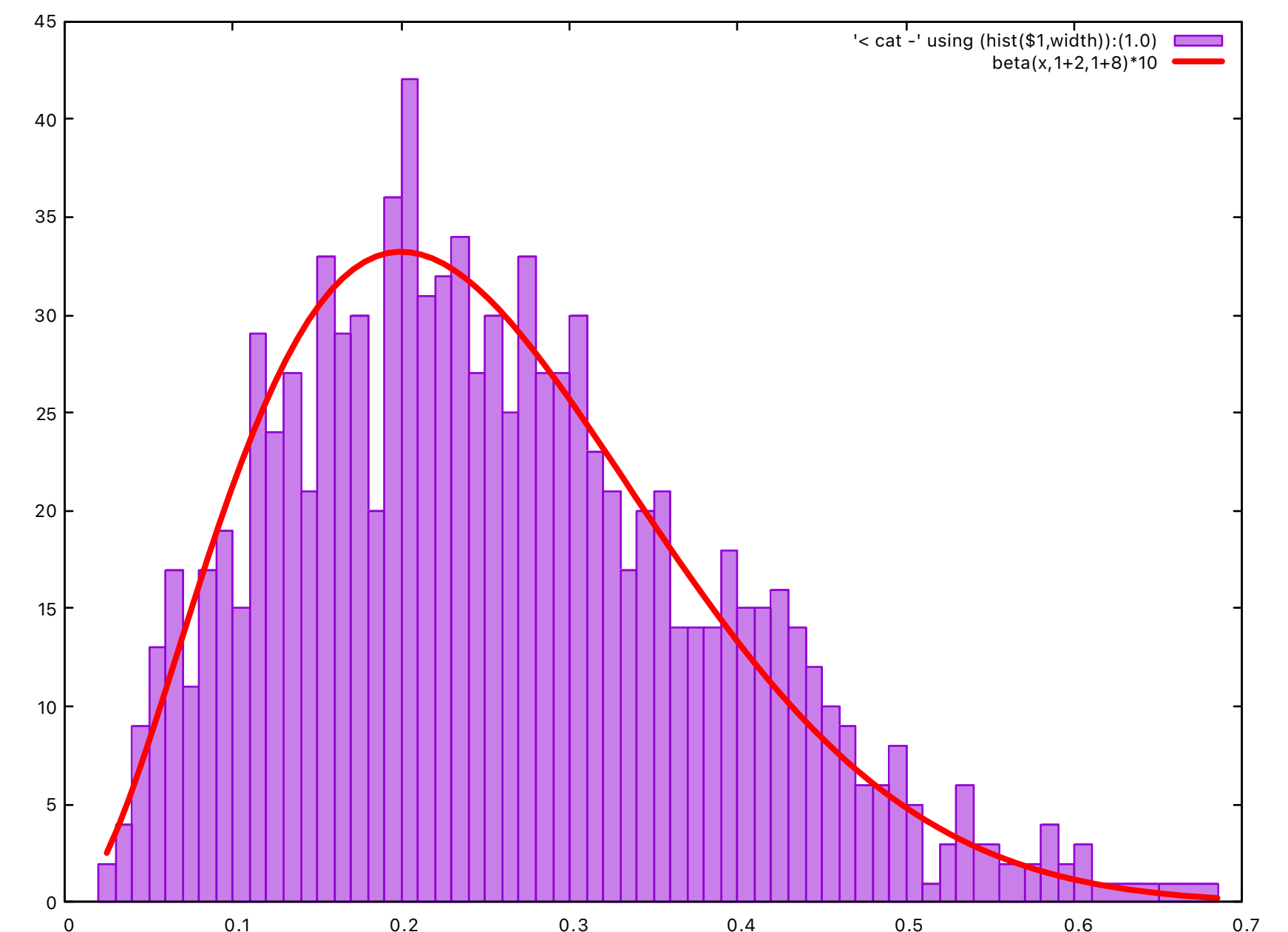
coin.ml

```
let coin prob data =  
  let z = sample prob (uniform ~a:0. ~b:1.) in  
  List.iter (observe prob (bernoulli ~p:z)) data;  
  z  
  
let _ =  
  let d = infer coin [ 1; 1; 0; 0; 0; 0; 0; 0; 0; 0 ] in  
  plot d
```

Executing the model generates one sample

- **sample**: draw from a distribution
- **assume/observe**: hard conditioning, reject invalid samples
- Terminates with *n* valid samples

1000 particles



Weighted rejection sampling

Adapt rejection sampling to soft conditioning

- Execute the sampler to get a pair (v_i, w_i)
- Suppose w_{\max} is known
- Accept the sample with probability w_i/w_{\max} or retry

But w_{\max} is not known...

Execution trace

Consider a program execution with

- $X_i = x_0, \dots, x_n$: set of random variables sampled at step i , i.e., the trace
- $Y_i = y_0, \dots, y_m$: set of random variables observed at step i .

Remarks

- Sets X_i and Y_i depend on the execution path
- We can only control X_i

```
let bimodal y =  
  let z = sample (bernoulli ~p:0.5) in  
  let mu =  
    if z then sample (gaussian ~mu:-1. ~sigma:1.)  
    else sample (gaussian ~mu:1. ~sigma:1.)  
  in  
  let () = observe (gaussian ~mu ~sigma:1.) y in  
  z
```

Multi-sites Metropolis Hastings

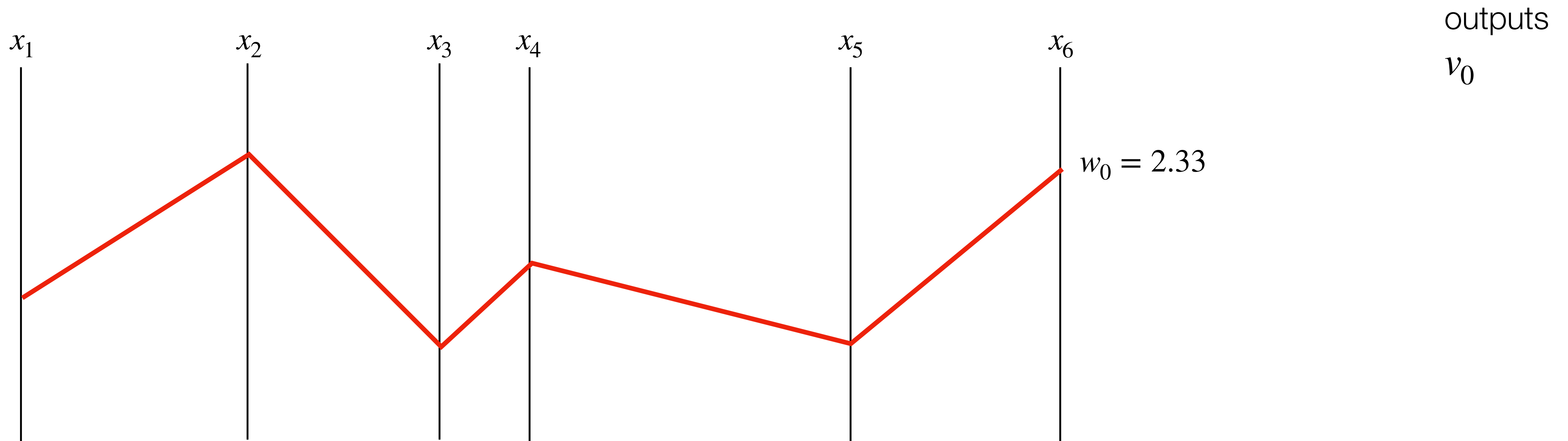
Markov chain on execution traces

- Execute the sampler to get a trace and associated score (X_i, y_i)
- If $w_i \geq w_{i-1}$ accept the trace (and the associated output)
- Else accept the trace with probability w_i/w_{i-1}
- Otherwise return the previous trace X_{i-1}

Multi-sites Metropolis Hastings

Markov chain on execution traces

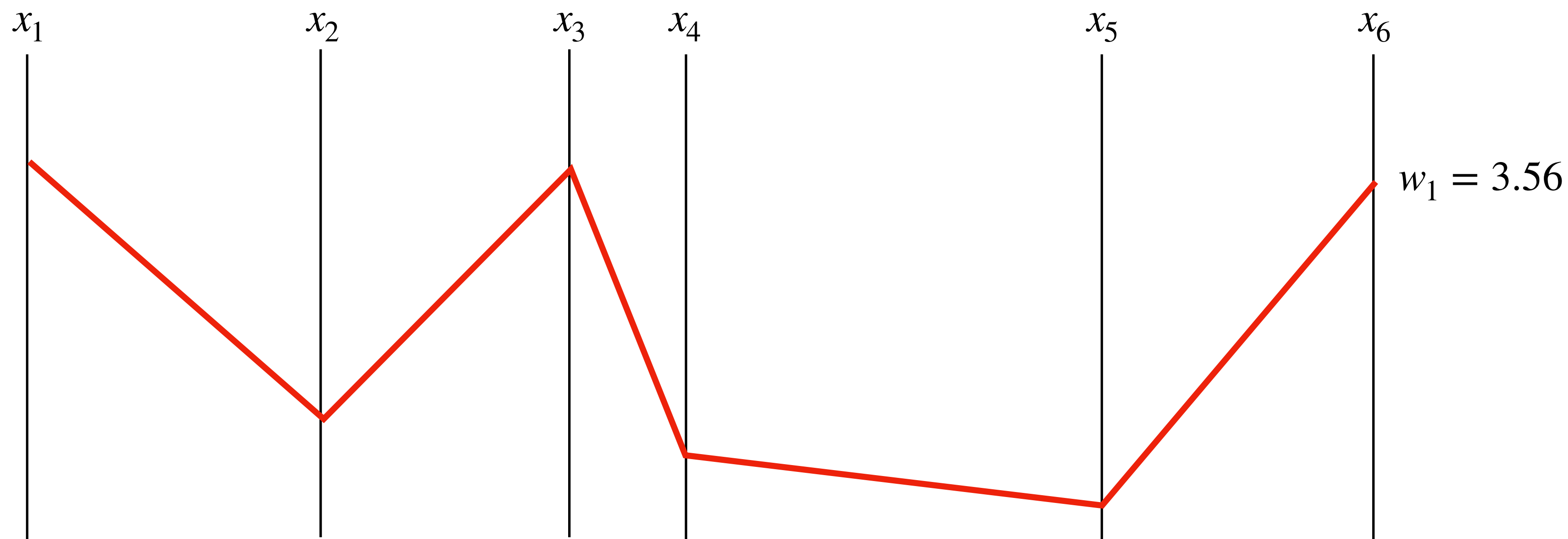
- Execute the sampler to get a trace and associated score (X_i, y_i)
- If $w_i \geq w_{i-1}$ accept the trace (and the associated output)
- Else accept the trace with probability w_i/w_{i-1}
- Otherwise return the previous trace X_{i-1}



Multi-sites Metropolis Hastings

Markov chain on execution traces

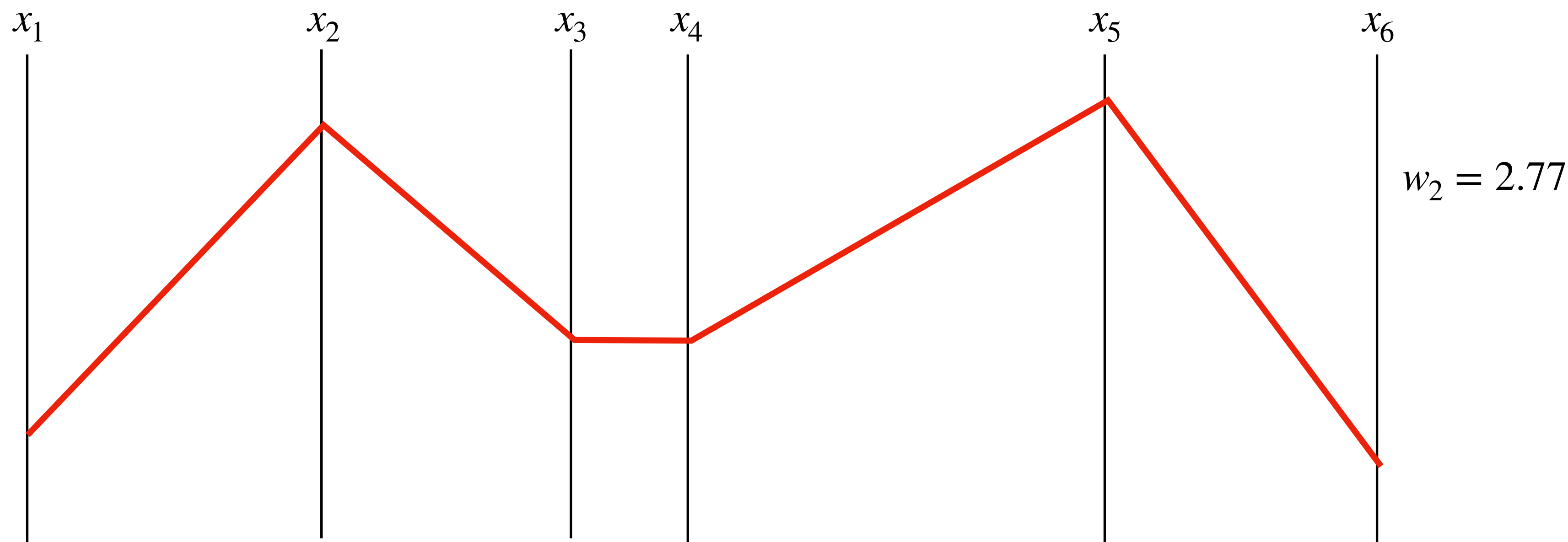
- Execute the sampler to get a trace and associated score (X_i, y_i)
- If $w_i \geq w_{i-1}$ accept the trace (and the associated output)
- Else accept the trace with probability w_i/w_{i-1}
- Otherwise return the previous trace X_{i-1}



Multi-sites Metropolis Hastings

Markov chain on execution traces

- Execute the sampler to get a trace and associated score (X_i, y_i)
- If $w_i \geq w_{i-1}$ accept the trace (and the associated output)
- Else accept the trace with probability w_i/w_{i-1}
- Otherwise return the previous trace X_{i-1}



outputs

v_0

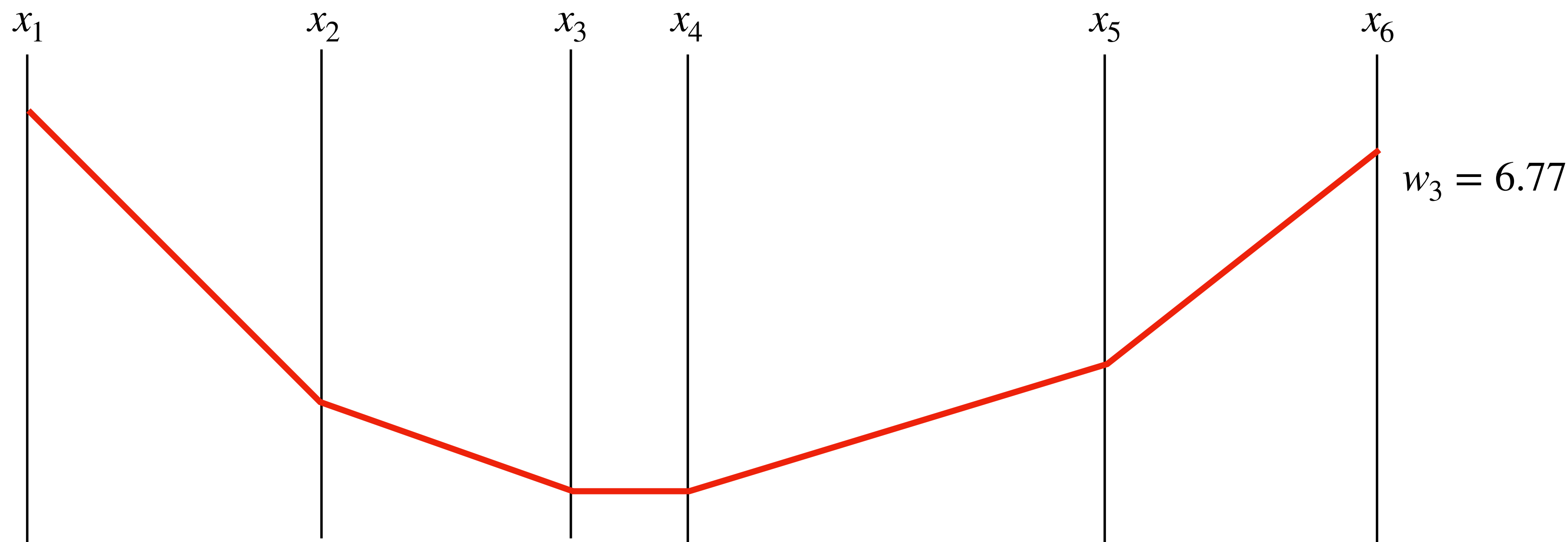
v_1

v_1

Multi-sites Metropolis Hastings

Markov chain on execution traces

- Execute the sampler to get a trace and associated score (X_i, y_i)
- If $w_i \geq w_{i-1}$ accept the trace (and the associated output)
- Else accept the trace with probability w_i/w_{i-1}
- Otherwise return the previous trace X_{i-1}



outputs

v_0

v_1

v_1

v_3

...

General Metropolis Hastings

More generally

- $X_i = x_0, \dots, x_n$: set of random variables sampled at step i , i.e., the trace
- $Y_i = y_0, \dots, y_m$: set of random variables observed at step i
- Propose a new trace from a proposal distribution $q(X_i \mid X_{i-1})$
- Accept the trace X_i with probability α , where

$$\alpha = \min \left(1, \frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} \mid X_i)}{q(X_i \mid X_{i-1})} \right)$$

- Otherwise return the previous trace X_{i-1}

Multi-sites Metropolis Hastings: acceptance

- Draw proposal from priors $q(X_i | X_{i-1}) = p(X_i)$
- Resample all variables in X_i at each iteration

$$\begin{aligned} \frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} | X_i)}{q(X_i | X_{i-1})} &= \frac{p(Y_i | X_i) p(X_i)}{p(Y_{i-1} | X_{i-1}) p(X_{i-1})} \frac{q(X_{i-1} | X_i)}{q(X_i | X_{i-1})} \\ &= \frac{p(Y_i | X_i) p(X_i)}{p(Y_{i-1} | X_{i-1}) p(X_{i-1})} \frac{p(X_{i-1})}{p(X_i)} \\ &= \frac{p(Y_i | X_i)}{p(Y_{i-1} | X_{i-1})} \\ &= \frac{w_i}{w_{i-1}} \end{aligned}$$

Markov chain on execution traces

- Execute the sampler to get a trace and associated score (X_i, y_i)
- If $w_i \geq w_{i-1}$ accept the trace (and the associated output)
- Else accept the trace with probability w_i/w_{i-1}
- Otherwise return the previous trace X_{i-1}

Single-site Metropolis Hastings

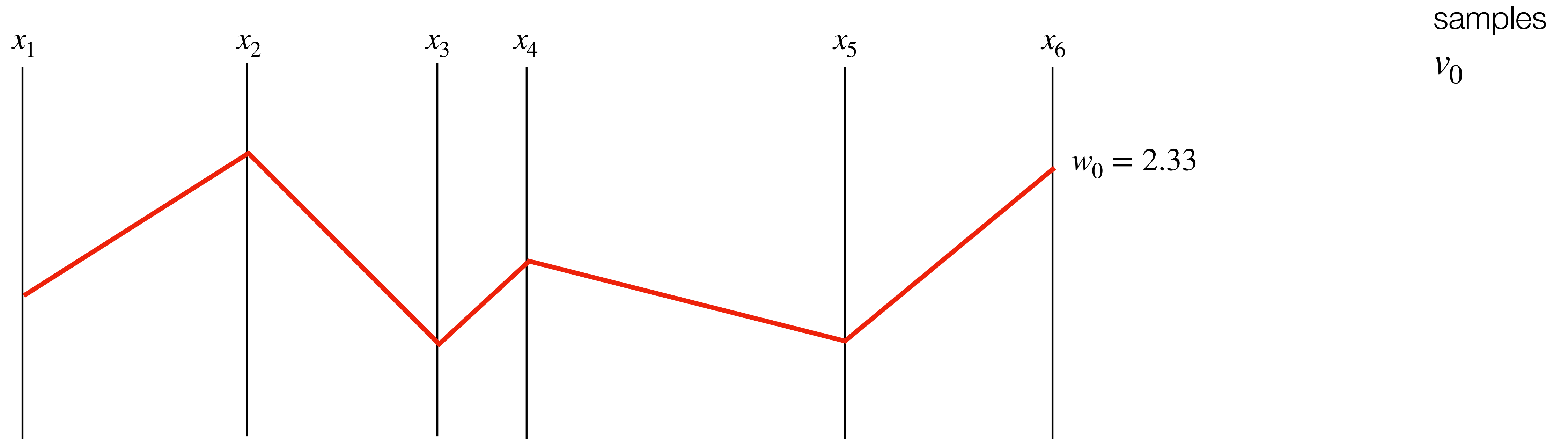
Reuse most of the previous trace (i.e., sampled values)

- Choose one random variable to resample to obtain a new execution
- Accept the trace with probability α
- Otherwise use the previous trace

Single-site Metropolis Hastings

Reuse most of the previous trace (i.e., sampled values)

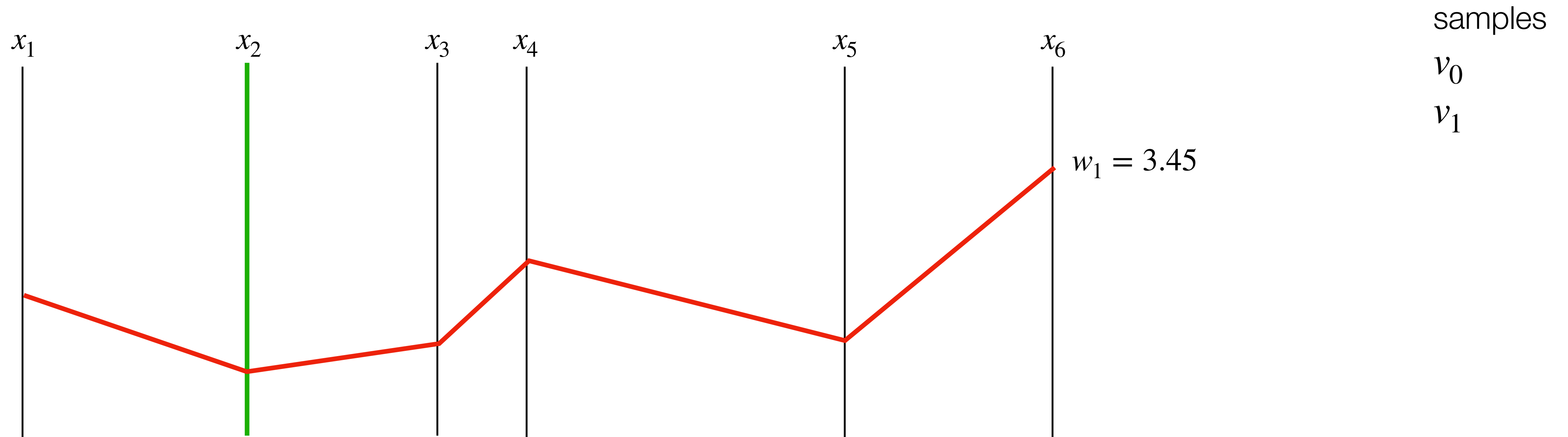
- Choose one random variable to resample to obtain a new execution
- Accept the trace with probability α
- Otherwise use the previous trace



Single-site Metropolis Hastings

Reuse most of the previous trace (i.e., sampled values)

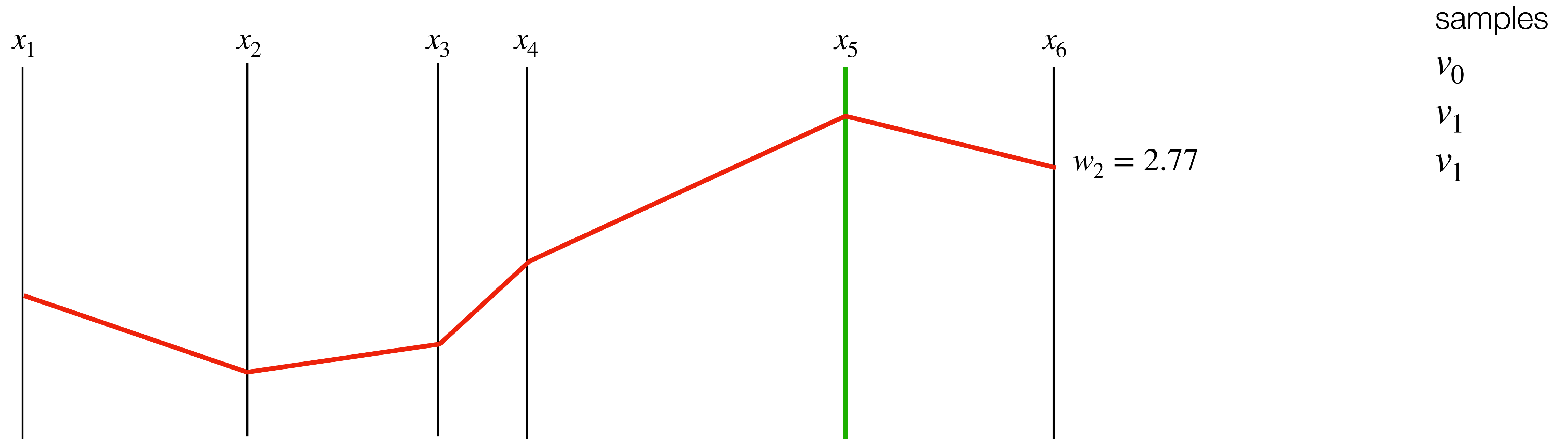
- Choose one random variable to resample to obtain a new execution
- Accept the trace with probability α
- Otherwise use the previous trace



Single-site Metropolis Hastings

Reuse most of the previous trace (i.e., sampled values)

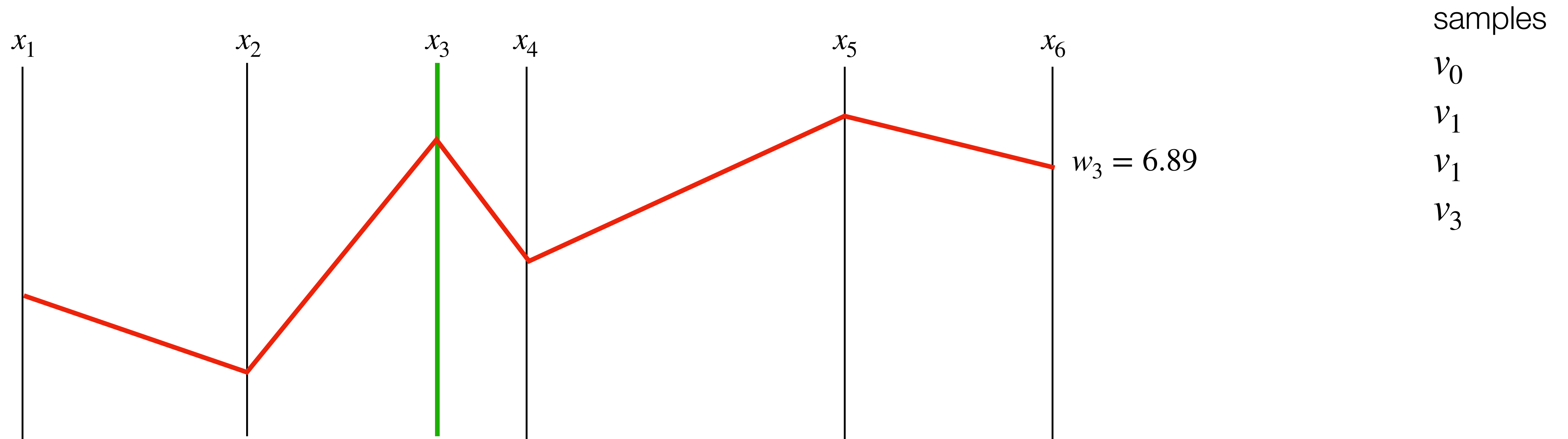
- Choose one random variable to resample to obtain a new execution
- Accept the trace with probability α
- Otherwise use the previous trace



Single-site Metropolis Hastings

Reuse most of the previous trace (i.e., sampled values)

- Choose one random variable to resample to obtain a new execution
- Accept the trace with probability α
- Otherwise use the previous trace



Single-site Metropolis Hastings: acceptance

Track the likelihood of all random variables during execution

- $x = \text{sample}(d) \rightarrow w(x) = (x, \text{pdf}(d)(x))$
- $\text{observe}(d, y) \rightarrow w(y) = (y, \text{pdf}(d)(y))$ score (as in importance sampling)

```
let bimodal y =  
  let z = sample (bernoulli ~p:0.5) in  
  let mu =  
    if z then sample (gaussian ~mu:-1. ~sigma:1.)  
    else sample (gaussian ~mu:1. ~sigma:1.)  
  in  
  let () = observe (gaussian ~mu ~sigma:1.) y in  
  z
```

$$w(z) = (0, 0.5)$$

$$w(\mu) = (1.2, 0.40)$$

$$w(y) = (2.0, 0.27)$$

Single-site Metropolis Hastings: acceptance

- Pick one variable x_0 at random in the trace
- Resample only this variable and all its dependencies $X = X^{\text{sample}} \cup X^{\text{reuse}}$

$$\frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} \mid X_i)}{q(X_i \mid X_{i-1})} = \frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} \mid X_i, x_0) q(x_0 \mid X_i)}{q(X_i \mid X_{i-1}, x_0) q(x_0 \mid X_{i-1})}$$

Single-site Metropolis Hastings: acceptance

- Pick one variable x_0 at random in the trace
- Resample only this variable and all its dependencies $X = X^{\text{sample}} \cup X^{\text{reuse}}$

$$\frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} \mid X_i)}{q(X_i \mid X_{i-1})} = \frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} \mid X_i, x_0) q(x_0 \mid X_i)}{q(X_i \mid X_{i-1}, x_0) q(x_0 \mid X_{i-1})}$$

$$p(X, Y) = \prod_{x \in X} w(x) \prod_{y \in Y} w(y) \quad \text{(not necessarily independent)}$$

$$q(x_0 \mid X) = \frac{1}{|X|} \quad \text{1/(n choices)}$$

$$q(X \mid X', x_0) = \prod_{x \in X^{\text{sample}}} w(x) \quad \text{resampled from the prior}$$

Single-site Metropolis Hastings: acceptance

- Pick one variable x_0 at random in the trace
- Resample only this variable and all its dependencies $X = X^{\text{sample}} \cup X^{\text{reuse}}$

$$\begin{aligned}
 \frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} | X_i)}{q(X_i | X_{i-1})} &= \frac{p(X_i, Y_i)}{p(X_{i-1}, Y_{i-1})} \frac{q(X_{i-1} | X_i, x_0) q(x_0 | X_i)}{q(X_i | X_{i-1}, x_0) q(x_0 | X_{i-1})} \\
 &= \frac{q(x_0 | X_i)}{q(x_0 | X_{i-1})} \frac{p(X_i, Y_i)}{q(X_i | X_{i-1}, x_0)} \frac{q(X_{i-1} | X_i, x_0)}{p(X_{i-1}, Y_{i-1})} \\
 &= \frac{|X_{i-1}|}{|X_i|} \frac{\prod_{x \in X_i} w(x) \prod_{y \in Y_i} w(y)}{\prod_{x \in X_i^{\text{sample}}} w(x)} \frac{\prod_{x \in X_{i-1}^{\text{sample}}} w(x)}{\prod_{x \in X_{i-1}} w(x) \prod_{y \in Y_{i-1}} w(y)} \\
 &= \frac{|X_{i-1}|}{|X_i|} \frac{\prod_{x \in X_i^{\text{reuse}}} w(x)}{\prod_{x \in X_{i-1}^{\text{reuse}}} w(x)} \frac{\prod_{y \in Y_i} w(y)}{\prod_{y \in Y_{i-1}} w(y)} \\
 &\quad \begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ \text{choice } x_0 & \text{reused} & \text{scores} \end{array}
 \end{aligned}$$

$$p(X, Y) = \prod_{x \in X} w(x) \prod_{y \in Y} w(y)$$

$$q(x_0 | X) = \frac{1}{|X|}$$

$$q(X | X', x_0) = \prod_{x \in X^{\text{sample}}} w(x)$$

Single-site Metropolis Hastings

infer.ml

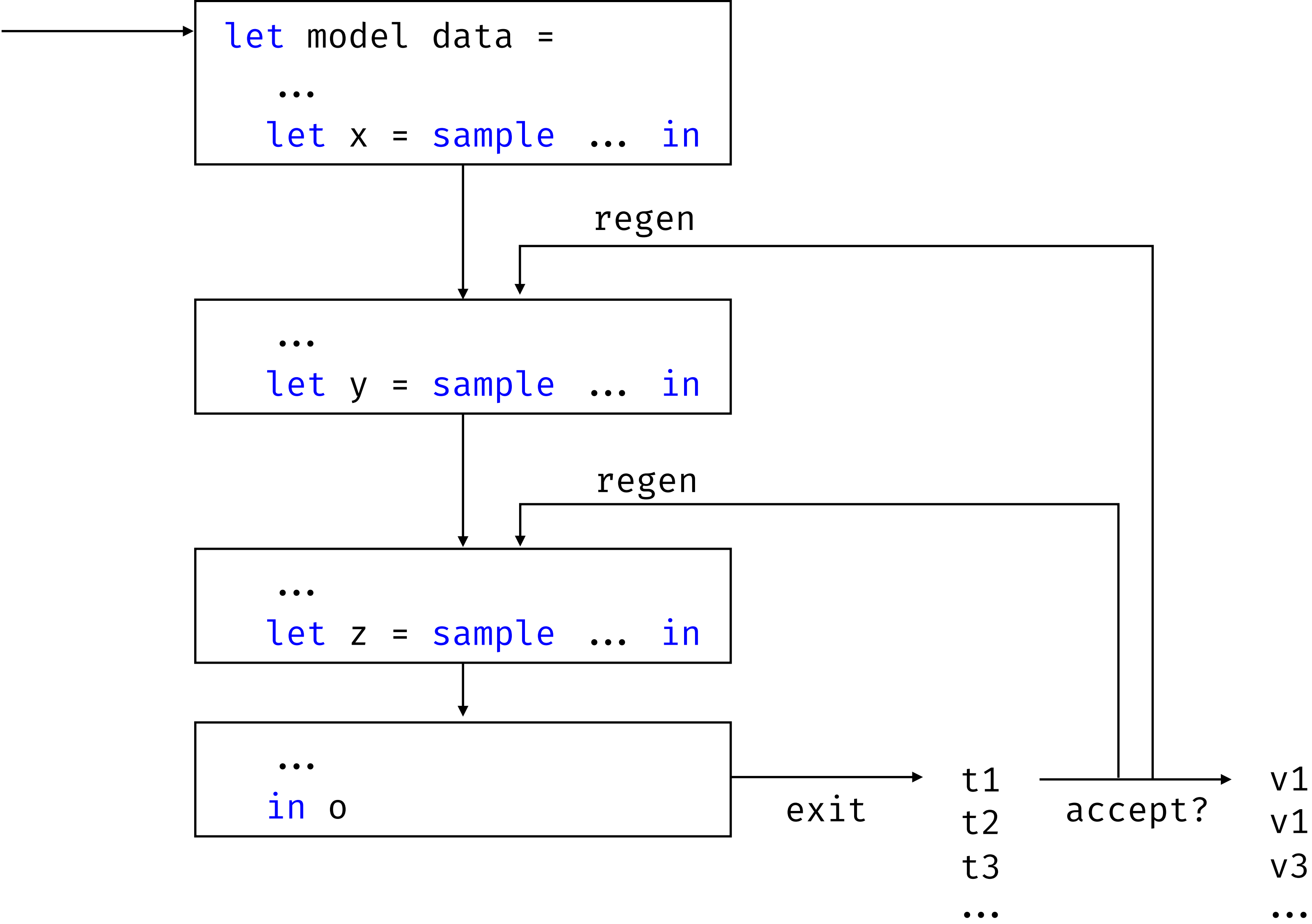
```
module Metropolis_hasting: sig
  type 'a prob
  and 'a next = 'a prob → 'a prob
  and ('a, 'b) model = 'a → ('b → 'b next) → 'b next

  val sample : 'a Distribution.t → ('a → 'b next) → 'b next
  val factor : float → (unit → 'b next) → 'b next
  val infer : ('a, 'b) model → 'a → 'b Distribution.t
end = struct ... end
```

Inference algorithm

- Generate an execution trace
- Compare the score to the previous trace
- Accept (keep the new trace) or reject (keep the previous trace)
- Draw a sample site at random in the execution path
- Restart execution from this site to generate a new trace
- Stop after n iterations

Single-site Metropolis Hastings



Single-site Metropolis Hastings

```
module Metropolis_hasting = struct
  type 'a prob = {
    score : float;
    trace : 'a sample_site list;
    value : 'a option;
  }

  and 'a sample_site =
    | Sample : {
      k : 'b → 'a next;
      score : float;
      dist : 'b Distribution.t;
    }
    → 'a sample_site

  and 'a next = 'a prob → 'a prob
  and ('a, 'b) model = 'a → ('b → 'b next) → 'b next
```

Single-site Metropolis Hastings

```
let sample dist k prob =  
  let value = Distribution.draw dist in  
  let sample_site = Sample { k; score = prob.score; dist } in  
  k value { prob with trace = sample_site :: prob.trace }  
  
let factor s k prob = k () { prob with score = prob.score +. s }  
  
let observe d x k prob = factor (Distribution.logpdf d x) k prob  
  
let exit v prob = { prob with value = Some v }
```

Single-site Metropolis Hastings

```
let mh prob prob' =  
  let fw = -.log (prob.trace ▷ List.length ▷ Float.of_int) in  
  let bw = -.log (prob'.trace ▷ List.length ▷ Float.of_int) in  
  min 1. (exp (prob'.score -. prob.score +. bw -. fw))
```

Single-site Metropolis Hastings

```
let rec gen n values prob =  
  if n = 0 then values  
  else  
    let regen_from = Random.int (List.length prob.trace) in  
    let (Sample regen) = List.nth prob.trace regen_from in  
    let prob' =  
      sample regen.dist regen.k  
      {  
        prob with  
        trace = Utils.slice prob.trace regen_from;  
        score = regen.score;  
      }  
    in  
    let next_prob = if Random.float 1. < mh prob prob' then prob' else prob in  
    gen (n - 1) (Option.get next_prob.value :: values) next_prob
```

Single-site Metropolis Hastings

```
let infer ?(n = 1000) m data =  
  let prob = (m data) exit { score = 0.; trace = []; value = None } in  
  let values = gen n [] prob ▷ Array.of_list in  
  Distribution.uniform_support ~values  
end
```

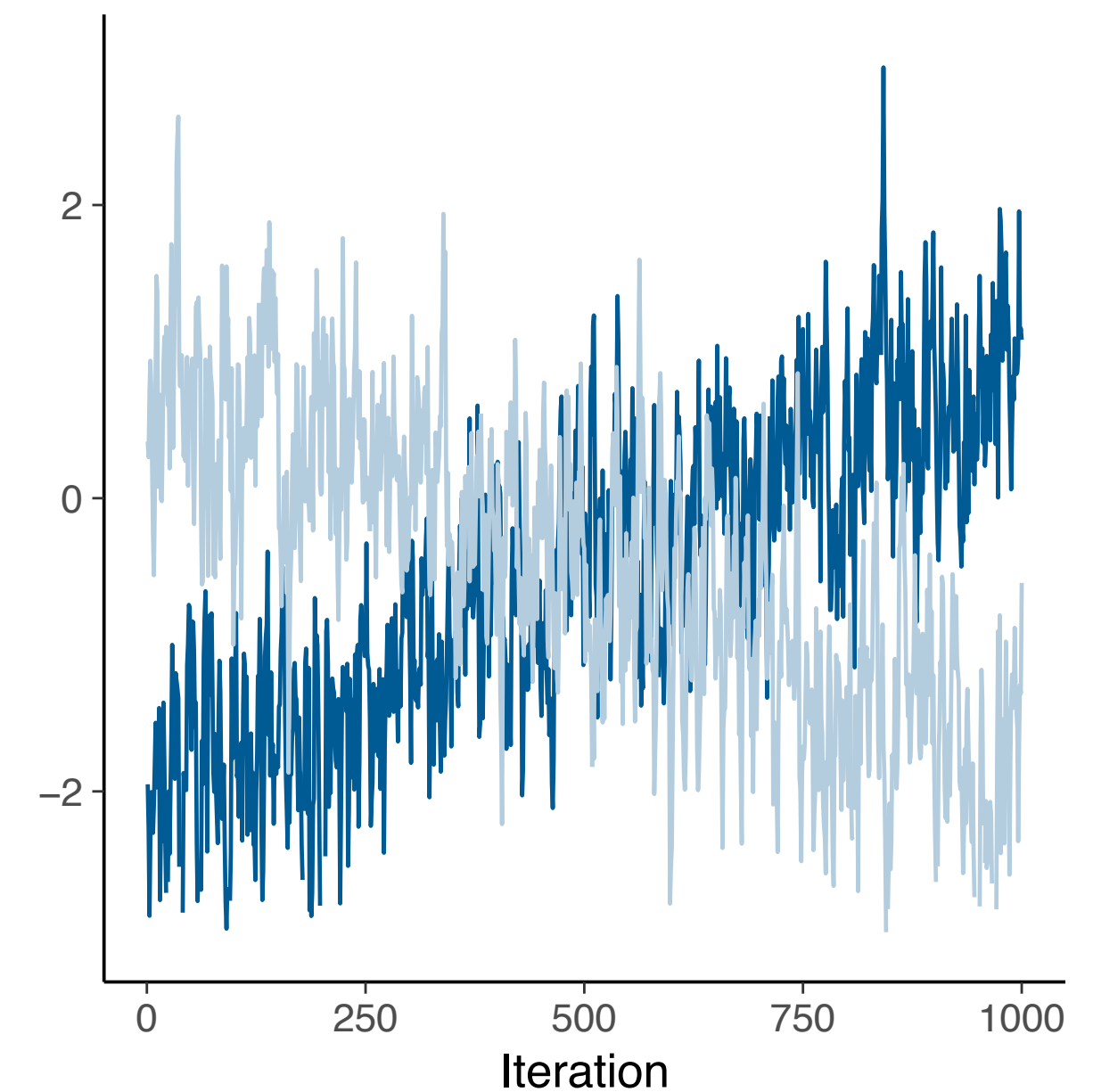
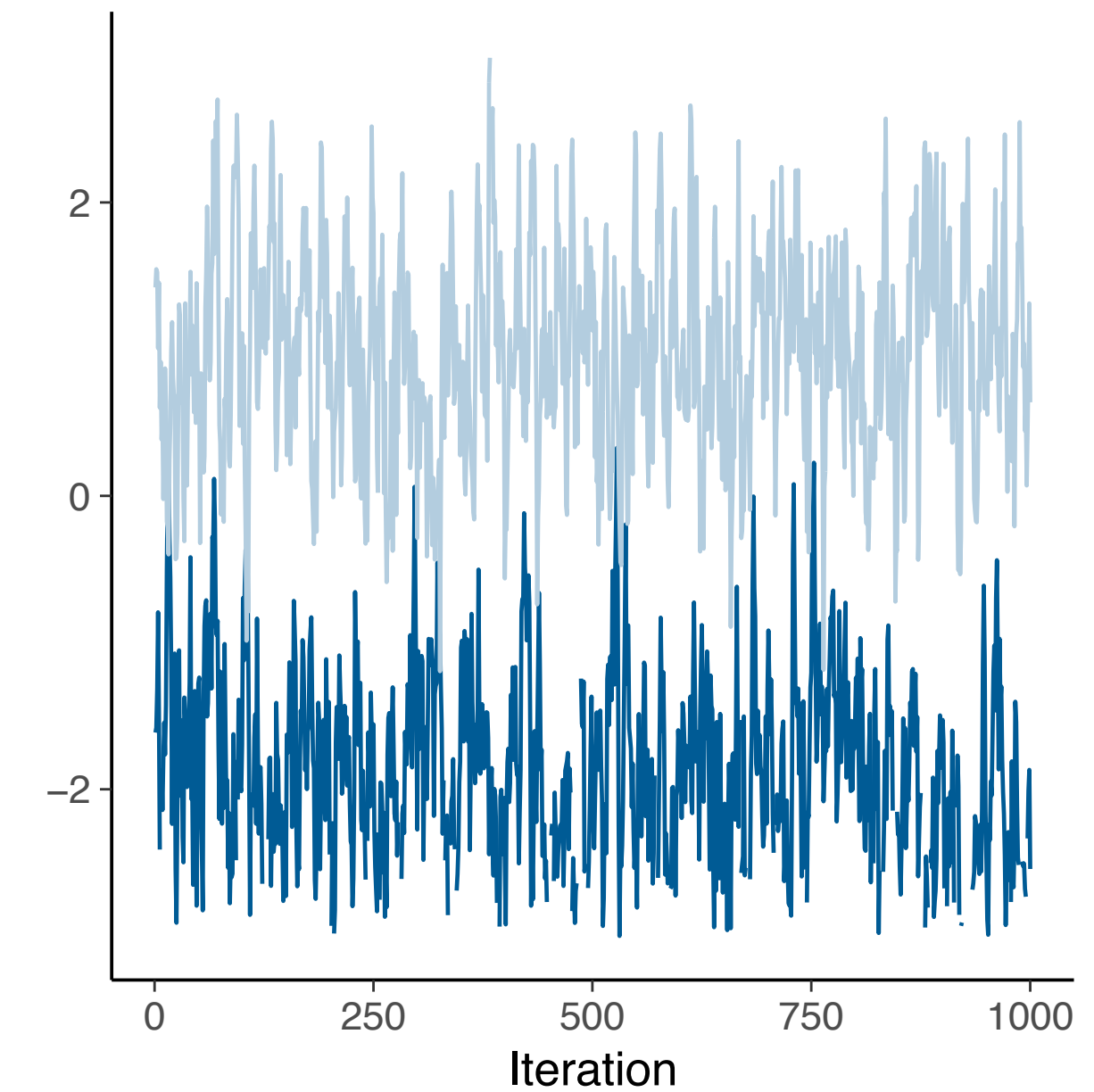

Limitations

Convergence

- Theoretical conditions are complex
- Must be checked experimentally
- Diagnostic tools: trace plot, R-hat (multi-chains)
- Solution: warmup, change initial conditions, reparameterization, ...

Sample correlation

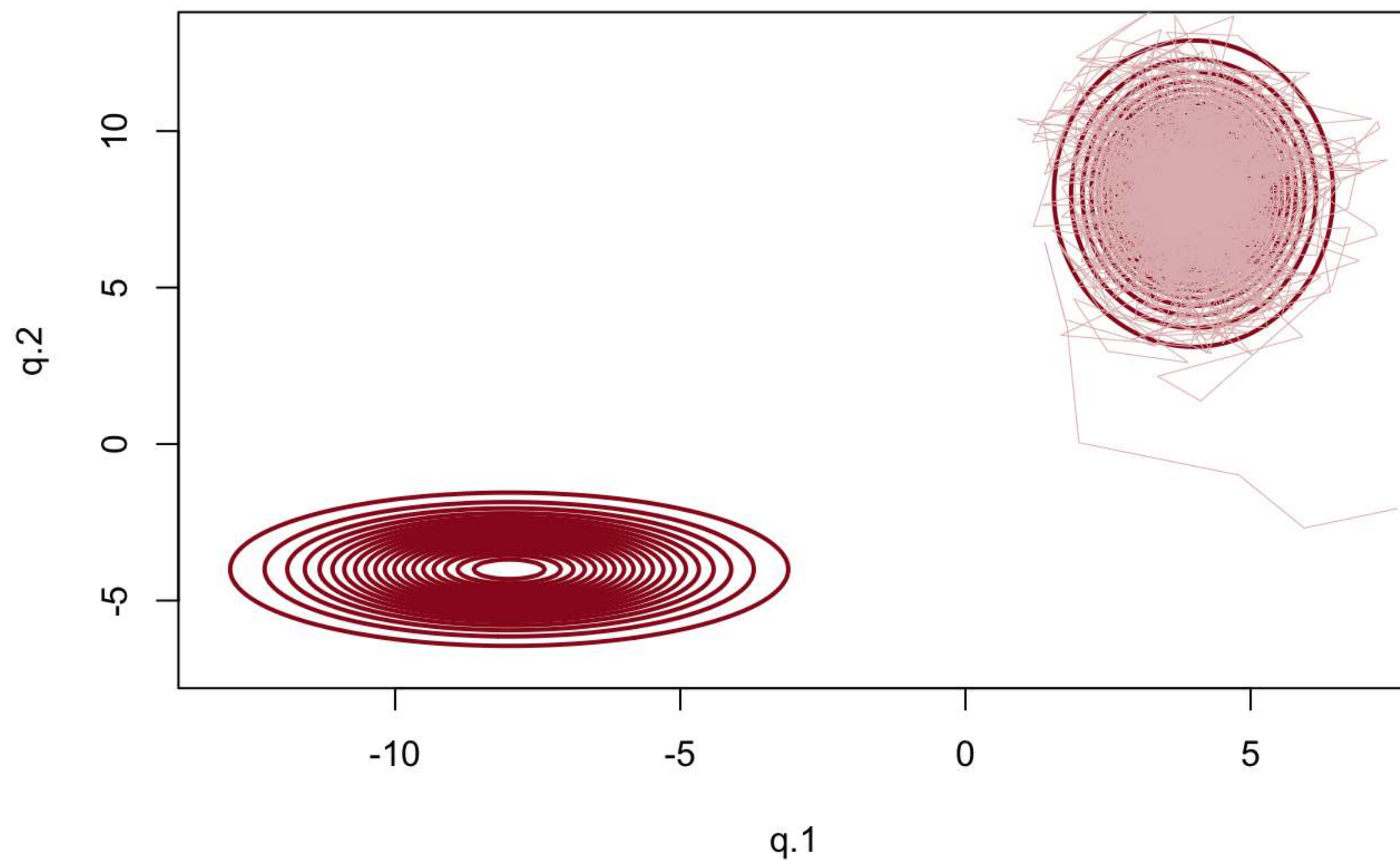
- Next sample depends on the previous one
- Diagnostic tools ESS (effective sample size)
- Solution: thinning, ...



Limitations

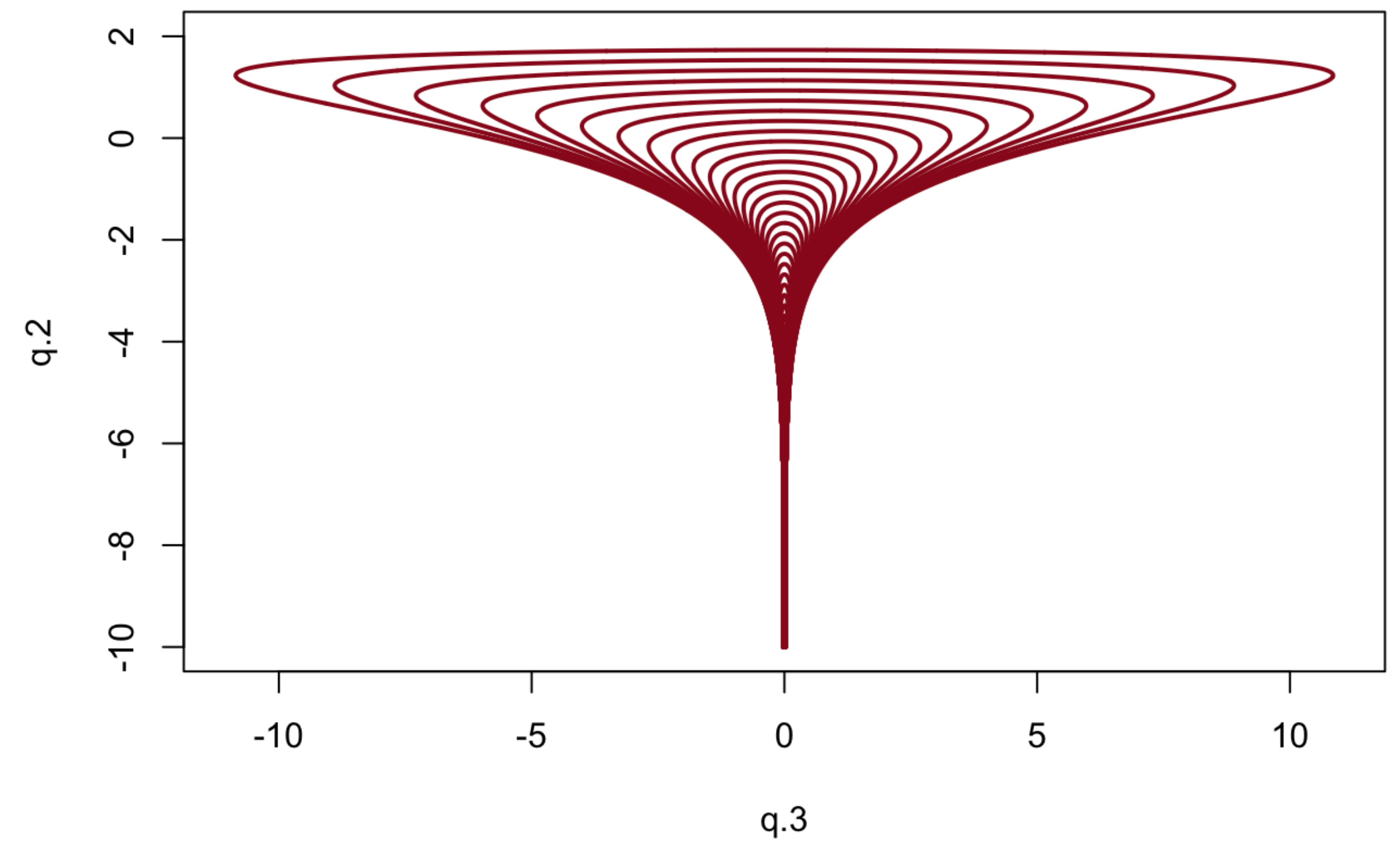
Pathological models

Metastable Target Density



Multimodal distribution

Funnel Target Density



Neal's funnel

Advanced inference

Probabilistic Programming Languages

Stan: A probabilistic programming language

“A Stan program imperatively defines a log probability function over parameters conditioned on specified data and constants”

Brief History:

- Named after Stanisław Marcin Ulam (Monte Carlo method)
- Initial release: August 30, 2012
- Today: version 2.33 (50+ iterations after 1.0.0)

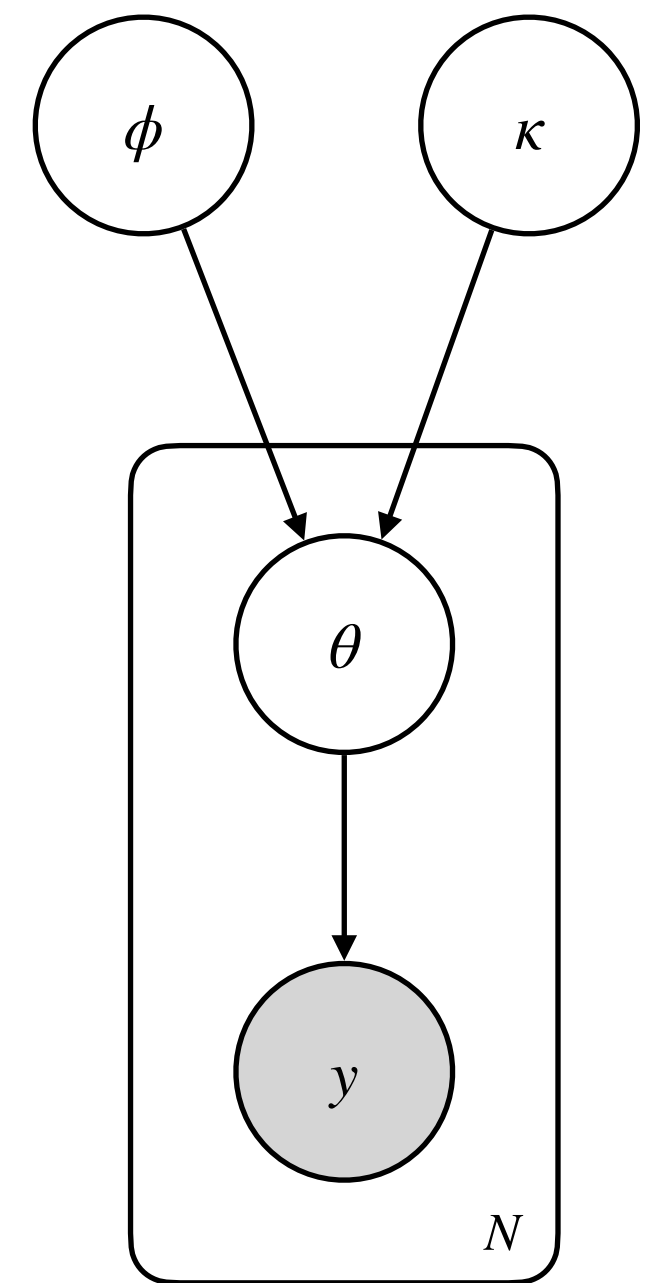
In a nutshell:

- Small imperative language to describe probabilistic models
- Bayesian inference on continuous latent variables
- No U-Turn Sampler (NUTS): an optimized Hamiltonian Monte Carlo (HMC) inference
- Interface with popular programming language: R, Python, etc...



Stan: reminders

```
data {  
  int<lower=0> N;           // players  
  int<lower=0> K[N];        // initial trials  
  int<lower=0> y[N];        // initial successes  
}  
parameters {  
  real<lower=0, upper=1> phi; // population chance of success  
  real<lower=1> kappa;        // population concentration  
  vector<lower=0, upper=1>[N] theta; // chance of success  
}  
model {  
  kappa ~ pareto(1, 1.5);    // hyperprior  
  theta ~ beta(phi * kappa, (1 - phi) * kappa); // prior  
  y ~ binomial(K, theta);    // likelihood  
}
```



Stan semantics

Parameters and Data blocks

- Declared in their respective blocks (before the model)
- Precisely identified before inference

Model block

- A small imperative deterministic language with a global accumulator `target`
- Two constructs to update this accumulator: `~` and `target +=`

The language is deterministic: statements update the environment

- $\llbracket s \rrbracket : \Gamma \rightarrow \Gamma$
- The model block defines an un-normalized log-density function

model is deterministic

Probabilistic semantics: turn the log-density into a measure over parameters

- $\{p\} : \mathcal{D} \rightarrow \Sigma_X \rightarrow [0, \infty)$, where X is the parameters domain
- $\{p\}_D = \lambda U. \int_U \exp(\llbracket \text{model}(p) \rrbracket_{D,\theta}(\text{target})) d\theta.$
- D is the initial environment containing the (observed) data.

function of parameter and data

Stan semantics

Statements: $\llbracket s \rrbracket : \Gamma \rightarrow \Gamma$

$$\llbracket \text{skip} \rrbracket_{\gamma} = \gamma$$

$$\llbracket x = e \rrbracket_{\gamma} = \gamma + [x \leftarrow \llbracket e \rrbracket_{\gamma}]$$

$$\llbracket x[e_1, \dots, e_n] = e \rrbracket_{\gamma} = \gamma + [x[\llbracket e_1 \rrbracket_{\gamma}, \dots, \llbracket e_n \rrbracket_{\gamma}] \leftarrow \llbracket e \rrbracket_{\gamma}]$$

$$\llbracket \text{if}(e) s_1 \text{ else } s_2 \rrbracket_{\gamma} = \text{if } \llbracket e \rrbracket_{\gamma} = 0 \text{ then } \llbracket s_1 \rrbracket_{\gamma} \text{ else } \llbracket s_2 \rrbracket_{\gamma}$$

$$\llbracket s_1 ; s_2 \rrbracket_{\gamma} = \llbracket s_2 \rrbracket_{\llbracket s_1 \rrbracket_{\gamma}}$$

$$\llbracket \text{for}(x \text{ in } e_1 : e_2) \{s\} \rrbracket_{\gamma} = \text{let } n_1 = \llbracket e_1 \rrbracket_{\gamma} \text{ in let } n_2 = \llbracket e_2 \rrbracket_{\gamma} \text{ in} \\ \text{if } n_1 > n_2 \text{ then } \gamma \text{ else } \llbracket \text{for}(x \text{ in } n_1 + 1 : n_2) \{s\} \rrbracket_{\llbracket s \rrbracket_{\gamma} + [x \leftarrow n_1]}$$

$$\llbracket \text{while}(e) \{s\} \rrbracket_{\gamma} = \text{if } \llbracket e \rrbracket_{\gamma} = 0 \text{ then } \gamma \text{ else } \llbracket \text{while}(e) \{s\} \rrbracket_{\llbracket s \rrbracket_{\gamma}}$$

$$\llbracket \text{target} += e \rrbracket_{\gamma} = \gamma + [\text{target} \leftarrow \gamma(\text{target}) + \llbracket e \rrbracket_{\gamma}]$$

$$\llbracket e_1 \sim e_2 \rrbracket = \text{let } \mu = \llbracket e_2 \rrbracket_{\gamma} \text{ in } \llbracket \text{target} += \text{pdf}(\mu)(e_1) \rrbracket_{\gamma}$$

imperative, deterministic

Density semantics beyond Stan

Remember: sampler semantics

- Expressions are interpreted as weighted samplers in log space
- Given an environment γ , $\llbracket e \rrbracket_\gamma = v, w$
- $\llbracket e \rrbracket : \Gamma \rightarrow V \times \mathbb{R}$
- Parameters are now inputs

$$\begin{aligned}\llbracket c \rrbracket_\gamma &= c, 0 \\ \llbracket x \rrbracket_\gamma &= \gamma(x), 0 \\ \llbracket \text{sample}(e) \rrbracket_\gamma &= ??? \\ \llbracket \text{factor}(e) \rrbracket_\gamma &= (), -\llbracket e \rrbracket_\gamma \\ \llbracket \text{observe}(e_1, e_2) \rrbracket_\gamma &= (), -\log\text{pdf}(\llbracket e_1 \rrbracket_\gamma)(\llbracket e_2 \rrbracket_\gamma) \\ \llbracket \text{let } x = e_1 \text{ in } e_2 \rrbracket &= \text{let } v_1, w_1 = \llbracket e_1 \rrbracket_\gamma \text{ in} \\ &\quad \text{let } v_2, w_2 = \llbracket e_2 \rrbracket_{\gamma + [x \leftarrow v_1]} \text{ in} \\ &\quad v_2, w_1 + w_2\end{aligned}$$

```
sample mu (* where mu is defined on [a, b] *)  
≡  
let x = sample (uniform (a, b)) in  
let () = observe (mu, x) in  
x
```

all parameters must be know statically

`let x = sample(e) in ... → observe(e, x) ; ...`

Hamiltonian Monte-Carlo (HMC)

Preferred inference algorithm for Stan

Analogy: Particle in an energy field

- Program define a density of the form $\exp(-U(X))$
- On continuous spaces U can be interpreted as an energy
- Low energy wells correspond to high probability regions
- HMC simulate the trajectory of a particle in this energy field

Hamiltonian Dynamics

- M : mass matrix
- P : momentum

$$K(P) = \frac{1}{2} P^T M^{-1} P$$

The diagram illustrates the components of the Hamiltonian. The central equation is $H(X, P) = K(P) + U(X)$. An upward arrow from $H(X, P)$ points to the label "hamiltonian". A downward arrow from $K(P)$ points to the label "kinetic energy". An upward arrow from $U(X)$ points to the label "potential energy".

$$\begin{array}{ccc} \text{hamiltonian} & & \text{potential energy} \\ \uparrow & & \uparrow \\ H(X, P) = K(P) + U(X) \\ \downarrow & & \\ \text{kinetic energy} & & \end{array}$$

Hamiltonian Monte-Carlo (HMC)

Energy conservation

$$\frac{dH}{dt} = (\nabla_P H)^T \frac{dP}{dt} + (\nabla_X H)^T \frac{dX}{dt}$$

Hamiltonian dynamics

$$\begin{cases} \frac{dX}{dt} = \nabla_P H(X, P) = M^{-1} P \\ \frac{dP}{dt} = -\nabla_X H(X, P) = -\nabla_X U(X) \end{cases}$$

The diagram illustrates the decomposition of the Hamiltonian function $H(X, P)$ into its constituent parts. The central equation is $H(X, P) = K(P) + U(X)$. An upward-pointing arrow from $H(X, P)$ is labeled "hamiltonian". A downward-pointing arrow from $K(P)$ is labeled "kinetic energy". An upward-pointing arrow from $U(X)$ is labeled "potential energy".

$$\begin{array}{ccc} \text{hamiltonian} & & \text{potential energy} \\ \uparrow & & \uparrow \\ H(X, P) = K(P) + U(X) & & \\ \downarrow & & \\ \text{kinetic energy} & & \end{array}$$

Hamiltonian Monte-Carlo (HMC)

Generate samples (X, P) from the density $\exp(-H(X, P))$

- At each iteration
- Sample an initial momentum $P_0 \sim \mathcal{N}(0, M)$
- Solve the Hamiltonian dynamics (discretized)
- Perform a Metropolis Hastings update with probability α

$$\alpha = \min \left(1, \frac{\exp(-H(X_i, P_i))}{\exp(-H(X_{i-1}, P_{i-1}))} \right)$$

momentum can then be marginalized

If the hamiltonian is preserved: accept with probability 1.

- Problem: numerical approximations
- Solution: leapfrog integrator and reject using MH acceptance probability

Hamiltonian Monte-Carlo (HMC)

```
let u x = let _, u = model data x in u
let k p = 0.5 * transpose p * inv m * p

let h x p = u x +. k p

let rec gen n values x =
  if n = 0 then values
  else
    let p = Distribution.draw mv_normal(0, m) in          autodiff magic!
    let x', p' = leapfrog (grad u) x p in
    let next_x = if Random.float 1. < exp(h x p -. h x' p') then x' else x in
    let next_value, _ = model data next_x in
    gen (n - 1) (next_value :: values) next_x
```

Warning: pseudo-code

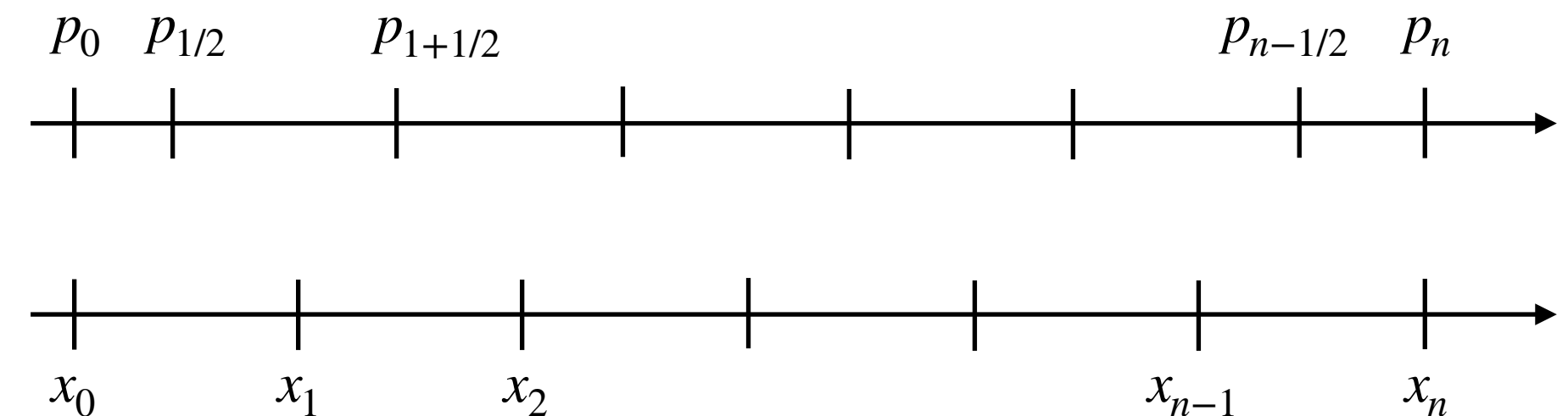
Leapfrog integration

```
let leapfrog u_grad x0 p0 =  
  let p = p0 - 0.5 * step_size * u_grad x0 in  
  let rec loop n x p =  
    if n = 0 then x, p  
    else  
      let x' = x + step_size * p in  
      let p' = p - step_size * u_grad x' in  
      loop (n-1) x' p'  
  in  
  let xt, pt = loop (path_len - 1) x0 p in  
  let x' = xt + step_size * pt in  
  let p' = pt - 0.5 * step_size * u_grad x' in  
  x', p'
```

Warning: pseudo-code

first half-step for the momentum

last half-step for the momentum



Stochastic Variational Inference (SVI)

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)} = \frac{p(x | z)p(z)}{\int_z p(x | z)p(z)dz}$$

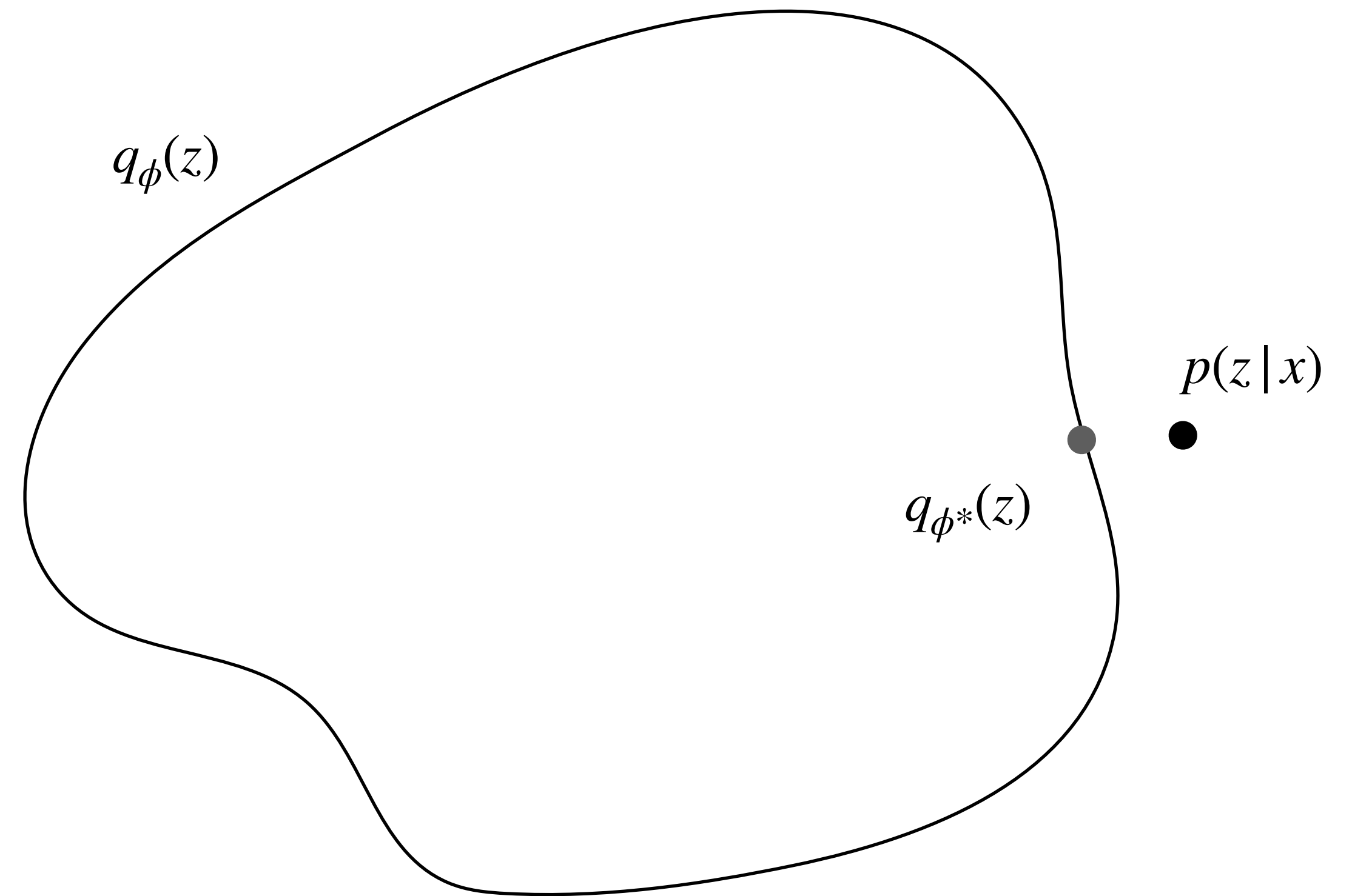
Variational family

- Parameterized by a parameter ϕ
- Find the closest member to the posterior $q_{\phi^*}(z)$
- Optimization problem

Metrics: Kullback-Leibler divergence

$$KL(q(x) \parallel p(x)) = - \int q(x) \log \frac{p(x)}{q(x)}$$

- $KL(q \parallel p) \geq 0$ positive
- $KL(q \parallel p) = 0 \iff |x| \neq 0 \implies p(x) = q(x)$, equal almost everywhere
- $KL(q \parallel p) \neq KL(p \parallel q)$ asymmetric
- No triangular inequality



Stochastic Variational Inference (SVI)

$$\begin{aligned} KL(q_\phi(z) \parallel p(z|x)) &= - \int q_\phi(z) \log \frac{p(z|x)}{q_\phi(z)} dz \\ &= - \int q_\phi(z) \log \frac{p(x, z)}{p(x)q_\phi(z)} dz \\ &= - \int q_\phi(z) \log \frac{p(x, z)}{q_\phi(z)} dz + \int q_\phi(z) \log p(x) dz \\ &= - \int q_\phi(z) \log \frac{p(x, z)}{q_\phi(z)} dz + \log p(x) \end{aligned}$$

Stochastic Variational Inference (SVI)

$$\begin{aligned} KL(q_\phi(z) \parallel p(z|x)) &= - \int q_\phi(z) \log \frac{p(z|x)}{q_\phi(z)} dz \\ &= - \int q_\phi(z) \log \frac{p(x, z)}{p(x)q_\phi(z)} dz \\ &= - \int q_\phi(z) \log \frac{p(x, z)}{q_\phi(z)} dz + \int q_\phi(z) \log p(x) dz \\ &= - \int q_\phi(z) \log \frac{p(x, z)}{q_\phi(z)} dz + \log p(x) \end{aligned}$$

$$\log p(x) = \underbrace{KL(q_\phi(z) \parallel p(x \mid z))}_{\text{minimize}} + \underbrace{\int q_\phi(z) \log \frac{p(x, z)}{q_\phi(z)} dz}_{\text{maximize ELBO}}$$

Stochastic Variational Inference (SVI)

How to solve the optimisation problem?

Program your own guide

- Pyro (first versions)
- Sample the same variables in the guide and the model

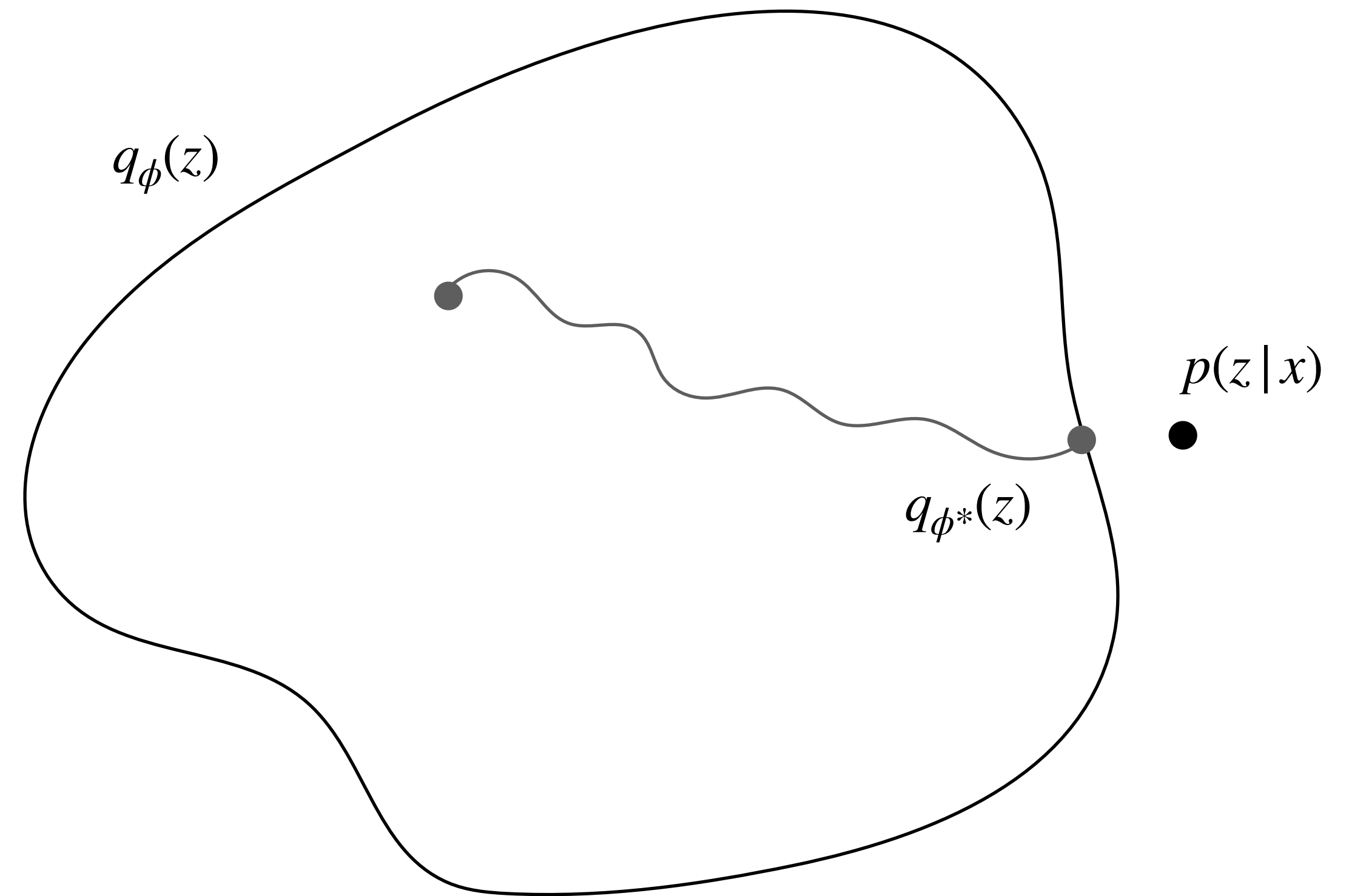
```
def model():  
    pyro.sample("z_1", ... )
```

```
def guide():  
    pyro.sample("z_1", ... )
```

Approximate gradient ascent.

$$\nabla_{\phi} KL(q_{\phi}(z) \parallel p(z \mid x)) \longrightarrow \nabla_{\phi} \mathcal{L} = \nabla_{\phi} \int q_{\phi}(z) \log \frac{p(x, z)}{q_{\phi}(z)} dz$$

autodiff magic!



Stochastic Variational Inference (SVI)

How to solve the optimisation problem?

Black-box variational inference

- Variational families with tractable solution
- Mean-field approximation $q_{\phi}(z) = \prod_{i=1}^n \mathcal{N}(z_i | \mu_i, \sigma_i)$ where $\phi = \{\mu_i, \sigma_i\}_{i \in [1, n]}$
- Full-rank approximation $q_{\phi}(z) = \mathcal{N}(z | \mu, \Sigma)$ where $\phi = (\mu, \Sigma)$

Assumptions

- Independences between random variables
- Only use Gaussians distributions

References

An Introduction to Probabilistic Programming

Jan-Willem van de Meent, Brooks Paige, Hongseok Yang, Frank Wood

<https://arxiv.org/abs/1809.10756>

Lightweight Implementations of Probabilistic Programming Languages Via Transformational Compilation

David Wingate Andreas Stuhlmüller Noah D. Goodman

AISTATS 2011

Markov Chain Monte Carlo in Practice

Michael Bettancourt

https://betanalpha.github.io/assets/case_studies/markov_chain_monte_carlo.html

Variational Inference: A Review for Statisticians

David M. Blei, Alp Kucukelbir, Jon D. McAuliffe

<https://arxiv.org/abs/1601.00670>