




# Sampling and MCMC Methods for Bayesian Inference

June, 2021



# Takeaways:

- If you have samples from a distribution, you know a lot about the distribution!
- It's not always straightforward to obtain samples, but there are techniques that can help
- There are accessible computational tools that can help
- We can think geometrically about MCMC sampling
- We can recognize some problems that can emerge from sampling

# Outline

Why sample?

Basic sampling algorithms

- Inverse cdf
- Rejection Sampling

MCMC Methods

- Idea and first example - Metropolis Algorithm
- Markov chains and stationary distributions
- Metropolis Hastings
- Gibbs Sampling
- Gradient-based sampling (HMC, NUTS)

Introduction to pymc3

Convergence diagnostics

# Samples - what are they good for?

Sometimes it is hard to compute a distribution of interest (like the posterior distribution  $p(\theta | y)$  in a non-conjugate Bayesian model BUT still possible to generate samples from that distribution

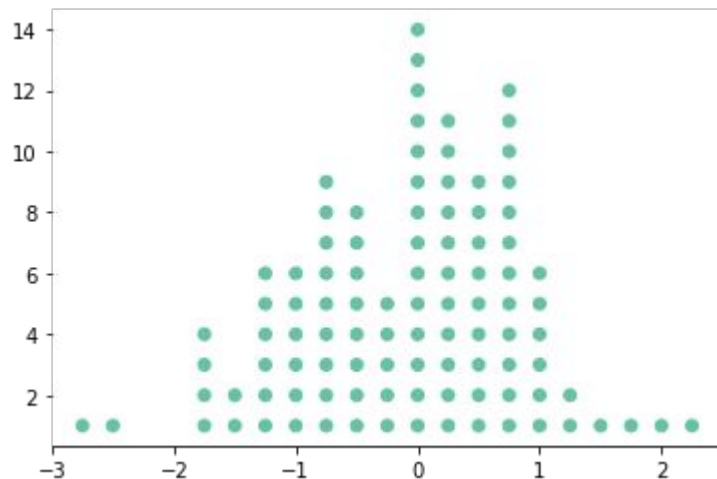
# Samples - what are they good for?

$$P(\theta > a)$$

$$\mathbb{E}[\theta]$$

$$\mathbb{E}[f(\theta)]$$

...

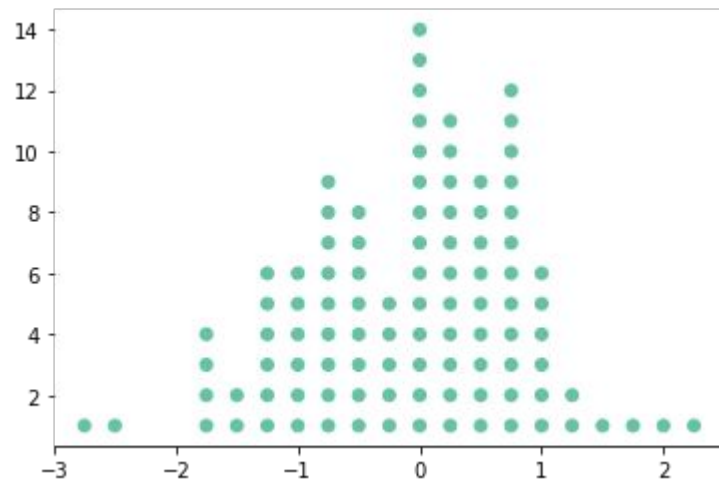


# Samples - what are they good for?

$$P(\theta > a) \approx \frac{\#\{\theta_i | \theta_i > a\}}{L}$$

$$\mathbb{E}[\theta] \approx \frac{1}{L} \sum_{i=0}^L \theta_i$$

$$\mathbb{E}[f(\theta)] \approx \frac{1}{L} \sum_{i=0}^L f(\theta_i)$$



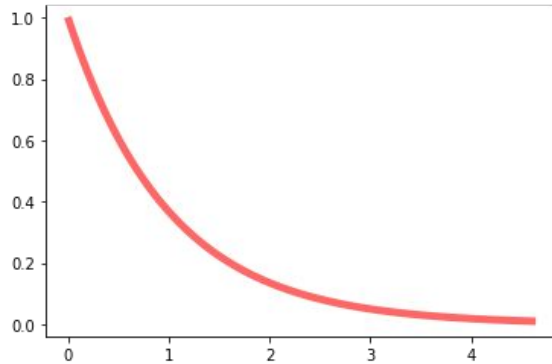
# So samples are useful.... How do we get them?

Many methods, from simple ones to various MCMC approaches.

We'll look at a bunch of methods in this workshop! Keep the big goal in mind.

# How would you sample from... an exponential distribution?

$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$

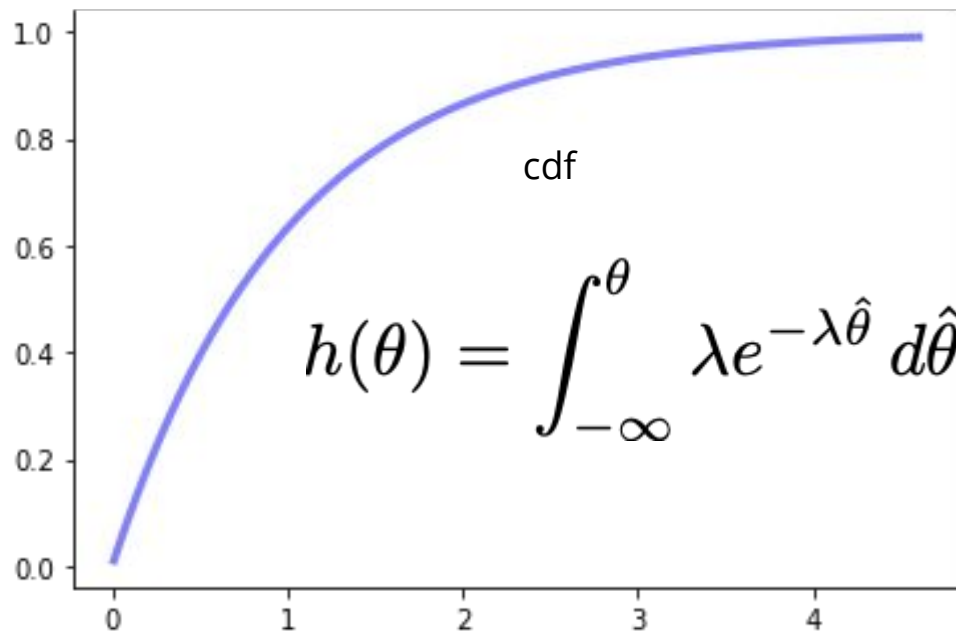
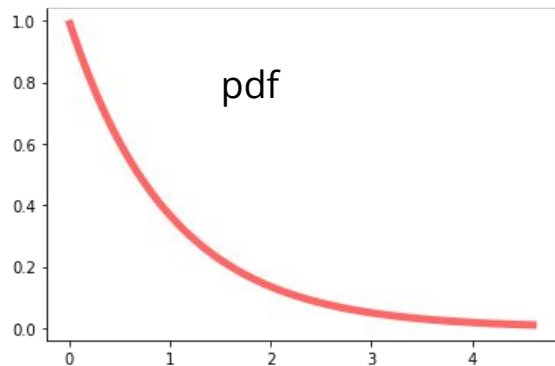


Assume we can generate (pseudo)-random samples from  $\text{Uniform}(0,1)$



# Sampling via the inverse cumulative density function

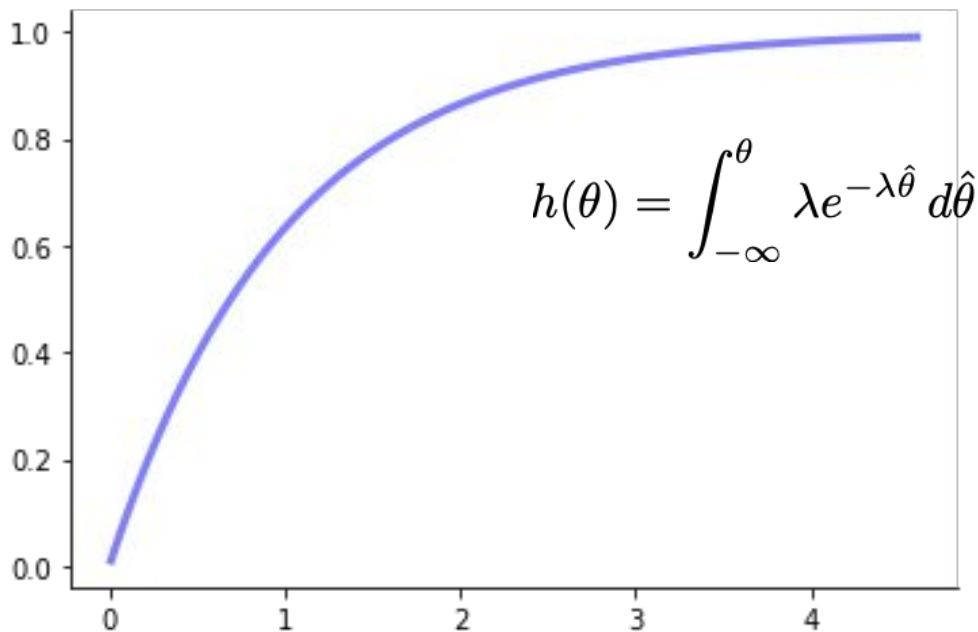
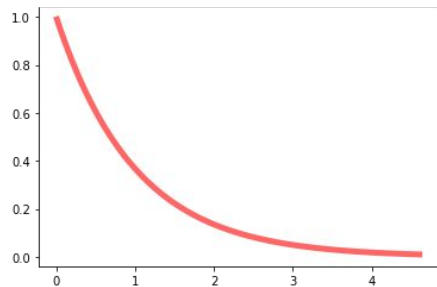
$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$



Assume we can generate (pseudo)-random samples from Uniform(0,1)

# Sampling via the inverse cumulative density function

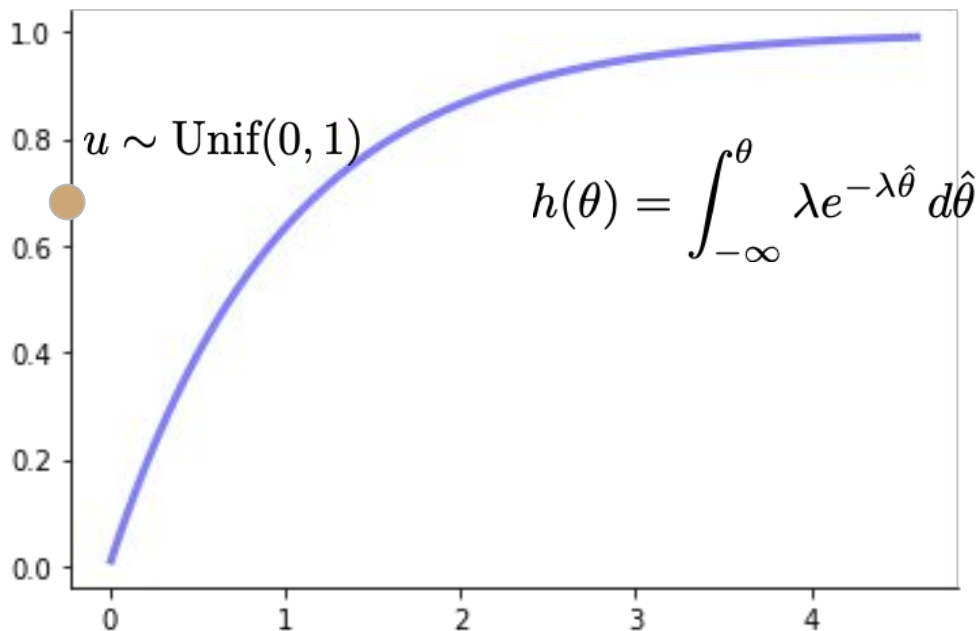
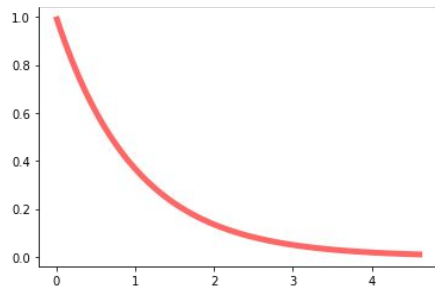
$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$



Assume we can generate (pseudo)-random samples from Uniform(0,1)

# Sampling via the inverse cumulative density function

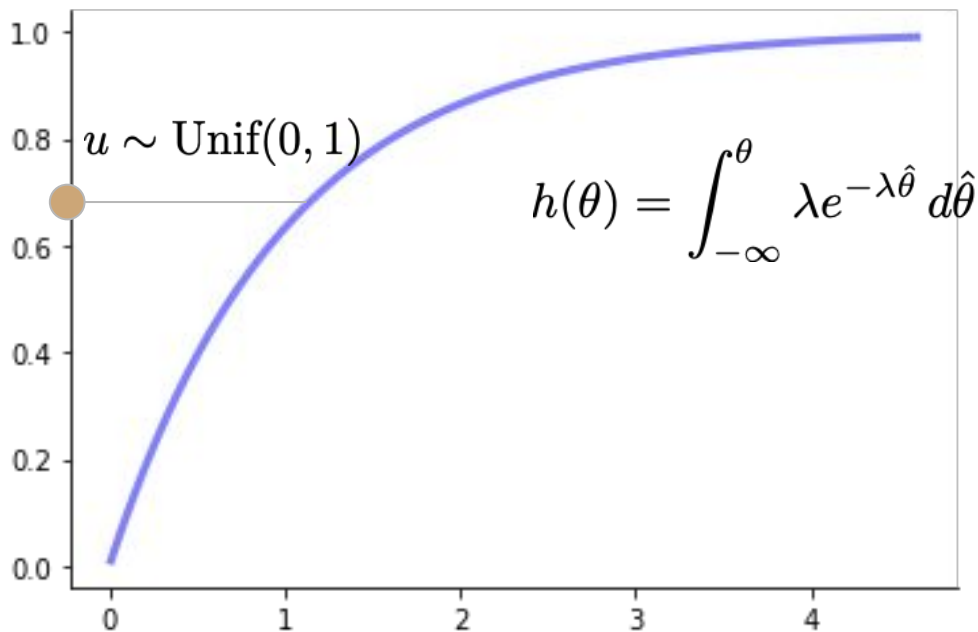
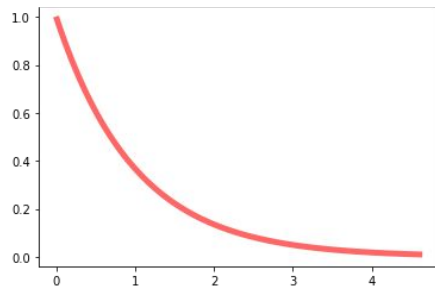
$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$



Assume we can generate (pseudo)-random samples from Uniform(0,1)

# Sampling via the inverse cumulative density function

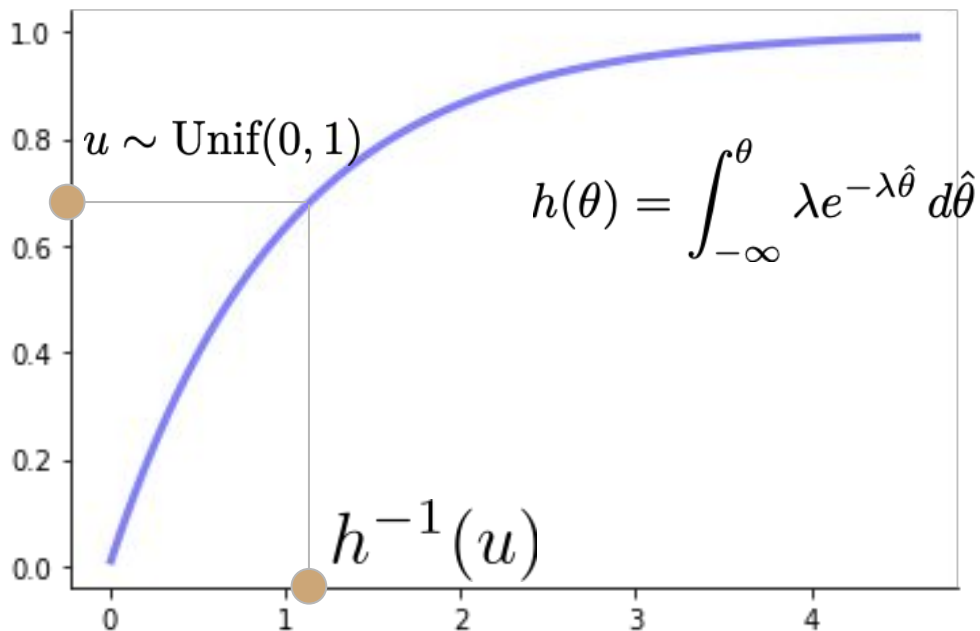
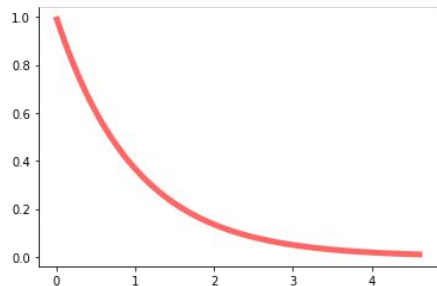
$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$



Assume we can generate (pseudo)-random samples from Uniform(0,1)

# Sampling via the inverse cumulative density function

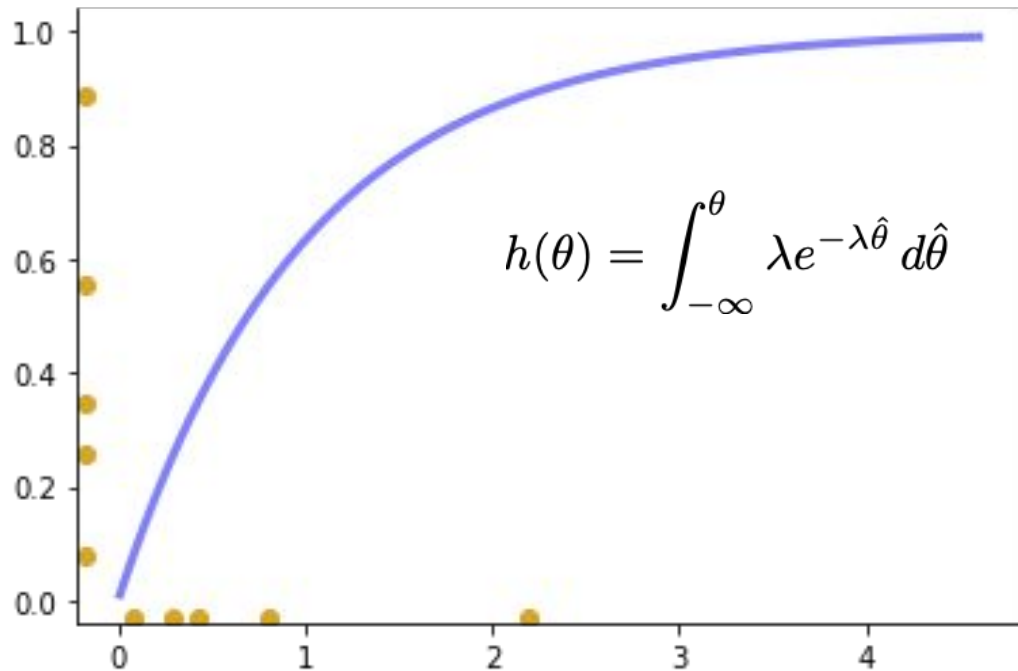
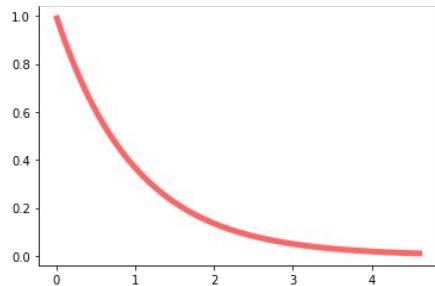
$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$



Assume we can generate (pseudo)-random samples from  $\text{Uniform}(0,1)$

# Sampling via the inverse cumulative density function

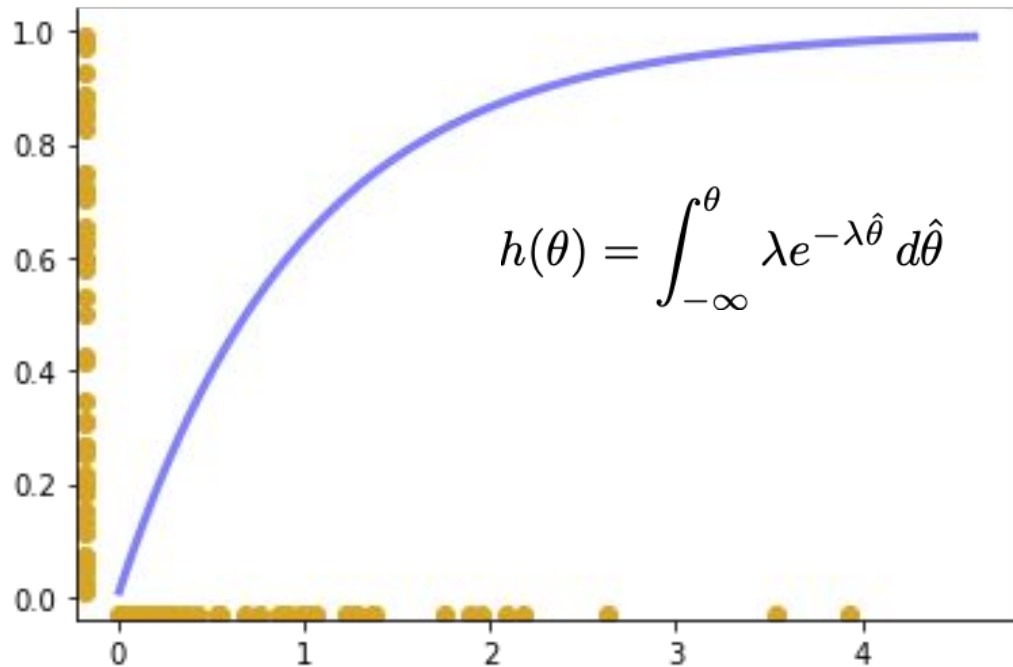
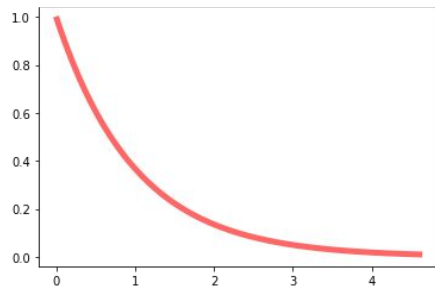
$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$



Assume we can generate (pseudo)-random samples from  $\text{Uniform}(0,1)$

# Sampling via the inverse cumulative density function

$$p(\theta) = \lambda e^{-\lambda\theta}, \theta > 0$$



Assume we can generate (pseudo)-random samples from  $\text{Uniform}(0,1)$

# Sampling via the inverse cumulative density function - summary:

Take random uniform samples  $u_i$

Calculate  $h^{-1}(u_i)$  where  $h(\theta) = \int_{-\infty}^{\theta} p(\hat{\theta}) d\hat{\theta}$ . These are the samples!

## For exponential distribution:

$$h(\theta) = \int_0^{\theta} \lambda e^{-\lambda \hat{\theta}} d\hat{\theta} = -e^{-\lambda \hat{\theta}} \Big|_0^{\theta} = 1 - e^{-\lambda \theta} \quad \text{cdf}$$

$$h^{-1}(u) = \lambda^{-1} \log(1 - u) \quad \text{Inverse cdf}$$



# Sampling via the inverse cumulative density function - justification

Take random uniform samples  $u_i$

Calculate  $h^{-1}(u_i)$  where  $h(\theta) = \int_{-\infty}^{\theta} p(\hat{\theta}) d\hat{\theta}$

# Sampling via the inverse cumulative density function - summary:

Take random uniform samples  $u_i$

Calculate  $h^{-1}(u_i)$  where  $h(\theta) = \int_{-\infty}^{\theta} p(\hat{\theta}) d\hat{\theta}$

But, we can't always calculate the inverse cdf!


# Rejection Sampling - set up

Suppose...

-Want to sample from a distribution with pdf  $p(z)$ ; don't have a simpler approach

-Can evaluate  $p(z)$  up to some normalizing constant  $Z_p$ . That, is can evaluate  $\tilde{p}(z)$ :

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

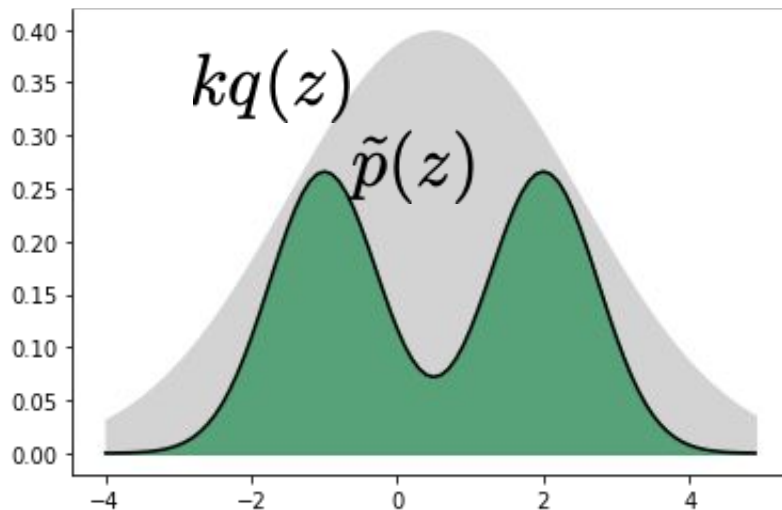
Recall:   $p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)} = \frac{p(y|\theta)p(\theta)}{\int p(y|\theta)p(\theta)}$  probably hard to evaluate

# Rejection Sampling - procedure

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Repeat until  
have sufficient  
samples

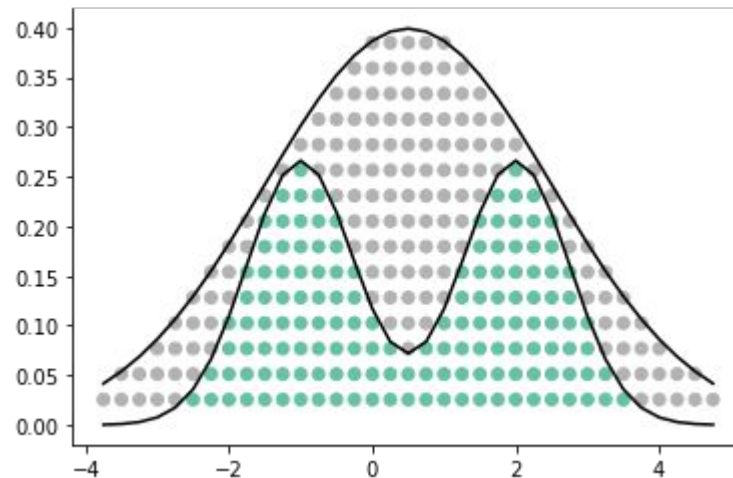
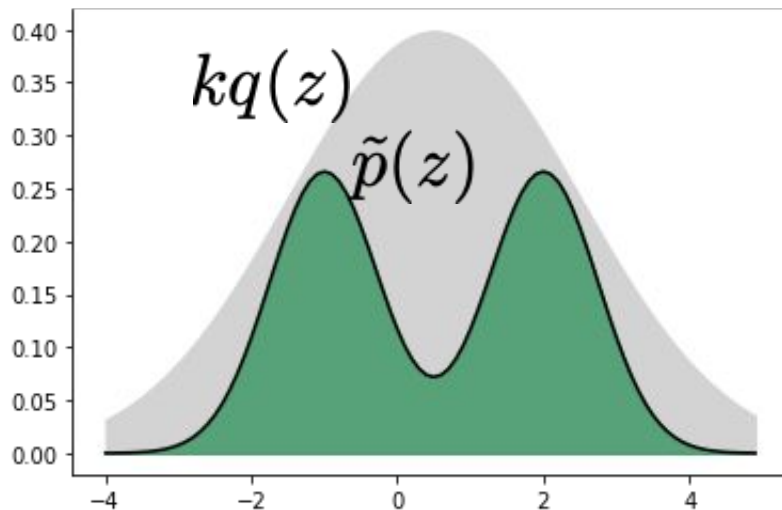


# Rejection Sampling - procedure

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

Repeat until  
have sufficient  
samples

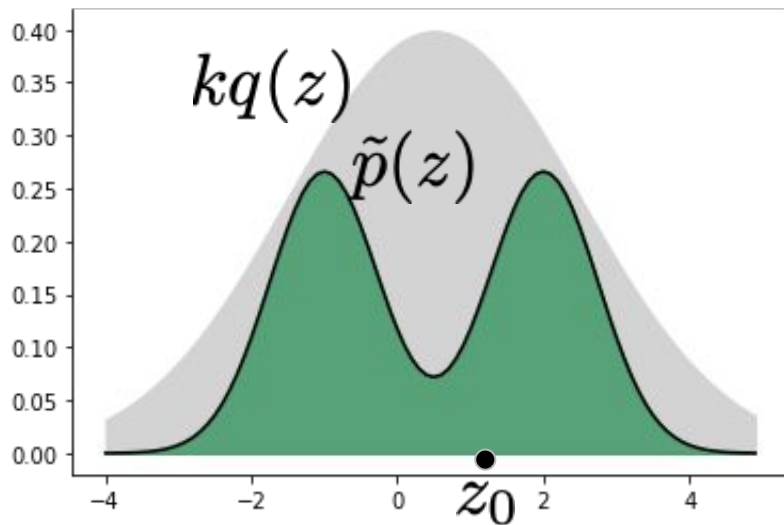


# Rejection Sampling - procedure

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Repeat until  
have sufficient  
samples

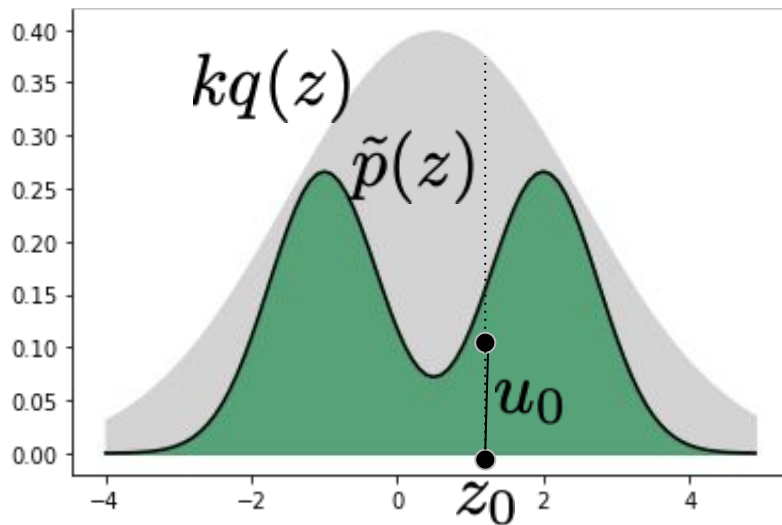


# Rejection Sampling - procedure

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Repeat until  
have sufficient  
samples

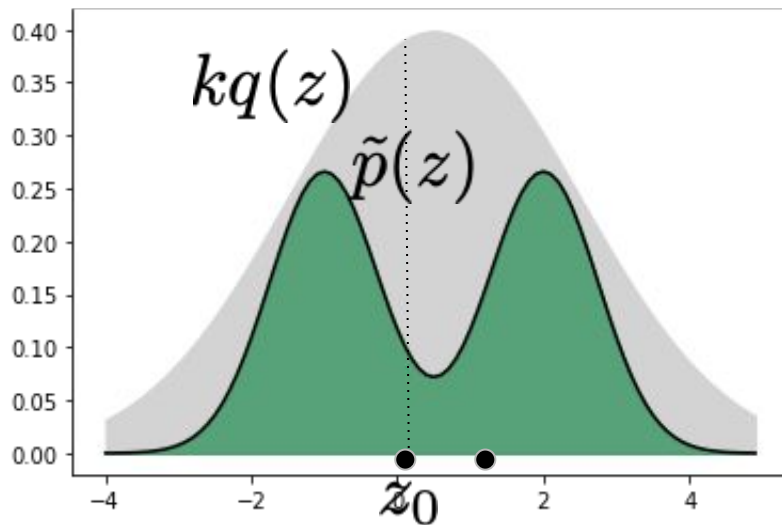


# Rejection Sampling - procedure

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Repeat until  
have sufficient  
samples



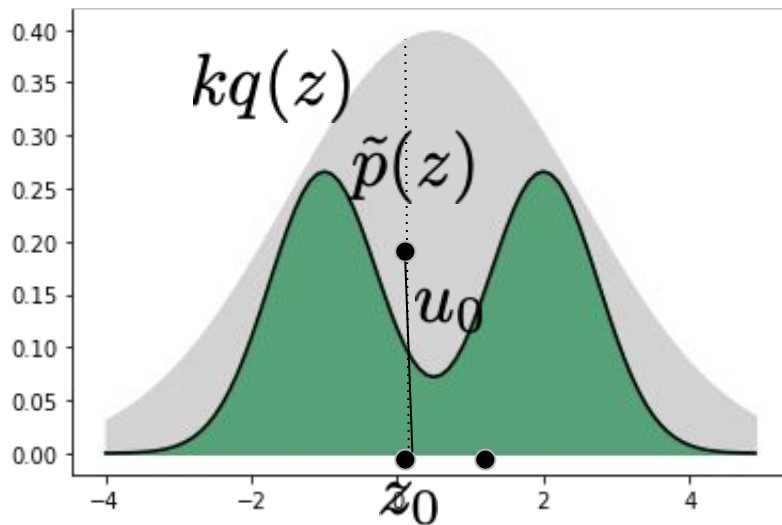


# Rejection Sampling - procedure

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Repeat until  
have sufficient  
samples

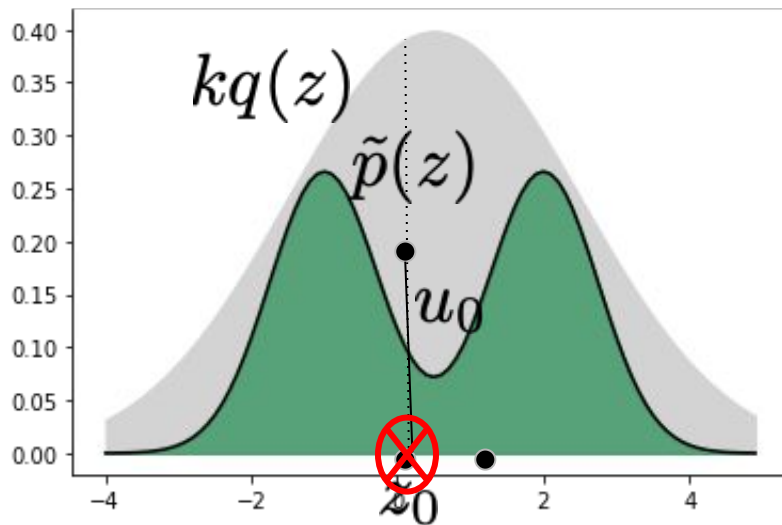


# Rejection Sampling - procedure

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Repeat until  
have sufficient  
samples

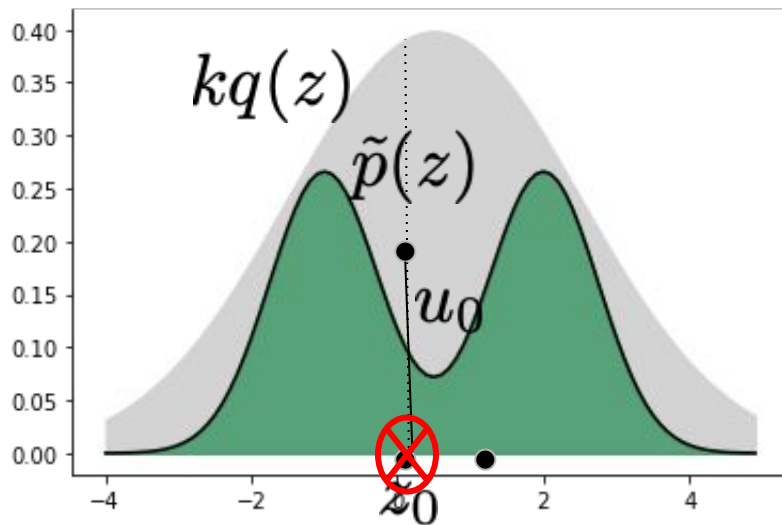


# Rejection Sampling - procedure

1. Choose a distribution  $q(z)$  that we can sample from and constant  $k$  such that  $kq(z) > \tilde{p}(z)$  for all  $z$
2. Sample  $z_0 \sim q(z)$
3. Sample  $u_0 \sim \text{Unif}(0, kq(z_0))$
4. Keep  $z_0$  as a sample if  $u_0 < \tilde{p}(z_0)$ , otherwise reject  $z_0$

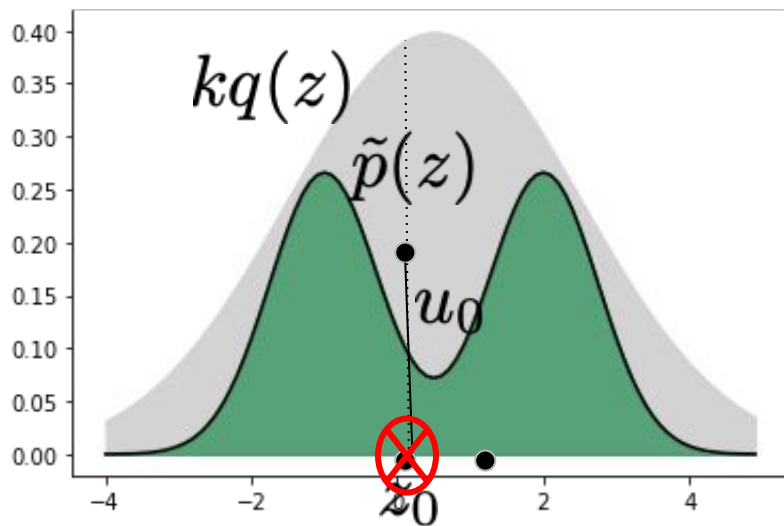
$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Repeat until  
have sufficient  
samples



# Rejection Sampling - Key limitations

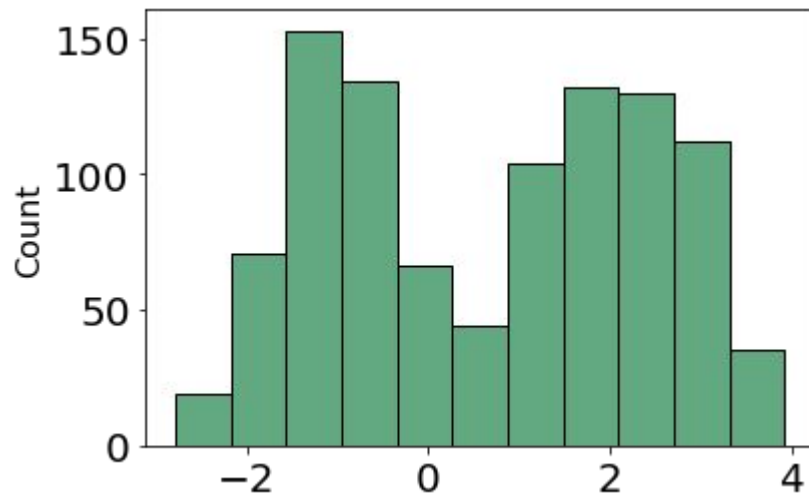
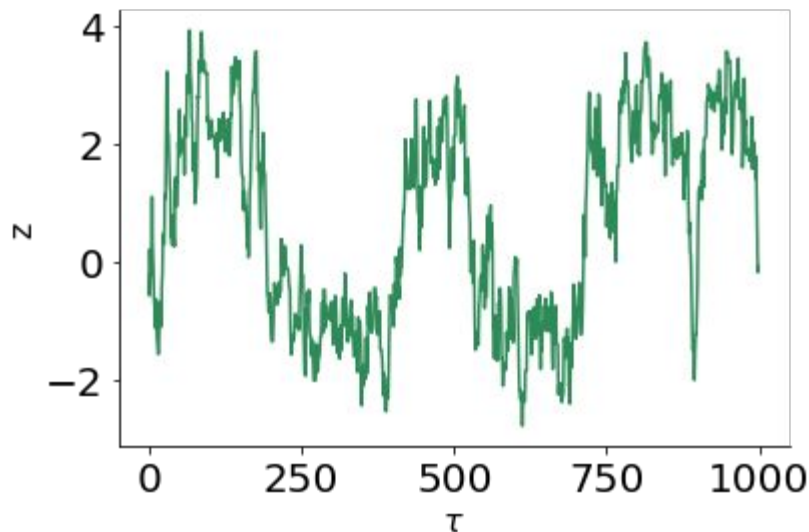
- Efficiency depends on good choice of  $k$  and  $q$
- Becomes inefficient very quickly as the dimensionality of  $z$  increases



# Markov Chain Monte Carlo

Idea:

- generate one sample based on the previous one (form a Markov Chain)
- do this *cleverly*, so that samples are from distribution of interest



# Markov Chain Monte Carlo - Metropolis Algorithm

Suppose we can evaluate  $\tilde{p}(z) = Z_p p(z)$ , the unnormalized pdf. Want samples.

Start with  $z^0$

Propose a new value for  $z^*$  according to some ***proposal distribution***

$q(z \mid z^\tau)$  that is symmetric in the sense that  $q(z_A \mid z_B) = q(z_B \mid z_A)$ .

Accept  $z^*$  with probability:  $A(z^*, z) = \min \left( 1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^\tau)} \right) = \frac{p(z^*)}{p(z^\tau)}$

If accept, set  $z^{\tau+1} = z^*$ . If not accept, set  $z^{\tau+1} = z^\tau$



# Markov Chain Monte Carlo - Metropolis Algorithm

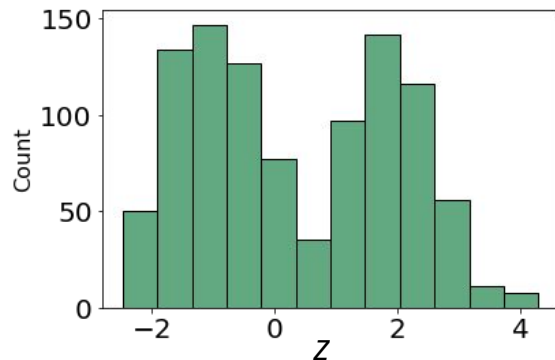
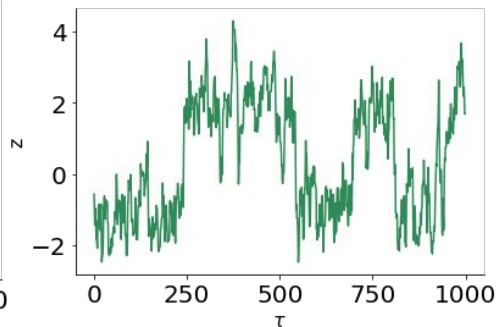
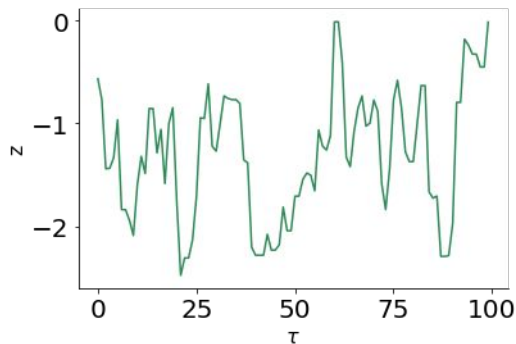
Start with  $z^0$

Propose a new value for  $z^*$  according to some **proposal distribution**

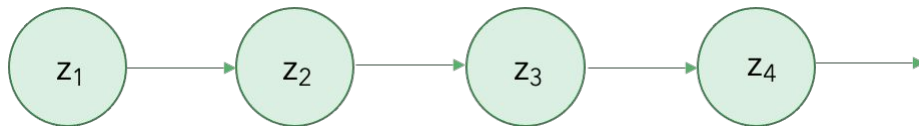
$q(z^* | z^\tau)$  that is symmetric in the sense that  $q(z_A | z_B) = q(z_B | z_A)$ .

Accept  $z^*$  with probability:  $A(z^*, z) = \min \left( 1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^\tau)} \right) = \frac{p(z^*)}{p(z^\tau)}$

If accept, set  $z^{\tau+1} = z^*$ . If not accept, set  $z^{\tau+1} = z^\tau$



# Markov chains



Sequence where future is independent from the past, given the present:

$$p(z^{\tau+1} | z^{\tau}, z^{\tau-1}, \dots, z_0) = p(z^{\tau+1} | z^{\tau})$$

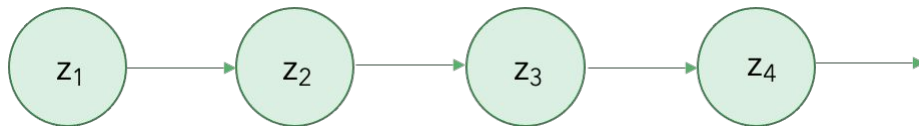
We call  $T(z^{\tau}, z^{\tau+1}) = p(z^{\tau+1} | z^{\tau})$  the **transition probabilities**

Example: random walk

$$z^{\tau+1} = z^{\tau} + s, s \sim N(0, \sigma^2)$$



# Markov chains



A Markov chain has a **stationary distribution**  $\pi$  if:

$$\pi(z) = \sum_{z'} T(z', z) \pi(z')$$

In matrix notation:

$$\pi = \pi P \quad \text{where}$$

$$P_{ij} = p(z^{\tau+1} = j | z^{\tau} = i)$$

or, in continuous state-space Markov chain, if:

$$\pi(z) = \int T(z', z) \pi(z') dz'$$

# Markov chains

In MCMC methods, we want to sample from some target distribution  $p$

We construct our Markov process so that it will asymptotically reach  $p$  as its stationary distribution.

What do we need to show to guarantee that this occurs?

1. That  $p$  is stationary distribution. A sufficient condition (that's often easier to work with than the definition of a stationary distribution) is that **detailed balance** holds:

$$\pi(z)T(z, z') = \pi(z')T(z', z)$$

2. The stationary distribution is unique. This holds if the Markov chain satisfies the is (roughly stated): *aperiodic* - gcd of number of steps to return to a state with positive probability is 1 and *positive recurrent* - expected return time to a state is finite. For continuous state space, replace state with set of states with positive probability )

# MCMC - Metropolis Hastings Algorithm

*Generalizes Metropolis algorithm by allowing non-symmetric proposal distributions*

Start with  $z^0$

Propose a new value for  $z^*$  according to some ***proposal distribution***

$q(z^* | z^\tau)$ , which is not necessarily symmetric

Accept  $z^*$  with probability:  $A(z^*, z) = \min \left( 1, \frac{\tilde{p}(z^*)q(z^\tau | z^*)}{\tilde{p}(z^\tau)q(z^* | z^\tau)} \right)$

If accept, set  $z^{\tau+1} = z^*$ . If not accept, set  $z^{\tau+1} = z^\tau$



# MCMC - Metropolis Hastings

recall:

$$A(z^*, z) = \min \left( 1, \frac{\tilde{p}(z^*)q(z^\tau|z^*)}{\tilde{p}(z^\tau)q(z^*|z^\tau)} \right)$$

Proving detailed balance: show  $p(z)T(z, z') = p(z')T(z', z)$

$$\begin{aligned} p(z)T(z, z') &= p(z)q(z'|z)A(z', z) \\ &= \min(p(z)q(z'|z), p(z')q(z|z')) \\ &= \min(p(z')q(z|z'), p(z)q(z'|z)) \\ &= p(z')q(z|z')A(z, z') \\ &= p(z')T(z', z) \end{aligned}$$

# MCMC - Metropolis Hastings

- Step size can be a challenge! If we choose it to be too small, we may explore too slowly to be efficient. If too large, we may reject too many samples to be efficient. Want acceptance ratio to be neither too large nor too small.
- Also, target distribution can have different standard deviations in different directions, another challenge if we use something like an isotropic Gaussian as proposal distribution
- We may want to discard the earliest samples in the chain (“burn-in” or “warm-up”)

# MCMC - Gibbs Sampling

-Iterate through each parameter, fixing all other parameters and sample from the conditional.

$$z_1^{\tau+1} \sim p(z_1 | z_2^\tau, z_3^\tau \dots z_m^\tau)$$

$$z_2^{\tau+1} \sim p(z_2 | z_1^{\tau+1}, z_3^\tau \dots z_m^\tau)$$

...

$$z_m^{\tau+1} \sim p(z_m | z_1^{\tau+1}, z_2^{\tau+1} \dots z_{m-1}^{\tau+1})$$

note that we're also typically conditioning on data

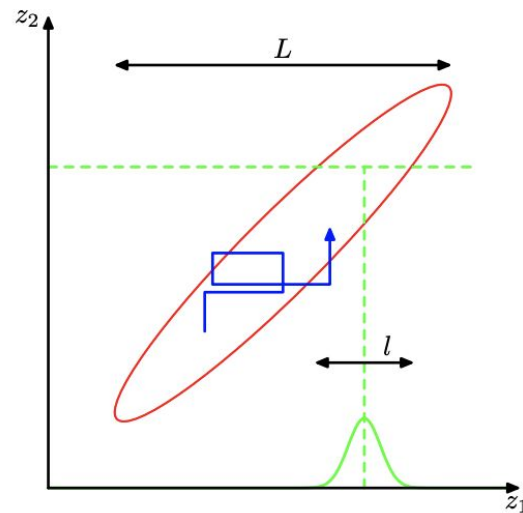
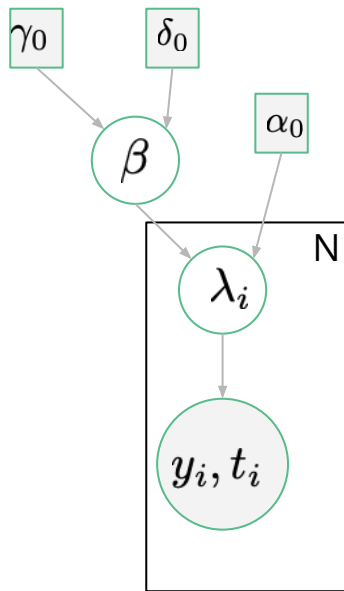


Fig 11.11 from *Pattern Recognition and Machine Learning*, Christopher Bishop

# MCMC - Gibbs Sampling

Conditionals are sometimes nice to calculate in hierarchical models.... e.g.:



$$\beta \sim \text{Gamma}(\gamma_0, \delta_0)$$

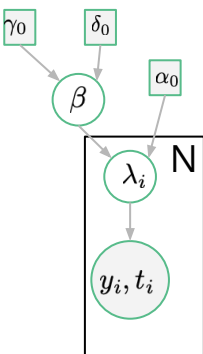
$$\lambda_i \sim \text{Gamma}(\alpha_0, \beta), i = 0, \dots, N$$

$$y_i \sim \text{Poisson}(\lambda_i t_i) i = 0, \dots, N$$

hyperparameters:  $\gamma_0, \delta_0, \alpha_0$

observations:  $y_i, t_i$

want to know about parameters:  $\beta, \lambda_i$



$$\begin{aligned}\beta &\sim \text{Gamma}(\gamma_0, \delta_0) \\ \lambda_i &\sim \text{Gamma}(\alpha_0, \beta), \quad i = 0, \dots, N \\ y_i &\sim \text{Poisson}(\lambda_i t_i) \quad i = 0, \dots, N\end{aligned}$$

hyperparameters:  $\gamma_0, \delta_0, \alpha_0$

observations:  $y_i, t_i$

want to know about parameters:  $\beta, \lambda_i$



$$\begin{aligned}\beta | \lambda, y, t &\sim \text{Gamma}(N\alpha_0 + \gamma_0, \sum_i \lambda_i + \delta_0) \\ \lambda_i | \beta, \lambda_{/i}, \mathbf{y}, \mathbf{t} &\sim \text{Gamma}(y_i + \alpha_0, t_i + \beta)\end{aligned}$$

Gamma and poisson pdfs

$$\text{Gamma}(x|a, b) = b^a x^{a-1} e^{-bx} / \Gamma(a)$$

$$\text{Poisson}(y|\lambda) = e^{-\lambda} \lambda^y / y!$$

$$\begin{aligned}p(\beta, \lambda | \mathbf{y}, \mathbf{t}) &\propto \overset{\text{likelihood}}{p(\mathbf{y}, \mathbf{t} | \lambda)} \overset{\text{prior}}{p(\lambda | \alpha_0, \beta)} p(\beta | \gamma_0, \delta_0) \\ &= \frac{\prod_i (e^{-\lambda_i t_i} (\lambda_i t_i)^{y_i}) \prod_i (\beta^\alpha \lambda_i^{\alpha_0-1} e^{-\beta \lambda_i}) \delta_0^{\gamma_0} \beta^{\delta_0-1} e^{-\delta_0 \beta}}{\Gamma(\alpha_0) \Gamma(\gamma_0) \prod_i y_i} \\ &\propto e^{-\sum_i (\lambda_i t_i + \beta \lambda_i) + \delta_0 \beta} \beta^{N\alpha_0 + \gamma_0 - 1} \prod_i \lambda_i^{y_i + \alpha_0 - 1}\end{aligned}$$

Combining terms and dropping factors that don't depend on beta or lambda

$$p(\beta | \lambda, y, t) \propto e^{-(\sum_i \lambda_i + \delta_0) \beta} \beta^{N\alpha_0 + \gamma_0 - 1}$$

$$\Rightarrow \beta | \lambda, y, t \sim \text{Gamma}(N\alpha_0 + \gamma_0, \sum_i \lambda_i + \delta_0)$$

$$p(\lambda_i | \lambda_{/i}, \beta, y, t) \propto e^{-(t_i + \beta) \lambda_i} \lambda_i^{y_i + \alpha_0 - 1}$$

$$\Rightarrow \lambda_i | \lambda_{/i}, \beta, y, t \sim \text{Gamma}(y_i + \alpha_0, t_i + \beta)$$

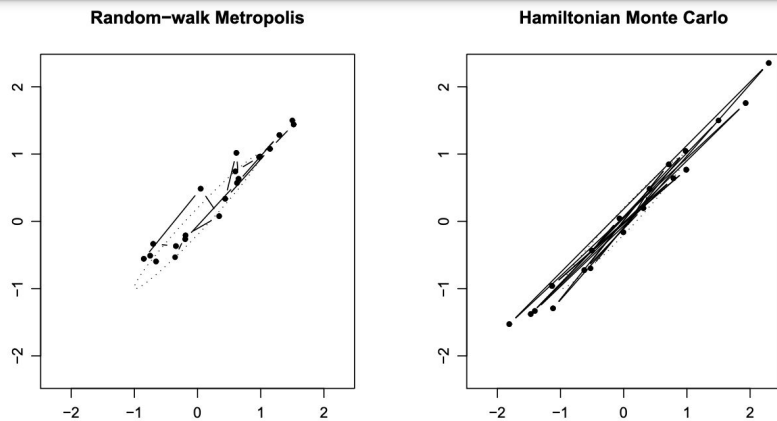


# MCMC - Gibbs Sampling - additional notes

- Can be an efficient choice if if you can actually get all the conditionals you need
- Gibbs sampling is actually a special case of Metropolis Hastings (see e.g PRML , Bishop, Section 11.3)
- Augmenting by including auxiliary variables can be (surprisingly?) useful, to make sampling easier
- Could use Metropolis-within-Gibbs if one or some of the conditionals are hard to sample from

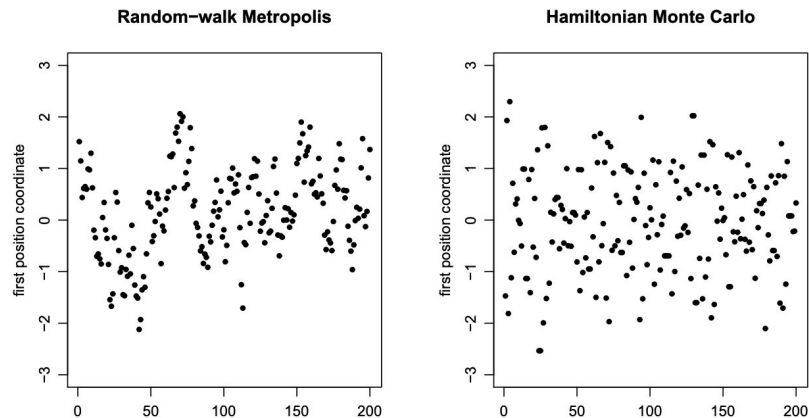
# Gradient-based MCMC methods - Hamiltonian MC

- Want a method that can propose samples that are far away from current, BUT still have a high probability of acceptance
- Will introduce some auxiliary variables and use Hamiltonian dynamics to generate proposals (parameters as **position**, auxiliary **momentum** variables)
- Flexible and efficient sampling method for continuous distributions. Need to be able to evaluate (log) density function (up to normalizing constant), and to compute partial derivatives of the log density



Figures 4, 5 from Neal, 2011

Figure 4: Twenty iterations of the random-walk Metropolis method (with 20 updates per iteration) and of the Hamiltonian Monte Carlo method (with 20 leapfrog steps per trajectory) for a 2D Gaussian distribution with marginal standard deviations of one and correlation 0.98. Only the two position coordinates are plotted, with ellipses drawn one standard deviation away from the mean.



# Gradient-based MCMC methods - HMC

-Position variables ( $\mathbf{q}$ ) are all the variables of interest (i.e. parameters in our Bayesian model)

-Momentum variables ( $\mathbf{p}$ ): auxiliary variables that we introduce; same dimensionality as  $\mathbf{q}$

Consider system described by Hamiltonian  $H(\mathbf{p}, \mathbf{q})$ , where  $\mathbf{p}$  and  $\mathbf{q}$  are each functions of time and *Hamilton's equations* hold:

$$\begin{aligned}\frac{dq_i}{dt} &= \frac{\partial H}{\partial p_i} \\ \frac{dp_i}{dt} &= -\frac{\partial H}{\partial q_i}\end{aligned}$$

Typically for HMC, assume:

$$H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + K(\mathbf{p})$$

Potential energy:  
negative log probability

Kinetic energy: usually use  
 $K(\mathbf{p}) = \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} / 2$   
(-log p of mean-0 Gaussian)

These equations tell us how the system evolves in time from a starting point along *trajectories* in  $\mathbf{p}, \mathbf{q}$  space

Will use points along these trajectories as proposals in a Metropolis step!

-Position variables ( $\mathbf{q}$ ) are all the variables of interest (i.e. parameters in our Bayesian model)

-Momentum variables ( $\mathbf{p}$ ): auxiliary variables that we introduce; same dimensionality as  $\mathbf{q}$

Consider system described by Hamiltonian  $H(\mathbf{p}, \mathbf{q})$ , where  $\mathbf{p}$  and  $\mathbf{q}$  are each functions of time and *Hamilton's equations* hold:

Typically for HMC, assume:

$$H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + K(\mathbf{p})$$

Potential energy:  
negative log probability

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i}$$

$$\frac{dp_i}{dt} = -\frac{\partial H}{\partial q_i}$$

Kinetic energy: usually use  
 $K(\mathbf{p}) = \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} / 2$

(-log p of mean-0 Gaussian)

Alternate:

1. Draw new momentum variables
2. Metropolis updates, using Hamiltonian dynamics to propose new state ( $\mathbf{p}^*$ ,  $\mathbf{q}^*$ )

-Position variables ( $\mathbf{q}$ ) are all the variables of interest (i.e. parameters in our Bayesian model)

-Momentum variables ( $\mathbf{p}$ ): auxiliary variables that we introduce; same dimensionality as  $\mathbf{q}$

Consider system described by Hamiltonian  $H(\mathbf{p}, \mathbf{q})$ , where  $\mathbf{p}$  and  $\mathbf{q}$  are each functions of time and *Hamilton's equations* hold:

Typically for HMC, assume:

$$H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + K(\mathbf{p})$$

Potential energy:  
negative log probability

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i}$$

$$\frac{dp_i}{dt} = -\frac{\partial H}{\partial q_i}$$

Kinetic energy: usually use  
 $K(\mathbf{p}) = \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} / 2$

(-log p of mean-0 Gaussian)

Alternate:

1. Draw new momentum variables
2. Metropolis updates, using Hamiltonian dynamics to propose new state ( $\mathbf{p}^*$ ,  $\mathbf{q}^*$ )

Can show:  $\frac{dH}{dt} = 0$

And Hamiltonian defines a joint probability:  $P(\mathbf{p}, \mathbf{q}) = \frac{1}{Z} \exp(-H(\mathbf{q}, \mathbf{p})/T)$

So, this joint probability is constant along a time trajectory

if  $H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + K(\mathbf{p})$

$$P(\mathbf{q}, \mathbf{p}) = \frac{1}{Z} \exp\left(\frac{-U(\mathbf{q})}{T}\right) \exp\left(\frac{-K(\mathbf{p})}{T}\right)$$

So, if we followed a trajectory exactly, our Metropolis step would accept with probability 1, since the proposed state would have same prob as initial.

In reality, we won't always accept, because we only approximate the trajectory

To compute (approx.) trajectories, need to integrate numerically. Will take  $L$  steps of size  $\epsilon$ . Easiest way could be Euler's method (alternate following the derivatives of  $p$  and  $q$  along a step of size  $\epsilon$ ):

$$p_i(t + \epsilon) = p_i(t) + \epsilon \frac{dp_i}{dt}(t) = p_i(t) - \epsilon \frac{\partial U}{\partial q_i}(q(t)), \quad q_i(t + \epsilon) = q_i(t) + \epsilon \frac{dq_i}{dt}(t) = q_i(t) + \epsilon \frac{p_i(t)}{m_i}$$

Euler's method **not** satisfactory here; instead use "leapfrog method" - half step momentum, full step position, another half step momentum.

Need to choose  $L, \epsilon$  - they matter

$$p_i(t + \epsilon/2) = p_i(t) - (\epsilon/2) \frac{\partial U}{\partial q_i}(q(t)),$$

$$q_i(t + \epsilon) = q_i(t) + \epsilon \frac{p_i(t + \epsilon/2)}{m_i},$$

$$p_i(t + \epsilon) = p_i(t + \epsilon/2) - (\epsilon/2) \frac{\partial U}{\partial q_i}(q(t + \epsilon)).$$

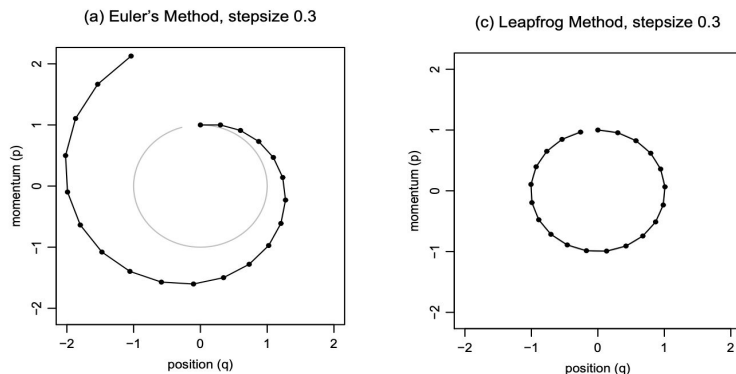


Figure 1 from Neal, 2011  
Position variable a zero mean gaussian

# Gradient-based MCMC methods - No U-Turn Sampler

- Extension of HMC

- Big idea: choose path lengths adaptively; don't need to choose  $L$  (number of steps in the integrator)

See:

Hoffman, Matthew D., and Andrew Gelman. "The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo." J. Mach. Learn. Res. 15, no. 1 (2014): 1593-1623.



# HMC - Two excellent resources for learning more

- Neal, Radford. *MCMC using Hamiltonian dynamics*, (2011). Chapter 5 of the Handbook of Markov Chain Monte Carlo. <https://arxiv.org/pdf/1206.1901.pdf>
- Betancourt, Michael. *A conceptual introduction to Hamiltonian Monte Carlo*. (2018) <https://arxiv.org/pdf/1701.02434.pdf>

# MCMC with PyMC3

- PyMC3 is an open source Python probabilistic programming library
- Makes it easy to specify a model and fit it using algorithms like NUTS
- Integrates nicely with tools for assessing your models and results
- Documentation here: <https://docs.pymc.io/> (including Quickstart: [https://docs.pymc.io/pymc-examples/examples/pymc3\\_howto/api\\_quickstart.html](https://docs.pymc.io/pymc-examples/examples/pymc3_howto/api_quickstart.html))
- Other probabilistic programming libraries exist. E.g. Stan is very popular and an excellent choice especially if you are an R user.

*See examples in colab*