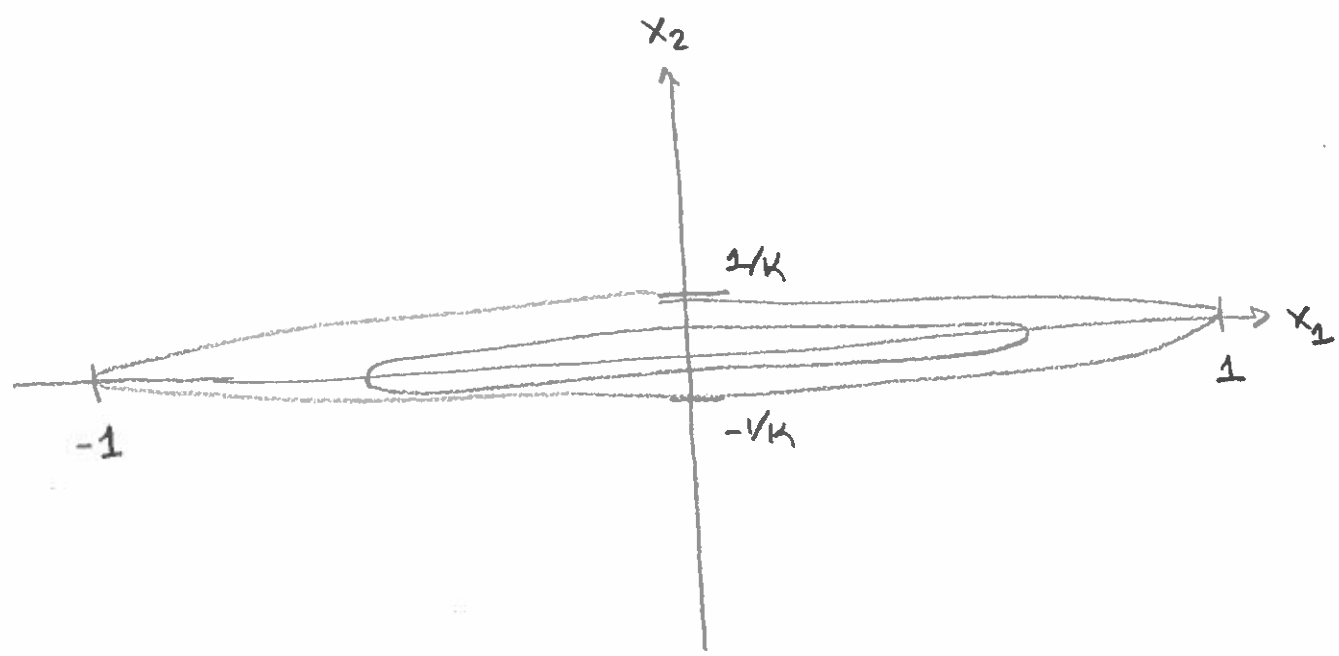


Lecture 8: All killer, no filler

- Review of MALA + HMC
- Goodman-Weare sampler
- Parallel tempering

MALA + HMC Review

- Imagine we are sampling a Gaussian  $\eta(0, (\frac{1}{k}))$  where  $k \gg 1$  is the "condition number".



- Target distribution is wide in  $x_1$  direction, narrow in  $x_2$  direction.

MALA uses the update

$$x \leftarrow x + \delta \nabla \log \pi(x) + \sqrt{2\delta} v,$$

where  $v$  is a random Gaussian velocity

$$v \sim \eta(0, I)$$

Q What is  $\nabla \log \pi(x)$  ?

Hint  $\pi(x) = \frac{\sqrt{k}}{2\pi} \exp\left(-\frac{x_1^2}{2} - \frac{k}{2}x_2^2\right)$

A  $\nabla \log \pi(x) = \begin{pmatrix} -x_1 \\ -kx_2 \end{pmatrix}$

- Update formula is

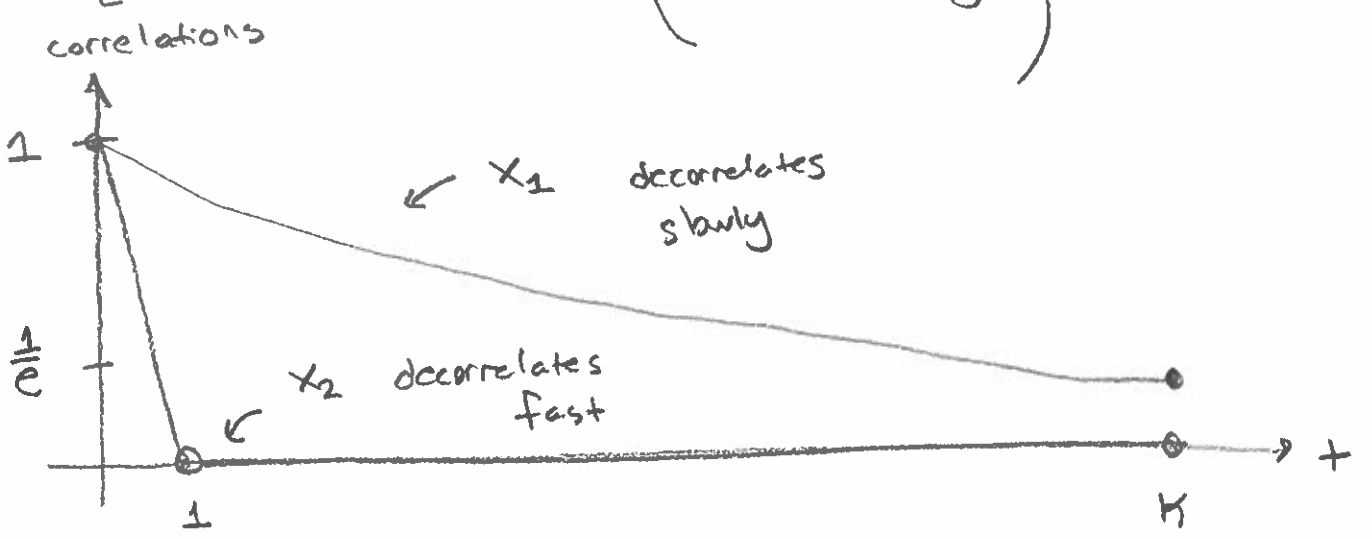
$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \leftarrow \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \delta \begin{pmatrix} -x_1 \\ -kx_2 \end{pmatrix} + \sqrt{2\delta} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

$$\Rightarrow \begin{cases} x_1 \leftarrow (1-\delta)x_1 + \sqrt{2\delta}v_1 \\ x_2 \leftarrow (1-k\delta)x_2 + \sqrt{2\delta}v_2 \end{cases}$$

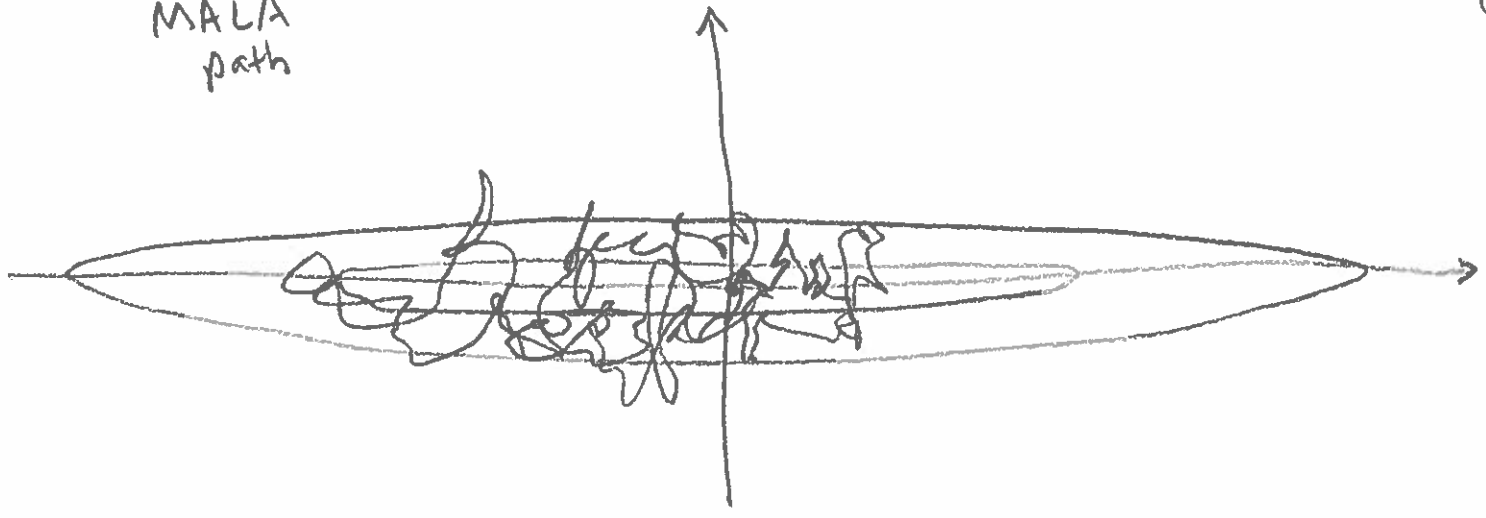
- For stability, we need  $\delta \in (0, \frac{2}{k})$ .

- If we set  $\delta = \frac{1}{k}$ , correlations decay as

$$\mathbb{E} \left[ x^{(t)} \mid x^{(0)} = x \right] = \begin{pmatrix} 1 - \frac{1}{k} \\ 0 \end{pmatrix}^t x$$



MALA path



⇒ It takes  $t = \mathcal{O}(K)$  timesteps for the MALA sampler to explore the distribution.

⇒  $t \approx 10^6$  if  $K \approx 10^6$  ;)

HMC uses the update

$$v \sim \mathcal{N}(0, I)$$

For  $t = 1, \dots, L$ :

$$v \leftarrow v + \frac{\delta}{2} \nabla \log \pi(x), \quad x \leftarrow x + \delta v, \quad v \leftarrow v + \frac{\delta}{2} \nabla \log \pi(x).$$

Q: What is the update formula for sampling a Gaussian  $\mathcal{N}(0, \begin{pmatrix} 1 & \\ & 1/K \end{pmatrix})$ ?

A:  $v_1 \leftarrow v_1 - \frac{\delta}{2} x_1, \quad x_1 \leftarrow x_1 + \delta v_1, \quad v_1 \leftarrow v_1 - \frac{\delta}{2} x_1$

$v_2 \leftarrow v_2 - \frac{\delta K}{2} x_2, \quad x_2 \leftarrow x_2 + \delta v_2, \quad v_2 \leftarrow v_2 - \frac{\delta K}{2} x_2$

- For stability, we need to take

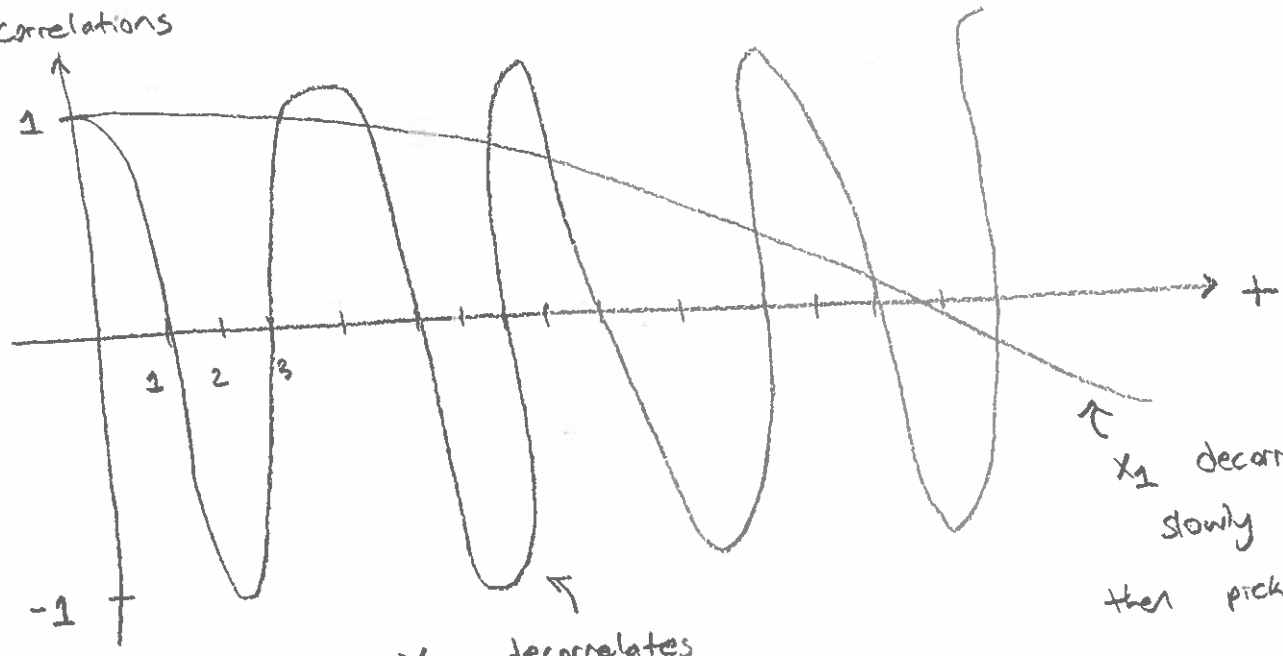
$\delta \in (0, \frac{2}{\sqrt{K}})$   $\Rightarrow$  Much larger steps!

- If we set  $\delta = \sqrt{\frac{2}{K}}$ , correlations decay like

$$\mathbb{E} [ X^{(t+1)} | X^{(t)} = x ] = \begin{pmatrix} \cos(\Theta + t) \\ \cos(\frac{\pi}{2} + t) \end{pmatrix} x$$

where  $\Theta = \arccos(1 - \frac{1}{K}) \approx \sqrt{\frac{2}{K}}$

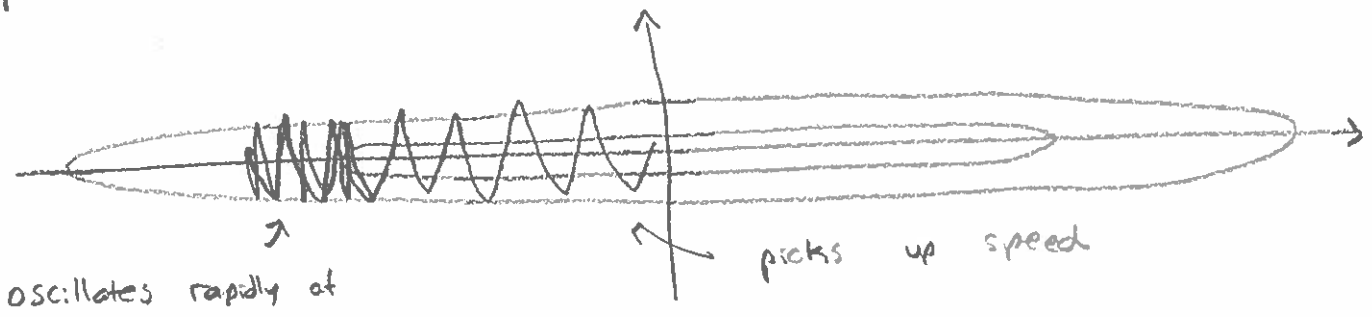
correlations



$x_1$  decorrelates slowly at first, then picks up speed

$x_2$  decorrelates completely at odd times  $t=1, 3, 5, \dots$

HMC path



oscillates rapidly at

picks up speed

⇒ It takes  $t = O(\sqrt{K})$  timesteps for the HMC sampler to explore the distribution

⇒  $\boxed{+ \approx 10^3 \text{ if } K \approx 10^6 \text{ ;}}$

⇒ "Heavy ball rolling downhill"

HMC algorithm

Input : target density  $\pi(x)$ , initial condition  $x_0$ , stepsize  $\delta$ , number of leapfrog steps  $L$ , terminal time  $T$

Output : approximate samples  $x_0, x_1, \dots, x_T$  from  $\pi$

For  $t = 1, 2, \dots, T$ :

- Generate  $v_t \sim \mathcal{N}(0, I)$

- Set  $x'_t = x_{t-1}$

- For  $i = 1, 2, \dots, L$ :

•  $v_t = v_t + \frac{\delta}{2} \nabla \log \pi(x'_t)$

•  $x'_t = x'_t + \delta v_t$

•  $v_t = v_t + \frac{\delta}{2} \nabla \log \pi(x'_t)$

- Generate  $u_t \sim \text{Unif}(0, 1)$

- Set  $x_t = \begin{cases} x'_t, & u_t \leq \pi(x'_t) / \pi(x_t) \\ x_{t-1}, & \text{otherwise.} \end{cases}$

Q: What can we do to remove the dependence on  $k$  entirely?

Goodman-Wear Sampler

Input: target density  $\pi(x)$ , number of ensemble members  $N$ , initial conditions  $X_0^{(1)}, X_0^{(2)}, \dots, X_0^{(N)} \in \mathbb{R}^d$ , terminal time  $T$

Output: approximate samples  $(X_t^{(1)}, \dots, X_t^{(N)})_{0 \leq t \leq T}$ , all from  $\pi$ .

For  $t = 1, 2, \dots, T$ :

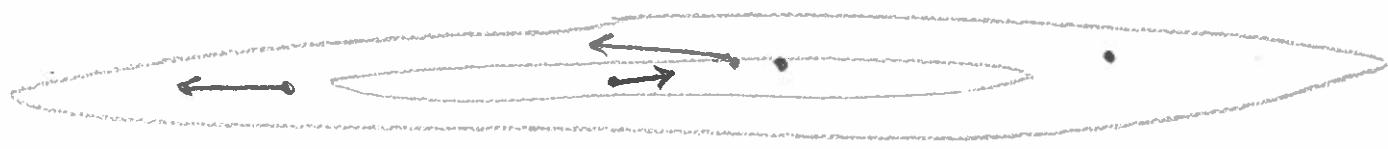
- Set  $(X_t^{(1)}, \dots, X_t^{(N)}) = (X_{t-1}^{(1)}, \dots, X_{t-1}^{(N)})$ .

- For  $i = 1, 2, \dots, N$ :

Stretch moves

- Select  $j \in \{1, \dots, N\} \setminus \{i\}$  uniformly at random.
- Sample  $U \sim \text{Unif}(2^{-3/2}, 2^{3/2})$
- Set  $Y_t^{(i)} = X_t^{(i)} + (1 - U^{2/3})(X_t^{(j)} - X_t^{(i)})$
- Sample  $V \sim \text{Unif}(0, 1)$
- Set  $X_t^{(i)} = Y_t^{(i)}$  if  $V \leq U^{\frac{2}{3}(d-1)} \frac{\pi(Y_t^{(i)})}{\pi(X_t^{(i)})}$ .

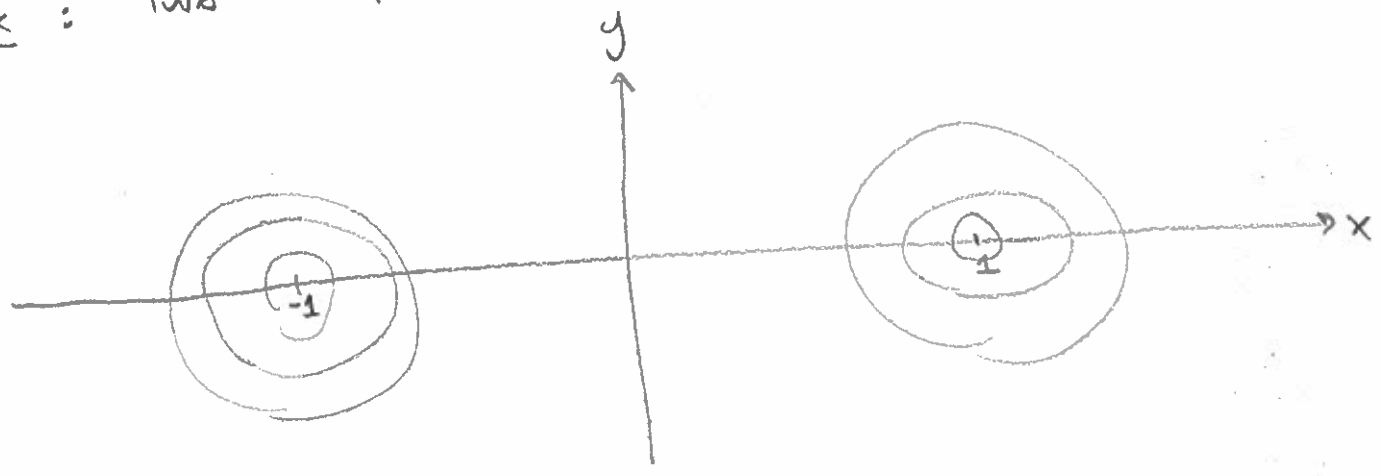
Intuition Randomly push one ensemble member toward or away from another ensemble member



- Goodman-Krichevsky is "affine-invariant"  $\Rightarrow$  correlations do not depend on  $K$  at all!
- Works well for  $d \leq 10-20$ .
- At each step,  $X_+^{(i)}$  either halves its distance from  $X_+^{(i)}$ , doubles its distance from  $X_+^{(i)}$ , or goes somewhere in between.

Q: What can we do to address multi-modality?

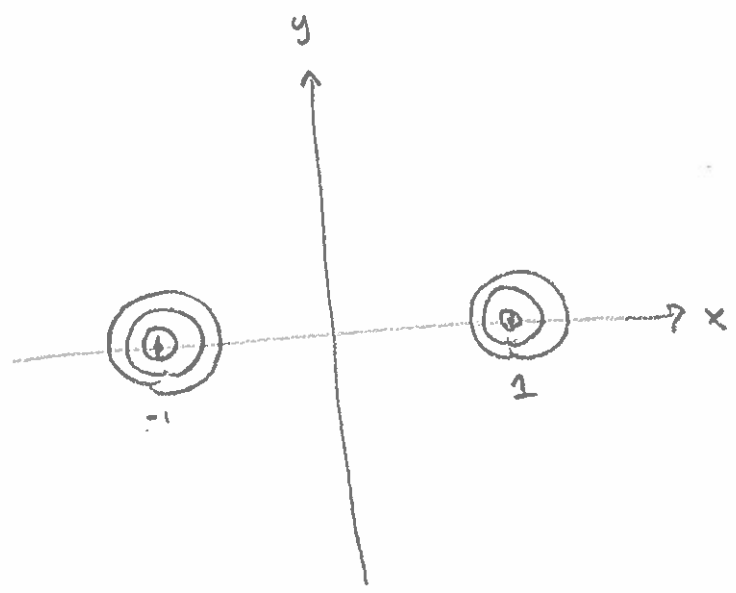
Ex: Two separated Gaussians



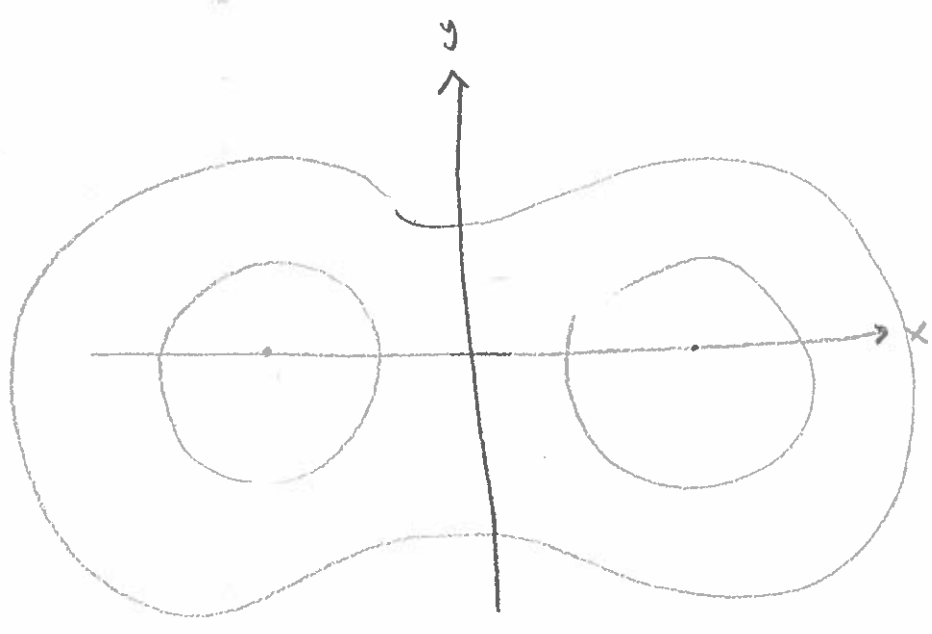
Idea: Why don't we increase the "temperature"?

$\pi(x) \propto \exp\left(-\frac{H(x)}{T}\right)$ ,  $T$  is the temperature.

Cold temperature  $T \ll 1$ :



Hot temperature  $T \gg 1$ :



Let's sample a range of temperatures

$$T_i = \frac{N}{i} \cdot T, \quad i = 1, 2, \dots, N$$

$$\Rightarrow \pi_i(x) \propto \exp\left(-\frac{H(x)}{T_i}\right) \propto \pi(x)^{i/N}, \quad i = 1, 2, \dots, N$$



# Parallel tempering

Input : target density  $\pi(x)$ , number of ensemble members  $N$ , initial conditions  $X_0^{(1)}, X_0^{(2)}, \dots, X_0^{(N)}$ , terminal time  $T$

Output : approximate samples  $X_0^{(n)}, \dots, X_T^{(n)}$  from  $\pi$ .

For  $t = 1, \dots, T$ :

Parallel moves { - For  $i = 1, \dots, N$  :  
 • Update  $X_{t-1}^{(i)}$  to  $X_t^{(i)}$  using an MCMC method that targets  $\pi_i(x) \propto \pi(x)^{i/N}$

Swap moves { - For  $i = 1, \dots, N-1$  :  
 • Sample  $U \sim \text{Unif}(0, 1)$   
 • Set  $(X_{t+}^{(i)}, X_{t+}^{(i+1)}) = (X_{t+}^{(i+1)}, X_{t+}^{(i)})$  if  $U \leq \left( \frac{\pi(X_{t+}^{(i)})}{\pi(X_{t+}^{(i+1)})} \right)^{1/N}$

Intuition We can jump between modes at high temperatures and information can filter down to low temperatures via a sequence of swaps.