

# LECTURE 5

---

## MULTI-MODAL DISTRIBUTIONS

**Saifuddin Syed**

# TUNING MCMC

# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$

# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$ 
  - ▶ e.g. step size, mass/covariance matrix, integration time, etc

# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$ 
  - ▶ e.g. step size, mass/covariance matrix, integration time, etc
- ▶ This defines a family of  $\pi$ -invariant kernels  $K_\theta$  such that give a valid MCMC algorithm

# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$ 
  - ▶ e.g. step size, mass/covariance matrix, integration time, etc
- ▶ This defines a family of  $\pi$ -invariant kernels  $K_\theta$  such that give a valid MCMC algorithm
  - ▶ For each  $\theta$  we can generate a Markov chain  $X_t$  that satisfies the ergodic theorem

$$X_t \sim K_\theta(X_{t-1}, dx) \quad \frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi[f]$$

# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$ 
  - ▶ e.g. step size, mass/covariance matrix, integration time, etc
- ▶ This defines a family of  $\pi$ -invariant kernels  $K_\theta$  such that give a valid MCMC algorithm
  - ▶ For each  $\theta$  we can generate a Markov chain  $X_t$  that satisfies the ergodic theorem

$$X_t \sim K_\theta(X_{t-1}, dx) \quad \frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi[f]$$

- ▶ The efficiency and rates of convergence depend on the parameters and must be tuned

# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$ 
  - ▶ e.g. step size, mass/covariance matrix, integration time, etc
- ▶ This defines a family of  $\pi$ -invariant kernels  $K_\theta$  such that give a valid MCMC algorithm
  - ▶ For each  $\theta$  we can generate a Markov chain  $X_t$  that satisfies the ergodic theorem

$$X_t \sim K_\theta(X_{t-1}, dx) \quad \frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi[f]$$

- ▶ The efficiency and rates of convergence depend on the parameters and must be tuned
- ▶ In practice, we do this adaptively and modify the parameters as we simulate the chain



# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$ 
  - ▶ e.g. step size, mass/covariance matrix, integration time, etc
- ▶ This defines a family of  $\pi$ -invariant kernels  $K_\theta$  such that give a valid MCMC algorithm
  - ▶ For each  $\theta$  we can generate a Markov chain  $X_t$  that satisfies the ergodic theorem

$$X_t \sim K_\theta(X_{t-1}, dx) \quad \frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi[f]$$

- ▶ The efficiency and rates of convergence depend on the parameters and must be tuned
- ▶ In practice, we do this adaptively and modify the parameters as we simulate the chain
  - ▶ E.g. modify the parameters until a desired step size is reached

# TUNING MCMC

- ▶ Given an MCMC algorithm there are often hyperparameters  $\theta \in \Theta$ 
  - ▶ e.g. step size, mass/covariance matrix, integration time, etc
- ▶ This defines a family of  $\pi$ -invariant kernels  $K_\theta$  such that give a valid MCMC algorithm
  - ▶ For each  $\theta$  we can generate a Markov chain  $X_t$  that satisfies the ergodic theorem

$$X_t \sim K_\theta(X_{t-1}, dx) \quad \frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi[f]$$

- ▶ The efficiency and rates of convergence depend on the parameters and must be tuned
- ▶ In practice, we do this adaptively and modify the parameters as we simulate the chain
  - ▶ E.g. modify the parameters until a desired step size is reached
- ▶ **Question:** Is this theoretically justified?

# PITFALLS OF ADAPTIVE TUNING

# PITFALLS OF ADAPTIVE TUNING

- ▶ Formally, we develop a  $\theta_t$  that depends on the history of  $X_1, \dots, X_t$

$$X_t \sim K_{\theta_t}(X_{t-1}, dx)$$

# PITFALLS OF ADAPTIVE TUNING

- ▶ Formally, we develop a  $\theta_t$  that depends on the history of  $X_1, \dots, X_t$

$$X_t \sim K_{\theta_t}(X_{t-1}, dx)$$

- ▶  $X_t$  is in general no longer a Markov chain,  $\pi$ -invariant or satisfies the ergodic theorem

# PITFALLS OF ADAPTIVE TUNING

- ▶ Formally, we develop a  $\theta_t$  that depends on the history of  $X_1, \dots, X_t$

$$X_t \sim K_{\theta_t}(X_{t-1}, dx)$$

- ▶  $X_t$  is in general no longer a Markov chain,  $\pi$ -invariant or satisfies the ergodic theorem
  - ▶ Intuitively, we can think of  $\theta_t(x_{0:t})$  as a control function governing a dynamical system  $X_t$



# PITFALLS OF ADAPTIVE TUNING

- ▶ Formally, we develop a  $\theta_t$  that depends on the history of  $X_1, \dots, X_t$

$$X_t \sim K_{\theta_t}(X_{t-1}, dx)$$

- ▶  $X_t$  is in general no longer a Markov chain,  $\pi$ -invariant or satisfies the ergodic theorem
  - ▶ Intuitively, we can think of  $\theta_t(x_{0:t})$  as a control function governing a dynamical system  $X_t$
- ▶ If not careful, we can nudge the law of the chain towards non-equilibrium states



# EXAMPLE: ADAPTIVE MARKOV CHAIN



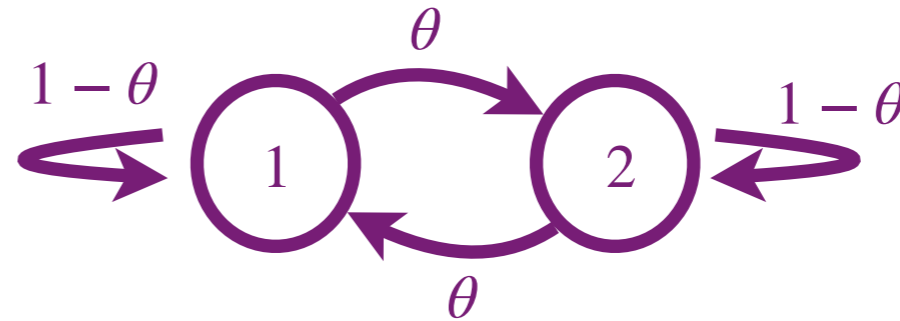
# EXAMPLE: ADAPTIVE MARKOV CHAIN

- ▶ Suppose  $\mathbb{X} = \{1,2\}$  and  $\Theta = (0,1)$

# EXAMPLE: ADAPTIVE MARKOV CHAIN

- ▶ Suppose  $\mathbb{X} = \{1,2\}$  and  $\Theta = (0,1)$
- ▶ Consider the family of Markov kernels  $K_\theta$  parametrised by  $\theta \in \Theta$

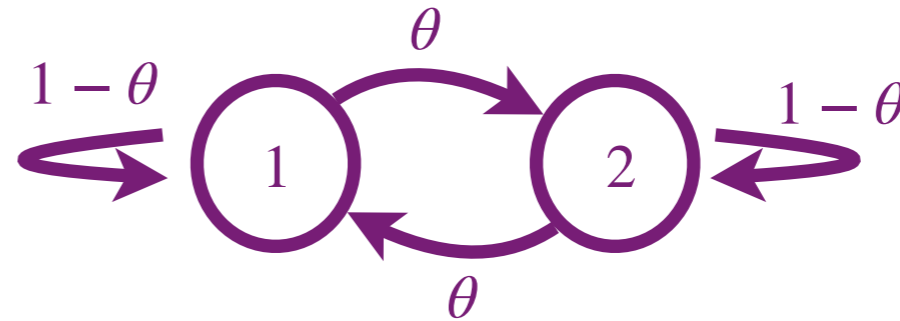
$$K_\theta = \begin{pmatrix} \theta & 1-\theta \\ 1-\theta & \theta \end{pmatrix}$$



# EXAMPLE: ADAPTIVE MARKOV CHAIN

- ▶ Suppose  $\mathbb{X} = \{1,2\}$  and  $\Theta = (0,1)$
- ▶ Consider the family of Markov kernels  $K_\theta$  parametrised by  $\theta \in \Theta$

$$K_\theta = \begin{pmatrix} \theta & 1-\theta \\ 1-\theta & \theta \end{pmatrix}$$

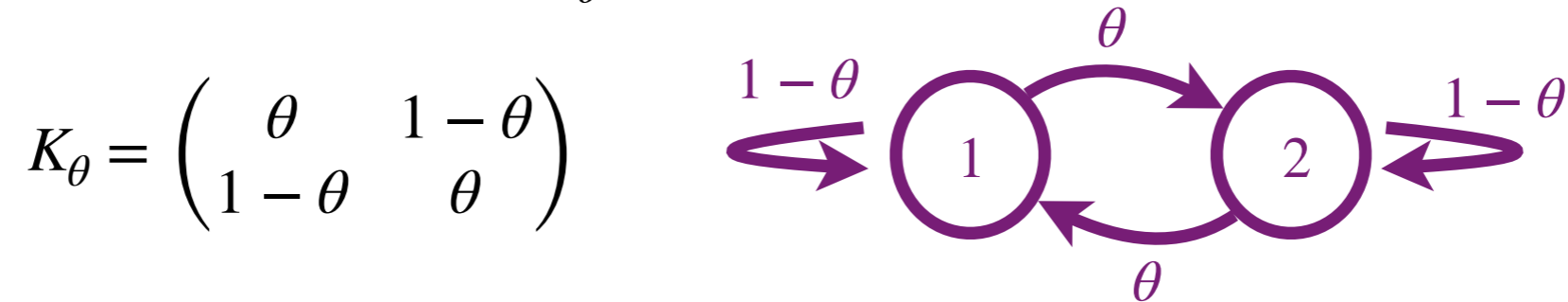


- ▶ For each  $\theta$ , the chain  $X_t$  generated by  $K_\theta$  is invariant to  $\pi(1) = \pi(2) = 0.5$

$$\pi K_\theta = \pi \quad X_t \sim K_\theta(X_{t-1}, dx)$$

# EXAMPLE: ADAPTIVE MARKOV CHAIN

- ▶ Suppose  $\mathbb{X} = \{1,2\}$  and  $\Theta = (0,1)$
- ▶ Consider the family of Markov kernels  $K_\theta$  parametrised by  $\theta \in \Theta$



- ▶ For each  $\theta$ , the chain  $X_t$  generated by  $K_\theta$  is invariant to  $\pi(1) = \pi(2) = 0.5$

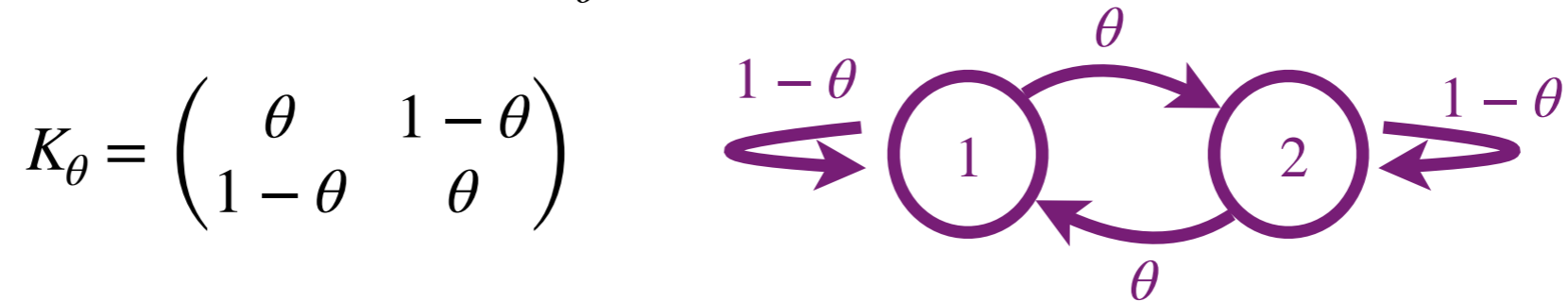
$$\pi K_\theta = \pi \quad X_t \sim K_\theta(X_{t-1}, dx)$$

- ▶ For  $\theta : \mathbb{X} \rightarrow [0,1]$  suppose we generate  $\hat{X}_t$  adaptively using the kernel at  $\theta(\hat{X}_t)$

$$\hat{X}_t \sim K_{\theta(\hat{X}_{t-1})}(\hat{X}_{t-1}, dx) = \hat{K}(\hat{X}_{t-1}, dx)$$

# EXAMPLE: ADAPTIVE MARKOV CHAIN

- ▶ Suppose  $\mathbb{X} = \{1,2\}$  and  $\Theta = (0,1)$
- ▶ Consider the family of Markov kernels  $K_\theta$  parametrised by  $\theta \in \Theta$



- ▶ For each  $\theta$ , the chain  $X_t$  generated by  $K_\theta$  is invariant to  $\pi(1) = \pi(2) = 0.5$

$$\pi K_\theta = \pi \quad X_t \sim K_\theta(X_{t-1}, dx)$$

- ▶ For  $\theta : \mathbb{X} \rightarrow [0,1]$  suppose we generate  $\hat{X}_t$  adaptively using the kernel at  $\theta(\hat{X}_t)$

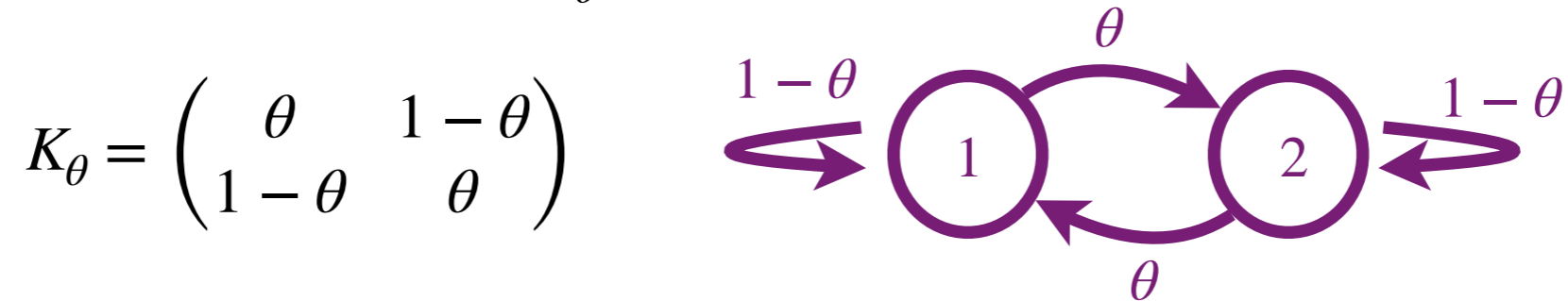
$$\hat{X}_t \sim K_{\theta(\hat{X}_{t-1})}(\hat{X}_{t-1}, dx) = \hat{K}(\hat{X}_{t-1}, dx)$$

- ▶ Defines a new Markov chain:



# EXAMPLE: ADAPTIVE MARKOV CHAIN

- ▶ Suppose  $\mathbb{X} = \{1,2\}$  and  $\Theta = (0,1)$
- ▶ Consider the family of Markov kernels  $K_\theta$  parametrised by  $\theta \in \Theta$



- ▶ For each  $\theta$ , the chain  $X_t$  generated by  $K_\theta$  is invariant to  $\pi(1) = \pi(2) = 0.5$

$$\pi K_\theta = \pi \quad X_t \sim K_\theta(X_{t-1}, dx)$$

- ▶ For  $\theta : \mathbb{X} \rightarrow [0,1]$  suppose we generate  $\hat{X}_t$  adaptively using the kernel at  $\theta(\hat{X}_t)$

$$\hat{X}_t \sim K_{\theta(\hat{X}_{t-1})}(\hat{X}_{t-1}, dx) = \hat{K}(\hat{X}_{t-1}, dx)$$

- ▶ Defines a new Markov chain:



- ▶ The adaptive chain  $\hat{X}_t$  is **not**  $\pi$ -invariant, i.e.  $\pi \hat{K} \neq \pi$

# PITFALLS OF ADAPTIVE TUNING

# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target



# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target
  - ▶ Doing this correctly is delicate to ensure equilibrium is reached

# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target
  - ▶ Doing this correctly is delicate to ensure equilibrium is reached
- ▶ Generally want to small or infrequent modification to  $\theta$

# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target
  - ▶ Doing this correctly is delicate to ensure equilibrium is reached
- ▶ Generally want to small or infrequent modification to  $\theta$ 
  - ▶ Can keep consistency by making **small** modification at each time  $\theta_t \approx \theta_{t-1}$

# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target
  - ▶ Doing this correctly is delicate to ensure equilibrium is reached
- ▶ Generally want to small or infrequent modification to  $\theta$ 
  - ▶ Can keep consistency by making **small** modification at each time  $\theta_t \approx \theta_{t-1}$
  - ▶ Or by making infrequent modifications, for example exponential time between updates

# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target
  - ▶ Doing this correctly is delicate to ensure equilibrium is reached
- ▶ Generally want to small or infrequent modification to  $\theta$ 
  - ▶ Can keep consistency by making **small** modification at each time  $\theta_t \approx \theta_{t-1}$
  - ▶ Or by making infrequent modifications, for example exponential time between updates
- ▶ In practice, we devote a portion of the compute budget for tuning:

# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target
  - ▶ Doing this correctly is delicate to ensure equilibrium is reached
- ▶ Generally want to small or infrequent modification to  $\theta$ 
  - ▶ Can keep consistency by making **small** modification at each time  $\theta_t \approx \theta_{t-1}$
  - ▶ Or by making infrequent modifications, for example exponential time between updates
- ▶ In practice, we devote a portion of the compute budget for tuning:
  - ▶ The goal is to learn a learn an optimal hyper paramaters  $\theta^*$  for the sampler

# PITFALLS OF ADAPTIVE TUNING

- ▶ **Problem:** Frequent adaptation never allows enough time for the chain to converge the target
  - ▶ Doing this correctly is delicate to ensure equilibrium is reached
- ▶ Generally want to small or infrequent modification to  $\theta$ 
  - ▶ Can keep consistency by making **small** modification at each time  $\theta_t \approx \theta_{t-1}$
  - ▶ Or by making infrequent modifications, for example exponential time between updates
- ▶ In practice, we devote a portion of the compute budget for tuning:
  - ▶ The goal is to learn a learn an optimal hyper parameters  $\theta^*$  for the sampler
  - ▶ Throw away the samples at the end then use the remaining budget to sample with  $\theta^*$

# LOCAL SAMPLERS



# LOCAL SAMPLERS

- ▶ Most MCMC methods are traditionally designed to be **local**

# LOCAL SAMPLERS

- ▶ Most MCMC methods are traditionally designed to be **local**
- ▶ We often assume the target distribution is continuous in the topological features of  $\mathbb{X}$

$$y \approx x \quad \implies \quad \pi(y) \approx \pi(x)$$

# LOCAL SAMPLERS

- ▶ Most MCMC methods are traditionally designed to be **local**
- ▶ We often assume the target distribution is continuous in the topological features of  $\mathbb{X}$

$$y \approx x \quad \implies \quad \pi(y) \approx \pi(x)$$

- ▶ The proposal  $Q$  appeals to the topology of the underlying statepace

# LOCAL SAMPLERS

- ▶ Most MCMC methods are traditionally designed to be **local**
- ▶ We often assume the target distribution is continuous in the topological features of  $\mathbb{X}$

$$y \approx x \quad \implies \quad \pi(y) \approx \pi(x)$$

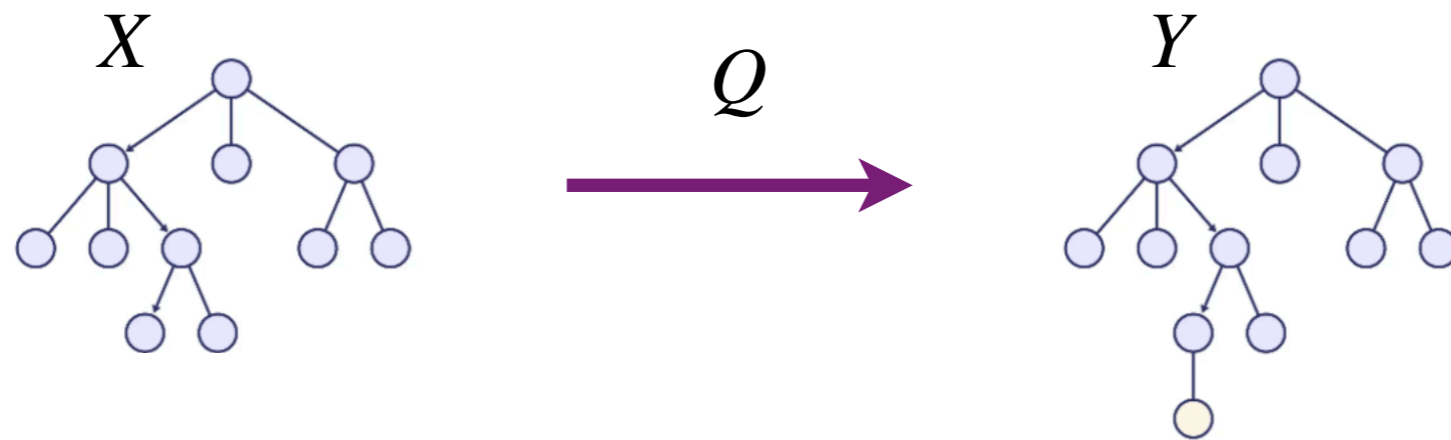
- ▶ The proposal  $Q$  appeals to the topology of the underlying statepace
  - ▶ Given  $X$  proposal samples  $Y \sim Q(X, \mathrm{d}y)$  within a neighbourhood of  $X$

# LOCAL SAMPLERS

- ▶ Most MCMC methods are traditionally designed to be **local**
- ▶ We often assume the target distribution is continuous in the topological features of  $\mathbb{X}$

$$y \approx x \implies \pi(y) \approx \pi(x)$$

- ▶ The proposal  $Q$  appeals to the topology of the underlying statepace
  - ▶ Given  $X$  proposal samples  $Y \sim Q(X, dy)$  within a neighbourhood of  $X$

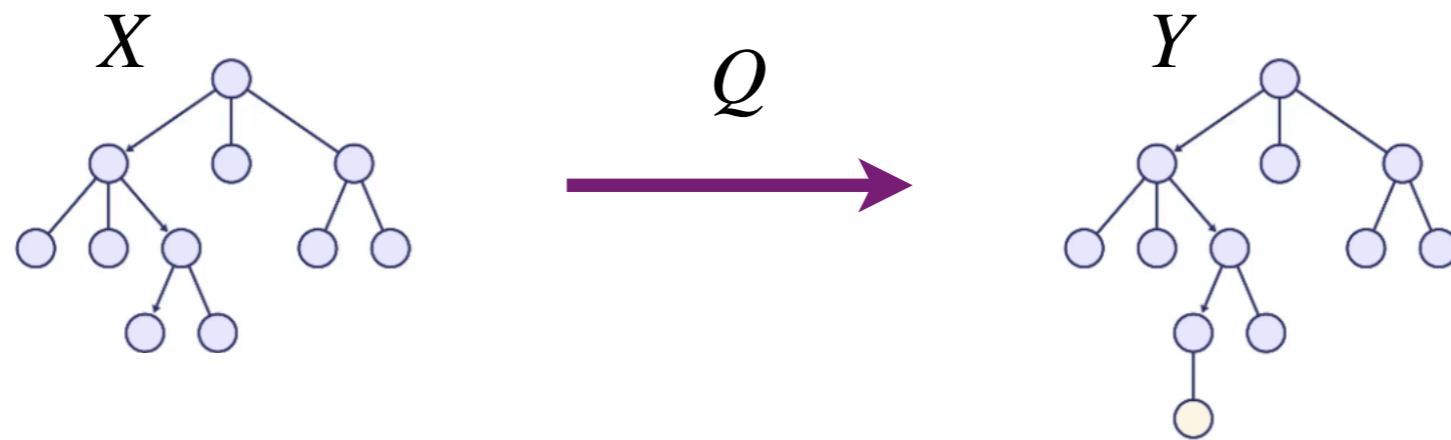


# LOCAL SAMPLERS

- ▶ Most MCMC methods are traditionally designed to be **local**
- ▶ We often assume the target distribution is continuous in the topological features of  $\mathbb{X}$

$$y \approx x \quad \implies \quad \pi(y) \approx \pi(x)$$

- ▶ The proposal  $Q$  appeals to the topology of the underlying statepace
  - ▶ Given  $X$  proposal samples  $Y \sim Q(X, dy)$  within a neighbourhood of  $X$

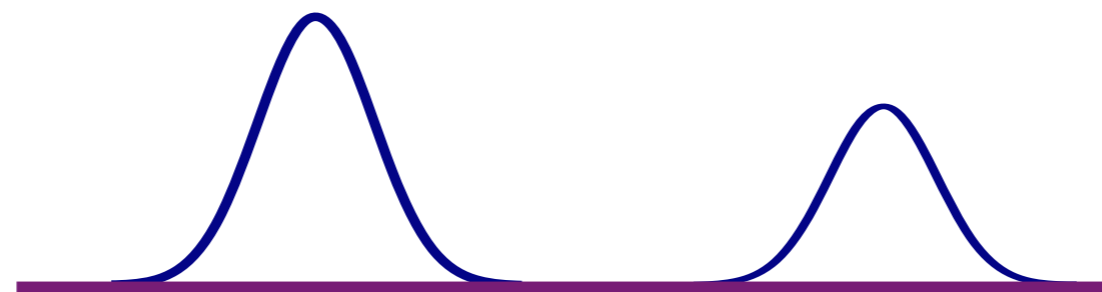


- ▶ Often appeals to the dynamics driven by a differential equation such as Langevin or Hamiltonian

# REDUCIBILITY OF LOCAL SAMPLERS

# REDUCIBILITY OF LOCAL SAMPLERS

- ▶ Local samplers prevent the chain from escaping local mode

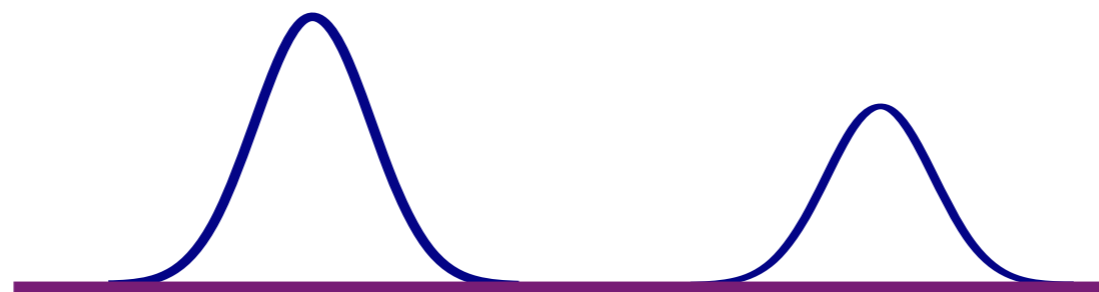




# REDUCIBILITY OF LOCAL SAMPLERS

- ▶ Local samplers prevent the chain from escaping local mode
- ▶ E.g. the continuous Langevin dynamics and Hamiltonian dynamics are not irreducible!

$$dY_\tau = -\nabla \log \pi(Y_\tau) d\tau + \sqrt{2} dW_\tau$$
$$dq_t = p_t dt \quad dp_t = -\log \nabla \pi(q_t) dt$$

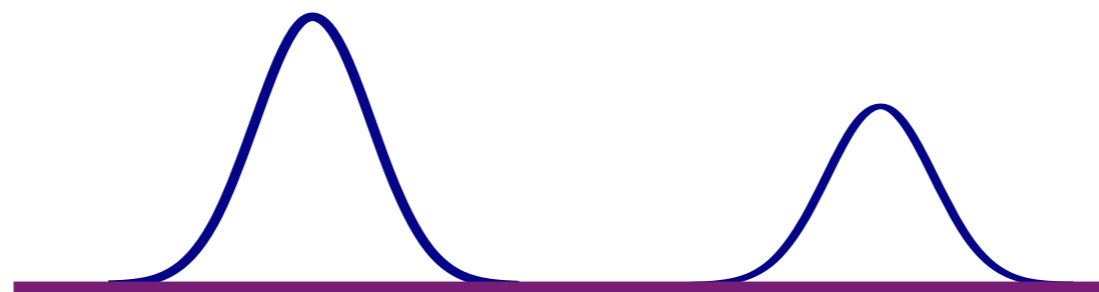


# REDUCIBILITY OF LOCAL SAMPLERS

- ▶ Local samplers prevent the chain from escaping local mode
- ▶ E.g. the continuous Langevin dynamics and Hamiltonian dynamics are not irreducible!

$$dY_\tau = -\nabla \log \pi(Y_\tau) d\tau + \sqrt{2} dW_\tau$$
$$dq_t = p_t dt \quad dp_t = -\log \nabla \pi(q_t) dt$$

- ▶ A local samplers are generally not irreducible if there exists there exists an open set  $U \subset \mathbb{X}$  such that  $\pi[\partial U] = 0$  and  $\pi[U] < 1$

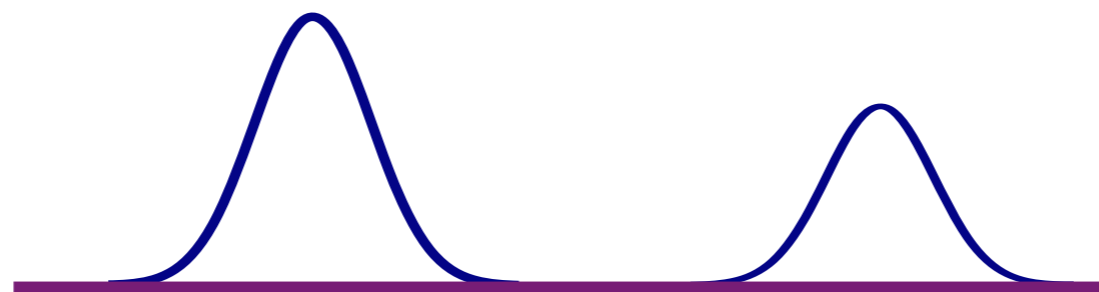


# REDUCIBILITY OF LOCAL SAMPLERS

- ▶ Local samplers prevent the chain from escaping local mode
- ▶ E.g. the continuous Langevin dynamics and Hamiltonian dynamics are not irreducible!

$$dY_\tau = -\nabla \log \pi(Y_\tau) d\tau + \sqrt{2} dW_\tau$$
$$dq_t = p_t dt \quad dp_t = -\log \nabla \pi(q_t) dt$$

- ▶ A local samplers are generally not irreducible if there exists there exists an open set  $U \subset \mathbb{X}$  such that  $\pi[\partial U] = 0$  and  $\pi[U] < 1$
- ▶ In theory, if  $\pi(x) > 0$ , but  $\pi(x) \approx 0$  we can technically have chains are uni-modal

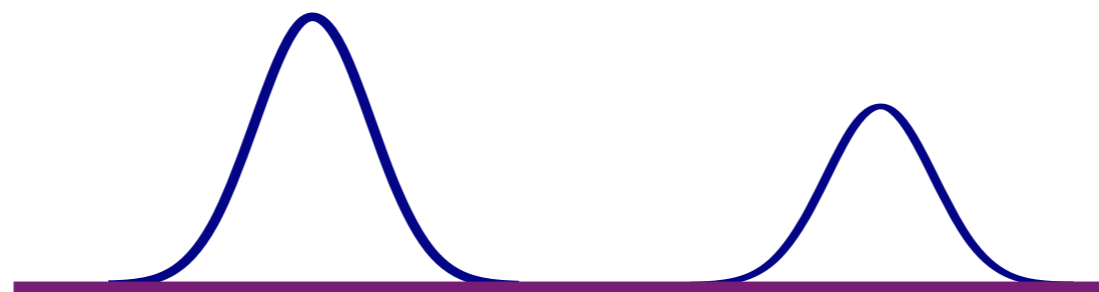


# REDUCIBILITY OF LOCAL SAMPLERS

- ▶ Local samplers prevent the chain from escaping local mode
- ▶ E.g. the continuous Langevin dynamics and Hamiltonian dynamics are not irreducible!

$$dY_\tau = -\nabla \log \pi(Y_\tau) d\tau + \sqrt{2} dW_\tau$$
$$dq_t = p_t dt \quad dp_t = -\log \nabla \pi(q_t) dt$$

- ▶ A local samplers are generally not irreducible if there exists there exists an open set  $U \subset \mathbb{X}$  such that  $\pi[\partial U] = 0$  and  $\pi[U] < 1$
- ▶ In theory, if  $\pi(x) > 0$ , but  $\pi(x) \approx 0$  we can technically have chains are uni-modal
- ▶ In practice, they are *effectively* irreducible



# WHAT IS A MODE

# WHAT IS A MODE

- ▶ What is a mode?

# WHAT IS A MODE

- ▶ What is a mode?
- ▶ Surprisingly ill-defined notion

# WHAT IS A MODE

- ▶ What is a mode?
- ▶ Surprisingly ill-defined notion
  - ▶ Intutively its a region of concentrated high probability



# WHAT IS A MODE

- ▶ What is a mode?
- ▶ Surprisingly ill-defined notion
  - ▶ Intuitively its a region of concentrated high probability
  - ▶ local optimum of the target







# WHAT IS A MODE

- ▶ What is a mode?
- ▶ Surprisingly ill-defined notion
  - ▶ Intuitively its a region of concentrated high probability
  - ▶ local optimum of the target

▶ **Question:** How many modes does this have:



- ▶ We want a definition that is philosophically aligned with intuitions
  - ▶ Should not overfit to not assuming extraneous structure of a specific problem
  - ▶ e.g. statespace, topology, geometry, smoothnes etc

# WHAT IS A MODE

- ▶ What is a mode?
- ▶ Surprisingly ill-defined notion
  - ▶ Intuitively its a region of concentrated high probability
  - ▶ local optimum of the target

▶ **Question:** How many modes does this have:



- ▶ We want a definition that is philosophically aligned with intuitions
  - ▶ Should not overfit to not assuming extraneous structure of a specific problem
  - ▶ e.g. statespace, topology, geometry, smoothnes etc
- ▶ Modes are a property of the sampler!

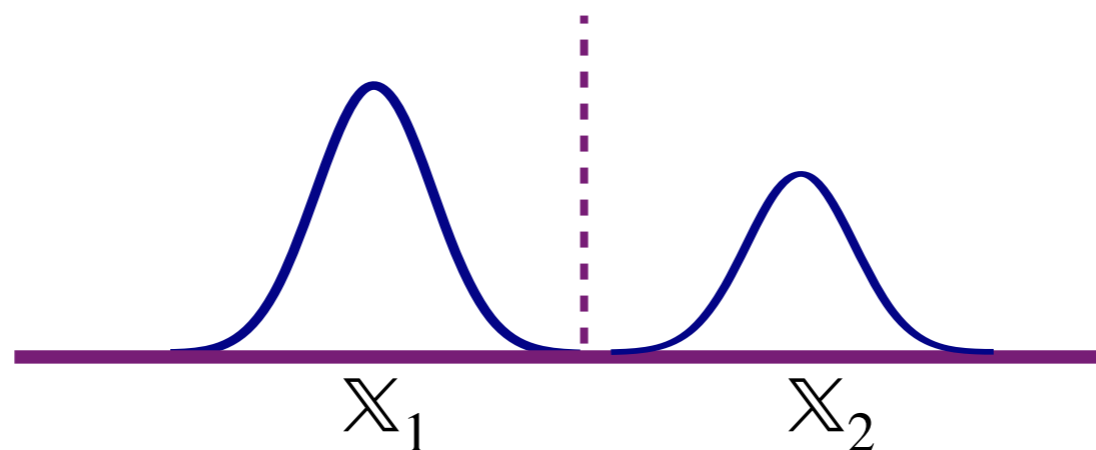






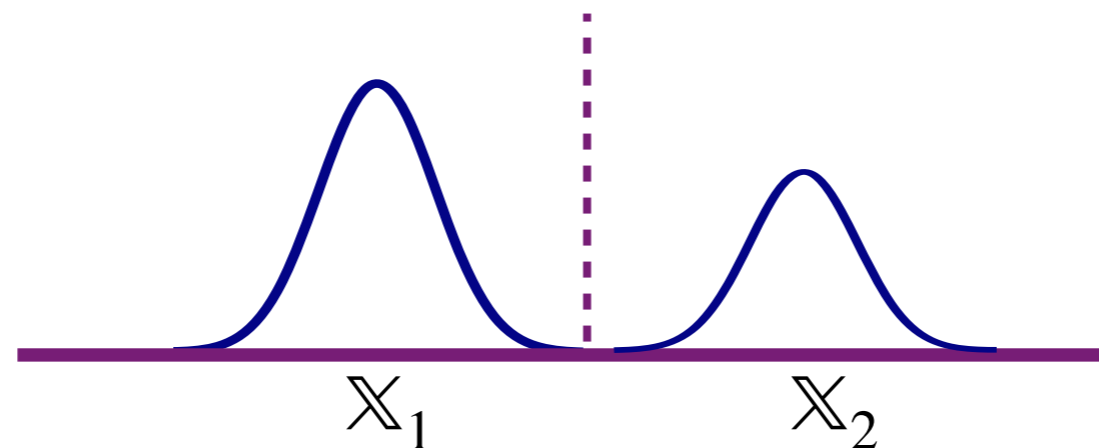
# MODES

- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if



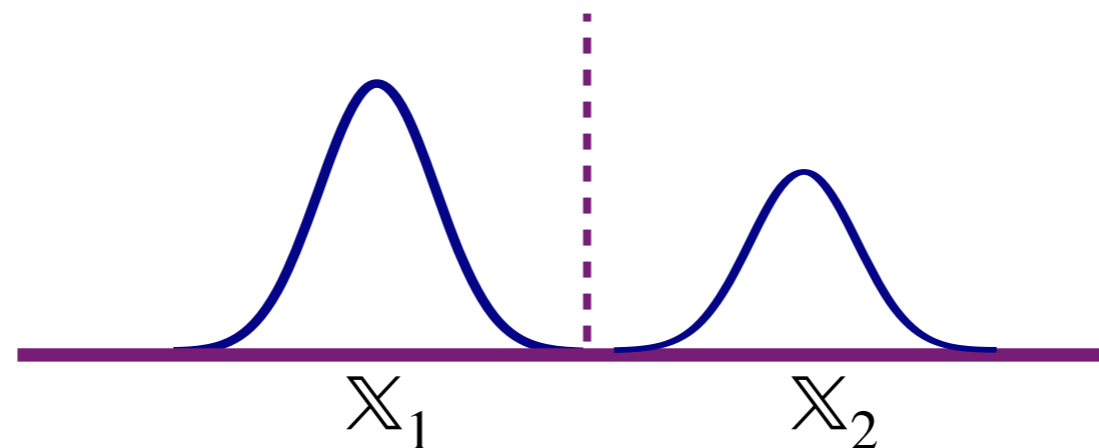
# MODES

- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if
  - ▶ It there is a possible probability to reach  $x'$  from  $x$  using  $K$  and vice-versa



# MODES

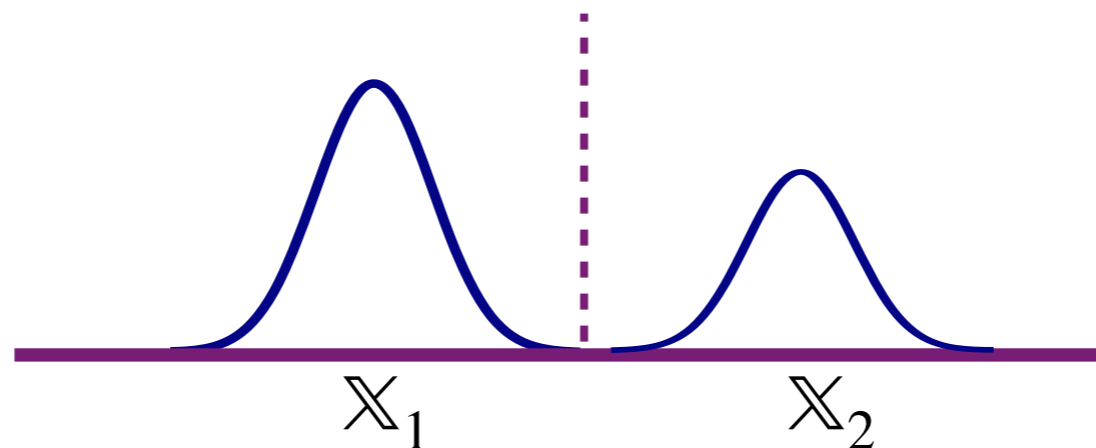
- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if
  - ▶ It there is a possible probability to reach  $x'$  from  $x$  using  $K$  and vice-versa
  - ▶ This defines an equivalent relation on the support of  $\pi$



# MODES

- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if
  - ▶ It there is a possible probability to reach  $x'$  from  $x$  using  $K$  and vice-versa
  - ▶ This defines an equivalent relation on the support of  $\pi$
- ▶ We will say the  $\mathbb{X}(x)$  is the mode of  $x$

$$\mathbb{X}(x) = \{x' \in \mathbb{X} : x \sim_K x'\}$$

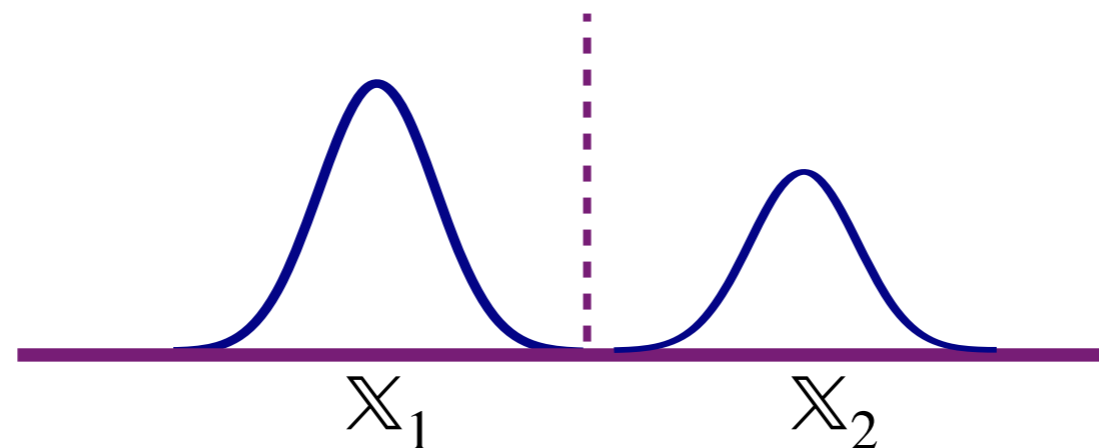


# MODES

- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if
  - ▶ It there is a possible probability to reach  $x'$  from  $x$  using  $K$  and vice-versa
  - ▶ This defines an equivalent relation on the support of  $\pi$
- ▶ We will say the  $\mathbb{X}(x)$  is the mode of  $x$

$$\mathbb{X}(x) = \{x' \in \mathbb{X} : x \sim_K x'\}$$

- ▶ I.e. the modes correspond to the irreducible components of the kernel

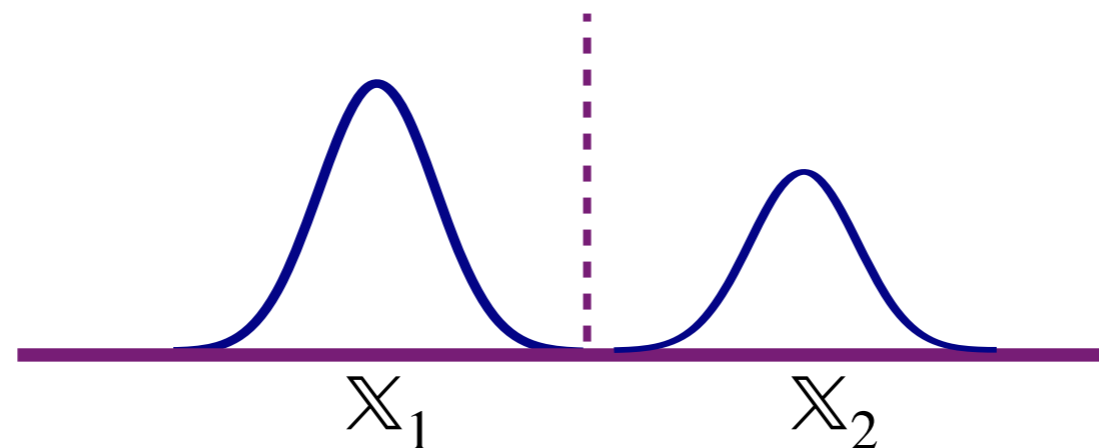


# MODES

- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if
  - ▶ It there is a possible probability to reach  $x'$  from  $x$  using  $K$  and vice-versa
  - ▶ This defines an equivalent relation on the support of  $\pi$
- ▶ We will say the  $\mathbb{X}(x)$  is the mode of  $x$

$$\mathbb{X}(x) = \{x' \in \mathbb{X} : x \sim_K x'\}$$

- ▶ I.e. the modes correspond to the irreducible components of the kernel
- ▶ The number of modes  $|\{\mathbb{X}(x) : x \in \mathbb{X}\}|$



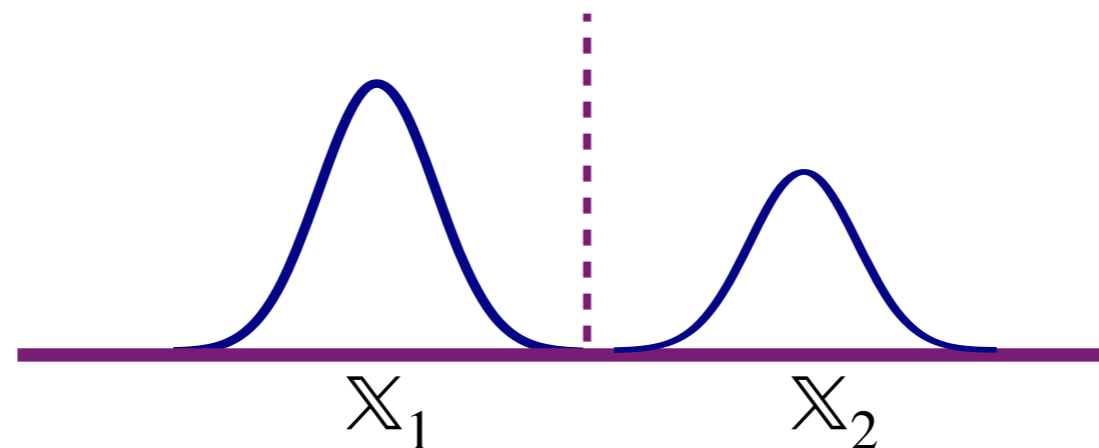
# MODES

- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if
  - ▶ It there is a possible probability to reach  $x'$  from  $x$  using  $K$  and vice-versa
  - ▶ This defines an equivalent relation on the support of  $\pi$

- ▶ We will say the  $\mathbb{X}(x)$  is the mode of  $x$

$$\mathbb{X}(x) = \{x' \in \mathbb{X} : x \sim_K x'\}$$

- ▶ I.e. the modes correspond to the irreducible components of the kernel
- ▶ The number of modes  $|\{\mathbb{X}(x) : x \in \mathbb{X}\}|$
- ▶ A distribution is uni-modal iff and only if  $K$  is  $\pi$ -irreducible



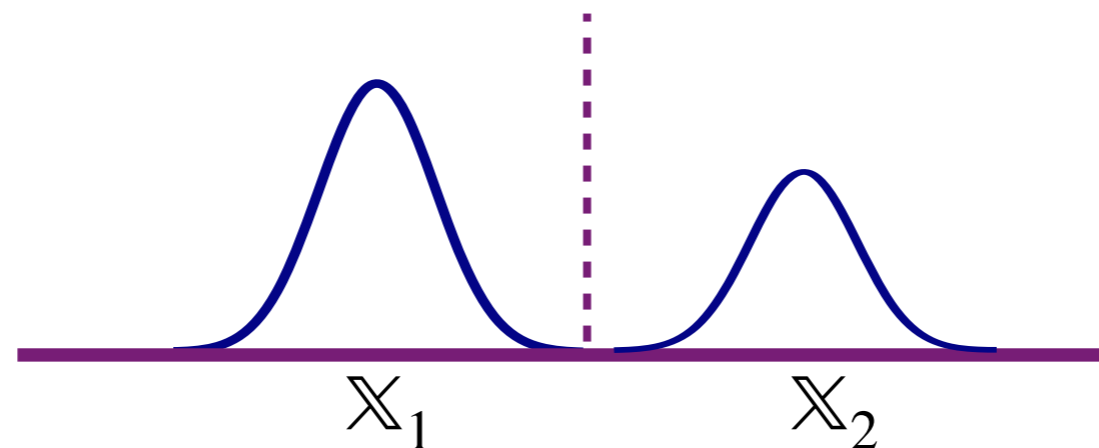
# MODES

- ▶ Let  $K$  be a  $\pi$ -invariant kernel and  $x, x' \in \text{supp}(\pi) = \mathbb{X}$ , we will say  $x \sim_K x'$  if
  - ▶ It there is a possible probability to reach  $x'$  from  $x$  using  $K$  and vice-versa
  - ▶ This defines an equivalent relation on the support of  $\pi$

- ▶ We will say the  $\mathbb{X}(x)$  is the mode of  $x$

$$\mathbb{X}(x) = \{x' \in \mathbb{X} : x \sim_K x'\}$$

- ▶ I.e. the modes correspond to the irreducible components of the kernel
- ▶ The number of modes  $|\{\mathbb{X}(x) : x \in \mathbb{X}\}|$
- ▶ A distribution is uni-modal iff and only if  $K$  is  $\pi$ -irreducible
- ▶ Notably  $\pi$  is uni-modal with respect to the independent kernel  $K(x, dx') = \pi(dx')$





# MULTI-MODAL DECOMPOSITION

# MULTI-MODAL DECOMPOSITION

- ▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$

# MULTI-MODAL DECOMPOSITION

- ▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$ 
  - ▶  $\mathbb{X}_i \cap \mathbb{X}_j = \mathbf{0}$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$

# MULTI-MODAL DECOMPOSITION

- ▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$ 
  - ▶  $\mathbb{X}_i \cap \mathbb{X}_j = \mathbf{0}$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$
- ▶ We have the following decomposition:

# MULTI-MODAL DECOMPOSITION

- ▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $\mathbf{K}$ 
  - ▶  $\mathbb{X}_i \cap \mathbb{X}_j = \mathbf{0}$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$
- ▶ We have the following decomposition:

$$\pi(x) = \sum_i \pi[\mathbb{X}_i] \pi[x | \mathbb{X}_i]$$

# MULTI-MODAL DECOMPOSITION

- ▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$ 
  - ▶  $\mathbb{X}_i \cap \mathbb{X}_j = \emptyset$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$
- ▶ We have the following decomposition:

$$\pi(x) = \sum_i \pi[\mathbb{X}_i] \pi[x | \mathbb{X}_i] = \sum_i w_i \pi_i(x)$$

# MULTI-MODAL DECOMPOSITION

- ▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$

- ▶  $\mathbb{X}_i \cap \mathbb{X}_j = \emptyset$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$

- ▶ We have the following decomposition:

$$\pi(x) = \sum_i \pi[\mathbb{X}_i] \pi[x | \mathbb{X}_i] = \sum_i w_i \pi_i(x)$$

- ▶ Where  $w_i$  is the weight of mode  $i$  and  $\pi_i$  is  $\pi$  conditional on the mode  $\mathbb{X}_i$

# MULTI-MODAL DECOMPOSITION

▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$

▶  $\mathbb{X}_i \cap \mathbb{X}_j = \emptyset$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$

▶ We have the following decomposition:

$$\pi(x) = \sum_i \pi[\mathbb{X}_i] \pi[x | \mathbb{X}_i] = \sum_i w_i \pi_i(x)$$

▶ Where  $w_i$  is the weight of mode  $i$  and  $\pi_i$  is  $\pi$  conditional on the mode  $\mathbb{X}_i$

$$w_i = \pi[\mathbb{X}_i]$$



# MULTI-MODAL DECOMPOSITION

▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$

▶  $\mathbb{X}_i \cap \mathbb{X}_j = \emptyset$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$

▶ We have the following decomposition:

$$\pi(x) = \sum_i \pi[\mathbb{X}_i] \pi[x | \mathbb{X}_i] = \sum_i w_i \pi_i(x)$$

▶ Where  $w_i$  is the weight of mode  $i$  and  $\pi_i$  is  $\pi$  conditional on the mode  $\mathbb{X}_i$

$$w_i = \pi[\mathbb{X}_i] \quad \pi_i(x) = \pi[x | \mathbb{X}_i] = \frac{\pi(x) 1[x \in \mathbb{X}_i]}{\pi[\mathbb{X}_i]}$$

# MULTI-MODAL DECOMPOSITION

▶ Let  $\mathbb{X}_i$  be the modes of the target  $\pi$  with respect to  $K$

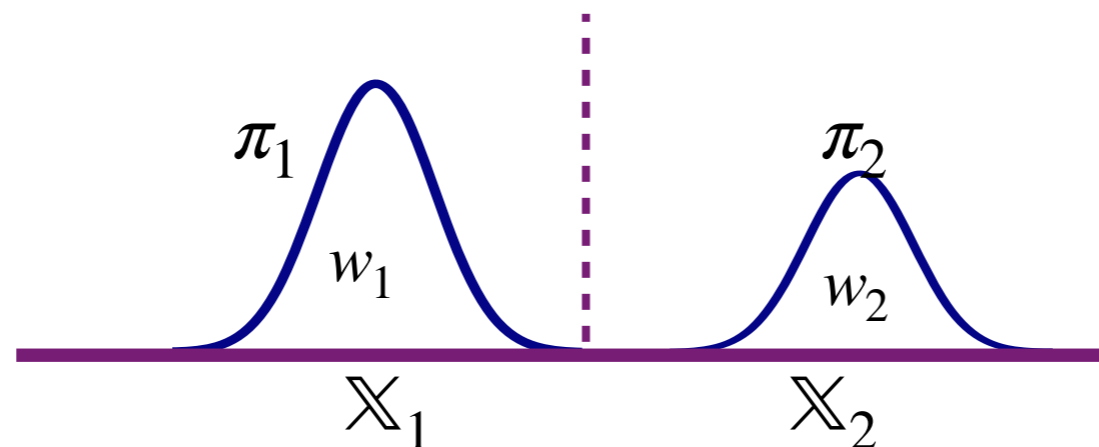
▶  $\mathbb{X}_i \cap \mathbb{X}_j = \emptyset$  for  $i \neq j$  and  $\mathbb{X} = \cup_i \mathbb{X}_i$

▶ We have the following decomposition:

$$\pi(x) = \sum_i \pi[\mathbb{X}_i] \pi[x | \mathbb{X}_i] = \sum_i w_i \pi_i(x)$$

▶ Where  $w_i$  is the weight of mode  $i$  and  $\pi_i$  is  $\pi$  conditional on the mode  $\mathbb{X}_i$

$$w_i = \pi[\mathbb{X}_i] \quad \pi_i(x) = \pi[x | \mathbb{X}_i] = \frac{\pi(x) 1[x \in \mathbb{X}_i]}{\pi[\mathbb{X}_i]}$$



# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

- ▶ Not that if  $K$  is  $\pi$ -invariant, then  $K$  is  $\pi_i$ -invariant and  $\pi_i$ -irreducible

# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

- ▶ Not that if  $K$  is  $\pi$ -invariant, then  $K$  is  $\pi_i$ -invariant and  $\pi_i$ -irreducible
- ▶ Let  $X_t$  be the Markov chain generated by  $K$ , assuming the ergodic theorem hold for each mode:

# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

- ▶ Not that if  $K$  is  $\pi$ -invariant, then  $K$  is  $\pi_i$ -invariant and  $\pi_i$ -irreducible
- ▶ Let  $X_t$  be the Markov chain generated by  $K$ , assuming the ergodic theorem hold for each mode:

- ▶ If  $X_0 \in \mathbb{X}_i$  then we have

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi_i[f]$$

# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

- ▶ Not that if  $K$  is  $\pi$ -invariant, then  $K$  is  $\pi_i$ -invariant and  $\pi_i$ -irreducible
- ▶ Let  $X_t$  be the Markov chain generated by  $K$ , assuming the ergodic theorem hold for each mode:

- ▶ If  $X_0 \in \mathbb{X}_i$  then we have

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi_i[f]$$

- ▶ They can forget the initial condition within a mode, and efficiently sample from the target conditional on the mode



# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

- ▶ Not that if  $K$  is  $\pi$ -invariant, then  $K$  is  $\pi_i$ -invariant and  $\pi_i$ -irreducible
- ▶ Let  $X_t$  be the Markov chain generated by  $K$ , assuming the ergodic theorem hold for each mode:

- ▶ If  $X_0 \in \mathbb{X}_i$  then we have

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi_i[f]$$

- ▶ They can forget the initial condition within a mode, and efficiently sample from the target conditional on the mode
  - ▶ **Problem:** They cannot forget the initial mode itself.

# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

- ▶ Not that if  $K$  is  $\pi$ -invariant, then  $K$  is  $\pi_i$ -invariant and  $\pi_i$ -irreducible
- ▶ Let  $X_t$  be the Markov chain generated by  $K$ , assuming the ergodic theorem hold for each mode:

- ▶ If  $X_0 \in \mathbb{X}_i$  then we have

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi_i[f]$$

- ▶ They can forget the initial condition within a mode, and efficiently sample from the target conditional on the mode
  - ▶ **Problem:** They cannot forget the initial mode itself.
    - ▶ What if we have gradients or do random initialisation or run multiple chains

# ERGODIC THEOREM FOR MULTI-MODAL DISTRIBUTIONS

- ▶ We have the following decomposition multi-modal decomposition  $\pi$  with respect to  $K$

$$\pi = \sum_i w_i \pi_i$$

- ▶ Not that if  $K$  is  $\pi$ -invariant, then  $K$  is  $\pi_i$ -invariant and  $\pi_i$ -irreducible
- ▶ Let  $X_t$  be the Markov chain generated by  $K$ , assuming the ergodic theorem hold for each mode:

- ▶ If  $X_0 \in \mathbb{X}_i$  then we have

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \pi_i[f]$$

- ▶ They can forget the initial condition within a mode, and efficiently sample from the target conditional on the mode
- ▶ **Problem:** They cannot forget the initial mode itself.
  - ▶ What if we have gradients or do random initialisation or run multiple chains
  - ▶ Unfortunately none of these work

# PITFALLS OF FIRST ORDER METHODS

- ▶ In general first-order methods don't help with multi-modal distributions:

# PITFALLS OF FIRST ORDER METHODS

- ▶ In general first-order methods don't help with multi-modal distributions:
- ▶ We we know each mode up to a normalising constant

$$\pi_i(x) = \frac{\gamma_i(x)}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x) dx$$

# PITFALLS OF FIRST ORDER METHODS

- ▶ In general first-order methods don't help with multi-modal distributions:
- ▶ We we know each mode up to a normalising constant

$$\pi_i(x) = \frac{\gamma_i(x)}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x) dx$$

- ▶ The normalising constant provides information about relative weighting of modes:

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

# PITFALLS OF FIRST ORDER METHODS

- ▶ In general first-order methods don't help with multi-modal distributions:
- ▶ We we know each mode up to a normalising constant

$$\pi_i(x) = \frac{\gamma_i(x)}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x) dx$$

- ▶ The normalising constant provides information about relative weighting of modes:

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ The score is agnostic to the normalising constant and relative weight of the modes  $w_i$

$$\nabla \log \pi(x) = \sum_i \nabla \log \gamma_i(x)$$

# PITFALLS OF FIRST ORDER METHODS

- ▶ In general first-order methods don't help with multi-modal distributions:
- ▶ We we know each mode up to a normalising constant

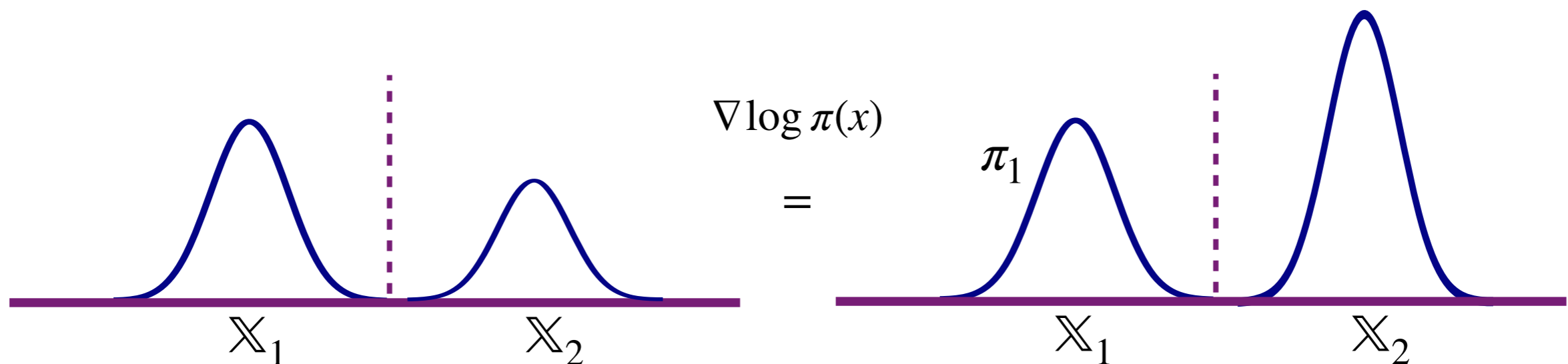
$$\pi_i(x) = \frac{\gamma_i(x)}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x) dx$$

- ▶ The normalising constant provides information about relative weighting of modes:

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ The score is agnostic to the normalising constant and relative weight of the modes  $w_i$

$$\nabla \log \pi(x) = \sum_i \nabla \log \gamma_i(x)$$





# PITFALL OF RANDOM INITIALISATION

# PITFALL OF RANDOM INITIALISATION

- ▶ Suppose we initialise  $X_0 \sim \mu$  for some distribution  $\mu$  over  $\mathbb{X}$

# PITFALL OF RANDOM INITIALISATION

- ▶ Suppose we initialise  $X_0 \sim \mu$  for some distribution  $\mu$  over  $\mathbb{X}$ 
  - ▶ If  $X_0 \sim \mu$  then let  $\alpha_i = \mu[\mathbb{X}_i]$

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \sum_i \alpha_i \pi_i[f]$$

# PITFALL OF RANDOM INITIALISATION

- ▶ Suppose we initialise  $X_0 \sim \mu$  for some distribution  $\mu$  over  $\mathbb{X}$ 
  - ▶ If  $X_0 \sim \mu$  then let  $\alpha_i = \mu[\mathbb{X}_i]$

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \sum_i \alpha_i \pi_i[f]$$

- ▶ Note that this implies that even if we initialise randomly, we can efficiently sample from

$$\pi_\mu = \sum_i \alpha_i \pi_i$$

# PITFALL OF RANDOM INITIALISATION

- ▶ Suppose we initialise  $X_0 \sim \mu$  for some distribution  $\mu$  over  $\mathbb{X}$ 
  - ▶ If  $X_0 \sim \mu$  then let  $\alpha_i = \mu[\mathbb{X}_i]$

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \sum_i \alpha_i \pi_i[f]$$

- ▶ Note that this implies that even if we initialise randomly, we can efficiently sample from

$$\pi_\mu = \sum_i \alpha_i \pi_i$$

- ▶  $\pi_\mu = \pi$  if and only if  $\mu[\mathbb{X}_i] = \pi[\mathbb{X}_i]$

# PITFALL OF RANDOM INITIALISATION

- ▶ Suppose we initialise  $X_0 \sim \mu$  for some distribution  $\mu$  over  $\mathbb{X}$ 
  - ▶ If  $X_0 \sim \mu$  then let  $\alpha_i = \mu[\mathbb{X}_i]$

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \sum_i \alpha_i \pi_i[f]$$

- ▶ Note that this implies that even if we initialise randomly, we can efficiently sample from

$$\pi_\mu = \sum_i \alpha_i \pi_i$$

- ▶  $\pi_\mu = \pi$  if and only if  $\mu[\mathbb{X}_i] = \pi[\mathbb{X}_i]$
- ▶ If we know the size and location of the modes of  $\pi$  then we can construct  $\mu$  such that the ergodic average is consistent

# PITFALL OF RANDOM INITIALISATION

- ▶ Suppose we initialise  $X_0 \sim \mu$  for some distribution  $\mu$  over  $\mathbb{X}$ 
  - ▶ If  $X_0 \sim \mu$  then let  $\alpha_i = \mu[\mathbb{X}_i]$

$$\frac{1}{T} \sum_{t=1}^T f(X_t) \xrightarrow[T \rightarrow \infty]{a.s.} \sum_i \alpha_i \pi_i[f]$$

- ▶ Note that this implies that even if we initialise randomly, we can efficiently sample from

$$\pi_\mu = \sum_i \alpha_i \pi_i$$

- ▶  $\pi_\mu = \pi$  if and only if  $\mu[\mathbb{X}_i] = \pi[\mathbb{X}_i]$
- ▶ If we know the size and location of the modes of  $\pi$  then we can construct  $\mu$  such that the ergodic average is consistent
- ▶ In general, we do not know this information, and therefore will get biased results.

# PITFALLS OF TARGETTING A DISTRIBUTION



# PITFALLS OF TARGETTING A DISTRIBUTION

- ▶ The normalising constant encode the information about the relative size of modes

$$\pi_i(x) = \frac{\gamma(x)1[x \in \mathbb{X}_i]}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x)dx$$

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

# PITFALLS OF TARGETTING A DISTRIBUTION

- ▶ The normalising constant encode the information about the relative size of modes

$$\pi_i(x) = \frac{\gamma(x)1[x \in \mathbb{X}_i]}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x)dx$$

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ For a valid MCMC algorithm should be invariant to the choice of normalising constant

# PITFALLS OF TARGETTING A DISTRIBUTION

- ▶ The normalising constant encode the information about the relative size of modes

$$\pi_i(x) = \frac{\gamma(x)1[x \in \mathbb{X}_i]}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x)dx$$

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ For a valid MCMC algorithm should be invariant to the choice of normalising constant
  - ▶ The choice of the mode is determined by the initial state  $X_0$

# PITFALLS OF TARGETTING A DISTRIBUTION

- ▶ The normalising constant encode the information about the relative size of modes

$$\pi_i(x) = \frac{\gamma(x)1[x \in \mathbb{X}_i]}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x)dx$$

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ For a valid MCMC algorithm should be invariant to the choice of normalising constant
  - ▶ The choice of the mode is determined by the initial state  $X_0$
  - ▶ The chain  $X_0, X_1, \dots, X_t, \dots$  generated by  $K$  is identical regardless of  $w_i$  or  $Z_i$

# PITFALLS OF TARGETTING A DISTRIBUTION

- ▶ The normalising constant encode the information about the relative size of modes

$$\pi_i(x) = \frac{\gamma(x)1[x \in \mathbb{X}_i]}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x)dx$$

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ For a valid MCMC algorithm should be invariant to the choice of normalising constant
  - ▶ The choice of the mode is determined by the initial state  $X_0$
  - ▶ The chain  $X_0, X_1, \dots, X_t, \dots$  generated by  $K$  is identical regardless of  $w_i$  or  $Z_i$

$$\pi \propto \sum_i w_i \gamma_i \propto \sum_i w'_i \gamma'_i$$

# PITFALLS OF TARGETTING A DISTRIBUTION

- ▶ The normalising constant encode the information about the relative size of modes

$$\pi_i(x) = \frac{\gamma(x)1[x \in \mathbb{X}_i]}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x)dx$$

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ For a valid MCMC algorithm should be invariant to the choice of normalising constant
  - ▶ The choice of the mode is determined by the initial state  $X_0$
  - ▶ The chain  $X_0, X_1, \dots, X_t, \dots$  generated by  $K$  is identical regardless of  $w_i$  or  $Z_i$

$$\pi \propto \sum_i w_i \gamma_i \propto \sum_i w'_i \gamma'_i$$

- ▶ It means we can not estimate  $Z_i$  even with a mode

# PITFALLS OF TARGETTING A DISTRIBUTION

- ▶ The normalising constant encode the information about the relative size of modes

$$\pi_i(x) = \frac{\gamma(x)1[x \in \mathbb{X}_i]}{Z_i}, \quad Z_i = \int_{\mathbb{X}_i} \gamma_i(x)dx$$

$$\pi(x) \propto \sum_i w_i \gamma_i(x), \quad Z = \sum_i w_i Z_i$$

- ▶ For a valid MCMC algorithm should be invariant to the choice of normalising constant
  - ▶ The choice of the mode is determined by the initial state  $X_0$
  - ▶ The chain  $X_0, X_1, \dots, X_t, \dots$  generated by  $K$  is identical regardless of  $w_i$  or  $Z_i$

$$\pi \propto \sum_i w_i \gamma_i \propto \sum_i w'_i \gamma'_i$$

- ▶ It means we can not estimate  $Z_i$  even with a mode
- ▶ To learn the normalising constant, we need something to compare against

# REFERENCE DISTRIBUTIONS

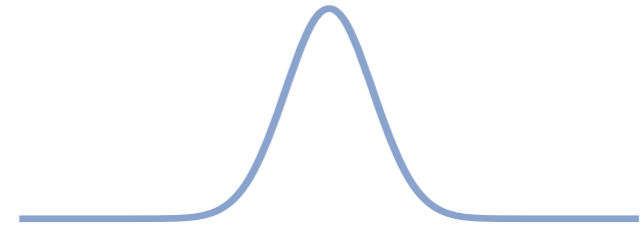
- ▶ Suppose we have a target distribution  $\pi$ :



# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$

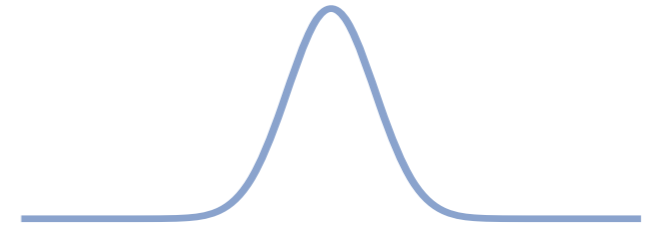


# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$

- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:

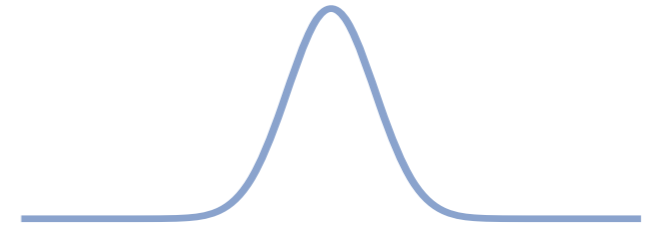


# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$

- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:
  - ▶ We can sample  $X \sim \eta$



# REFERENCE DISTRIBUTIONS

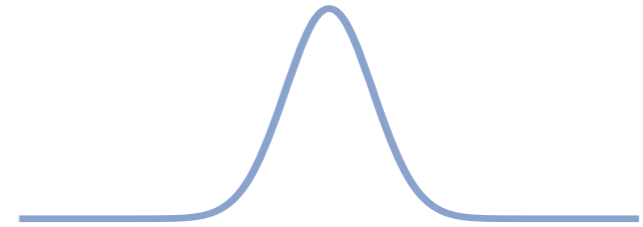
- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$

- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:

- ▶ We can sample  $X \sim \eta$

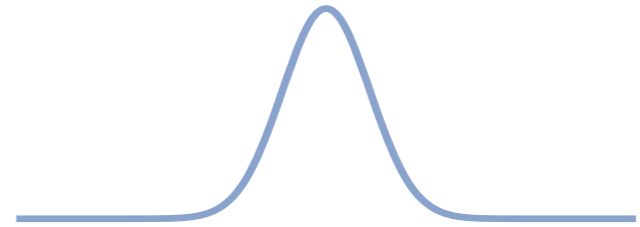
- ▶ We can evaluate the normalised density  $\eta(x)$ , where  $\eta(dx) = \eta(x) dx$



# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$

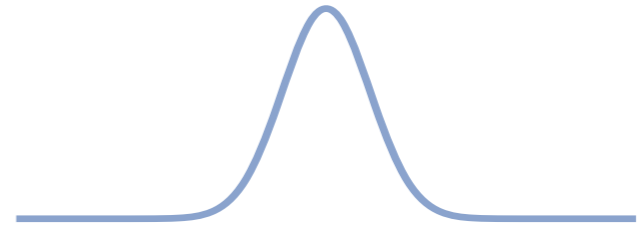


- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:
  - ▶ We can sample  $X \sim \eta$
  - ▶ We can evaluate the normalised density  $\eta(x)$ , where  $\eta(dx) = \eta(x) dx$
  - ▶ We have  $\pi \ll \eta$  and we can evaluate the weight function  $w : \mathbb{X} \rightarrow \mathbb{R}_+$  defined as the unnormalised likelihood ratio:

# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$



- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:

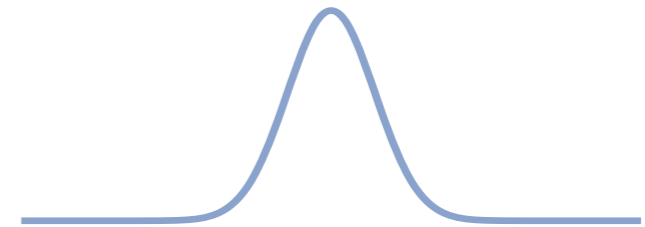
- ▶ We can sample  $X \sim \eta$
- ▶ We can evaluate the normalised density  $\eta(x)$ , where  $\eta(dx) = \eta(x) dx$
- ▶ We have  $\pi \ll \eta$  and we can evaluate the weight function  $w : \mathbb{X} \rightarrow \mathbb{R}_+$  defined as the unnormalised likelihood ratio:

$$w(x) = Z \frac{d\pi}{d\eta}(x) = \frac{\gamma(x)}{\eta(x)}$$

# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$



- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:

- ▶ We can sample  $X \sim \eta$
- ▶ We can evaluate the normalised density  $\eta(x)$ , where  $\eta(dx) = \eta(x) dx$
- ▶ We have  $\pi \ll \eta$  and we can evaluate the weight function  $w : \mathbb{X} \rightarrow \mathbb{R}_+$  defined as the unnormalised likelihood ratio:

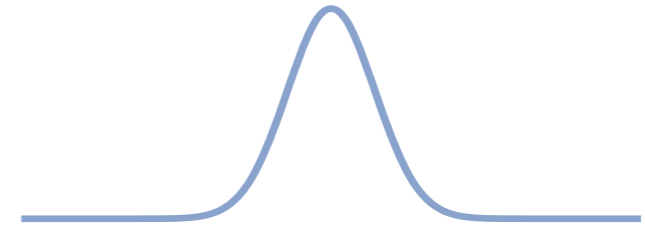
$$w(x) = Z \frac{d\pi}{d\eta}(x) = \frac{\gamma(x)}{\eta(x)}$$

- ▶ Reference distribution should be simple and tractable:

# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$



- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:

- ▶ We can sample  $X \sim \eta$
- ▶ We can evaluate the normalised density  $\eta(x)$ , where  $\eta(dx) = \eta(x) dx$
- ▶ We have  $\pi \ll \eta$  and we can evaluate the weight function  $w : \mathbb{X} \rightarrow \mathbb{R}_+$  defined as the unnormalised likelihood ratio:

$$w(x) = Z \frac{d\pi}{d\eta}(x) = \frac{\gamma(x)}{\eta(x)}$$

- ▶ Reference distribution should be simple and tractable:

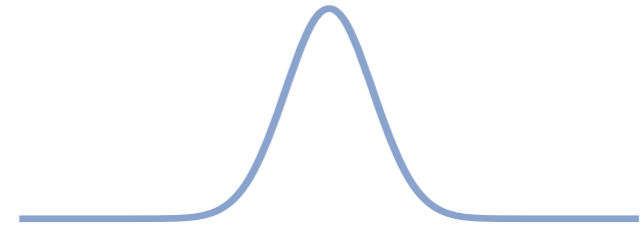
- ▶ Common choices will be Gaussians, uniform distribution, priors, etc



# REFERENCE DISTRIBUTIONS

- ▶ Suppose we have a target distribution  $\pi$ :

$$\pi(dx) = \frac{\gamma(x)}{Z} dx, \quad Z = \int_{\mathbb{X}} \gamma(x) dx$$



- ▶ We will say  $\eta$  is a **reference** distribution for  $\pi$  if:

- ▶ We can sample  $X \sim \eta$
- ▶ We can evaluate the normalised density  $\eta(x)$ , where  $\eta(dx) = \eta(x) dx$
- ▶ We have  $\pi \ll \eta$  and we can evaluate the weight function  $w : \mathbb{X} \rightarrow \mathbb{R}_+$  defined as the unnormalised likelihood ratio:

$$w(x) = Z \frac{d\pi}{d\eta}(x) = \frac{\gamma(x)}{\eta(x)}$$

- ▶ Reference distribution should be simple and tractable:

- ▶ Common choices will be Gaussians, uniform distribution, priors, etc
- ▶ Typically the intialisation for the MCMC chain

# REFERENCE DISTRIBUTIONS

# REFERENCE DISTRIBUTIONS

- ▶ Note that given a reference distributions we can approximate the normalising constant:

# REFERENCE DISTRIBUTIONS

- ▶ Note that given a reference distributions we can approximate the normalising constant:
- ▶ Recall that with a reference distribution:

$$Z = \int_{\mathbb{X}} \gamma(x) dx = \int_{\mathbb{X}} w(x) \eta(x) dx = \eta[w]$$

- ▶ Use Monte Carlo we can approximate:

$$X_1, \dots, X_N \sim \eta \quad \hat{Z} = \frac{1}{N} \sum_{n=1}^N w(X_n)$$

- ▶ It is unbiased and consistent estimator for  $Z$

$$\hat{Z} \xrightarrow[N \rightarrow \infty]{a.s.} Z$$

# REFERENCE DISTRIBUTIONS

- ▶ Note that given a reference distributions we can approximate the normalising constant:
- ▶ Recall that with a reference distribution:

$$Z = \int_{\mathbb{X}} \gamma(x) dx = \int_{\mathbb{X}} w(x) \eta(x) dx = \eta[w]$$

- ▶ Use Monte Carlo we can approximate:

$$X_1, \dots, X_N \sim \eta \quad \hat{Z} = \frac{1}{N} \sum_{n=1}^N w(X_n)$$

- ▶ It is unbiased and consistent estimator for  $Z$

$$\hat{Z} \xrightarrow[N \rightarrow \infty]{a.s.} Z$$

- ▶ If we have a local sampler  $K$ , we can learn a mode

# REFERENCE DISTRIBUTIONS

- ▶ Note that given a reference distributions we can approximate the normalising constant:
- ▶ Recall that with a reference distribution:

$$Z = \int_{\mathbb{X}} \gamma(x) dx = \int_{\mathbb{X}} w(x) \eta(x) dx = \eta[w]$$

- ▶ Use Monte Carlo we can approximate:

$$X_1, \dots, X_N \sim \eta \quad \hat{Z} = \frac{1}{N} \sum_{n=1}^N w(X_n)$$

- ▶ It is unbiased and consistent estimator for  $Z$

$$\hat{Z} \xrightarrow[N \rightarrow \infty]{a.s.} Z$$

- ▶ If we have a local sampler  $K$ , we can learn a mode
  - ▶ Can get consistency if we can initialise in the modes with the right weights

# REFERENCE DISTRIBUTIONS

- ▶ Note that given a reference distributions we can approximate the normalising constant:
- ▶ Recall that with a reference distribution:

$$Z = \int_{\mathbb{X}} \gamma(x) dx = \int_{\mathbb{X}} w(x) \eta(x) dx = \eta[w]$$

- ▶ Use Monte Carlo we can approximate:

$$X_1, \dots, X_N \sim \eta \quad \hat{Z} = \frac{1}{N} \sum_{n=1}^N w(X_n)$$

- ▶ It is unbiased and consistent estimator for  $Z$

$$\hat{Z} \xrightarrow[N \rightarrow \infty]{a.s.} Z$$

- ▶ If we have a local sampler  $K$ , we can learn a mode
  - ▶ Can get consistency if we can initialise in the modes with the right weights
- ▶ Reference is able to learn the normalising constant and hence the modes

# REFERENCE DISTRIBUTIONS

- ▶ Note that given a reference distributions we can approximate the normalising constant:
- ▶ Recall that with a reference distribution:

$$Z = \int_{\mathbb{X}} \gamma(x) dx = \int_{\mathbb{X}} w(x) \eta(x) dx = \eta[w]$$

- ▶ Use Monte Carlo we can approximate:

$$X_1, \dots, X_N \sim \eta \quad \hat{Z} = \frac{1}{N} \sum_{n=1}^N w(X_n)$$

- ▶ It is unbiased and consistent estimator for  $Z$

$$\hat{Z} \xrightarrow[N \rightarrow \infty]{a.s.} Z$$

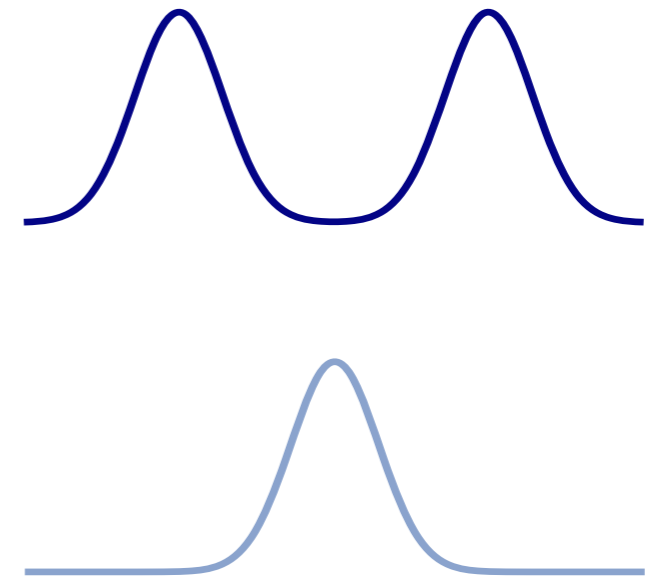
- ▶ If we have a local sampler  $K$ , we can learn a mode
  - ▶ Can get consistency if we can initialise in the modes with the right weights
- ▶ Reference is able to learn the normalising constant and hence the modes
- ▶ Can use the reference to stabilise generated by  $K$



# REFERENCE STABILISED MCMC

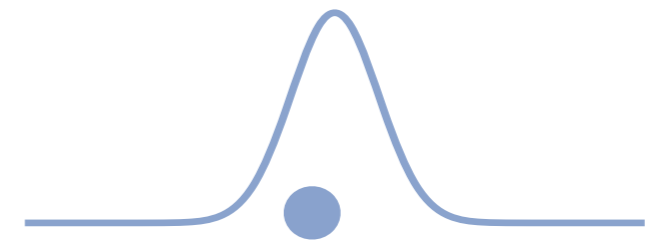
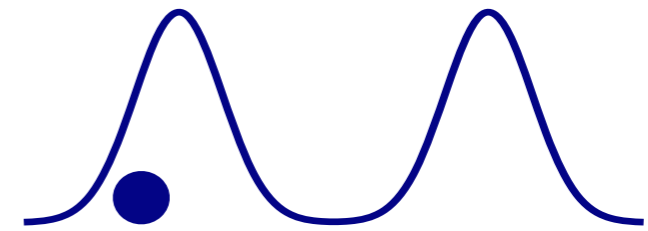
# REFERENCE STABILISED MCMC

- ▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$



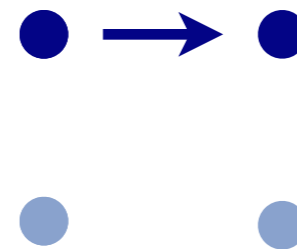
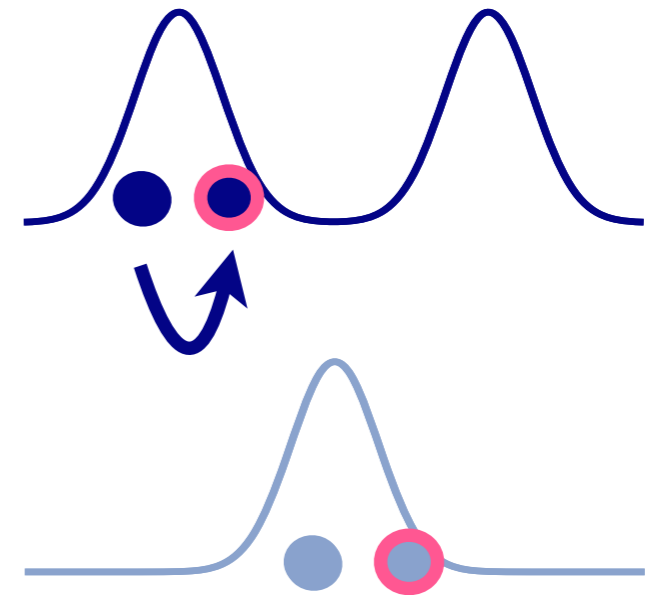
# REFERENCE STABILISED MCMC

- ▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$ 
  - ▶ **Intitialise:**  $X_0, X'_0 \sim \eta$



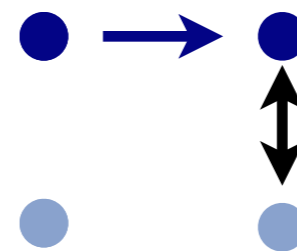
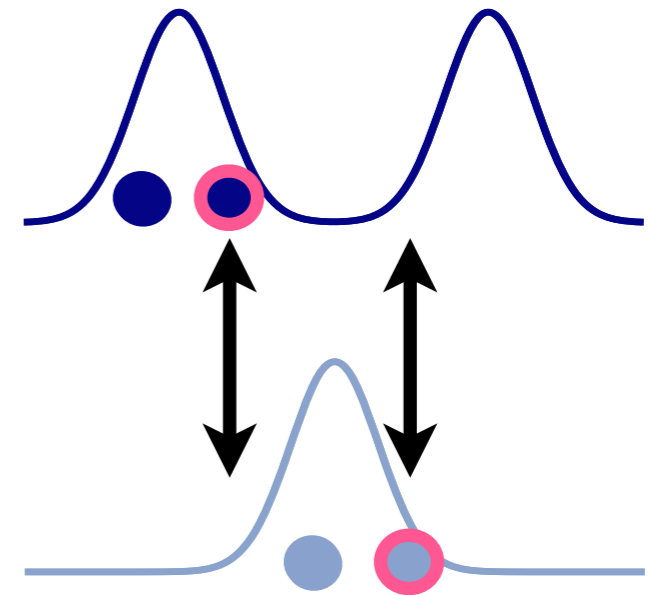
# REFERENCE STABILISED MCMC

- ▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$ 
  - ▶ **Intitialise:**  $X_0, X'_0 \sim \eta$
  - ▶ **Local move:**
    - ▶ Propogate target state  $X \sim K(X_{t-1}, dx)$
    - ▶ Generate new reference sample  $X' \sim \eta$



# REFERENCE STABILISED MCMC

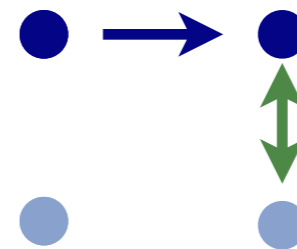
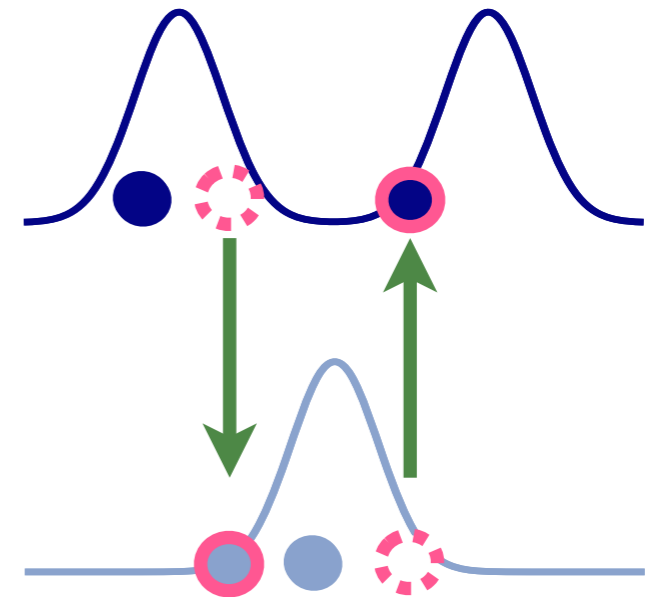
- ▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$ 
  - ▶ **Intitialise:**  $X_0, X'_0 \sim \eta$
  - ▶ **Local move:**
    - ▶ Propogate target state  $X \sim K(X_{t-1}, dx)$
    - ▶ Generate new reference sample  $X' \sim \eta$
  - ▶ **Communication:** propose a metropolised swap



# REFERENCE STABILISED MCMC

- ▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$ 
  - ▶ **Intitialise:**  $X_0, X'_0 \sim \eta$
  - ▶ **Local move:**
    - ▶ Propogate target state  $X \sim K(X_{t-1}, dx)$
    - ▶ Generate new reference sample  $X' \sim \eta$
  - ▶ **Communication:** propose a metropolised swap
    - ▶ **Accept** with probability  $\alpha(X, X')$

$$(X_t, X'_t) = (X', X)$$



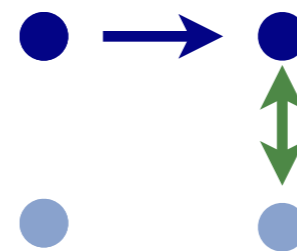
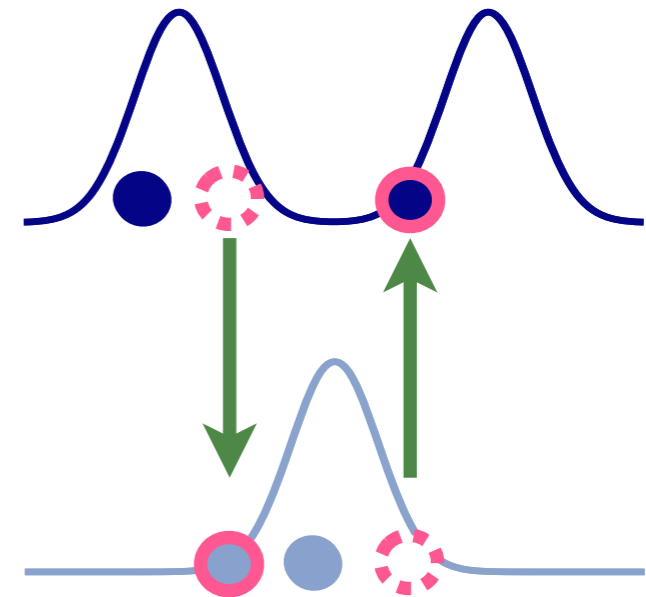
# REFERENCE STABILISED MCMC

- ▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$ 
  - ▶ **Intitialise:**  $X_0, X'_0 \sim \eta$
  - ▶ **Local move:**
    - ▶ Propogate target state  $X \sim K(X_{t-1}, dx)$
    - ▶ Generate new reference sample  $X' \sim \eta$
- ▶ **Communication:** propose a metropolised swap
  - ▶ **Accept** with probability  $\alpha(X, X')$

$$(X_t, X'_t) = (X', X)$$

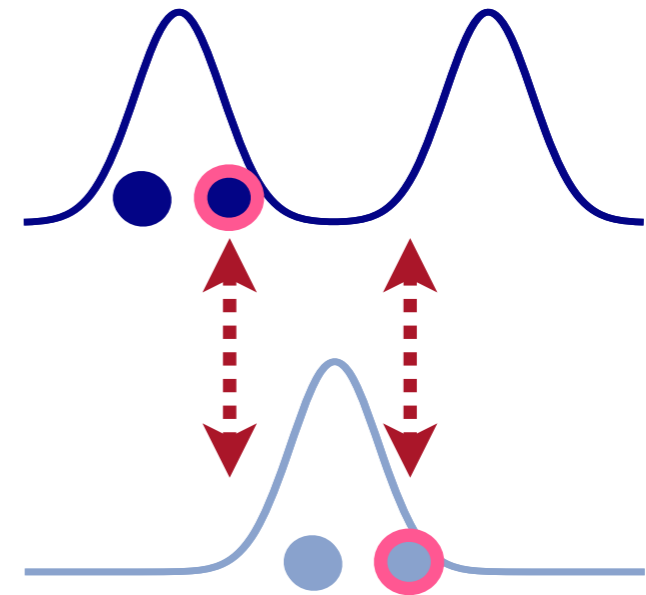
- ▶ Where  $\alpha(x, x')$  equals:

$$\alpha(x, x') = 1 \wedge \frac{\pi(x')\eta(x)}{\pi(x)\eta(x')} = 1 \wedge \frac{w(x')}{w(x)}$$



# RERERENCE STABILISED MCMC

- ▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$ 
  - ▶ **Intitialise:**  $X_0, X'_0 \sim \eta$
  - ▶ **Local move:**
    - ▶ Propogate target state  $X \sim K(X_{t-1}, dx)$
    - ▶ Generate new reference sample  $X' \sim \eta$



- ▶ **Communication:** propose a metropolised swap
  - ▶ **Accept** with probability  $\alpha(X, X')$

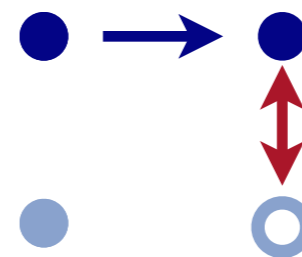
$$(X_t, X'_t) = (X', X)$$

- ▶ **Reject** otherwise and set

$$(X_t, X'_t) = (X, X')$$

- ▶ Where  $\alpha(x, x')$  equals:

$$\alpha(x, x') = 1 \wedge \frac{\pi(x')\eta(x)}{\pi(x)\eta(x')} = 1 \wedge \frac{w(x')}{w(x)}$$





# RERERENCE STABILISED MCMC

▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$

▶ **Intitialise:**  $X_0, X'_0 \sim \eta$

▶ **Local move:**

▶ Propogate target state  $X \sim K(X_{t-1}, dx)$

▶ Generate new reference sample  $X' \sim \eta$

▶ **Communication:** propose a metropolised swap

▶ **Accept** with probability  $\alpha(X, X')$

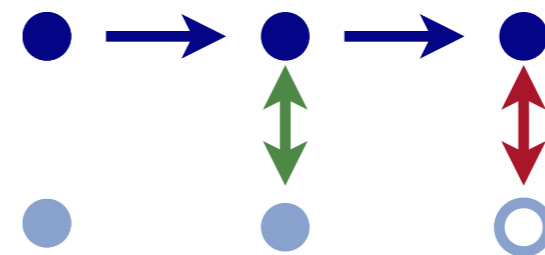
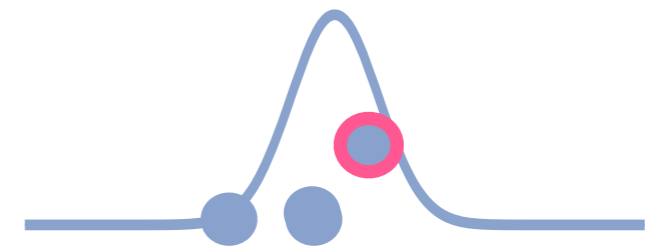
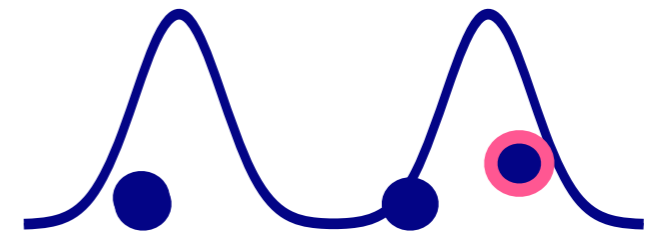
$$(X_t, X'_t) = (X', X)$$

▶ **Reject** otherwise and set

$$(X_t, X'_t) = (X, X')$$

▶ Where  $\alpha(x, x')$  equals:

$$\alpha(x, x') = 1 \wedge \frac{\pi(x')\eta(x)}{\pi(x)\eta(x')} = 1 \wedge \frac{w(x')}{w(x)}$$



# REFERENCE STABILISED MCMC

▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$

▶ **Intitialise:**  $X_0, X'_0 \sim \eta$

▶ **Local move:**

▶ Propogate target state  $X \sim K(X_{t-1}, dx)$

▶ Generate new reference sample  $X' \sim \eta$

▶ **Communication:** propose a metropolised swap

▶ **Accept** with probability  $\alpha(X, X')$

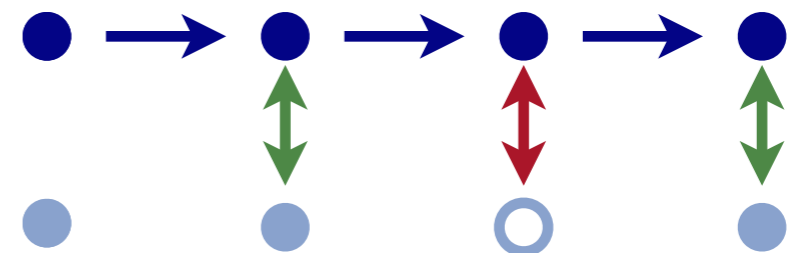
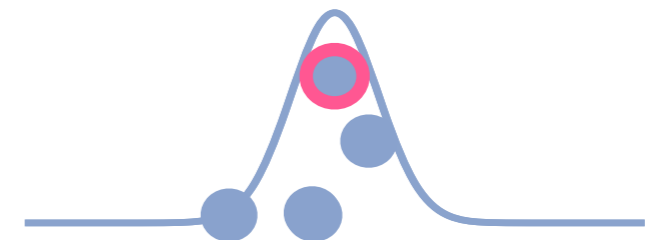
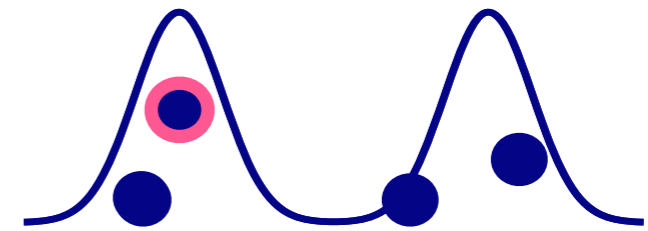
$$(X_t, X'_t) = (X', X)$$

▶ **Reject** otherwise and set

$$(X_t, X'_t) = (X, X')$$

▶ Where  $\alpha(x, x')$  equals:

$$\alpha(x, x') = 1 \wedge \frac{\pi(x')\eta(x)}{\pi(x)\eta(x')} = 1 \wedge \frac{w(x')}{w(x)}$$



# REFERENCE STABILISED MCMC

▶ Generate a chain in  $(X_t, X'_t) \in \mathbb{X}^2$  targeting  $\pi \otimes \eta$

▶ **Intialise:**  $X_0, X'_0 \sim \eta$

▶ **Local move:**

▶ Propogate target state  $X \sim K(X_{t-1}, dx)$

▶ Generate new reference sample  $X' \sim \eta$

▶ **Communication:** propose a metropolised swap

▶ **Accept** with probability  $\alpha(X, X')$

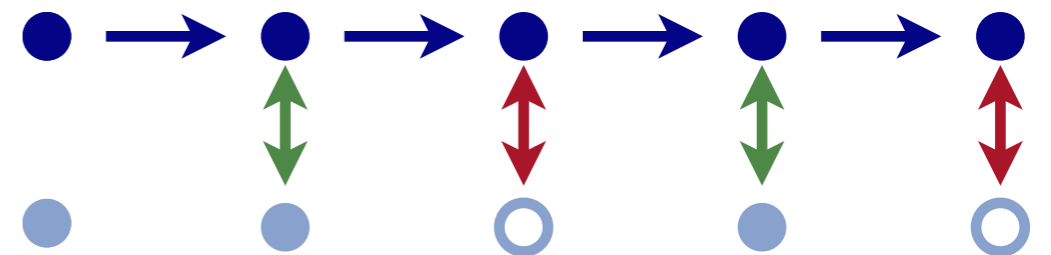
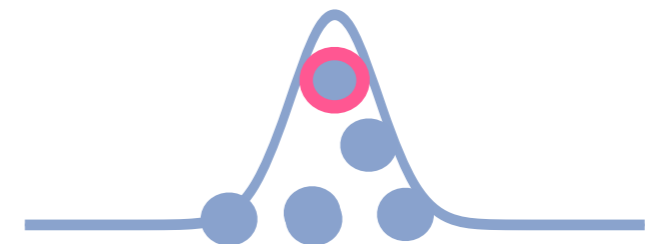
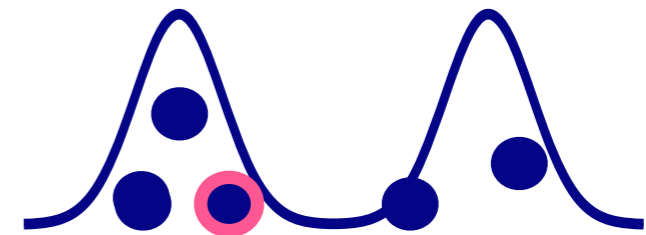
$$(X_t, X'_t) = (X', X)$$

▶ **Reject** otherwise and set

$$(X_t, X'_t) = (X, X')$$

▶ Where  $\alpha(x, x')$  equals:

$$\alpha(x, x') = 1 \wedge \frac{\pi(x')\eta(x)}{\pi(x)\eta(x')} = 1 \wedge \frac{w(x')}{w(x)}$$



# REFERENCE DISTRIBUTIONS

# REFERENCE DISTRIBUTIONS

- ▶ A successful swap acts as a random initialisation that matches the proportionality of the target!

# REFERENCE DISTRIBUTIONS

- ▶ A successful swap acts as a random initialisation that matches the proportionality of the target!
- ▶ Can use the reference samples to also obtain normalising constant estimates

$$\hat{Z} = \frac{1}{T} \sum_{t=1}^T w(X'_t), \quad \mathbb{V} \left[ \frac{\hat{Z}}{Z} \right] = \frac{\exp(D(\pi||\eta)) - 1}{N}$$

$$D(\pi||\eta) = \log(1 + \chi^2(\pi||\eta)), \quad \chi^2(\pi||\eta) = \mathbb{V}_\eta \left[ \frac{d\pi}{d\eta} \right]$$

# REFERENCE DISTRIBUTIONS

- ▶ A successful swap acts as a random initialisation that matches the proportionality of the target!
- ▶ Can use the reference samples to also obtain normalising constant estimates

$$\hat{Z} = \frac{1}{T} \sum_{t=1}^T w(X'_t), \quad \mathbb{V} \left[ \frac{\hat{Z}}{Z} \right] = \frac{\exp(D(\pi||\eta)) - 1}{N}$$

$$D(\pi||\eta) = \log(1 + \chi^2(\pi||\eta)), \quad \chi^2(\pi||\eta) = \mathbb{V}_\eta \left[ \frac{d\pi}{d\eta} \right]$$

- ▶ The swaps and normalising constants are stable if  $\eta \approx \pi$  rapidly deteriorates otherwise

# REFERENCE DISTRIBUTIONS

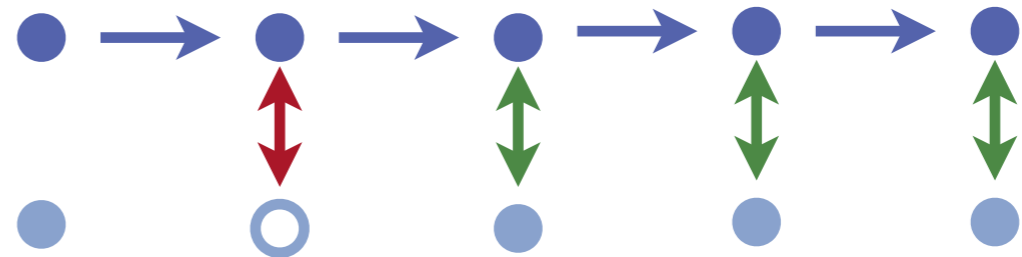
- ▶ A successful swap acts as a random initialisation that matches the proportionality of the target!
- ▶ Can use the reference samples to also obtain normalising constant estimates

$$\hat{Z} = \frac{1}{T} \sum_{t=1}^T w(X'_t), \quad \mathbb{V} \left[ \frac{\hat{Z}}{Z} \right] = \frac{\exp(D(\pi||\eta)) - 1}{N}$$

$$D(\pi||\eta) = \log(1 + \chi^2(\pi||\eta)), \quad \chi^2(\pi||\eta) = \mathbb{V}_\eta \left[ \frac{d\pi}{d\eta} \right]$$

- ▶ The swaps and normalising constants are stable if  $\eta \approx \pi$  rapidly deteriorates otherwise

$$\lim_{D \rightarrow 0^+} \mathbb{V} \left[ \frac{\hat{Z}}{Z} \right] \sim \frac{D}{T}$$





# REFERENCE DISTRIBUTIONS

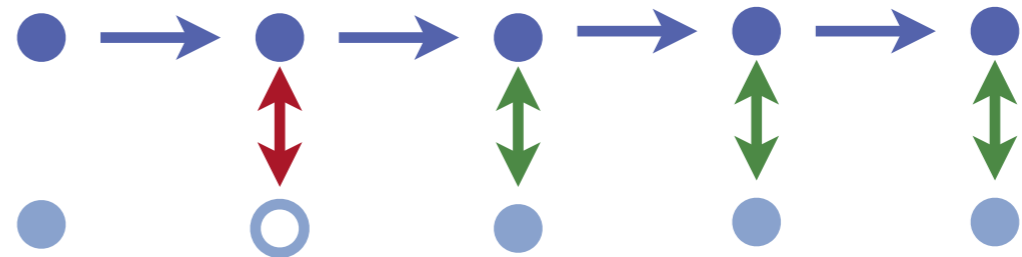
- ▶ A successful swap acts as a random initialisation that matches the proportionality of the target!
- ▶ Can use the reference samples to also obtain normalising constant estimates

$$\hat{Z} = \frac{1}{T} \sum_{t=1}^T w(X'_t), \quad \mathbb{V} \left[ \frac{\hat{Z}}{Z} \right] = \frac{\exp(D(\pi||\eta)) - 1}{N}$$

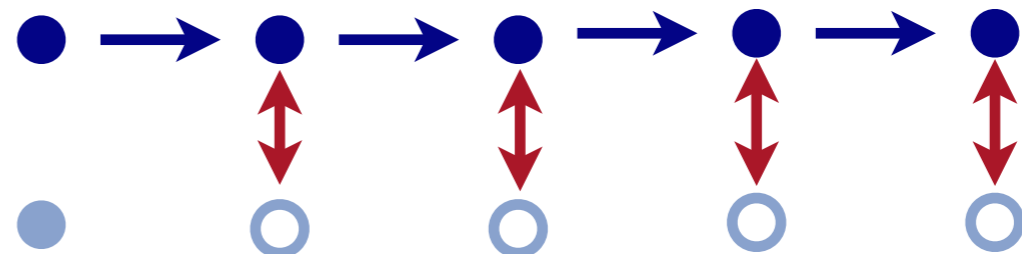
$$D(\pi||\eta) = \log(1 + \chi^2(\pi||\eta)), \quad \chi^2(\pi||\eta) = \mathbb{V}_\eta \left[ \frac{d\pi}{d\eta} \right]$$

- ▶ The swaps and normalising constants are stable if  $\eta \approx \pi$  rapidly deteriorates otherwise

$$\lim_{D \rightarrow 0^+} \mathbb{V} \left[ \frac{\hat{Z}}{Z} \right] \sim \frac{D}{T}$$



$$\lim_{D \rightarrow \infty} \mathbb{V} \left[ \frac{\hat{Z}}{Z} \right] \sim \frac{\exp(D)}{T}$$



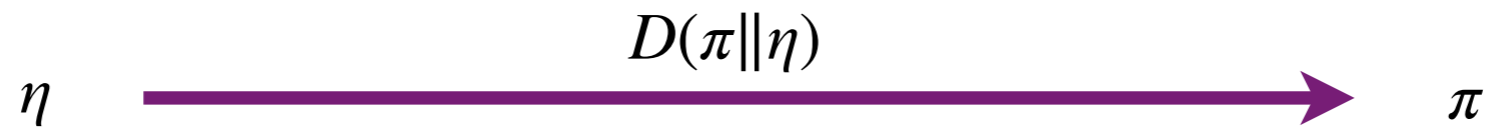
# VARIATIONAL INFERENCE

# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$

# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



# VARIATIONAL INFERENCE

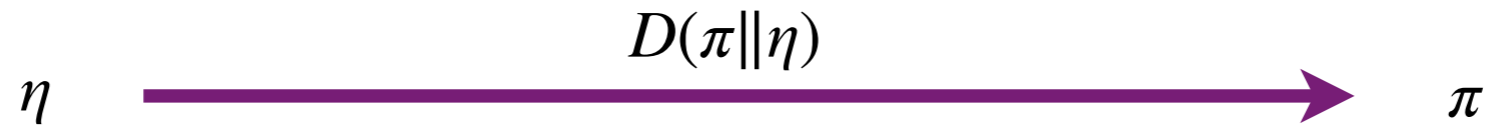
- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial

# VARIATIONAL INFERENCE

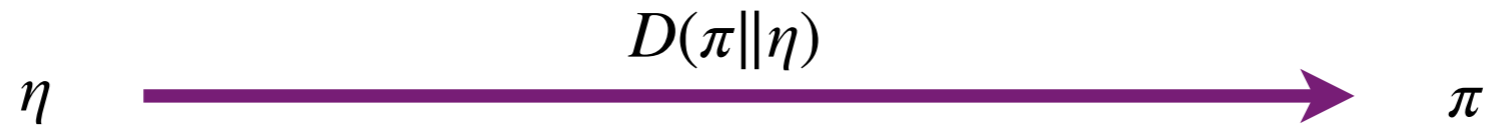
- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial
  - ▶ We can efficiently estimate  $\pi$  using the reference  $\eta$  if and only if  $D(\pi||\eta)$  is small

# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial
  - ▶ We can efficiently estimate  $\pi$  using the reference  $\eta$  if and only if  $D(\pi||\eta)$  is small
  - ▶ We want a reference that is as close to the target as possible

# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial
  - ▶ We can efficiently estimate  $\pi$  using the reference  $\eta$  if and only if  $D(\pi||\eta)$  is small
  - ▶ We want a reference that is as close to the target as possible
- ▶ **Solution 1:** Keep the target fixed and modify the reference (i.e. Variational inference)

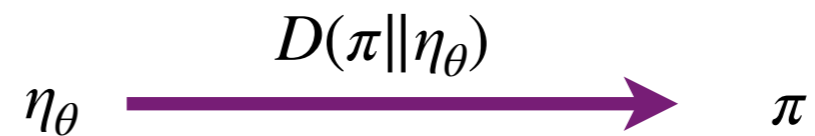


# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial
  - ▶ We can efficiently estimate  $\pi$  using the reference  $\eta$  if and only if  $D(\pi||\eta)$  is small
  - ▶ We want a reference that is as close to the target as possible
- ▶ **Solution 1:** Keep the target fixed and modify the reference (i.e. Variational inference)

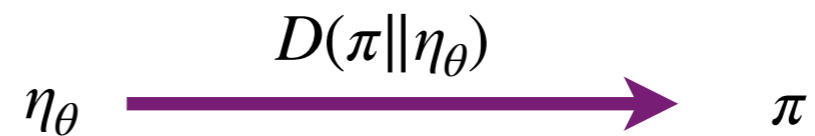


# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial
  - ▶ We can efficiently estimate  $\pi$  using the reference  $\eta$  if and only if  $D(\pi||\eta)$  is small
  - ▶ We want a reference that is as close to the target as possible
- ▶ **Solution 1:** Keep the target fixed and modify the reference (i.e. Variational inference)



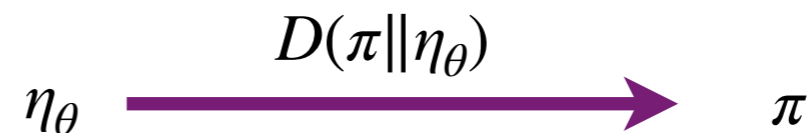
- ▶ Let  $\mathcal{Q} = \{\eta_\theta\}_{\theta \in \Theta}$  be a **variational** family of reference distribution

# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial
  - ▶ We can efficiently estimate  $\pi$  using the reference  $\eta$  if and only if  $D(\pi||\eta)$  is small
  - ▶ We want a reference that is as close to the target as possible
- ▶ **Solution 1:** Keep the target fixed and modify the reference (i.e. Variational inference)



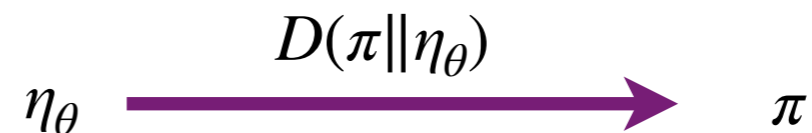
- ▶ Let  $\mathcal{Q} = \{\eta_\theta\}_{\theta \in \Theta}$  be a **variational** family of reference distribution
- ▶ **Goal:** Find the best reference in  $\mathcal{Q}$  that approximates the target

# VARIATIONAL INFERENCE

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ The choice of reference is crucial
  - ▶ We can efficiently estimate  $\pi$  using the reference  $\eta$  if and only if  $D(\pi||\eta)$  is small
  - ▶ We want a reference that is as close to the target as possible
- ▶ **Solution 1:** Keep the target fixed and modify the reference (i.e. Variational inference)



- ▶ Let  $\mathcal{Q} = \{\eta_\theta\}_{\theta \in \Theta}$  be a **variational** family of reference distribution
- ▶ **Goal:** Find the best reference in  $\mathcal{Q}$  that approximates the target

$$\eta \in \arg \min_{\theta} D(\pi||\eta_\theta)$$

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi \| \eta_{\theta})$$

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi \| \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi || \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi || \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn



# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi \| \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi || \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over
- ▶ **Challenge 2:** How do you choose the discrepancy  $D$ ?

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi || \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over
- ▶ **Challenge 2:** How do you choose the discrepancy  $D$ ?
  - ▶ Eg. Total variation, forward/reverse KL, Wasserstein, MMD,  $\alpha$ -divergence,  $f$ -divergence, etc

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi \| \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over
- ▶ **Challenge 2:** How do you choose the discrepancy  $D$ ?
  - ▶ Eg. Total variation, forward/reverse KL, Wasserstein, MMD,  $\alpha$ -divergence,  $f$ -divergence, etc
  - ▶ Each have their own pros and cons and measure different things

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi \| \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over
- ▶ **Challenge 2:** How do you choose the discrepancy  $D$ ?
  - ▶ Eg. Total variation, forward/reverse KL, Wasserstein, MMD,  $\alpha$ -divergence,  $f$ -divergence, etc
  - ▶ Each have their own pros and cons and measure different things
- ▶ **Challenge 3:** In general, the optimisations can be fragile with limited target information

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi || \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over
- ▶ **Challenge 2:** How do you choose the discrepancy  $D$ ?
  - ▶ Eg. Total variation, forward/reverse KL, Wasserstein, MMD,  $\alpha$ -divergence,  $f$ -divergence, etc
  - ▶ Each have their own pros and cons and measure different things
- ▶ **Challenge 3:** In general, the optimisations can be fragile with limited target information
  - ▶ The more flexible the harder it is to optimise

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi || \eta_{\theta})$$

- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over
- ▶ **Challenge 2:** How do you choose the discrepancy  $D$ ?
  - ▶ Eg. Total variation, forward/reverse KL, Wasserstein, MMD,  $\alpha$ -divergence,  $f$ -divergence, etc
  - ▶ Each have their own pros and cons and measure different things
- ▶ **Challenge 3:** In general, the optimisations can be fragile with limited target information
  - ▶ The more flexible the harder it is to optimise
  - ▶ Design choices traditionally appeal towards stability over flexibility

# CHALLENGES WITH VARIATIONAL INFERENCE

$$\eta \in \arg \min_{\theta} D(\pi || \eta_{\theta})$$

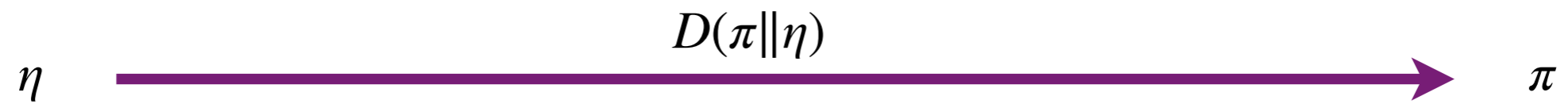
- ▶ **Challenge 1:** How do you design the variational family?
  - ▶ Extremely problem-specific
  - ▶ There is a tradeoff between how flexible your family is and how easy it is to learn
  - ▶ The more flexible the harder it is to optimise over
- ▶ **Challenge 2:** How do you choose the discrepancy  $D$ ?
  - ▶ Eg. Total variation, forward/reverse KL, Wasserstein, MMD,  $\alpha$ -divergence,  $f$ -divergence, etc
  - ▶ Each have their own pros and cons and measure different things
- ▶ **Challenge 3:** In general, the optimisations can be fragile with limited target information
  - ▶ The more flexible the harder it is to optimise
  - ▶ Design choices traditionally appeal towards stability over flexibility
- ▶ **Challenge 4:** even the best reference may not be close enough to stabilise the acceptance



# ANNEALING

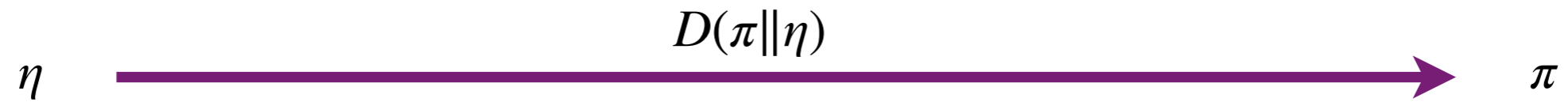
# ANNEALING

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



# ANNEALING

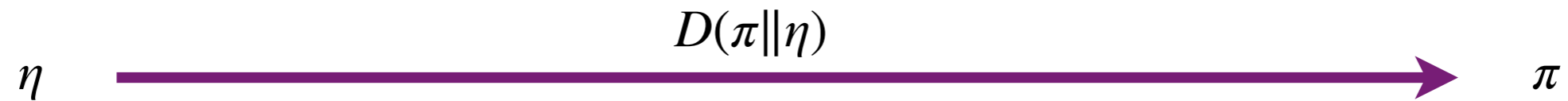
- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ **Solution 2:** Keep the reference fixed and modify the target

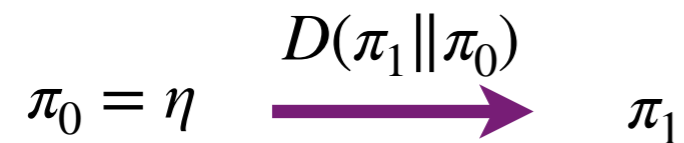
# ANNEALING

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



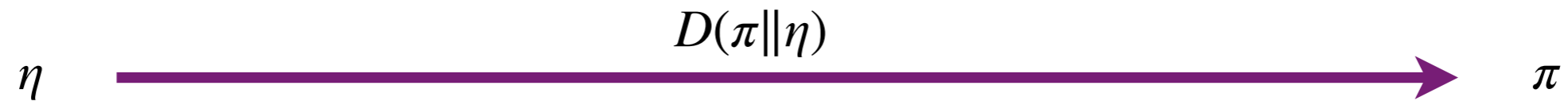
- ▶ **Solution 2:** Keep the reference fixed and modify the target

- ▶ We can efficiently propagate inferences from  $\eta = \pi_0$  to  $\pi_1$  if  $D(\pi_1||\pi_0)$  is small



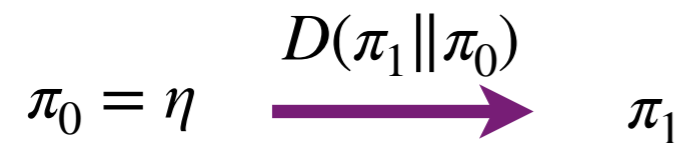
# ANNEALING

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ **Solution 2:** Keep the reference fixed and modify the target

- ▶ We can efficiently propagate inferences from  $\eta = \pi_0$  to  $\pi_1$  if  $D(\pi_1||\pi_0)$  is small

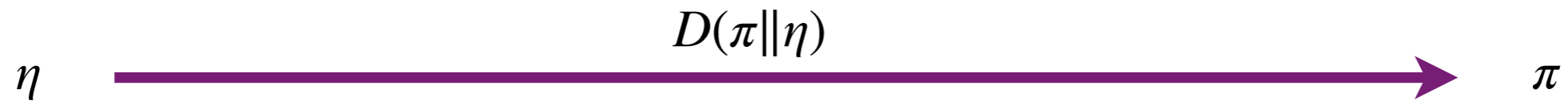


- ▶ We can then efficiently propagate inference from  $\pi_1$  to  $\pi_2$  close to  $\pi_1$



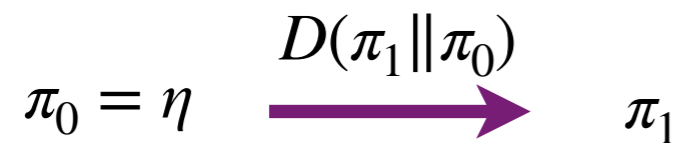
# ANNEALING

- ▶ Given a divergence  $D$  and a reference  $\eta$  and target  $\pi$



- ▶ **Solution 2:** Keep the reference fixed and modify the target

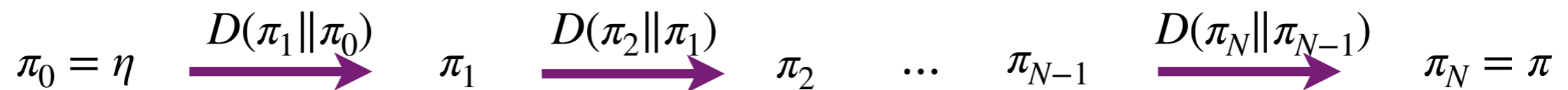
- ▶ We can efficiently propagate inferences from  $\eta = \pi_0$  to  $\pi_1$  if  $D(\pi_1||\pi_0)$  is small



- ▶ We can then efficiently propagate inference from  $\pi_1$  to  $\pi_2$  close to  $\pi_1$



- ▶ Can repeat until we reach the target

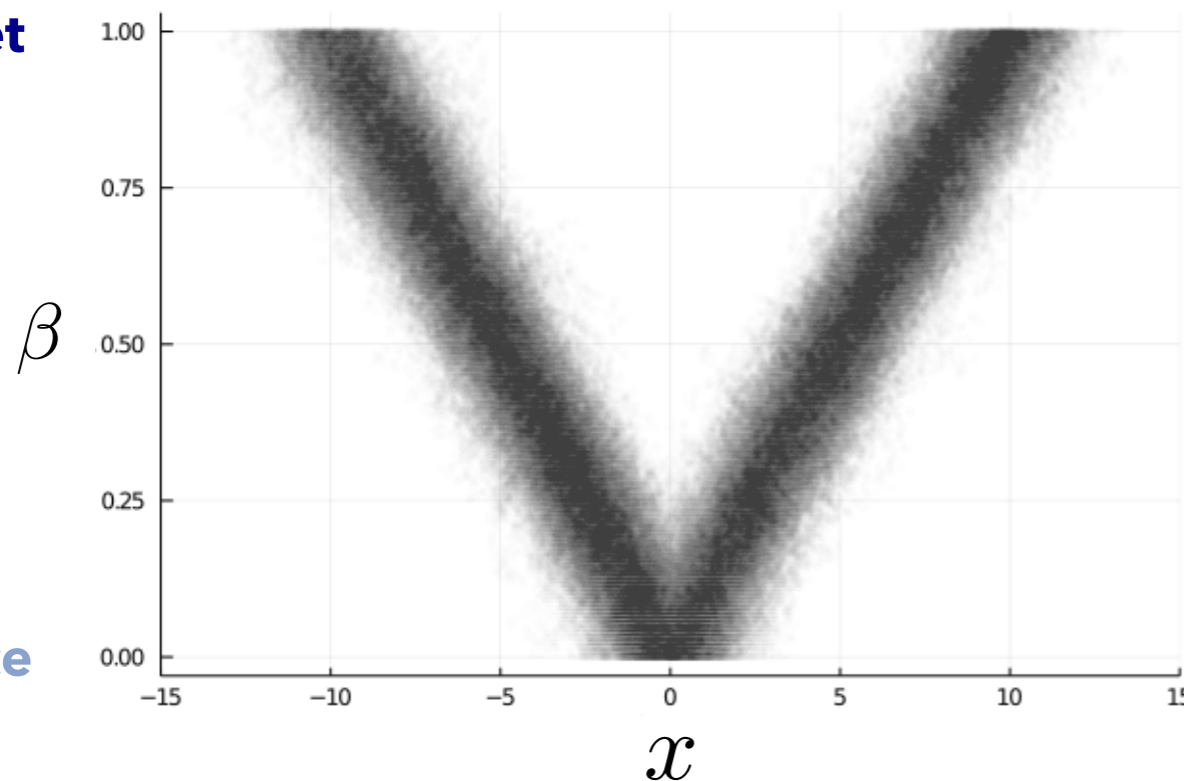


# ANNEALING

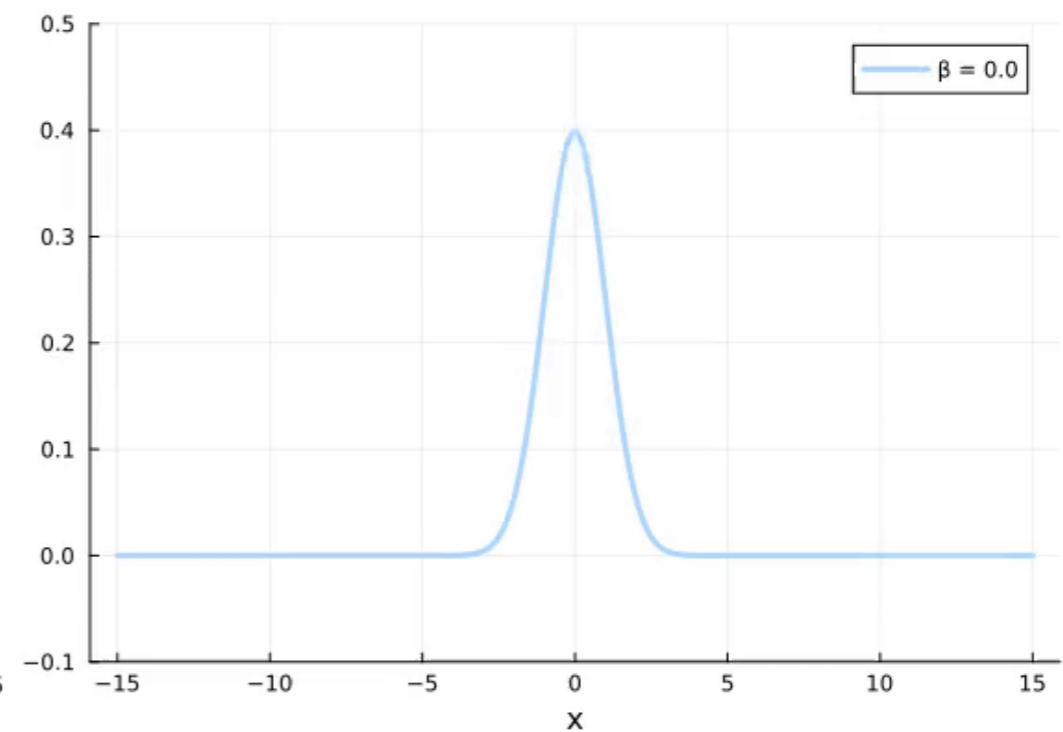
# ANNEALING

- ▶ For  $\beta \in [0,1]$  suppose we have distributions  $\pi_\beta$  such that  $\pi_0 = \eta$  and  $\pi_1 = \pi$ 
  - ▶  $\pi_\beta$  is the **annealing distribution** corresponding to the **annealing parameter**  $\beta$

Target



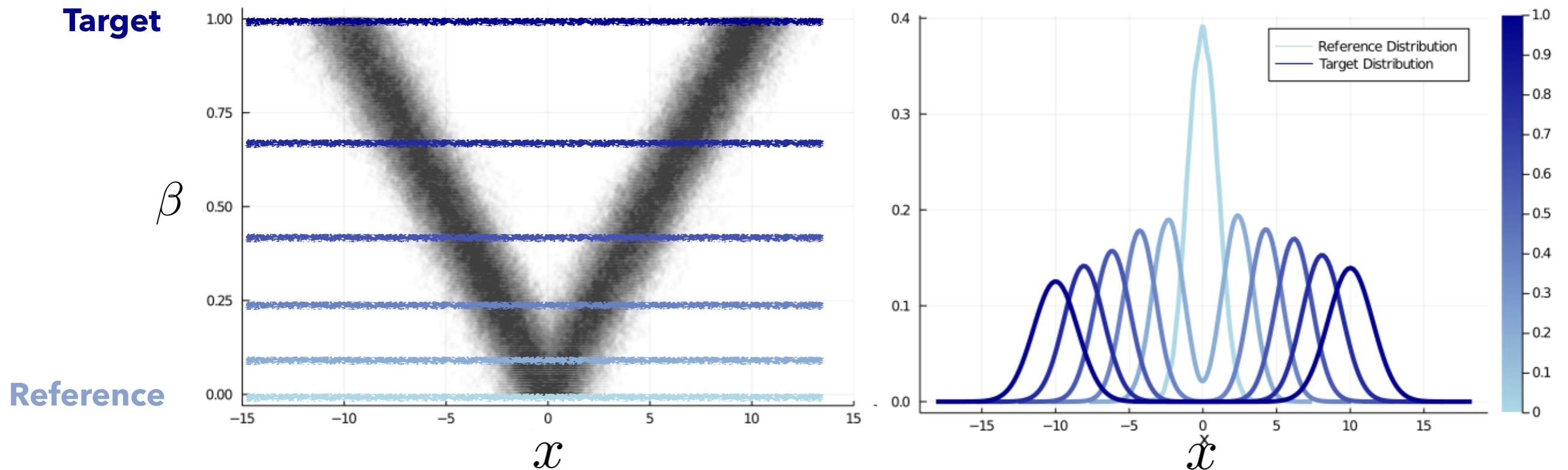
Reference





# ANNEALING

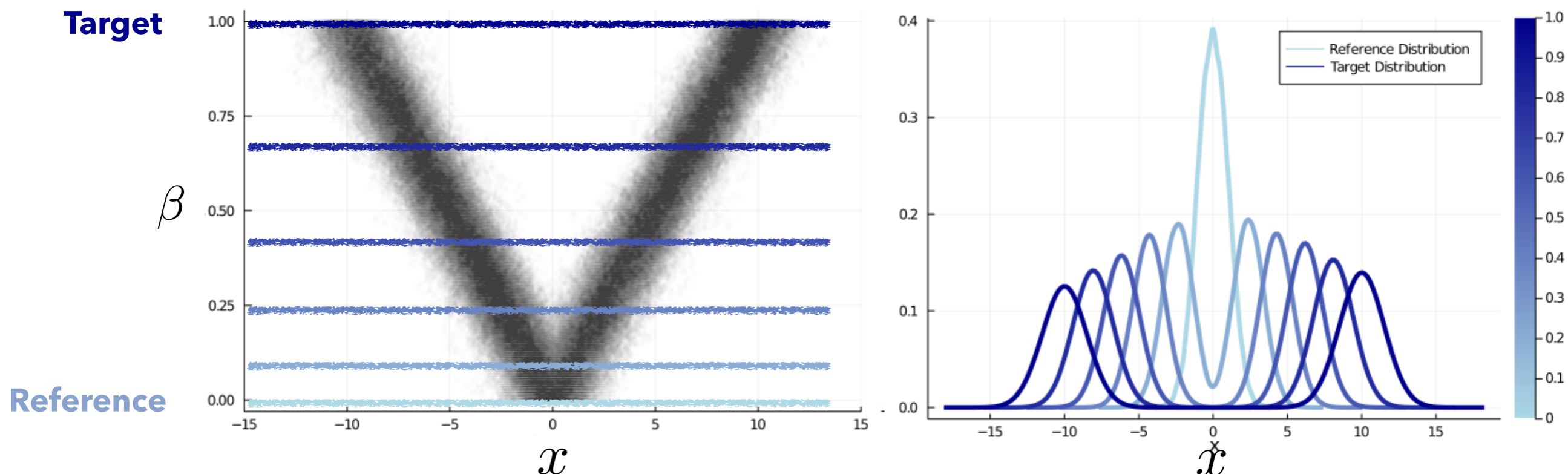
- ▶ For  $\beta \in [0,1]$  suppose we have distributions  $\pi_\beta$  such that  $\pi_0 = \eta$  and  $\pi_1 = \pi$ 
  - ▶  $\pi_\beta$  is the **annealing distribution** corresponding to the **annealing parameter**  $\beta$
- ▶ We can construct our sequence of annealing distributions by setting  $\pi_n = \pi_{\beta_n}$



# ANNEALING

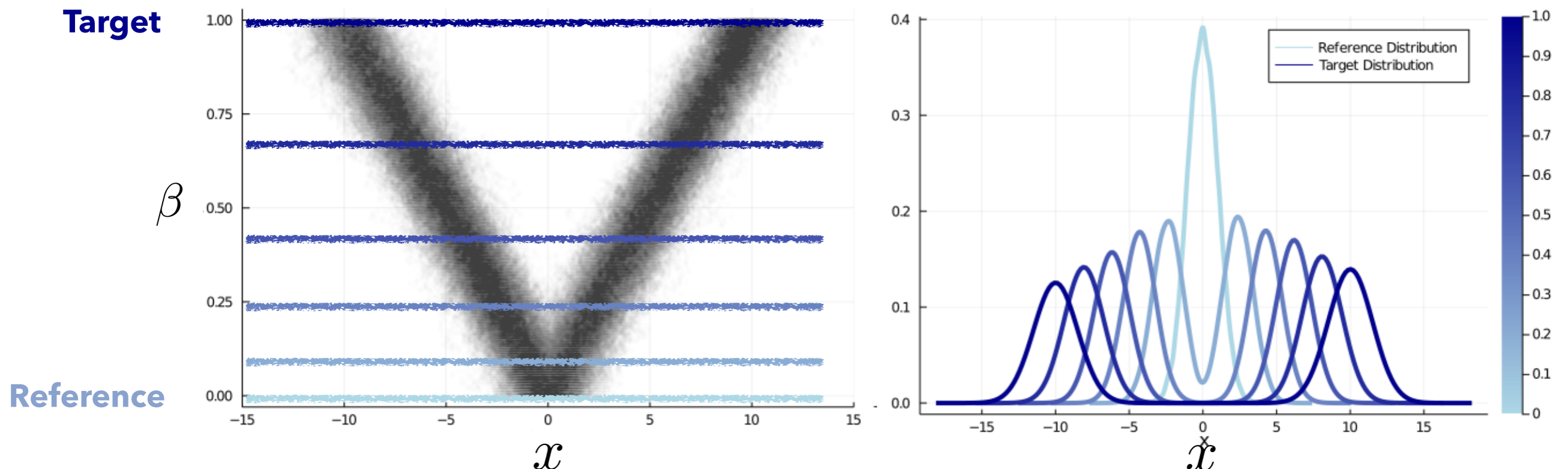
- ▶ For  $\beta \in [0,1]$  suppose we have distributions  $\pi_\beta$  such that  $\pi_0 = \eta$  and  $\pi_1 = \pi$ 
  - ▶  $\pi_\beta$  is the **annealing distribution** corresponding to the **annealing parameter**  $\beta$
- ▶ We can construct our sequence of annealing distributions by setting  $\pi_n = \pi_{\beta_n}$ 
  - ▶ Where  $\mathcal{B} = \beta_{0:N}$  is the **annealing schedule** satisfying:

$$0 = \beta_0 < \beta_1 < \dots < \beta_N = 1$$



# ANNEALING

- ▶ For  $\beta \in [0,1]$  suppose we have distributions  $\pi_\beta$  such that  $\pi_0 = \eta$  and  $\pi_1 = \pi$ 
  - ▶  $\pi_\beta$  is the **annealing distribution** corresponding to the **annealing parameter**  $\beta$
- ▶ We can construct our sequence of annealing distributions by setting  $\pi_n = \pi_{\beta_n}$ 
  - ▶ Where  $\mathcal{B} = \beta_{0:N}$  is the **annealing schedule** satisfying:
$$0 = \beta_0 < \beta_1 < \dots < \beta_N = 1$$
- ▶ Note that  $\pi_{n-1}$  and  $\pi_n$  are close the schedule if fine and  $\beta, \beta' \mapsto D(\pi_{\beta'} || \pi_\beta)$  is continuous

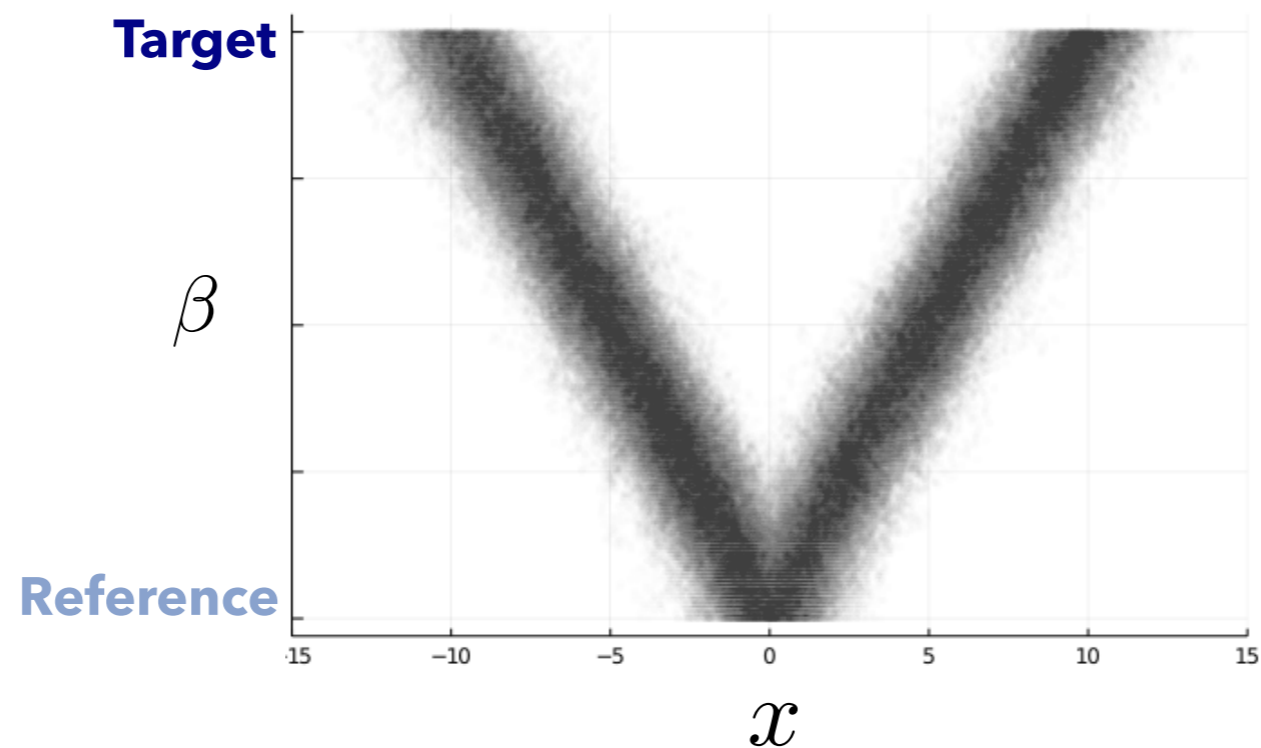


# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of

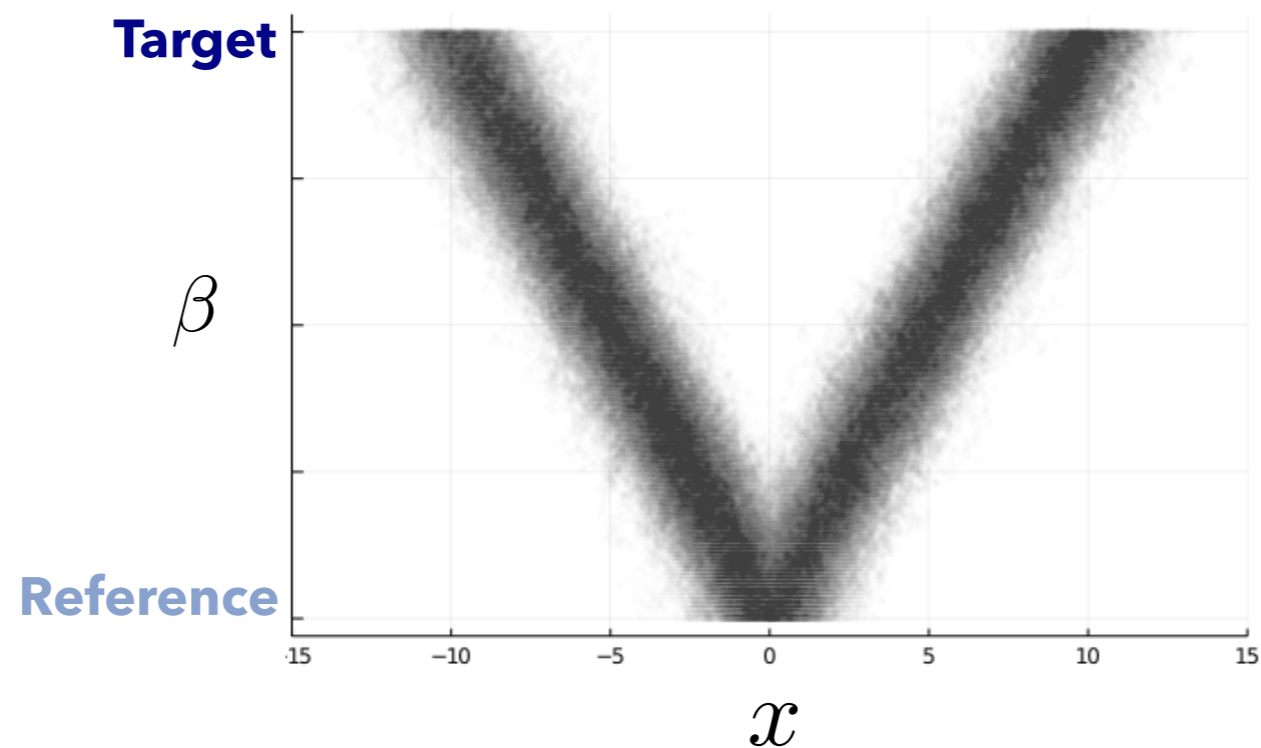
# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of



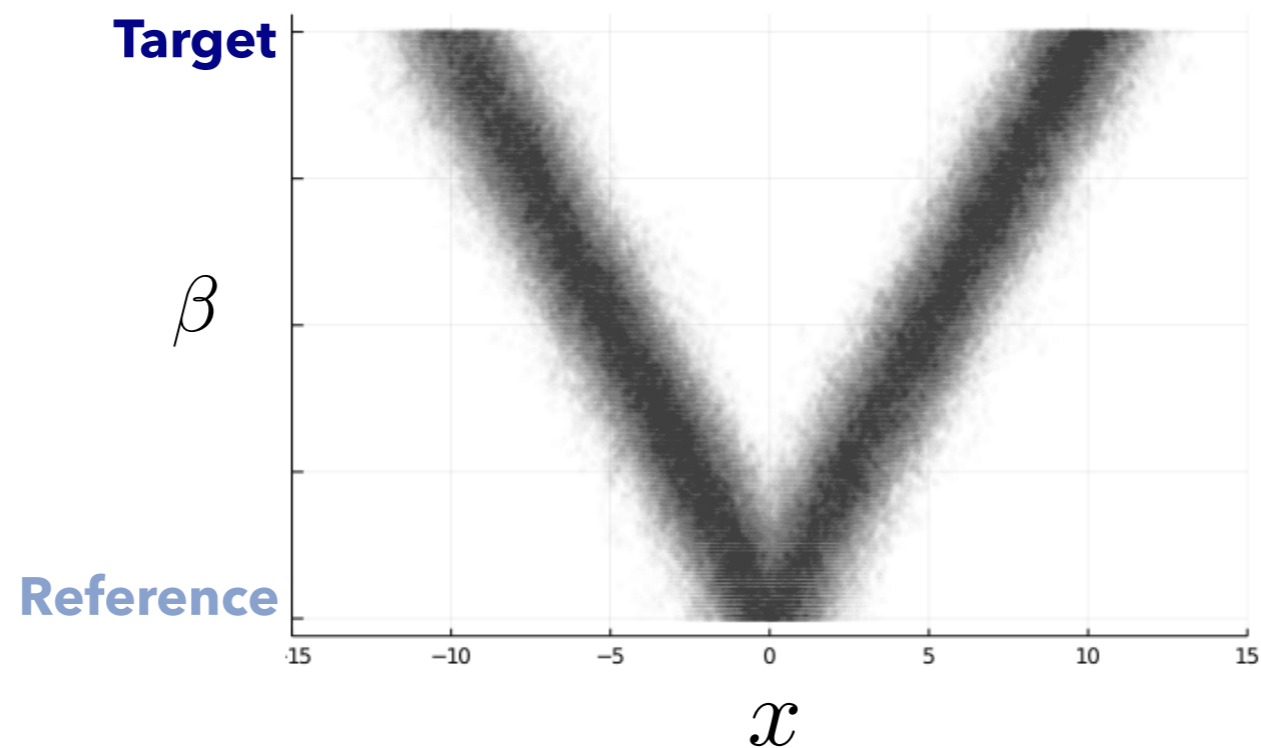
# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of
- ▶ **Annealing algorithms** draw inference from the entire path not just the target



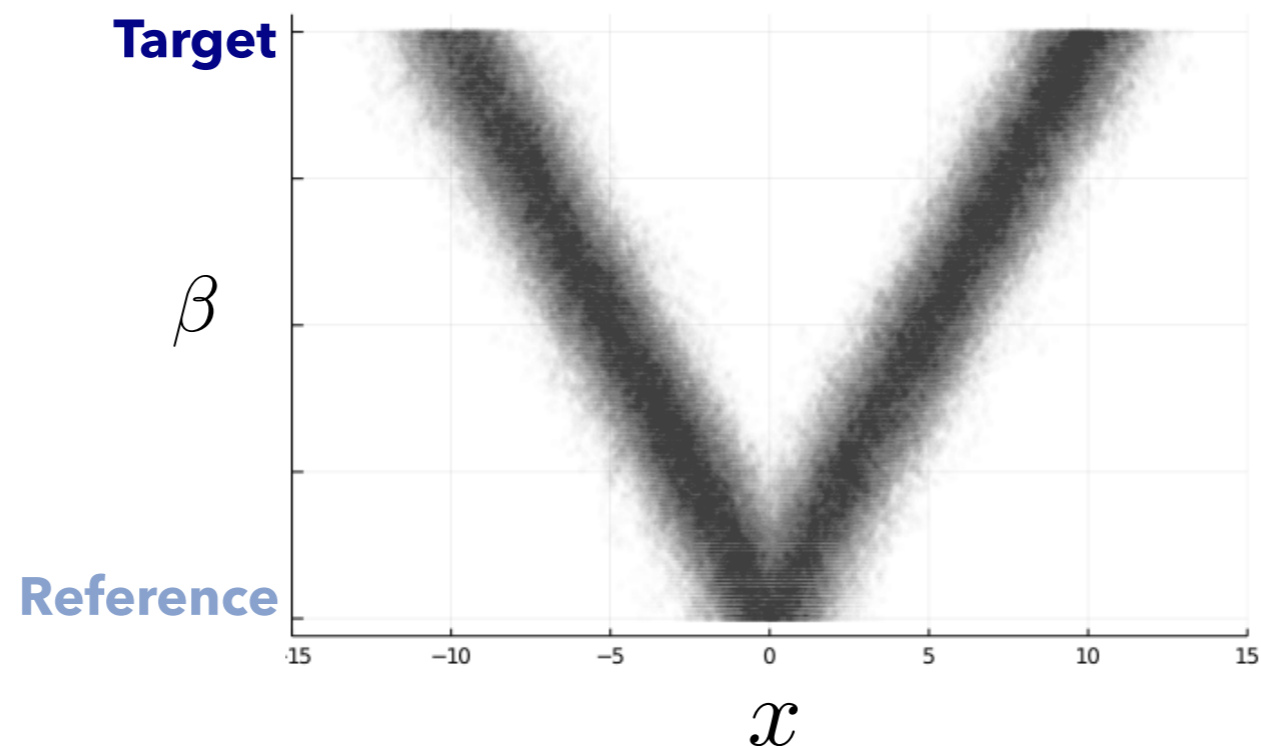
# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of
- ▶ **Annealing algorithms** draw inference from the entire path not just the target
- ▶ Annealing algorithms are meta algorithms:



# ANNEALING ALGORITHMS

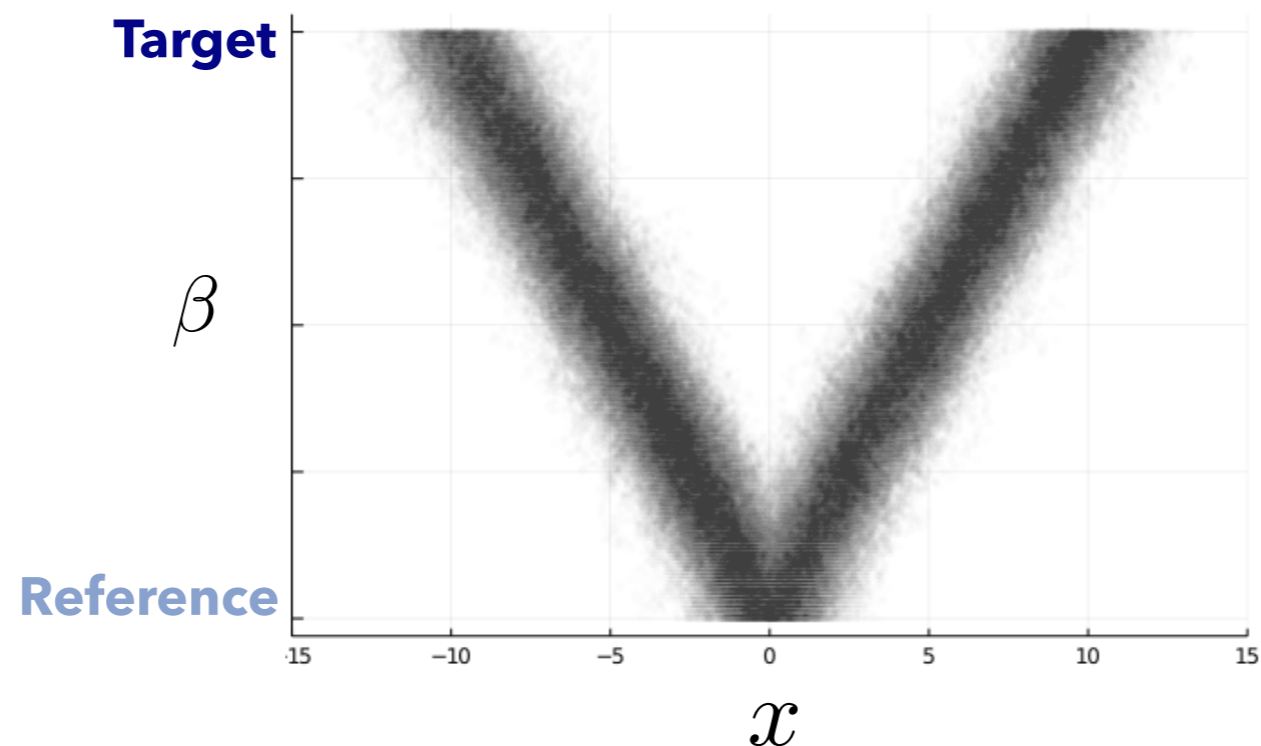
- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of
- ▶ **Annealing algorithms** draw inference from the entire path not just the target
- ▶ Annealing algorithms are meta algorithms:
  - ▶ **Input:** An annealing path  $\pi_\beta$  + a local efficient inference algorithm for each  $\pi_\beta$





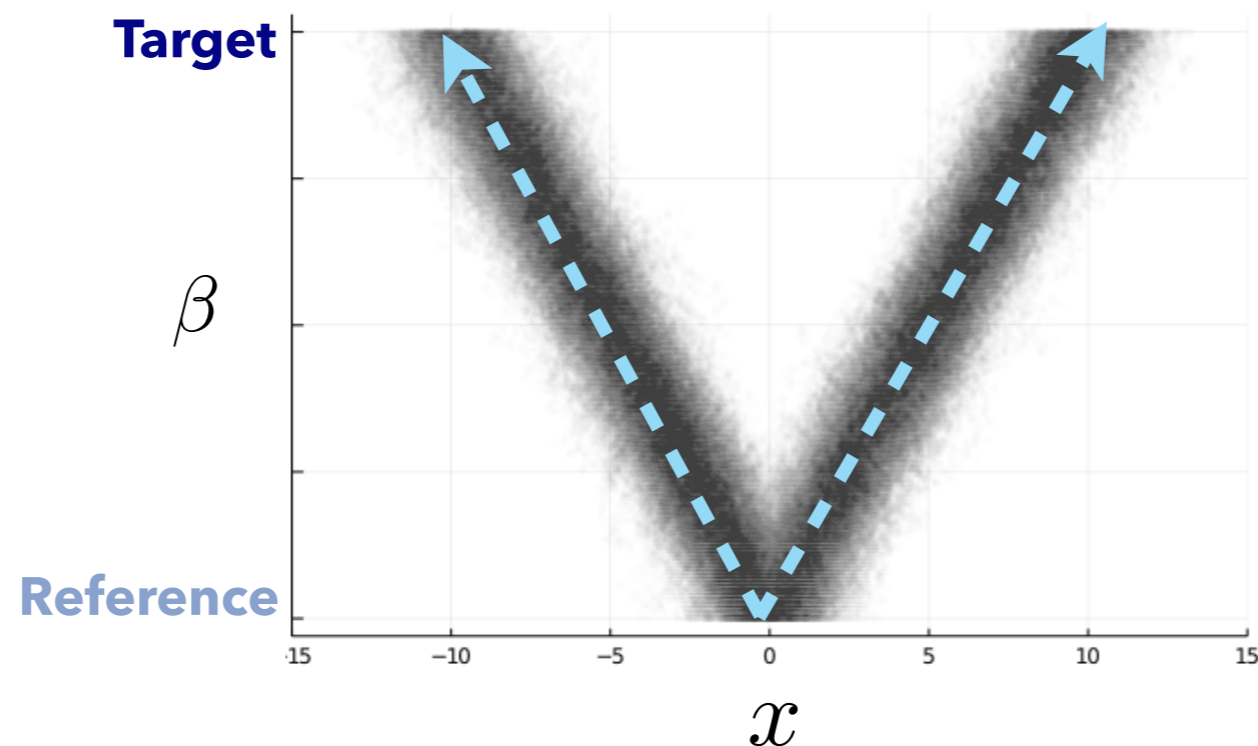
# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of
- ▶ **Annealing algorithms** draw inference from the entire path not just the target
- ▶ Annealing algorithms are meta algorithms:
  - ▶ **Input:** An annealing path  $\pi_\beta$  + a local efficient inference algorithm for each  $\pi_\beta$
  - ▶ **Output:** a globally efficient inference algorithm for the path  $\beta \mapsto \pi_\beta$



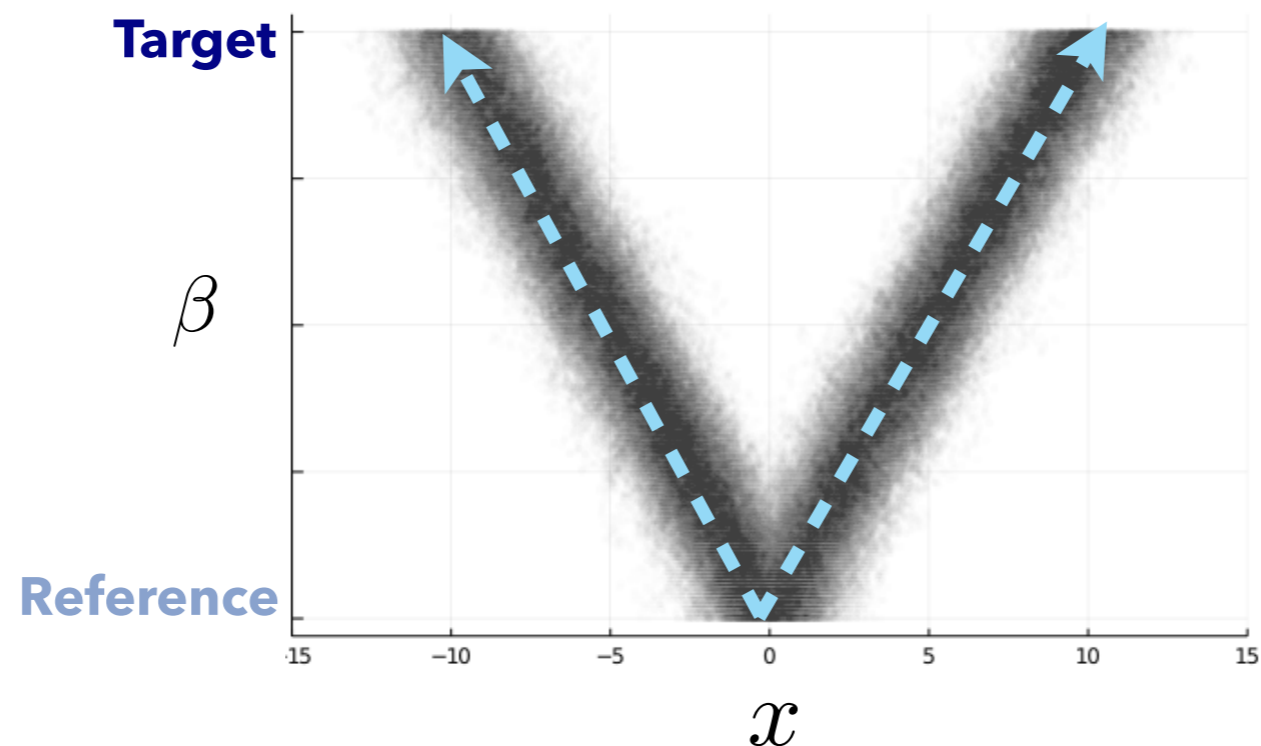
# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of
- ▶ **Annealing algorithms** draw inference from the entire path not just the target
- ▶ Annealing algorithms are meta algorithms:
  - ▶ **Input:** An annealing path  $\pi_\beta$  + a local efficient inference algorithm for each  $\pi_\beta$
  - ▶ **Output:** a globally efficient inference algorithm for the path  $\beta \mapsto \pi_\beta$
- ▶ Transform a  $d$ -dimensional multi-modal target into a  $d + 1$ -dimensional unimodal one



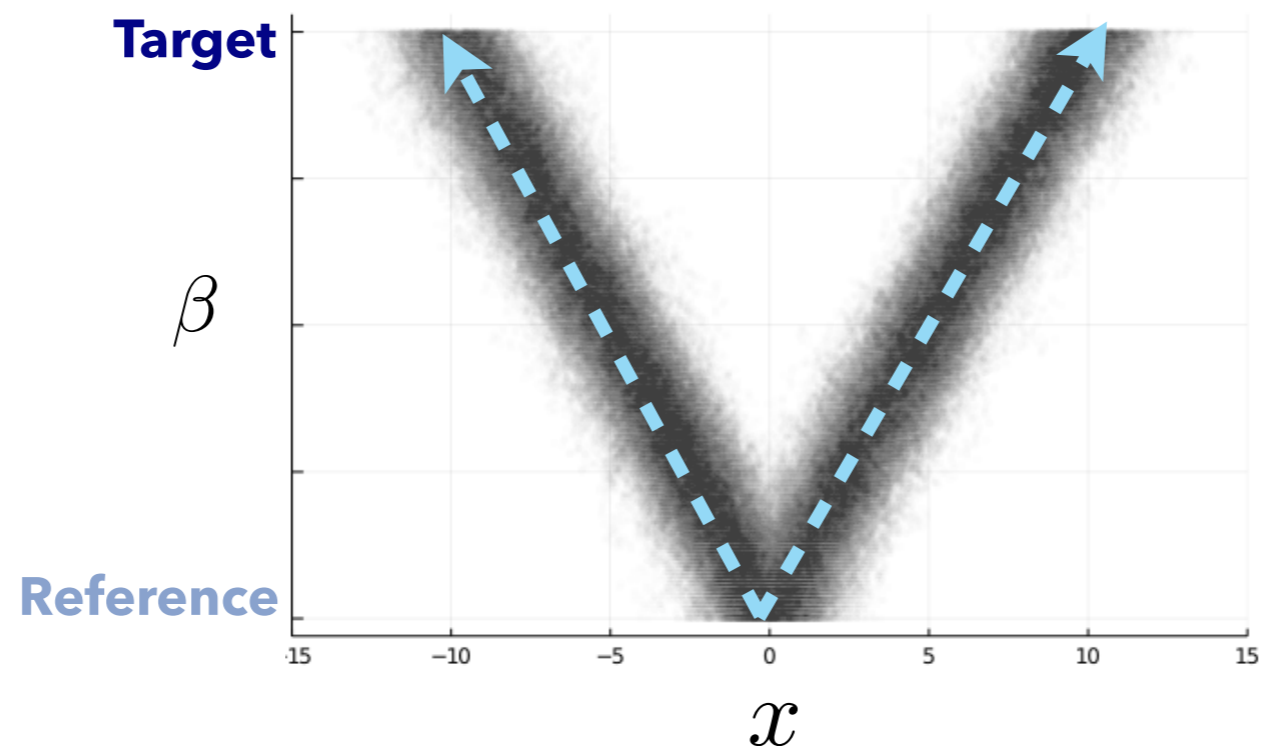
# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of
- ▶ **Annealing algorithms** draw inference from the entire path not just the target
- ▶ Annealing algorithms are meta algorithms:
  - ▶ **Input:** An annealing path  $\pi_\beta$  + a local efficient inference algorithm for each  $\pi_\beta$
  - ▶ **Output:** a globally efficient inference algorithm for the path  $\beta \mapsto \pi_\beta$
- ▶ Transform a  $d$ -dimensional multi-modal target into a  $d + 1$ -dimensional unimodal one
- ▶ They achieve state-of-the-art performance for hard multi-modal distributions



# ANNEALING ALGORITHMS

- ▶ Annealing is **not** an algorithm but a platform to build algorithms on top of
- ▶ **Annealing algorithms** draw inference from the entire path not just the target
- ▶ Annealing algorithms are meta algorithms:
  - ▶ **Input:** An annealing path  $\pi_\beta$  + a local efficient inference algorithm for each  $\pi_\beta$
  - ▶ **Output:** a globally efficient inference algorithm for the path  $\beta \mapsto \pi_\beta$
- ▶ Transform a  $d$ -dimensional multi-modal target into a  $d + 1$ -dimensional unimodal one
- ▶ They achieve state-of-the-art performance for hard multi-modal distributions
  - ▶ Can use the reference to assess the quality and ensure robustness



# HOW TO CONSTRUCT AN ANNEALING PATH

# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$

# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$
- ▶ In the context of sampling, we are given unnormalised densities:

$$\pi(\mathrm{d}x) = \frac{\gamma(x)}{Z} \mathrm{d}x$$

# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$
- ▶ In the context of sampling, we are given unnormalised densities:

$$\pi(dx) = \frac{\gamma(x)}{Z} dx$$

- ▶ We build annealing distribution by corrupting the densities:

$$\pi_\beta = \frac{\gamma_\beta(x)}{Z(\beta)}, \quad Z(\beta) = \int_{\mathbb{X}} \gamma_\beta(x) dx$$



# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$
- ▶ In the context of sampling, we are given unnormalised densities:

$$\pi(\mathrm{d}x) = \frac{\gamma(x)}{Z} \mathrm{d}x$$

- ▶ We build annealing distribution by corrupting the densities:

$$\pi_\beta = \frac{\gamma_\beta(x)}{Z(\beta)}, \quad Z(\beta) = \int_{\mathbb{X}} \gamma_\beta(x) \mathrm{d}x$$

- ▶  $\gamma_\beta : \mathbb{X} \rightarrow \mathbb{R}$  is the un-normalised density satisfying  $\gamma_0 = \eta$  and  $\gamma_1 = \gamma$

# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$
- ▶ In the context of sampling, we are given unnormalised densities:

$$\pi(\mathrm{d}x) = \frac{\gamma(x)}{Z} \mathrm{d}x$$

- ▶ We build annealing distribution by corrupting the densities:

$$\pi_\beta = \frac{\gamma_\beta(x)}{Z(\beta)}, \quad Z(\beta) = \int_{\mathbb{X}} \gamma_\beta(x) \mathrm{d}x$$

- ▶  $\gamma_\beta : \mathbb{X} \rightarrow \mathbb{R}$  is the un-normalised density satisfying  $\gamma_0 = \eta$  and  $\gamma_1 = \gamma$
- ▶ Normalizing constants  $Z(\beta)$  satisfies  $Z(0) = 1$  and  $Z(1) = Z$

# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$
- ▶ In the context of sampling, we are given unnormalised densities:

$$\pi(\mathrm{d}x) = \frac{\gamma(x)}{Z} \mathrm{d}x$$

- ▶ We build annealing distribution by corrupting the densities:

$$\pi_\beta = \frac{\gamma_\beta(x)}{Z(\beta)}, \quad Z(\beta) = \int_{\mathbb{X}} \gamma_\beta(x) \mathrm{d}x$$

- ▶  $\gamma_\beta : \mathbb{X} \rightarrow \mathbb{R}$  is the un-normalised density satisfying  $\gamma_0 = \eta$  and  $\gamma_1 = \gamma$
- ▶ Normalizing constants  $Z(\beta)$  satisfies  $Z(0) = 1$  and  $Z(1) = Z$
- ▶ There is a lot of flexibility in choosing a path. We will study what makes for a good path later.

# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$
- ▶ In the context of sampling, we are given unnormalised densities:

$$\pi(\mathrm{d}x) = \frac{\gamma(x)}{Z} \mathrm{d}x$$

- ▶ We build annealing distribution by corrupting the densities:

$$\pi_\beta = \frac{\gamma_\beta(x)}{Z(\beta)}, \quad Z(\beta) = \int_{\mathbb{X}} \gamma_\beta(x) \mathrm{d}x$$

- ▶  $\gamma_\beta : \mathbb{X} \rightarrow \mathbb{R}$  is the un-normalised density satisfying  $\gamma_0 = \eta$  and  $\gamma_1 = \gamma$
- ▶ Normalizing constants  $Z(\beta)$  satisfies  $Z(0) = 1$  and  $Z(1) = Z$
- ▶ There is a lot of flexibility in choosing a path. We will study what makes for a good path later.
- ▶ The canonical choice is the linear path

$$\gamma_\beta(x) \propto \eta(x)^{1-\beta} \gamma(x)^\beta$$

# HOW TO CONSTRUCT AN ANNEALING PATH

- ▶ Can think of the annealing path as a continuous corruption of the target  $\pi$  to noise  $\eta$
- ▶ In the context of sampling, we are given unnormalised densities:

$$\pi(\mathrm{d}x) = \frac{\gamma(x)}{Z} \mathrm{d}x$$

- ▶ We build annealing distribution by corrupting the densities:

$$\pi_\beta = \frac{\gamma_\beta(x)}{Z(\beta)}, \quad Z(\beta) = \int_{\mathbb{X}} \gamma_\beta(x) \mathrm{d}x$$

- ▶  $\gamma_\beta : \mathbb{X} \rightarrow \mathbb{R}$  is the un-normalised density satisfying  $\gamma_0 = \eta$  and  $\gamma_1 = \gamma$
- ▶ Normalizing constants  $Z(\beta)$  satisfies  $Z(0) = 1$  and  $Z(1) = Z$
- ▶ There is a lot of flexibility in choosing a path. We will study what makes for a good path later.
- ▶ The canonical choice is the linear path

$$\gamma_\beta(x) \propto \eta(x)^{1-\beta} \gamma(x)^\beta$$

- ▶ Before designing algorithms, let's build some intuition about annealing paths.