

# Sampling 1

A self-study materials for PRML [1]

Jisung Lim<sup>1</sup>

<sup>1</sup>B.S. Candidate of Industrial Engineering  
Yonsei University, South Korea.

8th February, 2017



Christopher M. Bishop. **Pattern Recognition and Machine Learning**. Springer, 2006.



Stephen P. Brooks. “Markov Chain Monte Carlo Method and Its Application”.  
In: **Journal of the Royal Statistical Society. Series D (The Statistician)** 47:1  
(1998), pp. 69–100. ISSN: 00390526, 14679884. URL: <http://www.jstor.org/stable/2988428>.

# Summary

## 1 Introduction

- Bayesian Approach and Intractability Relaxation
- Introduction to Sampling

## 2 Basic Sampling Methods

- Basic Sampling
- Generating Random Samples
- General Strategy for Generating Samples
- Importance Sampling
- SIP
- Sampling and the EM algorithm

## 3 Markov Chain

- Markov Chain
- Markov Chain for MCMC

## 4 Markov Chain Monte Carlo

- MCMC: Introduction
- MCMC Algorithms

# Bayesian and Intractability

## Bayesian and Intractability

### ■ Yes, it is the Bayesian.

As we've seen in the previous chapters, the Bayesian approach have been more strongly focused than the frequentist veiwpoint, reflecting the huge growth in the practical importance of Bayesian methods.

### ■ Marginalization and Intractability

Full Bayesian approach requires intensive marginalization over the posterior distribution  $p(\theta|\mathcal{D})$  of parameters  $\theta$ . For instance, when we solve regression problem with a full bayesian procedure, we encountered the equation such as

$$p(t|\mathbf{x}, \mathcal{D}) = \int_{\theta} \int_{\mathbf{w}} p(t|\mathbf{x}, \mathbf{w}, \theta) p(\mathbf{w}|\mathcal{D}, \theta) p(\theta|\mathcal{D}) d\mathbf{w} d\theta$$

which is analytically intractable.

# Intractability Relaxation

## Two Main Paradigm of Intractability Relaxation

### ■ Intractability is around us.

Throughout other chapters, We've also seen this kind of intractable integrals in other chapters such as: 'Bayesian Logistic Regression', 'Bayesian Neural Network', 'Gaussian Process', 'Mixture Model', etc.

### ■ Sampling

The development of sampling methods, such as Markov chain Monte Carlo (discussed in Chapter 11) along with dramatic improvements in the speed and memory capacity of computers, opened the door to the practical use of Bayesian techniques in an impressive range of problem domains. Monte Carlo methods are very flexible and can be applied to a wide range of models. However, they are computationally intensive and have mainly been used for small-scale problems.

### ■ Variational Inference

More recently, highly efficient deterministic approximation schemes such as variational Bayes and expectation propagation (discussed in Chapter 10) have been developed. These offer a complementary alternative to sampling methods and have allowed Bayesian techniques to be used in large-scale applications.

# Introduction to Sampling

## Our Goal

In many cases, the posterior distribution is required primarily for the purpose of evaluating expectations. Hence our fundamental goal is to evaluate the expectation of some function  $f(\mathbf{z})$  with respect to a probability distribution  $p(\mathbf{z})$  i.e.,

$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z}, \quad (1)$$

or in the case of discrete variables,

$$\mathbb{E}[f] = \sum_{\mathbf{z}} f(\mathbf{z})p(\mathbf{z}). \quad (2)$$

## Analytical intractability and approximation

Let us suppose that such expectations  $\mathbb{E}[f]$  are too complex to be evaluated exactly using analytical techniques.

For random variable  $\mathbf{z}$ , iid random samples  $\{\mathbf{z}^{(l)}\}_{l=1}^L$  allow us to approximate the expectation by a statistic given by

$$\bar{f} = \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}^{(l)}). \quad (3)$$



# Introduction to Sampling

## Nice Properties of Statistic $\bar{f}$

- $\bar{f}$  is unbiasedness estimator for  $\mathbb{E}[f]$

$$\mathbb{E}[f] = \mathbb{E}[\bar{f}]$$

i.e., the estimator  $\bar{f}$  has the correct mean.

- The variance of  $\bar{f}$  is given by

$$\text{Var}[\bar{f}] = \frac{1}{L^2} (\text{Var}[f(\mathbf{z}^{(1)})] + \cdots + \text{Var}[f(\mathbf{z}^{(L)})]) = \frac{1}{L} \text{Var}[f] = \frac{1}{L} \mathbb{E}[(f - \mathbb{E}[f])^2]$$

i.e., the accuracy of the estimator does not depend on the dimensionality of  $\mathbf{z}$  but depends on the number of samples  $L$ . That is, in principle, high accuracy may be achievable with a relative small samples.

## Common Problems

- The samples  $\{\mathbf{z}^{(l)}\}_{l=1}^L$  might not be independent. That is, the effective sample size might be much smaller than the apparent sample size.
- Some functions  $f$  have large value with low  $p$  at some point  $\mathbf{z}$ , which can cause biased result evaluating expectation. That is, the approximation of expectations is undesirably dominated by some specific tendency.
- In practical, we therefore need relatively large sample size than we expected to achieve sufficient accuracy.

# Ancestral Sampling

## Sampling from the joint probability distribution

- The joint distribution can be specified by

$$p(\mathbf{z}) = \prod_{i=1}^M p(\mathbf{z}_i | \text{pa}_i) \quad (4)$$

- Let the graph has partial order between parent(low) and child(high).
- Get sample  $\hat{\mathbf{z}}_i$  from the pdf  $p(\mathbf{z}_i | \text{pa}_i)$  in topological order, from lowest to highest (i.e., parent to child).
- This is always possible because at each step all of the parent values will have been instantiated.



# Logic Sampling (w/ evidence)

## What if some data (evidence) is already taken?

- Now we consider the case of a directed graph in which some of the nodes are instantiated with observed values.
- The most simple, straight forward algorithm is *logic sampling*
- Do *Ancestral Sampling* repetitively until the sample satisfies the evidence.
- However, this algorithm is impractical because the sample is highly likely to be disagreed so as to be discarded.
- There are better algorithms such as *rejection sampling* and *adaptive rejection sampling* and *importance sampling*.

# Logic Sampling (w/ evidence)

## How to deal with indirected graph model?

- In the case of probability distributions defined by an undirected graph, there is no one-pass sampling strategy that will sample even from the prior distribution with no observed variables.
- Instead, computationally more expensive techniques must be employed, such as Gibbs sampling, which is discussed in Section 11.3.

## How to deal with the sampling problem from marginal distribution

- As well as sampling from conditional distributions, we may also require samples from a marginal distribution.
- If we already have a strategy for sampling from a joint distribution  $p(u, v)$ , then it is straightforward to obtain samples from the marginal distribution  $p(u)$  simply by ignoring the values for  $v$  in each sample.

# Generating Random Samples

- **Goal:** How to generate random numbers from simple non-uniform distributions?

random variable:  $y$  (5)

- **Given:** A source of uniformly distributed random numbers and a function maps  $z$  to  $y$ .

$$z \sim U(0, 1) \quad \text{and} \quad y = f(z) \quad (6)$$

- **find  $f$ :** By change of variables

$$p(y) = p(z) \left| \frac{dz}{dy} \right| \quad (7)$$

and integrating each side

$$h(y) := \int_{-\infty}^y p(y) \, dy = z \quad (8)$$

thus

$$y = h^{-1}(z) \quad \text{i.e.} \quad f = h^{-1} \quad (9)$$

# Generating Random Samples

## Generalization to Multiple Variables

The generalized form for multiple variables is started with

$$p(y_1, y_2, \dots, y_M) = p(z_1, z_2, \dots, z_M) \left| \frac{\partial(z_1, z_2, \dots, z_M)}{\partial(y_1, y_2, \dots, y_M)} \right| \quad (10)$$

cf) The Box-Muller method.

## Limitation

- Obviously, the transformation technique depends for its success on the ability to calculate and then invert the indefinite integral of the required distribution.
- That is, such operations will only be feasible for a limited number of simple distributions.
- and so we must turn to alternative approaches in search of a more general strategy: **rejection sampling** and **importance sampling**.

# Rejection Sampling

- **Goal:** Sample from relatively complex distribution, subject to certain constraints.
- **Given 1:** The distribution  $p(z)$  is not a simple, standard distributions considered in the previous slides. That is, it is too complex to find

$$h(z) := \int_{-\infty}^z p(t) dt = u \quad (11)$$

inverse function of which is the function  $f : u \rightarrow z$  where  $u \sim U(0, 1)$ .

- **Given 2:** We can easily evaluate evaluate  $p(z)$  for any given  $z$  up to some normalizing constant  $Z_p$  as follows

$$p(z) = \frac{1}{Z_p} \tilde{p}(z) \quad (12)$$

where  $\tilde{p}(z)$  can readily be evaluated but  $Z_p$  is unknown.

- **Main concept:** The main concept of rejection sampling is using distribution cover  $q(z)$  which is much simpler than  $p(z)$ . That envelope distribution is called 'proposal distribution'.



# Rejection Sampling

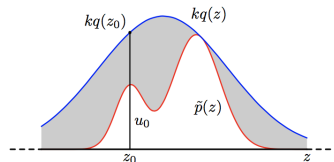


Figure 1: Rejection Sampling.

- **Proposal distribution  $q(z)$ :** We first introduce *proposal distribution*, from which we can readily draw samples.

$$q(z) \quad (13)$$

- **Bound Constant  $k$ :** We next introduce a constant  $k$  whose value is chosen such that

$$kq(z) \geq \tilde{p}(z) \quad \forall z \quad (14)$$

# Rejection Sampling

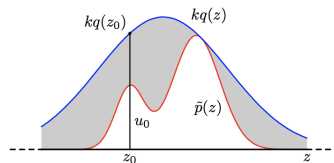


Figure 2: Rejection Sampling.

## Drawing sample

- 1 Generate a number  $z_0$  from the distribution  $q(z)$ .
- 2 Generate a number  $u_0$  from the uniform distribution over  $[0, kq(z_0)]$ .
- 3 If  $u_0 > \tilde{p}(z_0)$  then the sample is rejected, otherwise accepted.

## The acceptance rate

$$\begin{aligned}
 p(\text{accept}) &= \int \{\tilde{p}(z)/kq(z)\}q(z)dz \\
 &= \frac{1}{k} \int \tilde{p}(z)dz
 \end{aligned}$$

# Adaptive Rejection Sampling

## Motivation

- In many instances where we might wish to apply rejection sampling, it proves difficult to determine a suitable analytical form for the envelope distribution  $q(z)$ .
- How about constructing the envelope function  $q(z)$  on the fly based on measured values of the distribution  $p(z)$ .

## Settings

- $p(z)$  is log-concave.
- The function  $\ln p(z)$  and its gradient are evaluated at some initial set of grid points, and the intersections of the resulting tangent lines are used to construct the envelope function.

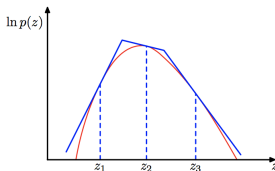


Figure 3: log envelope function.



# Adaptive Rejection Sampling

## envelope function

- Because the log envelope distribution is a succession of linear functions, and hence the envelope distribution itself comprises a piecewise, smooth exponential distribution of the form

$$q(z) = k_i \lambda_i \exp \{-\lambda_i(z - z_i)\} \quad \text{where} \quad \text{interval}(z_i) \quad (16)$$

## Sampling step

- 1 A sample valuse is drawn from the envelope distribution, let the sample be drawn at  $z'$ .
- 2 If accepted, it will be a sample drawn from the desired distribution.
- 3 If the sample is rejected, a new line is evaluated at the corresponding position  $z'$ , and the log envelope function (and eventually envelope function) is thereby refined.

# Conclusion to Rejection Sampling

## Conclusions

- The log concave constraint can be vanished simply by following each rejection sampling step with a Metropolis-Hastings step (section 11.2), giving rise to *adaptive rejection metropolis* sampling.
- That is, clearly, for rejection sampling to be of practical value, we required that the comparison function be close to the required distribution so that the rate of rejection is kept to minimum.

## Limitations

- For more practical examples, where the desired distribution may be multimodal and sharply peaked, it will be extremely difficult to find a good proposal distribution and comparison function. (e.g. Gaussian cover)
- Furthermore, the exponential decrease of acceptance rate with dimensionality is a generic feature of rejection sampling.

# Importance Sampling

## Motivation

- Only wish to evaluate the expectation, not the samples themselves.
- Discretize whole  $z$ -space into a uniform grid and evaluate the integration. (c.f. Riemannian sum)

$$\mathbb{E}[f] \approx \sum_{l=1}^L p(\mathbf{z}^{(l)}) f(\mathbf{z}^{(l)}). \quad (17)$$

## Concerns

- An obvious problem, however, is that the number of terms in the summation grows exponentially with the dimensionality of  $z$ .
- Only relatively small regions of  $z$  space has probability.
- So this first motivation is very inefficient.
- We therefore introduce a proposal distribution  $q(z)$ .

# Importance Sampling

## Approximation

- We can approximate  $\mathbb{E}[f]$  as follows

$$\begin{aligned}\mathbb{E}[f] &= \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int f(\mathbf{z})\frac{p(\mathbf{z})}{q(\mathbf{z})}q(\mathbf{z})d\mathbf{z} \\ &\approx \frac{1}{L} \sum_{l=1}^L \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})} f(\mathbf{z}^{(l)}) = \frac{1}{L} \sum_{l=1}^L r_l f(\mathbf{z}^{(l)})\end{aligned}\tag{18}$$

where  $r_l = \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})}$  are known as *importance weights*.

## Normalization constant

- In practical sense,

$$\begin{aligned}p(\mathbf{z}) &= \tilde{p}(\mathbf{z})/Z_p \\ q(\mathbf{z}) &= \tilde{q}(\mathbf{z})/Z_q\end{aligned}\tag{19}$$

where  $Z_p$  and  $Z_q$  are unknown.

- The approximation is then given by

$$\mathbb{E}[f] \approx \frac{Z_q}{Z_p} \frac{1}{L} \sum_{l=1}^L \frac{\tilde{p}(\mathbf{z}^{(l)})}{\tilde{q}(\mathbf{z}^{(l)})} f(\mathbf{z}^{(l)}) = \frac{Z_q}{Z_p} \frac{1}{L} \sum_{l=1}^L \tilde{r}_l f(\mathbf{z}^{(l)})\tag{20}$$



# Importance Sampling

- The ratio of normalization constants  $Z_p$  and  $Z_q$  can be approximated with same sample set

$$\frac{Z_p}{Z_q} = \frac{1}{Z_q} \int \tilde{p}(\mathbf{z}) d\mathbf{z} = \int \frac{\tilde{p}(\mathbf{z})}{\tilde{q}(\mathbf{z})} q(\mathbf{z}) d\mathbf{z} \approx \frac{1}{L} \sum_{l=1}^L \tilde{r}_l \quad (21)$$

- The approximation of  $\mathbb{E}[f]$  is therefore given by

$$\mathbb{E}[f] \approx \sum_{l=1}^L w_l f(\mathbf{z}^{(l)}) \quad (22)$$

where

$$w_l = \frac{\tilde{r}_l}{\sum_m \tilde{r}_m} = \frac{\tilde{q}(\mathbf{z}^{(l)})/q(\mathbf{z}^{(l)})}{\sum_m \tilde{q}(\mathbf{z}^{(m)})/q(\mathbf{z}^{(m)})} \quad (23)$$

which implies that the success of the importance sampling approach depends crucially on how well the sampling distribution  $q(\mathbf{z})$  matches the desired distribution  $p(\mathbf{z})$ .



# Importance Sampling

## Limitations and Hazards

- If, as is often the case,
  - 1  $p(z)f(z)$  is strongly varying i.e.,  $q(z)$  is not tight as much as we desired.
  - 2 Significant proportion of its mass concentrated over relatively small regions of  $z$  space.
  - 3 The set  $\{r_l\}$  may be therefore dominated by a few weights having large values.
  - 4 Thus the effective sample size can be much smaller than the apparent sample size  $L$ .
- The extreme case is where none of the samples falls in the regions where  $p(z)f(z)$  is large. Just consider the case where the function  $p(z)f(z)$  is sharply peaked in some small regions of  $z$  space.

## Conclusion

- Hence a major drawback of the importance sampling method is the potential to produce results that are arbitrarily in error and with no diagnostic indication.
- This also highlights a key requirement for the sampling distribution  $q(z)$ , namely that it should not be small or zero in regions where  $p(z)$  may be significant.

## Further study

- Likelihood weighted sampling
- Self-importance sampling



# Sampling-importance-resampling

## Motivation

- Importance sampling approach only gives approximation to expectation, not samples.
- Rejection sampling approach gives samples but it requires to determine the constant  $k$  with desirable acceptance rate.

**Concepts:** There are two stages to the scheme.

- 1  $L$  samples  $\mathbf{z}^{(l)}$  are drawn from  $q(\mathbf{z})$ .
- 2  $w_l$  are constructed using (23).
- 3 A second set of  $L$  samples is drawn from the discrete distribution  $(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)})$  with probability given by the weights  $(w_1, \dots, w_L)$ .

## Intuition

- Let us consider the approximation

$$\int f(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})} q(\mathbf{z}) d\mathbf{z} \approx \frac{1}{L} \sum_{l=1}^L \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})} f(\mathbf{z}^{(l)}) \approx \int f(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z} \quad (24)$$

where  $r_l = \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})}$ .

- Consider then what is the role of  $w_l$

$$w_l = \frac{\tilde{r}_l}{\sum_m \tilde{r}_m}$$

# Monte Carlo EM algorithm

- Between E step and M step we must evaluate the complete log likelihood given by

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \int p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{Z}, \mathbf{X}|\boldsymbol{\theta}) d\mathbf{Z} \quad (26)$$

which can be approximated by sampling methods as follows

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \frac{1}{L} \sum_{l=1}^L \ln p(\mathbf{Z}^{(l)}, \mathbf{X}|\boldsymbol{\theta}) \quad (27)$$

- This procedure is called Monte Carlo EM algorithm.
- A particular instance of the Monte Carlo EM algorithm, called stochastic EM, arises if we consider a finite mixture model, and draw just one sample at each E step.



# IP algorithm

- When we move from a maximum likelihood approach to a full Bayesian treatment, marginalization is required over the posterior distribution  $p(\boldsymbol{\theta}|\mathcal{D})$  of the parameter.
- Furthermore, we shall suppose that drawing samples from the posterior distribution  $p(\boldsymbol{\theta}|\mathcal{D})$  is computationally difficult.
- **I-step:** For some case where there exists latent variable  $\mathbf{Z}$  as follows

$$p(\mathbf{Z}|\mathcal{D}) = \int p(\mathbf{Z}|\boldsymbol{\theta}, \mathcal{D})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta}, \quad (28)$$

let first draw a sample  $\boldsymbol{\theta}^{(l)}$  from the current estimate for  $p(\boldsymbol{\theta}|\mathcal{D})$ , and then draw a sample  $\mathbf{Z}^{(l)}$  from  $p(\mathbf{Z}|\boldsymbol{\theta}^{(l)}, \mathcal{D})$ .

- **P-step:** Using the sample  $\mathbf{Z}^{(l)}$ , approximate  $p(\boldsymbol{\theta}|\mathcal{D})$  as follows

$$p(\boldsymbol{\theta}|\mathcal{D}) \approx \frac{1}{L} \sum_{l=1}^L p(\boldsymbol{\theta}|\mathbf{Z}^{(l)}, \mathcal{D}). \quad (29)$$

- Note that we are making a (somewhat artificial) distinction between parameters  $\boldsymbol{\theta}$  and hidden variables  $\mathbf{Z}$ . From now on, we blur this distinction and focus simply on the problem of drawing samples from a given posterior distribution.

# Summary and Limitations

- At the heart of the Bayesian inference is “marginalization.” Therefore, the ability to integrate analytically complex and high dimensional functions is extremely important in Bayesian statistics.
- Sampling methods allows us to sample from posterior distribution directly, and obtain sample estimates of the quantities of interests.
- With high dimensionality, however, the sampling methods what we’ve discussed so far still have severe limitations.
- e.g., gaussian enveloping function in rejection sampling.
- e.g., Scattered clusters in importance sampling.

# Markov Chain; Definition

Before discussing MCMC methods, it is fair to study some general and useful properties of Markov chains in some detail. We only deal with first-order Markov chain in this chapter which is defined by

## Definition

A first-order Markov Chain is defined to be a series of **random variables**

$$\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(M)} \quad (30)$$

satisfying the Markov property given by

$$p(\mathbf{z}^{(m+1)} | \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}) = p(\mathbf{z}^{(m+1)} | \mathbf{z}^{(m)})$$

which also can be represented as a directed graph in the form of a chain. And the finite-dimensional (here, it stands for time dimension) is given by

$$P(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(M)}) \quad (31)$$

# Markov Chain; Examining the state space

## General definition of Stochastic process

$$X : (\Omega, \mathcal{F}) \rightarrow (S^\infty, \mathcal{E}^\infty) \quad (32)$$

where  $(S, \mathcal{E})$  is state space on which  $X^{(n)}$ 's are defined as follows

$$X^{(n)} : (\Omega, \mathcal{F}) \rightarrow (S, \mathcal{E}) \quad (33)$$

*The random variable from State space to Sampling space* The stochastic process, or discrete-time Markov chain, have countable state space  $S$ . To simplify the notation and focus on the sample itself we introduce the random vector maps original state  $S$  to  $\mathbf{z}$  which is equiprobable,  $D$ -dimensional space.

$$\mathbf{z}^{(n)} : (S, \mathcal{E}) \rightarrow (\mathbf{z}, \mathcal{A}) \quad (34)$$

where  $\mathcal{A}$  is a  $\sigma$ -algebra of  $\mathbf{z}$ .

# Markov Chain; Properties

- *Initial state:*  $z^{(0)}$
- *Transition probability:*

$$T_m(z^{(m)}, z^{(m+1)}) := p(z^{(m+1)} | z^{(m)}) \quad (35)$$

- *Time homogeneity:* Time-homogeneous Markov chains (or stationary Markov chains) are processes where

$$T_m(z^{(m)}, z^{(m+1)}) = T_n(z^{(n)}, z^{(n+1)}) \quad \forall m, n \quad (36)$$

- *Invariant distribution:*

$$p(z) = \sum_{z'} T(z', z) p(z') \quad (37)$$

- *Detailed balance:*

$$p(z) T(z, z') = p(z') T(z', z) \quad (38)$$

# Markov Chain; properties

**Properties** From the definition and the specification of the first-order Markov chain, we can examine some useful properties.

- *Reducibility*: A Markov chain is said to be irreducible if it is possible to get to any state from any state. The accessibility can be defined by

$$P(\mathbf{z}_n | \mathbf{z}_m) > 0 \quad \forall n, m \quad (39)$$

- *Periodicity*: A state  $n$  has period  $p$  if any return to state  $n$  must occur in multiples of  $p$  time steps. If  $p = 1$ , then the state is said to be *aperiodic*.
- *Recurrent*: A state  $n$  is said to be *transient* if, given that we start in state  $n$ , there is a non-zero probability that we will never return to  $n$ . State  $n$  is *recurrent* (or persistent) if it is not transient.
- *Ergodicity*: The Markov chain is called ergodic if it is irreducible, and its states are positive recurrent and aperiodic.

# Markov Chain for MCMC

**Our goal** is still to generate a sample from  $p(\mathbf{z})$ .

- A standard Markov chain Monte Carlo approach is to construct an **ergodic** Markov chain with a stationary distribution.

The main properties underpinning the MCMC method, to accomplish our goal, are followings.

- *Time reversible property*

$$\sum_{\mathbf{z}'} p(\mathbf{z}') T(\mathbf{z}', \mathbf{z}) = \sum_{\mathbf{z}'} p(\mathbf{z}) T(\mathbf{z}, \mathbf{z}') = p(\mathbf{z}) \sum_{\mathbf{z}'} T(\mathbf{z}, \mathbf{z}') = p(\mathbf{z}) \quad (40)$$

where  $p(\mathbf{z})$  is invariant distribution.

- *Existing and unique stationary distribution*

Any chain which is *irreducible* and *aperiodic* will have a unique stationary distribution.

- *Ergodic property*

$$p(\mathbf{z}^{(m)}) \rightarrow p(\mathbf{z}) \quad \text{as } m \rightarrow \infty \quad (41)$$

# Markov Chain for MCMC[2]

- 1 A Markov chain with stationary distribution  $p(\mathbf{z})$  is called *time reversible* (or simply *reversible*) if its transition kernel  $T_m(\mathbf{z}^{(m)}, \mathbf{z}^{(m+1)})$  is such that it exhibits *detailed balance*, i.e.

$$p(\mathbf{z})T(\mathbf{z}, \mathbf{z}') = p(\mathbf{z}')T(\mathbf{z}', \mathbf{z}). \quad (42)$$

This essentially means that the chain would look the same whether you ran it forwards in time or backwards. The behaviour of reversible chains is well understood, and it is a desirable property for any MCMC transition kernel to have, since any transition kernel for which *detailed balance* holds will have stationary distribution  $p(\mathbf{z})$ .

- 2 Any chain which is *irreducible* and *aperiodic* will have a unique stationary distribution, and that the  $t$ -step transition kernel will ‘converge’ to that stationary distribution as  $t \rightarrow \infty$ . See Meyn and Tweedie (1993).
- 3 A Markov chain is *ergodic* if, given realizations  $\{\mathbf{z}^{(m)} : m = 0, 1, \dots\}$  from the chain, typical asymptotic results include the distributional convergence of the realizations. That is, the distribution of the state of the chain at time  $m$  converges to the stationary distribution  $p(\mathbf{z})$  as  $m \rightarrow \infty$ .

$$p(\mathbf{z}^{(m)}) \rightarrow p(\mathbf{z}) \quad \text{as } m \rightarrow \infty$$



# MCMC: Introduction

## Motivation: Limitation of Basic Samplings

- In the previous section, we discussed the rejection sampling and importance sampling strategies for evaluating expectations of functions, and we saw that they suffer from severe limitations particularly in spaces of high dimensionality.

## Intuition: Markov Chain Monte Carlo

- General strategy which allows sampling from a large class of distribution (based on the mechanism of Markov chains)
- It also uses a proposal distribution to generate samples from another distribution.
- Introduces state, denoted by  $\tau$ , and remember the previous information, a sample,  $\mathbf{z}^{(\tau)}$
- The proposal distribution then depends on the current state:  $q(\mathbf{z}|\mathbf{z}^{(\tau)})$

# MCMC Algorithms

## Goal and Requirement

- **Goal:** to generate a sample from  $p(z)$
- **Requirement 1: stationary distribution**  
A markov chain where  $p(z)$  is invariant.
- **Requirement 2: ergodicity property**  
With requirement 1,  $p(z^{(m)}) \rightarrow p(z)$  as  $m \rightarrow \infty$

## Sampling Step

- 1 Remembering the current sample  $z^{(\tau)}$ , generate a candidate sample  $z^*$  from a proposal distribution  $q(z|z^{(\tau)})$ .
- 2 Accept the sample according to the criterion.
- 3 If the candidate sample is accepted, then  $z^{(\tau+1)} \leftarrow z^*$ , otherwise  $z^{(\tau+1)} \leftarrow z^{(\tau)}$ .

## Note that

- The sequence of samples  $z^{(1)}, z^{(1)}, \dots$  are strongly dependent to previous sample due to its sequential sampling. We therefore chose only  $M$ th samples to ensure the independence.



# Metropolis-Hastings Algorithm

## Acceptance function:

$$A_k(z^*, z^{(\tau)}) = \min_k \left( 1, \frac{\tilde{p}(z^*)q_k(z^{(\tau)}|z^*)}{\tilde{p}(z^{(\tau)})q_k(z^*|z^{(\tau)})} \right) \quad (44)$$

Here  $k$  labels the members of the set of possible transitions being considered.

**Rejection criteria:** For  $z^*$  from the proposal distribution  $q(z^*|z^{(\tau)})$ , Consider the function  $A_k(z^*, z^{(\tau)})$ . Here the transition probability is given by

$$q(z^*|z^{(\tau)})A_k(z^*, z^{(\tau)}) \quad (45)$$

## Validation of limiting distribution

$$\begin{aligned} p(z)q_k(z'|z)A_k(z', z) &= \min(p(z)q_k(z'|z), p(z')q_k(z|z')) \\ &= \min(p(z')q_k(z|z'), p(z)q_k(z'|z)) \\ &= p(z')q_k(z|z')A_k(z, z') \end{aligned} \quad (46)$$

# Metropolis Algorithm

**Symmetry specification:**  $q_k(z|z') = q_k(z'|z)$  **Acceptance function:**

$$A(z^*, z^{(\tau)}) = \min \left( 1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^{(\tau)})} \right) \quad (47)$$

where  $q(z_A|z_B) = q(z_B|z_A)$  for all  $A$  and  $B$ .

**Rejection criterion**

$$\text{Reject: } A(z^*, z^{(\tau)}) \leq u \quad (48)$$

$$\text{Accept: } A(z^*, z^{(\tau)}) > u$$

where  $u$  is randomly picked from uniform distribution  $U(0, 1)$ .

**Validation of limiting distribution**

$$\begin{aligned} p(z)q_k(z'|z)A_k(z', z) &= q_k(z'|z) \min \{p(z), p(z')\} \\ &= q_k(z|z') \min \{p(z'), p(z)\} \\ &= p(z')q_k(z|z')A_k(z, z') \end{aligned} \quad (49)$$



# Gibbs Sampling

## Gibbs Sampling

1. Initialize  $\{z_i : i = 1, \dots, M\}$
2. For  $\tau = 1, \dots, T$ :
  - Sample  $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$ .
  - Sample  $z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau+1)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$ .
  - $\vdots$
  - Sample  $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \dots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \dots, z_M^{(\tau)})$ .
  - $\vdots$
  - Sample  $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \dots, z_{M-1}^{(\tau+1)})$ .

Figure 4: Gibbs Sampling.

# Gibbs Sampling

## Acceptance function:

$$A(\mathbf{z}^*, \mathbf{z}) = \frac{p(\mathbf{z}^*)q_k(\mathbf{z}|\mathbf{z}^*)}{p(\mathbf{z}^{(\tau)})q_k(\mathbf{z}^*|\mathbf{z})} = \frac{p(z_k^*|\mathbf{z}_{-k}^*)p(\mathbf{z}_{-k}^*)p(z_k|\mathbf{z}_{-k}^*)}{p(z_k|\mathbf{z}_{-k})p(\mathbf{z}_{-k})p(z_k^*|\mathbf{z}_{-k})} \quad (50)$$

where  $\mathbf{z}_{-k} = \mathbf{z}_{-k}^*$

## Validation of equilibrium distribution $p$

$$p(z_i = z_i^{(\tau+1)}) = T(z_i = z_i^{(\tau)}, z_i = z_i^{(\tau+1)})p(z_i = z_i^{(\tau)}) \quad (51)$$

where  $T = p(z_i|z_{-i})$ .

# Slice Sampling

## Motivation

- In Metropolis-Hastings Algorithm, If the step size is too small, the result is slow decorrelation due to random walk behaviour, whereas if it is too large the result is inefficiency due to a high rejection rate.

## Concept

- adaptive step size that is automatically adjusted to match the characteristics of the distribution.

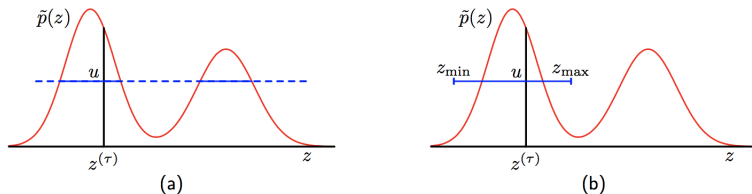


Figure 5: Slice Sampling.