

# Advanced Simulation - Lecture 7

Patrick Rebeschini

February 5th, 2018

- Given a target  $\pi(x) = \pi(x_1, x_2, \dots, x_d)$ , Gibbs sampling works by sampling from  $\pi_{X_j|X_{-j}}(x_j|x_{-j})$  for  $j = 1, \dots, d$ .
- Sampling exactly from one of these conditionals might be a hard problem itself.
- Even if it is possible, the Gibbs sampler might converge slowly if components are highly correlated.
- If the components are not highly correlated then Gibbs sampling performs well, even when  $d \rightarrow \infty$ , e.g. with an error increasing “only” polynomially with  $d$ .
- Metropolis–Hastings algorithm (1953, 1970) is a more general algorithm that can bypass these problems.
- Additionally Gibbs can be recovered as a special case.

# Metropolis–Hastings algorithm

- Target distribution on  $\mathbb{X} = \mathbb{R}^d$  of density  $\pi(x)$ .
- Proposal distribution: for any  $x, x' \in \mathbb{X}$ , we have  $q(x'|x) \geq 0$  and  $\int_{\mathbb{X}} q(x'|x) dx' = 1$ .
- Starting with  $X^{(1)}$ , for  $t = 2, 3, \dots$

**1** Sample  $X^* \sim q(\cdot | X^{(t-1)})$ .

**2** Compute

$$\alpha(X^* | X^{(t-1)}) = \min \left( 1, \frac{\pi(X^*) q(X^{(t-1)} | X^*)}{\pi(X^{(t-1)}) q(X^* | X^{(t-1)})} \right).$$

**3** Sample  $U \sim \mathcal{U}_{[0,1]}$ . If  $U \leq \alpha(X^* | X^{(t-1)})$ , set  $X^{(t)} = X^*$ , otherwise set  $X^{(t)} = X^{(t-1)}$ .

# Metropolis–Hastings algorithm

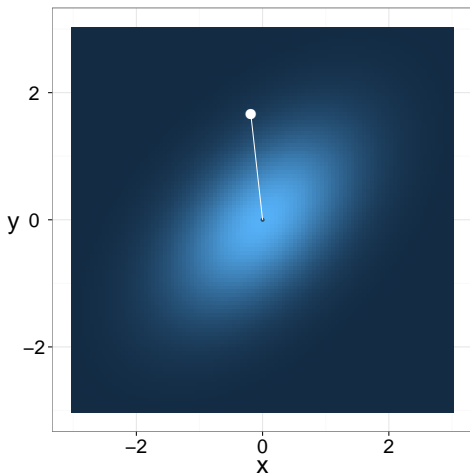


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

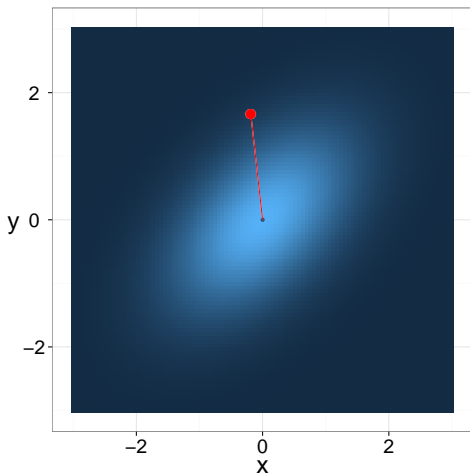


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

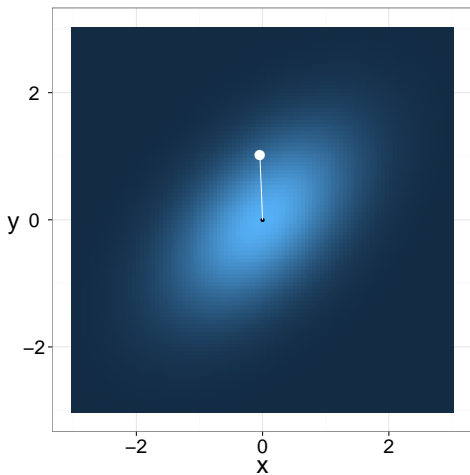


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

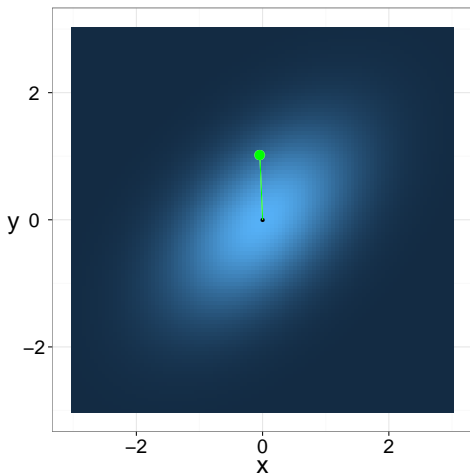


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

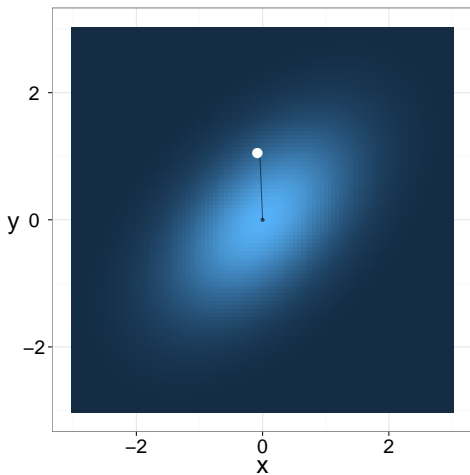


Figure: Metropolis–Hastings on a bivariate Gaussian target.



# Metropolis–Hastings algorithm

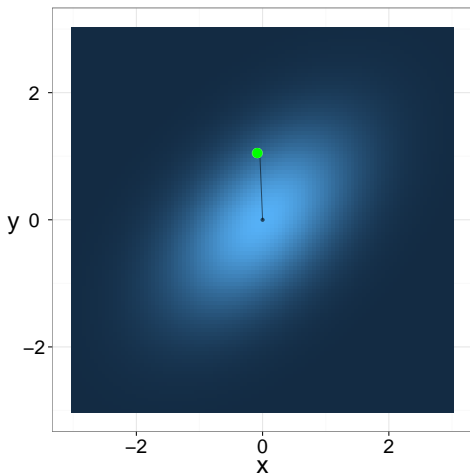


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

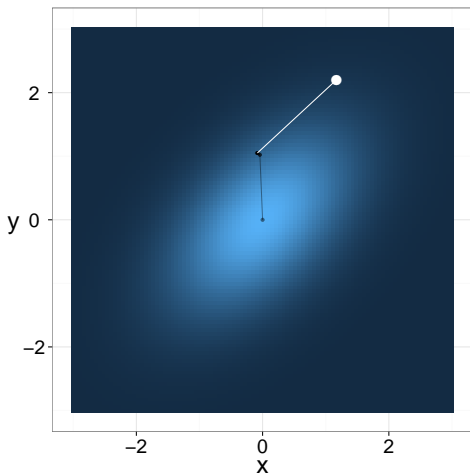


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

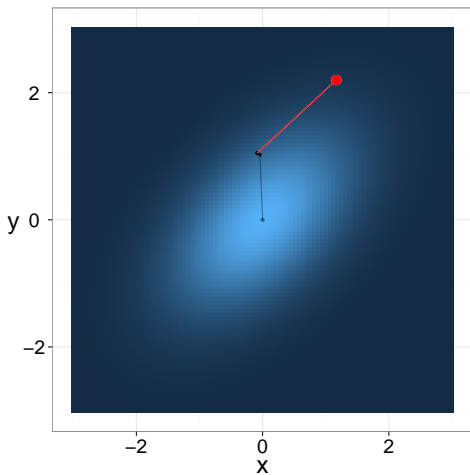


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

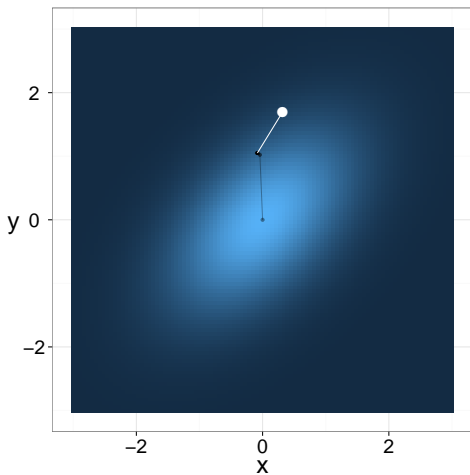


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

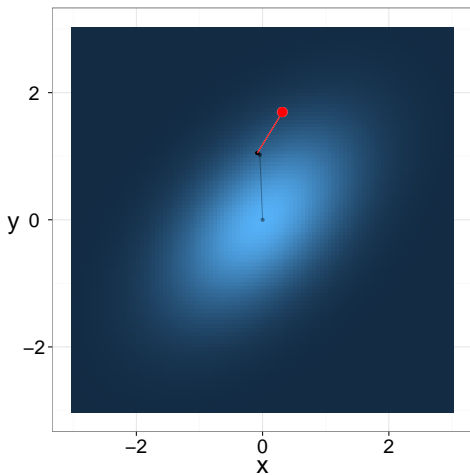


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

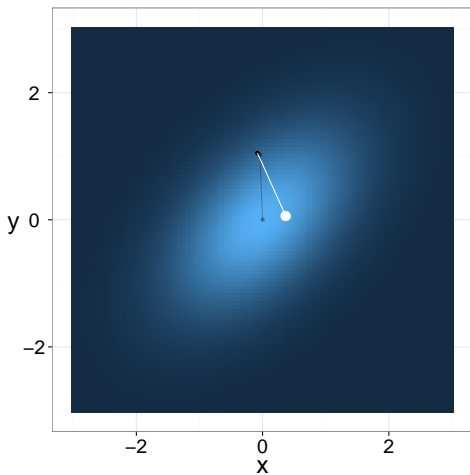


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

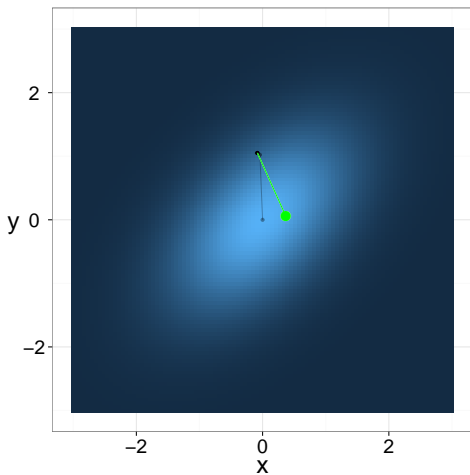


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

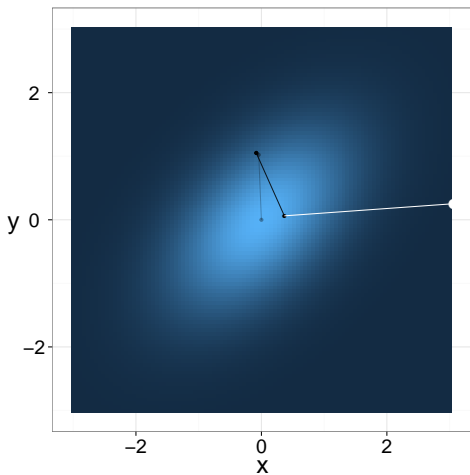


Figure: Metropolis–Hastings on a bivariate Gaussian target.



# Metropolis–Hastings algorithm

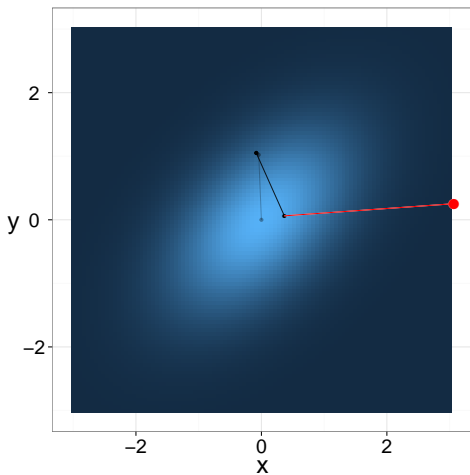


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

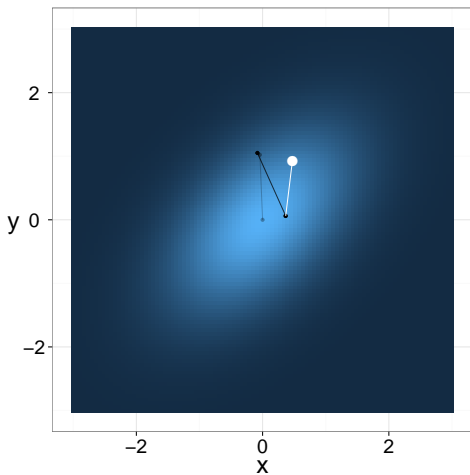


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

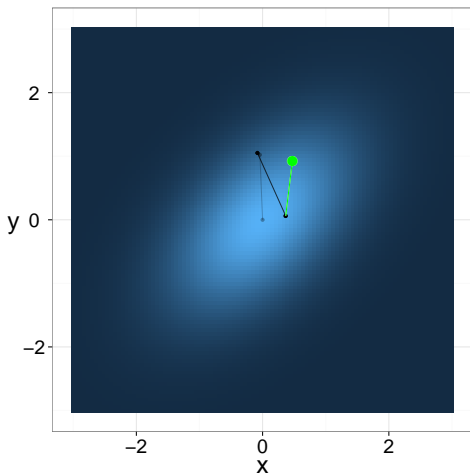


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

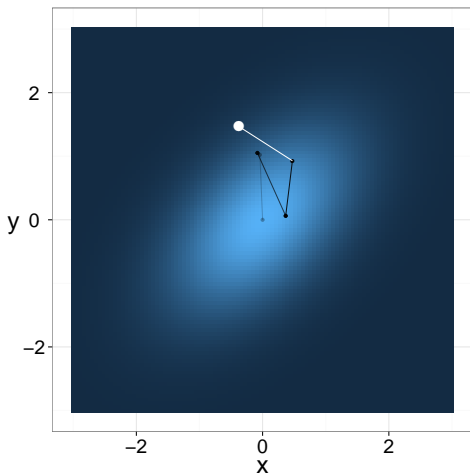


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

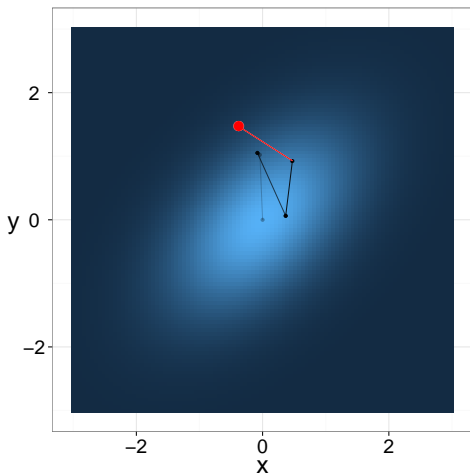


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

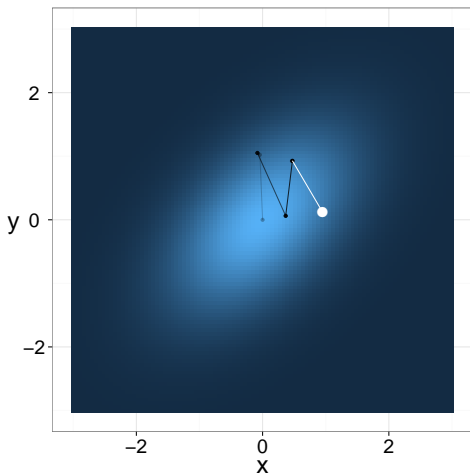


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

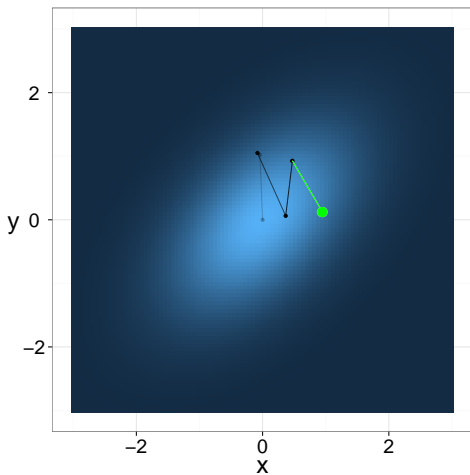


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

- Metropolis–Hastings only requires point-wise evaluations of  $\pi(x)$  up to a normalizing constant; indeed if  $\tilde{\pi}(x) \propto \pi(x)$  then

$$\frac{\pi(x^*) q(x^{(t-1)} | x^*)}{\pi(x^{(t-1)}) q(x^* | x^{(t-1)})} = \frac{\tilde{\pi}(x^*) q(x^{(t-1)} | x^*)}{\tilde{\pi}(x^{(t-1)}) q(x^* | x^{(t-1)})}.$$

- At each iteration  $t$ , a candidate is proposed. The probability of a candidate being accepted is given by

$$a(x^{(t-1)}) = \int_{\mathbb{X}} \alpha(x | x^{(t-1)}) q(x | x^{(t-1)}) dx$$

in which case  $X^{(t)} = X$ , otherwise  $X^{(t)} = X^{(t-1)}$ .

- This algorithm clearly defines a Markov chain  $(X^{(t)})_{t \geq 1}$ .



# Transition Kernel and Reversibility

- **Lemma.** The transition kernel of the Metropolis–Hastings algorithm is given by

$$K(y | x) \equiv K(x, y) = \alpha(y | x)q(y | x) + (1 - a(x))\delta_x(y)$$

where  $\delta_x$  denotes the Dirac mass at  $x$ .

- *Proof.* We have

$$\begin{aligned} K(x, y) &= \int q(x^* | x) \{ \alpha(x^* | x) \delta_{x^*}(y) + (1 - \alpha(x^* | x)) \delta_x(y) \} dx^* \\ &= q(y | x) \alpha(y | x) + \left\{ \int q(x^* | x) (1 - \alpha(x^* | x)) dx^* \right\} \delta_x(y) \\ &= q(y | x) \alpha(y | x) + \left\{ 1 - \int q(x^* | x) \alpha(x^* | x) dx^* \right\} \delta_x(y) \\ &= q(y | x) \alpha(y | x) + \{1 - a(x)\} \delta_x(y). \end{aligned}$$

- **Proposition.** The Metropolis–Hastings kernel  $K$  is  $\pi$ –reversible and thus admit  $\pi$  as invariant distribution.
- *Proof.* For any  $x, y \in \mathbb{X}$ , with  $x \neq y$

$$\begin{aligned}\pi(x)K(x, y) &= \pi(x)q(y | x)\alpha(y | x) \\ &= \pi(x)q(y | x)\min\left(1, \frac{\pi(y)q(x | y)}{\pi(x)q(y | x)}\right) \\ &= \min(\pi(x)q(y | x), \pi(y)q(x | y)) \\ &= \pi(y)q(x | y)\min\left(\frac{\pi(x)q(y | x)}{\pi(y)q(x | y)}, 1\right) \\ &= \pi(y)K(y, x).\end{aligned}$$

If  $x = y$ , then obviously  $\pi(x)K(x, y) = \pi(y)K(y, x)$ .

- Consider the target distribution

$$\pi(x) = \left( \mathcal{U}_{[0,1]}(x) + \mathcal{U}_{[2,3]}(x) \right) / 2$$

and the proposal distribution

$$q(x^* | x) = \mathcal{U}_{(x-\delta, x+\delta)}(x^*).$$

- The MH chain is reducible if  $\delta \leq 1$ : the chain stays either in  $[0, 1]$  or  $[2, 3]$ .
- Note that the MH chain is aperiodic if it always has a non-zero chance of staying where it is.

- The MH chain  $(X^{(t)})_{t \geq 1}$  is irreducible if  $q(x^*|x) > 0$  for any  $x, x^* \in \text{supp}(\pi)$ : every state can be reached in a single step (strongly irreducible). Less strict conditions in (Roberts & Rosenthal, 2004).
- The MH chain is Harris recurrent if it is irreducible (see Tierney, 1994).
- **Theorem.** If the Markov chain generated by the Metropolis–Hastings sampler is  $\pi$ –irreducible, then we have for any integrable function  $\varphi : \mathbb{X} \rightarrow \mathbb{R}$ :

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^t \varphi(X^{(i)}) = \int_{\mathbb{X}} \varphi(x) \pi(x) dx$$

for *every* starting value  $X^{(1)}$ .

# Random Walk Metropolis–Hastings

- In the Metropolis–Hastings, pick  $q(x^* | x) = g(x^* - x)$  with  $g$  being a *symmetric* distribution, thus

$$X^* = X + \varepsilon, \quad \varepsilon \sim g;$$

e.g.  $g$  is a zero-mean multivariate normal or t-student.

- Acceptance probability becomes

$$\alpha(x^* | x) = \min \left( 1, \frac{\pi(x^*)}{\pi(x)} \right).$$

- We accept...
  - a move to a more probable state with probability 1;
  - a move to a less probable state with probability

$$\pi(x^*)/\pi(x) < 1.$$

# Independent Metropolis–Hastings

- If the proposal distribution  $q(x^* | x)$  does not depend on  $x$ , we call it an independent proposal.
- Acceptance probability becomes

$$\alpha(x^* | x) = \min \left( 1, \frac{\pi(x^*)q(x)}{\pi(x)q(x^*)} \right).$$

- For instance, multivariate normal or t-student distribution.
- If  $\pi(x)/q(x) < M$  for all  $x$  and some  $M < \infty$ , then the chain is uniformly ergodic.
- It can be shown that the acceptance probability at stationarity is then at least  $1/M$  (Lemma 7.9 of Robert & Casella).
- On the other hand, if such an  $M$  does not exist, the chain is not even geometrically ergodic!

# Choosing a good proposal distribution

- Goal: to design a Markov chain with small correlation  $\rho\left(X^{(t-1)}, X^{(t)}\right)$  between subsequent values (why?).
- Two sources of correlation:
  - between the current state  $X^{(t-1)}$  and proposed value  $X \sim q\left(\cdot | X^{(t-1)}\right)$ ,
  - correlation induced if  $X^{(t)} = X^{(t-1)}$ , if proposal is rejected.
- Trade-off: there is a compromise between
  - proposing large moves,
  - obtaining a decent acceptance probability.
- For multivariate distributions: covariance of proposal should reflect the covariance structure of the target.

- Target distribution, we want to sample from

$$\pi(x) = \mathcal{N}\left(x; \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}\right).$$

- We use a random walk Metropolis—Hastings algorithm with

$$g(\varepsilon) = \mathcal{N}\left(\varepsilon; 0, \sigma^2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right).$$

- What is the optimal choice of  $\sigma^2$ ?
- We consider three choices:  $\sigma^2 = 0.1^2, 1, 10^2$ .



# Metropolis–Hastings algorithm

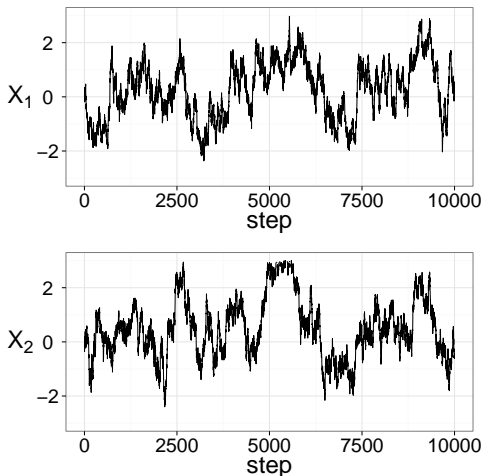


Figure: Metropolis–Hastings on a bivariate Gaussian target. With  $\sigma^2 = 0.1^2$ , the acceptance rate is  $\approx 94\%$ .

# Metropolis–Hastings algorithm

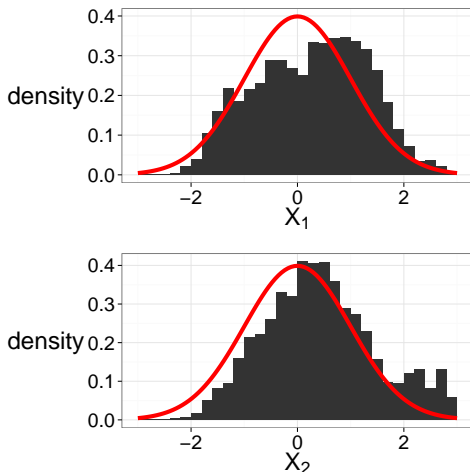


Figure: Metropolis–Hastings on a bivariate Gaussian target. With  $\sigma^2 = 0.1^2$ , the acceptance rate is  $\approx 94\%$ .

# Metropolis–Hastings algorithm

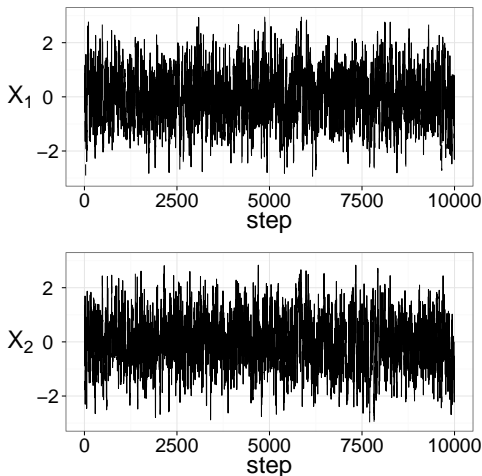


Figure: Metropolis–Hastings on a bivariate Gaussian target. With  $\sigma^2 = 1$ , the acceptance rate is  $\approx 52\%$ .

# Metropolis–Hastings algorithm

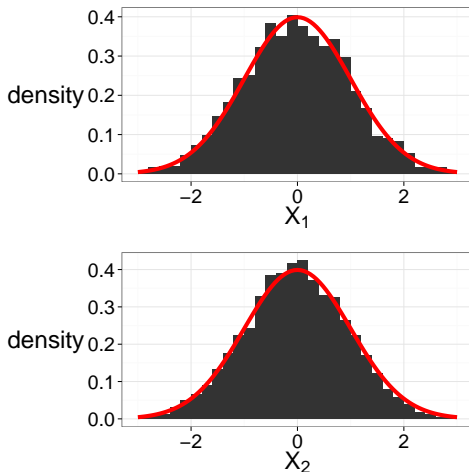


Figure: Metropolis–Hastings on a bivariate Gaussian target. With  $\sigma^2 = 1$ , the acceptance rate is  $\approx 52\%$ .

# Metropolis–Hastings algorithm

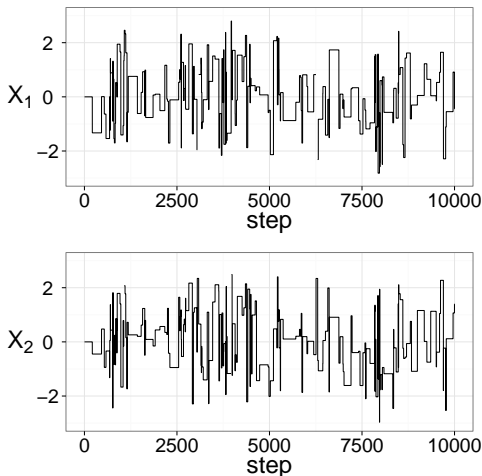


Figure: Metropolis–Hastings on a bivariate Gaussian target. With  $\sigma^2 = 10$ , the acceptance rate is  $\approx 1.5\%$ .

# Metropolis–Hastings algorithm

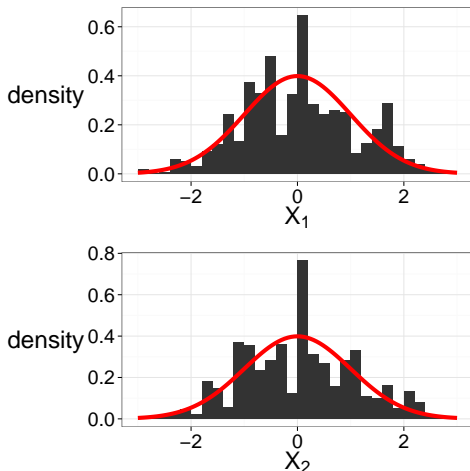


Figure: Metropolis–Hastings on a bivariate Gaussian target. With  $\sigma^2 = 10$ , the acceptance rate is  $\approx 1.5\%$ .

# Choice of proposal

- Aim at some intermediate acceptance ratio: 20%? 40%?  
Some hints come from the literature on “optimal scaling”.
- Maximize the expected square jumping distance:

$$\mathbb{E} \left[ \|X_{t+1} - X_t\|^2 \right]$$

- In multivariate cases, try to mimick the covariance structure of the target distribution.

Cooking recipe: run the algorithm for  $T$  iterations, check some criterion, tune the proposal distribution accordingly, run the algorithm for  $T$  iterations again . . .

“Constructing a chain that mixes well is somewhat of an art.”  
*All of Statistics*, L. Wasserman.

# The adaptive MCMC approach

- One can make the transition kernel  $K$  adaptive, i.e. use  $K_t$  at iteration  $t$  and choose  $K_t$  using the past sample  $(X_1, \dots, X_{t-1})$ .
- The Markov chain is not homogeneous anymore: the mathematical study of the algorithm is much more complicated.
- Adaptation can be counterproductive in some cases (see Atchadé & Rosenthal, 2005)!
- Adaptive Gibbs samplers also exist.