

Lecture 1: Basics of Monte Carlo Methods

Tianpei Xie

Sep. 28th., 2022

Contents

1	Basic Concepts	2
1.1	Random variable generation	2
1.2	Why Monte Carlo ?	3
2	Rejection and Weighting	5
2.1	Reject Sampling	5
2.2	Variance Reduction	7
2.3	Importance Sampling	8
3	Sequential Importance Sampling	10
3.1	Sequential Importance Sampling (SIS)	10
3.2	Nonlinear filtering	12

1 Basic Concepts

1.1 Random variable generation

- The most fundamental methods for random variable generation is the ***uniform pseudo-random number generator***.

Definition A *uniform pseudo-random number generator* is an algorithm which, starting from an initial value u_0 and a transformation D , produces a sequence $(u_i) = (D^i(u_0))$ of values in $[0, 1]$. For all n , the values (u_1, \dots, u_n) reproduce the behavior of an *iid* sample (V_1, \dots, V_n) of uniform random variables when compared through a usual set of tests.

- This definition is clearly restricted to testable aspects of the random variable generation, which are connected through the deterministic transformation $u_i = D^i(u_{i-1})$. Thus, the validity of the algorithm consists in the verification that the sequence U_1, \dots, U_n leads to acceptance of the hypothesis

$$H_0 : U_1, \dots, U_n \text{ i.i.d. } \mathcal{U}[0, 1]$$

- Definition above is therefore ***functional***: An algorithm that generates uniform numbers is acceptable if it is not rejected by a set of tests.
- This methodology is not without problems, however. Consider, for example, particular applications that might demand a large number of iterations, as the theory of large deviations, or particle physics, where algorithms resistant to standard tests may exhibit fatal faults. In particular, algorithms having hidden periodicities (e.g. Ising models) or which are not uniform for the smaller digits may be difficult to detect.
- Given uniform distribution $\mathcal{U}[0, 1]$, and a known c.d.f. $F(x)$, we can generate samples via the ***inverse transform*** method. [Robert and Casella, 1999]

Definition For a ***non-decreasing*** function F on \mathbb{R} , the ***generalized inverse*** of F , F^- , is the function defined by

$$F^-(u) = \inf\{x \in \mathcal{R}(F) : F(x) \geq u\}. \quad (1)$$

where $\mathcal{R}(F)$ is the range of F . If F is the cumulative distribution function $F_X(x) = P(X \leq x)$, then F_X^{-1} is called ***quantile function***.

- We then have the following lemma, sometimes known as the ***probability integral transform***, which gives us a representation of any random variable as a transform of a uniform random variable.

Lemma 1.1 [Robert and Casella, 1999]

If $U \sim \mathcal{U}[0, 1]$, then the random variable $F^-(U)$ has the distribution F .

- Examples for some quantile functions [Robert and Casella, 1999]
 - **Exponential distribution**: $F^-(u) = -\log(1-u)$. Since U and $1-U$ are both uniform in $[0, 1]$, $X = \log(U) \sim \exp(1)$ is exponentially distributed with $\lambda = 1$;
 - **Gamma distribution** and χ^2 distribution can be generated from the exponential dis-

tribution $\exp(1)$

$$Y = \beta \sum_{j=1}^a X_j \sim \Gamma(a, \beta), ; a \in \mathbb{N} +$$

$$Y = 2 \sum_{j=1}^{\nu} X_j \sim \chi_{2\nu}^2; \nu \in \mathbb{N} +$$

– **Normal distribution:**

$$X_1 = \sqrt{-2 \log(U_1)} \cos(2\pi U_2)$$

$$X_2 = \sqrt{-2 \log(U_1)} \sin(2\pi U_2)$$

where U_1, U_2 are i.i.d from $\mathcal{U}[0, 1]$. Then X_1, X_2 are i.i.d from $\mathcal{N}(0, 1)$.

- Theoretically, any random variable can be generated using the inverse transform. But in practice, it is only available when we can compute the c.d.f and its inverse F^{-} explicitly.

1.2 Why Monte Carlo ?

- The most fundamental operations in probabilistic modeling and statistics involve integration on high dimensional spaces:

$$I := \int_{\mathcal{X}} g(\mathbf{x}) d\mathbf{x}$$

For example,

– **Expectation:**

$$\mathbb{E}_{\mathbf{X} \sim p} [h(\mathbf{X})] = \int_{\mathcal{X}} h(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \quad (2)$$

where $\mathbf{x} := [x_1, \dots, x_d]$

– **Marginalization:**

$$P(x_1, \dots, x_{i-1}, x_{i+1} \dots, x_d) = \int_{x_i \in \mathcal{X}_i} P(x_1, \dots, x_{i-1}, x_i, x_{i+1} \dots, x_d) dx_i$$

– **Conditioning:**

$$P(x_1, \dots, x_{i-1}, x_{i+1} \dots, x_d | X_i = x_i) = \frac{P(x_1, \dots, x_{i-1}, x_i, x_{i+1} \dots, x_d)}{\int_{\mathbf{x}_{-i} \in \mathbf{X}_{-i}} P(x_1, \dots, x_{i-1}, x_i, x_{i+1} \dots, x_d) d\mathbf{x}_{-i}}$$

where $\mathbf{x}_{-i} = [x_1, \dots, x_{i-1}, x_{i+1} \dots, x_d]$.

- This is especially important for **Bayesian inference** and analysis. Since it relies on computing *posterior distribution* given data using **the Bayes rule**

$$P(\Theta | \mathcal{D}) = \frac{P(\mathcal{D} | \Theta) P(\Theta)}{P(\mathcal{D})}$$

$$\begin{aligned}
&= \frac{P(\mathcal{D}|\Theta)P(\Theta)}{\int_{\Theta} P(\mathcal{D}|\Theta)P(\Theta)d\Theta} \\
&\Rightarrow \log P(\Theta|\mathcal{D}) = \log P(\mathcal{D}|\Theta) + \log P(\Theta) - \log \int_{\Theta} P(\mathcal{D}|\Theta)P(\Theta)d\Theta
\end{aligned}$$

where $\mathcal{D} := \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ are i.i.d. samples and Θ is a set of parameters for model $p(\mathcal{D}|\Theta)$.

- The main challenge for Bayesian methods in high dimensional space is to compute the normalization factor

$$\int_{\Theta} P(\mathcal{D}|\Theta)P(\Theta)d\Theta.$$

- The **log-partition function** in **exponential families** play a critical role since it defines a *bijjective mapping* from *natural parameters* to *mean parameters* and is related to **negative entropy** via *variational principles*. The log-partition function for exponential families is a convex function as well.

$$A(\eta) = \log \int_{\Omega} \exp(\langle \eta, \phi(\mathbf{x}) \rangle) d\mathbf{x}$$

- It is natural to propose using a sample $(\mathbf{X}_1, \dots, \mathbf{X}_m)$ generated from the density p to *approximate integration* by the *empirical average*. This approach is referred to as the **Monte Carlo methods**. Specifically from (2),

$$\bar{h}_m = \frac{1}{m} \sum_{i=1}^m h(\mathbf{X}_i). \quad (3)$$

since \bar{h}_m converges almost surely to $\mathbb{E}_{\mathbf{X} \sim p}[h(\mathbf{X})]$ by the *Strong Law of Large Numbers*. Moreover, when \bar{h}_m has a finite expectation under p , the speed of convergence of \bar{h}_m can be assessed since the variance

$$\text{Var}(\bar{h}_m) = \frac{1}{m} \int_{\mathcal{X}} (h(\mathbf{x}) - \mathbb{E}_{\mathbf{X} \sim p}[h(\mathbf{X})])^2 p(\mathbf{x}) d\mathbf{x}$$

can also be estimated from the sample $(\mathbf{X}_1, \dots, \mathbf{X}_m)$ through

$$v_m = \frac{1}{m} \sum_{i=1}^m (h(\mathbf{X}_i) - \bar{h}_m)^2$$

For m large,

$$\frac{h(\mathbf{X}) - \bar{h}_m}{\sqrt{v_m}} \rightarrow \mathcal{N}(0, 1).$$

This leads to the construction of a **convergence test** and of **confidence bounds** on the approximation of $\mathbb{E}_{\mathbf{X} \sim p}[h(\mathbf{X})]$.

2 Rejection and Weighting

2.1 Reject Sampling

- There exists a fundamental idea underlying the **Accept-Reject sampling method**:

Theorem 2.1 (Fundamental Theorem of Simulation) [Robert and Casella, 1999]
Simulating $X \sim F(x)$ is equivalent to simulating

$$(X, U) \sim \mathcal{U}\{(x, u) : 0 < u < F(x)\}.$$

In essence, If f is the density of interest, on an arbitrary space, we can write this theorem as

$$f(x) = \int_0^{f(x)} du.$$

Thus, f appears as the marginal density (in X) of the joint distribution, $(X, U) \sim \mathcal{U}\{(x, u) : 0 < u < f(x)\}$.

- Suppose $l(\mathbf{x}) = c f(\mathbf{x})$ is computable, where f is the probability density and c is unknown. If we can find a sampling distribution $g(\mathbf{x})$ and a **covering constant** M so that $Mg(\mathbf{x}) \geq l(\mathbf{x})$ (i.e. the envelop property) is satisfied for all \mathbf{x} . We can apply the following procedure:

Accept-Reject sampling sampling [Liu, 2001]

1. Draw sample \mathbf{x} from $g(\cdot)$ and compute the **ratio**

$$r = \frac{l(\mathbf{x})}{M g(\mathbf{x})} \quad (\leq 1)$$

2. Flip a coin with **success probability** r . If the head turns up, accept \mathbf{x} ; otherwise, reject \mathbf{x}

Then the accepted samples follow the distribution $f(\mathbf{x})$.

Proof: Let $I \in \{0, 1\}$ be the event when head turns up ($= 1$) vs. tail turns up ($= 0$).

$$P(I = 1) = \int P(I = 1|\mathbf{x})g(\mathbf{x})d\mathbf{x} = \int \frac{c f(\mathbf{x})}{M g(\mathbf{x})}g(\mathbf{x})d\mathbf{x} = \frac{c}{M}$$

Then

$$P(\mathbf{x}|I = 1) = \frac{g(\mathbf{x}) \frac{c f(\mathbf{x})}{M g(\mathbf{x})}}{P(I = 1)} = f(\mathbf{x}). \quad \blacksquare$$

- This is equivalent to say

1. Draw sample \mathbf{X} from $g(\mathbf{x})$ and U from $\mathcal{U}[0, 1]$
2. Accept $\mathbf{Y} = \mathbf{X}$ if $U \leq f(\mathbf{x})/(M g(\mathbf{x}))$;
3. Otherwise, return to step 1.

- The key for a successful application of rejection sampling method is to choose a good trial distribution $g(\cdot)$ so that M is **small**. In high dimensional settings, it is usually hard to find a good trial distribution.

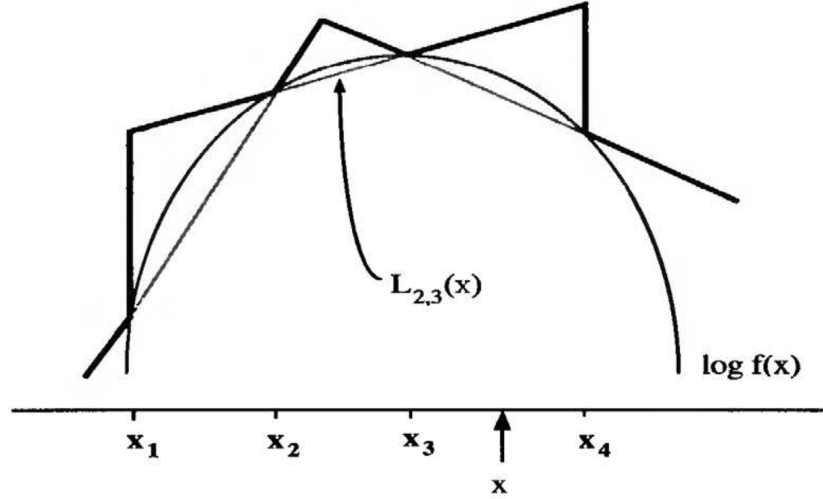


Fig. 2.5. Lower and upper envelopes of $h(x) = \log f(x)$, f a log-concave density (Source: Gilks et al. 1995).

Figure 1: The log-concave function h and its upper and lower bounds [Robert and Casella, 1999]

- **Lemma 2.2** *If there exist a density $g_m(\mathbf{x})$, a function $g_l(\mathbf{x})$ and a constant M such that*

$$g_l(\mathbf{x}) \leq f(\mathbf{x}) \leq M g_m(\mathbf{x})$$

then the algorithm **Envelope Accept-Reject**

1. Draw sample \mathbf{X} from $g_m(\mathbf{x})$ and U from $\mathcal{U}[0, 1]$
2. Accept \mathbf{X} if $U \leq g_l(\mathbf{x}) / (M g_m(\mathbf{x}))$;
3. Otherwise, accept \mathbf{X} if $U \leq f(\mathbf{x}) / (M g_m(\mathbf{x}))$ else return to step 1.

produces random variables that are distributed according to f .

- By the construction of a lower envelope on f , based on the function g_l , the number of evaluations of f is potentially decreased by a factor

$$\frac{1}{M} \int g_l(\mathbf{x}) d\mathbf{x},$$

which is the probability that f is not evaluated. This method is called the **squeeze principle**.

- For the exponential families, the density function

$$f(\mathbf{x}; \boldsymbol{\theta}) = h(\mathbf{x}) \exp(\langle \boldsymbol{\theta}, \mathbf{x} \rangle - A(\boldsymbol{\theta}))$$

is *log-concave* in terms of \mathbf{x} . An envelop accept-reject sampling method is discussed in [Robert and Casella, 1999]. The method is called **adaptive rejection sampling (ARS)** and it provides a sequential evaluation of lower and upper envelopes of the density f when $h = \log f$ is concave. Consider a set of n samples $\mathcal{D}_n = \{\mathbf{x}_i\}$. For each sample \mathbf{x}_i , $h(\mathbf{x}_i) = \log f(\mathbf{x}_i)$ is concave, so the following upper bound and lower bound can be found using linear segments $L_{i,i+1}$ through $(\mathbf{x}_i, h(\mathbf{x}_i))$ and $(\mathbf{x}_{i+1}, h(\mathbf{x}_{i+1}))$. Define

$$\bar{h}_n(\mathbf{x}) := \min \{L_{i-1,i}(\mathbf{x}), L_{i+1,i+2}(\mathbf{x})\};$$

$$\underline{h}_n(\mathbf{x}) := L_{i,i+1}(\mathbf{x}).$$

Then by concavity of h , for $\mathbf{x} \in [\mathbf{x}_i, \mathbf{x}_{i+1}]$

$$\underline{h}_n(\mathbf{x}) \leq h(\mathbf{x}) \leq \bar{h}_n(\mathbf{x})$$

Therefore, for $\underline{f}_n(\mathbf{x}) = \exp \underline{h}_n(\mathbf{x})$ and $\bar{f}_n(\mathbf{x}) = \exp \bar{h}_n(\mathbf{x})$ implies that

$$\underline{f}_n(\mathbf{x}) \leq f(\mathbf{x}) \leq \bar{f}_n(\mathbf{x}) = \omega_n g_n(\mathbf{x})$$

where ω_n is a normalization factor so that $g_n(\mathbf{x})$ is a density.

Then the **adaptive rejection sampling (ARS)** is

1. Initiate $\mathcal{D}_n = \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$. Compute the line segments $L_{i,i+1}$, $i = 0, \dots, n-1$.
2. Draw sample \mathbf{X} from $g_n(\mathbf{x})$ and U from $\mathcal{U}[0, 1]$
3. Accept \mathbf{X} if $U \leq \underline{f}_n(\mathbf{x})/(\omega_n g_n(\mathbf{x}))$;
4. Otherwise, accept $\mathbf{X}_n = \mathbf{X}$ if $U \leq f(\mathbf{x})/(\omega_n g_n(\mathbf{x}))$ and update $\mathcal{D}_{n+1} = \mathcal{D}_n \cup \{\mathbf{X}_n\}$.

An interesting feature of this algorithm is that the set \mathcal{D}_n is only updated when $f(\mathbf{x})$ has been previously computed. As the algorithm produces variables the lower bound $\underline{f}_n(\mathbf{x})$ and upper bound $\bar{f}_n(\mathbf{x})$ become increasingly accurate and, therefore, we progressively reduce the number of evaluations of f .

2.2 Variance Reduction

We introduce several techniques that reduce variance while maintain the unbiasedness of estimator.

- **Stratified sampling:** We partition the region \mathcal{X} into k sub-regions \mathcal{X}^i and suppose that the probability distribution f in each sub-regions is **relative homogeneous**. Then we can generate samples $\mathbf{X}_1^i, \dots, \mathbf{X}_{m_i}^i$ from each region \mathcal{X}^i and compute the region-wise sample average

$$\bar{\mu}_{m_i}^i = \frac{1}{m_i} \sum_{j=1}^{m_i} f(\mathbf{X}_j^i)$$

Then the final result is the average of region-wise sample average $\bar{\mu} = \sum_{i=1}^k \bar{\mu}_{m_i}^i$ and the variance $\bar{\sigma}^2 = \sum_{i=1}^k \frac{\sigma_i^2}{m_i}$, which is reduced from original.

Clearly, if the probability distribution is not homogeneous within each region, the stratified sampling would increase the bias and makes the estimate less accurate.

- **Rao-Blackwellization:** An approach to reduce the variance of an estimator is to use the *conditioning inequality*

$$\text{var}(\mathbb{E}[\delta(\mathbf{X})|\mathbf{Z}]) \leq \text{var}(\mathbb{E}[\delta(\mathbf{X})]) \quad (4)$$

sometimes called **Rao-Blackwellization** [Liu, 2001] because the inequality is associated with the *Rao-Blackwell Theorem* [Lehmann and Casella, 2006], although the conditioning is

not always in terms of sufficient statistics. This method reflects a **basic principle**: *one should carry out analytical computation as much as possible*. The more information available for sampler, the less variance it would have.

Suppose the sample can be decomposed into two parts (\mathbf{X}, \mathbf{Z}) and the conditional expectation can be carried out explicitly $\mathbb{E}[h(\mathbf{X})|\mathbf{Z}]$ for each sub-population $(\mathbf{X}, \mathbf{Z}_i)$, then the estimator

$$\hat{I}_m = \frac{1}{m} \sum_{i=1}^m \mathbb{E}[h(\mathbf{X})|\mathbf{Z}_i]$$

is unbiased and has lower variance than the direct sample average of $h(\mathbf{X}_j)$. In statistic, \hat{I}_m is often called **mixture estimator** since it combines a mixture of distributions for each sub-population.

- **Control Variates Methods**: In this method, one uses a **control variate** C that is highly correlated with sample \mathbf{X} to reduce the variance. Suppose that $\mu_C = \mathbb{E}[C]$ is known and the sample estimator

$$h(X) := X(b) = X + b(C - \mu_C),$$

which has the same mean as \mathbf{X} . Suppose that the optimal weight $b^* = \text{Cov}(X, C)/\text{var}(C)$ can be computed explicitly, the variance of $\mathbf{X}(b)$

$$\begin{aligned} \text{var}(X(b)) &= \text{var}(X) - 2b \text{Cov}(X, C) + b^2 \text{var}(C) \\ &= (1 - \rho_{X,C}^2) \text{var}(X) \leq \text{var}(X). \end{aligned}$$

Note that the technique of control variates is manageable only in very specific cases: the control function $\mu_C = \mathbb{E}[C]$ must be available, as well as the optimal weight b^* .

- **Antithetic Variates Methods**: This method describes a way of **creating negative correlated samples**. Suppose $X_1 = F^{-1}(U)$ for some random variable $U \sim \mathcal{U}[0, 1]$. Note that $X_2 = F^{-1}(1 - U)$ follows the same distribution as X_1 . X_2 is called **antithetic variables**. More generally, let g be a monotone function on $[0, 1]$, so that for any $u_1, u_2 \in [0, 1]$

$$(g(u_1) - g(u_2))(g(1 - u_1) - g(1 - u_2)) \leq 0.$$

Then for two i.i.d random variables $U_1, U_2 \sim \mathcal{U}[0, 1]$, $X_1 = g(U)$ and $X_2 = g(1 - U)$ we have

$$\mathbb{E}[(g(U_1) - g(U_2))(g(1 - U_1) - g(1 - U_2))] = \text{Cov}(X_1, X_2) \leq 0$$

Thus the variance of $\frac{1}{2}(X_1 + X_2)$ is less than the variance of two independent samples from Monte Carlo simulator.

2.3 Importance Sampling

- The vanilla rejection sampling method suffers from **low sample efficiency and slow convergence** since it wastes a lot of effort evaluating random samples located in regions where the target function value $h(\mathbf{x})p(\mathbf{x})$ is almost zero.
- The **importance sampling** idea suggests that one should **focus on regions of "importance"** so as to save computational resources. It is important to note that in high dimensional space, the support of the target distribution is exponentially small as compared to the entire region \mathcal{X} . In high dimensional setting, the **vanilla Monte Carlo methods are bound to fail**.

- Consider the expectation estimation

$$\begin{aligned}
\mu &= \mathbb{E}_{\mathbf{X} \sim p} [h(\mathbf{X})] = \int_{\mathcal{X}} h(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\
&= \mathbb{E}_{\mathbf{X} \sim g} \left[\frac{p(\mathbf{X})}{g(\mathbf{X})} h(\mathbf{X}) \right] := \mathbb{E}_{\mathbf{X} \sim g} [w(\mathbf{X}) h(\mathbf{X})] \\
&= \int_{\mathcal{X}} w(\mathbf{x}) h(\mathbf{x}) g(\mathbf{x}) d\mathbf{x} \\
&\text{where } w(\mathbf{x}) := \frac{p(\mathbf{x})}{g(\mathbf{x})} \text{ is the } \textit{importance weight}.
\end{aligned} \tag{5}$$

- The ***Importance Sampling Algorithm*** [Liu, 2001, Robert and Casella, 1999] is as below:

1. Draw $\mathbf{X}_1, \dots, \mathbf{X}_m \sim g(\mathbf{x})$ where g is trial distribution;
2. Calculate the *importance weight* $w(\mathbf{x})$:

$$w_i := w(\mathbf{X}_i) = \frac{p(\mathbf{X}_i)}{g(\mathbf{X}_i)}, \quad i = 1, \dots, m$$

3. Approximate μ by

$$\hat{\mu} = \frac{\sum_{i=1}^m w_i h(\mathbf{X}_i)}{\sum_{i=1}^m w_i} \tag{6}$$

Another way to approximate μ is via the unbiased estimator

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m w_i h(\mathbf{X}_i) \tag{7}$$

The normalized estimator (6) is ***biased*** but as $m \rightarrow \infty$ $\hat{\mu} \rightarrow \mu$, i.e. it is ***asymptotically unbiased***.

- While the (7) converges to $\mathbb{E}_p[h]$ almost surely, its variance

$$\mathbb{E}_g \left[\frac{p^2(\mathbf{X})}{g^2(\mathbf{X})} h^2(\mathbf{X}) \right] = \mathbb{E}_p \left[\frac{p(\mathbf{X})}{g(\mathbf{X})} h^2(\mathbf{X}) \right] = \int_{\mathcal{X}} \frac{p^2(\mathbf{x})}{g(\mathbf{x})} h^2(\mathbf{x}) d\mathbf{x} < \infty \tag{8}$$

need to be finite. Thus $\text{supp}(g) \supset \text{supp}(p)$. In practice, the importance sampling works ***poorly*** (high-amplitude ***jumps***, ***instability*** of the path of the average, ***slow convergence***) when

$$\int_{\mathcal{X}} \frac{p^2(\mathbf{x})}{g(\mathbf{x})} d\mathbf{x} = \infty.$$

A "Rule of thumb" estimate of the variance is $\text{var}(\hat{\mu}) \approx \text{var}_g(h) \mathbb{E}_g[w^2]$. Geweke (1989) mentions two types of *sufficient conditions* for finite variance estimator:

1. $p(\mathbf{x})/g(\mathbf{x}) < M$, $\forall \mathbf{x} \in \mathcal{X}$ and $\text{var}(h) < \infty$;
2. \mathcal{X} is compact, $p(\mathbf{x}) < F$ and $g(\mathbf{x}) > \epsilon$, $\forall \mathbf{x} \in \mathcal{X}$.

Note that the first condition implies that the rejection sampling also applies.

On the other hand, (6) can achieve finite variance and converges to $\mathbb{E}_p[h]$ when $m \rightarrow \infty$.

- The **relative efficiency** is defined as ratio between variance of vanilla Monte Carlo sampler and importance sampler:

$$\text{RE} = \frac{\frac{1}{m} \sum_{i=1}^m h(\mathbf{Y}_i)}{\frac{1}{m} \sum_{i=1}^m w_i h(\mathbf{X}_i)} \rightarrow 1,$$

if $w_i \rightarrow 1, \forall i$ or $\text{var}_g(w) \rightarrow 0$. Note that $\mathbb{E}_g[w] = \mathbb{E}_g\left[\frac{p}{g}\right] = 1$. Therefore $\text{RE} \leq 1$. This is the price we pay for sampling via simpler instrumental distribution g . If we use the normalized weight estimator (6), the relative efficiency is about

$$\text{RE} \approx \frac{1}{1 + \text{var}_g(w)} \quad (9)$$

The **effective sample size (ESS)** is thus

$$\text{ESS} \approx \frac{m}{1 + \text{var}_g(w)}.$$

These approximations do not involve h . Here we assume that

$$\mathbb{E}_p[(w(\mathbf{X}) - \mathbb{E}_p[w(\mathbf{X})])(h(\mathbf{X}) - \mu)] \text{ is small.}$$

- For optimal choice of g in (7), we have the following theorem

Theorem 2.3 [Robert and Casella, 1999] *The choice of g that minimizes the variance of the estimator (7) is*

$$g^*(z) = \frac{|h(z)| p(z)}{\int_{\mathcal{X}} |h(z)| p(z) dz}.$$

- There are several **advantages** using importance sampling:
 - For (6), one only need to know target distribution $p(\cdot)$ **up to a constant** when calculating the importance weight. This is very convenient for distributions such as the exponential families. Also (6) has lower variance with higher bias as compared to (7).
 - Importance sampling allows us to **approximate** the expectation of an unknown complex target distribution using **simple trial distribution** and then **correcting the bias** via **reweighting**. Similar to rejection sampling, the performance of importance sampling depends on the **closeness** of $g(\mathbf{x})$ to the target $h(\mathbf{x})p(\mathbf{x})$. In particular, the trial density $g(\mathbf{x})$ should have longer tail than $p(\mathbf{x})$ (i.e. $\text{supp}(g) \supset \text{supp}(p)$) so that the importance weight is well-defined in all support of target distribution. In high dimensional setting, it would still be challenging to find such a good trial distribution g .

3 Sequential Importance Sampling

3.1 Sequential Importance Sampling (SIS)

- It is nontrivial to find a good trial distribution g in high dimensional space. One of the most useful strategies in these problems is to build up the trial density **sequentially**.

- Denote $\mathbf{x} := [x_1, \dots, x_d]$. The trial distribution g can be *factorized* as

$$g(\mathbf{x}) = g(x_1) \prod_{i=2}^d g(x_i | x_1, \dots, x_{i-1})$$

by which we hope to obtain some information on target density while building up the trial density. Note that the target density $p(\mathbf{x})$ can also be factorized as

$$p(\mathbf{x}) = p(x_1) \prod_{i=2}^d p(x_i | x_1, \dots, x_{i-1}).$$

So the importance weight is calculated as

$$w(\mathbf{x}) = \frac{g(x_1) \prod_{i=2}^d g(x_i | x_1, \dots, x_{i-1})}{p(x_1) \prod_{i=2}^d p(x_i | x_1, \dots, x_{i-1})}. \quad (10)$$

From (10), we can reformulate the importance weight **recursively**

$$\begin{aligned} w_t(\mathbf{x}_t) &= w_{t-1}(\mathbf{x}_{t-1}) \frac{p(x_t | \mathbf{x}_{t-1})}{g(x_t | \mathbf{x}_{t-1})}, \quad t = 2, \dots, d \\ w_1(\mathbf{x}_1) &= \frac{p(x_1)}{g(x_1)} \end{aligned} \quad (11)$$

where $\mathbf{x}_t = (x_1, \dots, x_t)$.

- Note that computing $p(x_t | \mathbf{x}_{t-1})$ may be nontrivial since it require marginalization to find $p(\mathbf{x}_{t-1})$. Suppose we can find a sequence of ***auxiliary distributions*** $p_t(\mathbf{x}_t), t = 1, \dots, d$ so that $p_t(\mathbf{x}_t)$ is a reasonable approximation of the *marginal* $p(\mathbf{x}_t)$ and $p_d(\mathbf{x}) = p(\mathbf{x})$. Note that $p_t(\mathbf{x}_t)$ only need to be known up to a constant and only serves as a "guide" to construction of whole samples.

The ***Sequential Importance Sampling (SIS)*** algorithm is described as below:

1. Draw $X_t = x_t$ from $g(x_t | \mathbf{x}_{t-1})$ and let $\mathbf{x}_t \leftarrow [x_t, \mathbf{x}_{t-1}]$.
2. Compute the ***incremental weight***:

$$u_t = \frac{p_t(\mathbf{x}_t)}{p_{t-1}(\mathbf{x}_{t-1}) g(x_t | \mathbf{x}_{t-1})} \quad (12)$$

3. Compute $w_t \leftarrow w_{t-1} u_t$

It is easy to show that \mathbf{x}_t is properly weighted by w_t given that \mathbf{x}_{t-1} is properly weighted by w_{t-1} . Thus the whole sample \mathbf{x} obtained sequentially is properly weighted by the final importance w_d w.r.t. to target $p(\mathbf{x})$.

- The **benefits** for using sequential importance sampling include:
 - We can make use of the **characteristics of local factor** $p(x_t | \mathbf{x}_{t-1})$ in target density when designing trial density $g(x_t | \mathbf{x}_{t-1})$. An example is the *local Markov property* $p(x_t | \mathbf{x}_{t-1}) = p(x_t | x_{t-1})$ for probabilistic graphical models. By sequential sampling, we break the complex problem into smaller pieces.

- We can **stop** generating further components of \mathbf{x} if the *partial weights* w_k that derived from the sequentially generated the *partial samples* \mathbf{x}_k are *too small*. We can also **reject sample with small weight** and restart again. This way we avoid wasting effort generating samples with little effect on final estimation. This rejection process would introduce additional bias which should be corrected [Robert and Casella, 1999].
- The SIS algorithm is **attractive** since we can use a sequence of auxiliary distributions to construct more efficient sampling algorithm.

3.2 Nonlinear filtering

- Assume $\boldsymbol{\xi}_k \in \mathbb{R}^d$ is the state at time k and \mathbf{y}_k is the observation at time k . $\boldsymbol{\xi}_{1:t} = [\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_t] \in \mathbb{R}^{d \times t}$ is the trajectory of states up to t . The *state-space model* has the following equations

$$\begin{aligned}\boldsymbol{\xi}_t &\sim q_t(\boldsymbol{\xi}_t | \boldsymbol{\xi}_{t-1}, \theta) \quad (\text{state equation}) \\ \mathbf{y}_t &\sim f_t(\mathbf{y}_t | \boldsymbol{\xi}_t, \phi) \quad (\text{observation equation})\end{aligned}$$

- A main challenge is to find the efficient methods for *online estimation and prediction (filtering)* of the state $\boldsymbol{\xi}_t$ given sequential observations $\mathbf{y}_{1:t} = [\mathbf{y}_1, \dots, \mathbf{y}_t]$.
- By Bayes rule, the optimal online estimation of $\boldsymbol{\xi}_t$ given $\mathbf{y}_{1:t}$ is the *conditional mean estimator*

$$\begin{aligned}\hat{\boldsymbol{\xi}}_t &= \mathbb{E} [\boldsymbol{\xi}_t | \mathbf{y}_{1:t}] \\ &= \int \boldsymbol{\xi}_t p(\boldsymbol{\xi}_t | \mathbf{y}_{1:t}) d\boldsymbol{\xi}_t\end{aligned} \tag{13}$$

where the **posterior distribution** of state

$$\begin{aligned}p_t(\boldsymbol{\xi}_t) &:= p(\boldsymbol{\xi}_t | \mathbf{y}_{1:t}) = \frac{f_t(\mathbf{y}_t | \boldsymbol{\xi}_t, \phi) p(\boldsymbol{\xi}_t | \mathbf{y}_{1:(t-1)})}{p(\mathbf{y}_t | \mathbf{y}_{1:(t-1)})} \\ &= \frac{\int [f_t(\mathbf{y}_t | \boldsymbol{\xi}_t, \phi) q_t(\boldsymbol{\xi}_t | \boldsymbol{\xi}_{t-1}, \theta)] p(\boldsymbol{\xi}_{t-1} | \mathbf{y}_{1:(t-1)}) d\boldsymbol{\xi}_{t-1}}{p(\mathbf{y}_t | \mathbf{y}_{1:(t-1)})} \\ &\propto \int [q_t(\boldsymbol{\xi}_t | \boldsymbol{\xi}_{t-1}, \theta) f_t(\mathbf{y}_t | \boldsymbol{\xi}_t, \phi)] p_{t-1}(\boldsymbol{\xi}_{t-1}) d\boldsymbol{\xi}_{t-1}\end{aligned} \tag{14}$$

where $p_{t-1}(\boldsymbol{\xi}_{t-1}) := p(\boldsymbol{\xi}_{t-1} | \mathbf{y}_{1:(t-1)})$ is the "current" posterior distribution on state $\boldsymbol{\xi}_{t-1}$ and $p(\mathbf{y}_t | \mathbf{y}_{1:(t-1)})$ is a constant.

- There are several special cases
 - If q_t and f_t are both *linear Gaussian condition distributions*, the state-space model is called **linear state-space model** or **dynamic linear model**. The recursive update (14) can be computed analytically since the current posterior distribution $p(\boldsymbol{\xi}_{t-1} | \mathbf{y}_{1:(t-1)})$ is also Gaussian. The resulting algorithm is the basis of **Kalman filter** [Liu, 2001].
 - If $\boldsymbol{\xi} \in \Xi$ is discrete with finite element, the state-space model is called **Hidden Markov Model (HMM)**. The recursive update in (14) is close to the dynamic programming algorithm (*forward-backward algorithm* or *Baum-Welch*).
- Besides these two cases, exact computation the optimal online estimation $\hat{\boldsymbol{\xi}}_t = \mathbb{E} [\boldsymbol{\xi}_t | \mathbf{y}_{1:t}]$ is impossible since the computation of normalization constant becomes infeasible when t grows.

- Instead, we can approximate $\hat{\xi}_t = \mathbb{E}[\xi_t | \mathbf{y}_{1:t}]$ using sequential importance sampling (SIS). Define the trial distribution $g(\xi_t | \xi_{1:(t-1)}, \mathbf{y}_{1:t}) = q_t(\xi_t | \xi_{t-1})$. We choose the **auxiliary distribution** $p_t(\xi_{1:t}) = p(\xi_{1:t} | \mathbf{y}_{1:t})$ and **importance weight** is updated via

$$w_t = \frac{p(\xi_{1:t} | \mathbf{y}_{1:t})}{g(\xi_{1:t} | \mathbf{y}_{1:t})} \quad (15)$$

$$\begin{aligned} &= \frac{p(\xi_t, \xi_{1:(t-1)} | \mathbf{y}_t, \mathbf{y}_{1:(t-1)})}{g(\xi_t, \xi_{1:(t-1)} | \mathbf{y}_t, \mathbf{y}_{1:(t-1)})} = \frac{p(\xi_t, \mathbf{y}_t, \xi_{1:(t-1)} | \mathbf{y}_{1:(t-1)})}{g(\xi_t, \xi_{1:(t-1)} | \mathbf{y}_{1:t})} \frac{1}{p(\mathbf{y}_t | \mathbf{y}_{1:(t-1)})} \\ &\propto \frac{p(\xi_t, \mathbf{y}_t | \xi_{1:(t-1)}, \mathbf{y}_{1:(t-1)})}{g(\xi_t | \xi_{1:(t-1)}, \mathbf{y}_{1:t})} \frac{p(\xi_{1:(t-1)} | \mathbf{y}_{1:(t-1)})}{g(\xi_{1:(t-1)} | \mathbf{y}_{1:(t-1)})} \\ &= \frac{f_t(\mathbf{y}_t | \xi_t, \phi) q_t(\xi_t | \xi_{t-1})}{g(\xi_t | \xi_{1:(t-1)}, \mathbf{y}_{1:t})} w_{t-1} \\ &= \frac{f_t(\mathbf{y}_t | \xi_t, \phi) q_t(\xi_t | \xi_{t-1})}{q_t(\xi_t | \xi_{t-1})} w_{t-1} = f_t(\mathbf{y}_t | \xi_t, \phi) w_{t-1} := u_t w_{t-1}. \end{aligned} \quad (16)$$

The *incremental weight* $u_t = f_t(\mathbf{y}_t | \xi_t, \phi)$.

- Note that normalization on $w_{t+1}^{(j)} = w_1^{(j)} \prod_{k=1}^t f_k(\mathbf{y}_{k+1} | \xi_{k+1}^{(j)}, \phi)$ is equivalent to applying **softmax** operation on $\left\{ \prod_{k=1}^t f_k(\mathbf{y}_{k+1} | \xi_{k+1}^{(j)}, \phi) \right\}_j$. As $k \rightarrow \infty$, we will see $w_{t+1}^{(j)} \rightarrow 1$ but all other $w_{t+1}^{(i)} \rightarrow 0 \forall i \neq j$. This phenomenon is called **particle degeneracy**.
- A simple method called **particle filter** (or **bootstrap filter**) is proposed to fix the particle degeneracy by **weight resampling** [Liu, 2001].

The **Sampling-Importance-Resampling (SIR)**

1. Draw $\xi_{t+1}^{(*,j)}$ from the state equation $q_t(\xi_{t+1} | \xi_t^{(j)}, \theta)$ for $j = 1, \dots, m$
2. Weight each draw by $w_{t+1}^{(j)} = f_t(\mathbf{y}_{t+1} | \xi_{t+1}^{(*,j)}, \phi) \tilde{w}_t^{(j)}$
3. Normalize $w_{t+1}^{(j)}$ as $\tilde{w}_{t+1}^{(j)}$
4. **Resample** from $\left\{ \xi_{t+1}^{(*,1)}, \dots, \xi_{t+1}^{(*,m)} \right\}$ according to multinomial distribution with probability $\left\{ \tilde{w}_{t+1}^{(j)} \right\}_{j=1}^m$ to produce a random sample $\left\{ \xi_{t+1}^{(1)}, \dots, \xi_{t+1}^{(m)} \right\}$ for time $t+1$

Averaging $\left\{ \xi_{t+1}^{(1)}, \dots, \xi_{t+1}^{(m)} \right\}$ will obtain the approximate conditional posterior mean estimator at $t+1$.

- The sequence $(\xi_1^{(j)}, \dots, \xi_t^{(j)})$ is called j -th **particle trajectory**. The particle filter maintains m particles in the system at each time stamp.
- It is easy to see that if the $\xi_t^{(j)}$ follow the "current" posterior distribution $p_t(\xi_t)$ and if m is large enough, then the new random samples $\left\{ \xi_{t+1}^{(1)}, \dots, \xi_{t+1}^{(m)} \right\}$ follows the updated posterior $p_{t+1}(\xi_{t+1})$ **approximately**.

References

- Erich L Lehmann and George Casella. *Theory of point estimation*. Springer Science & Business Media, 2006.
- Jun S Liu. *Monte Carlo strategies in scientific computing*, volume 10. Springer, 2001.
- Christian P Robert and George Casella. *Monte Carlo statistical methods*, volume 2. Springer, 1999.