

## Sequential Monte Carlo methods

### Lecture 3 – Monte Carlo and importance sampling

Thomas Schön, Uppsala University  
2017-08-24

## Outline – Lecture 3

**Aim:** Motivate and introduce the Monte Carlo idea and derive importance sampling.

### Outline:

1. Why do we need Monte Carlo?
2. The Monte Carlo idea
3. Importance sampling
4. Ex. joint filtering using importance sampling

1/16

## Why do we need Monte Carlo methods?

Probabilistic modelling often produce intractable optimization and/or integration problems.

Recall the nonlinear filtering problem or consider the computation of a point estimate via **expectation**, e.g. the conditional mean

$$\hat{x}_{t|t} = \mathbb{E}[X_t | y_{1:t}] = \int x_t p(x_t | y_{1:t}) dx_t.$$

Monte Carlo methods are **computational solutions** where the distributions of interest are approximated by a large number of  $N$  random samples called particles.

2/16

## Common test functions

Hence, Monte Carlo methods can be used to solve integrals like

$$\mathbb{E}[\varphi(X) | y_{1:t}] = \int \varphi(x) p(x | y_{1:T}) dx$$

Common test functions  $\varphi(x)$  include:

- Conditional mean  $\varphi(x) = x$  (previous slide)
- Indicator function  $\varphi(x) = I(x > \vartheta)$  for some threshold  $\vartheta$ , which provides an estimate of tail probabilities (modelling e.g. extreme events).
- Covariances and other higher order moments.
- ...

3/16

## The Monte Carlo idea (I/II)

Let  $X \sim \pi(x)$ , where we refer to  $\pi(x)$  as the **target density**.

**(Very) restrictive assumption:** Assume that we have  $N$  samples  $\{x^i\}_{i=1}^N$  from the target density  $\pi(x)$ , making up an **empirical approximation**

$$\hat{\pi}^N(x) = \sum_{i=1}^N \frac{1}{N} \delta_{x^i}(x).$$

Allows for the following approximation of the integral,

$$\mathbb{E}_{\pi}[\varphi(X)] = \int \varphi(x) \pi(x) dx \approx \int \varphi(x) \sum_{i=1}^N \frac{1}{N} \delta_{x^i}(x) dx = \frac{1}{N} \sum_{i=1}^N \varphi(x^i)$$

$$\text{" } \int + \delta \rightarrow \sum \text{"}$$

4/16

## The Monte Carlo idea (II/II)

The integral

$$I(\varphi) = \mathbb{E}_{\pi}[\varphi(X)] = \int \varphi(x) \pi(x) dx$$

is approximated by

$$\hat{I}_N(\varphi) = \frac{1}{N} \sum_{i=1}^N \varphi(x^i).$$

The strong law of large numbers tells us that

$$\hat{I}_N(\varphi) \xrightarrow{\text{a.s.}} I(\varphi), \quad N \rightarrow \infty,$$

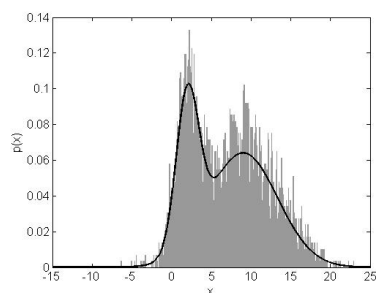
and the central limit theorem states that

$$\frac{\sqrt{N} (\hat{I}_N(\varphi) - I(\varphi))}{\sigma_{\varphi}} \xrightarrow{d} \mathcal{N}(0, 1), \quad N \rightarrow \infty.$$

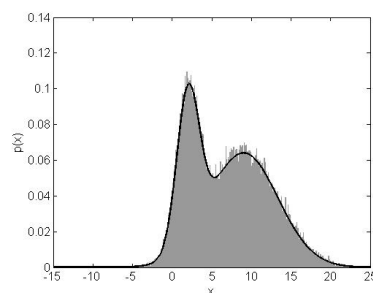
5/16

## The Monte Carlo idea – toy illustration

$$\pi(x) = 0.3\mathcal{N}(x | 2, 2) + 0.7\mathcal{N}(x | 9, 19)$$



5 000 samples



50 000 samples

**Obvious problem:** In general we are **not** able to directly sample from the density we are interested in.

6/16

## Importance sampling

## Importance sampling – proposal distribution

The **proposal distribution**<sup>1</sup> is chosen by the user:

1. It should be simple to sample from and
2. we require  $q(\mathbf{x}) > 0$  for all  $\mathbf{x}$  where  $\pi(\mathbf{x}) > 0$

**Idea:** Choose the **proposal** density  $q(\mathbf{x})$  such that it is easy to generate samples from it and somehow **compensate** for the mismatch between the target and the proposal.

<sup>1</sup>a.k.a. importance distribution or instrumental distribution.

7/16

## Point-wise evaluation of the target

It is often the case that the target density  $\pi(\mathbf{x})$  can only be evaluated "up to an unknown normalization constant  $Z$ ",

$$\pi(\mathbf{x}) = \frac{\tilde{\pi}(\mathbf{x})}{Z}$$

where  $\tilde{\pi}(\mathbf{x})$  can be evaluated for any  $\mathbf{x}$ , but the constant  $Z$  is unknown.

Ex. (nonlinear joint filtering problem) The target density given by  $\pi(\mathbf{x}) = p(\mathbf{x}_{0:t} | y_{1:t})$  and we have

$$\underbrace{p(\mathbf{x}_{0:t} | y_{1:t})}_{\pi(\mathbf{x})} = \frac{\overbrace{p(\mathbf{x}_{0:t}, y_{1:t})}^{\tilde{\pi}(\mathbf{x})}}{\underbrace{p(y_{1:t})}_Z},$$

where we can evaluate  $\tilde{\pi}(\mathbf{x}) = p(\mathbf{x}_{0:t}, y_{1:t})$  point-wise, but  $Z = p(y_{1:t})$  is intractable in general.

8/16

## Self-normalized importance sampling

### Self-normalized importance sampling

Insert

$$\pi(\mathbf{x}) = \frac{\tilde{\pi}(\mathbf{x})}{Z}$$

into the importance sampling integral results in

$$I(\varphi) = \mathbb{E}[\varphi(\mathbf{X})] = \int \varphi(\mathbf{x}) \frac{\tilde{\pi}(\mathbf{x})}{Z q(\mathbf{x})} q(\mathbf{x}) d\mathbf{x} = \frac{1}{Z} \int \varphi(\mathbf{z}) \underbrace{\frac{\tilde{\pi}(\mathbf{z})}{q(\mathbf{z})}}_{=\omega(\mathbf{z})} q(\mathbf{z}) d\mathbf{z}$$

Hence, the importance sampling estimator is

$$\hat{I}^N(\varphi) = \frac{1}{NZ} \sum_{i=1}^N \tilde{w}^i \varphi(\mathbf{x}^i),$$

where  $\tilde{w}^i = \omega(\mathbf{x}^i)$ .

The normalization constant  $Z$  is still problematic.

9/16

## Self-normalized importance sampling

The normalization constant is given by the following integral

$$Z = \int \tilde{\pi}(\mathbf{x}) d\mathbf{x},$$

which we can approximate using our samples  $\{\mathbf{x}^i\}_{i=1}^N$  from  $q(\mathbf{x})$ .

The result is

$$Z = \int \frac{\tilde{\pi}(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) d\mathbf{x} \approx \frac{1}{N} \sum_{i=1}^N \tilde{w}^i$$

The **self-normalized** importance sampling estimate is obtained by inserting this into  $\tilde{I}^N(\varphi)$ ,

$$\tilde{I}^N(\varphi) = \sum_{i=1}^N w^i \varphi(\mathbf{x}^i), \quad w^i = \frac{\tilde{w}^i}{\sum_{j=1}^N \tilde{w}^j}$$

10/16

## Self-normalized importance sampling

### Algorithm 1 Importance sampler

1. Sample  $\mathbf{x}^i \sim q(\mathbf{x})$ .
2. Compute the weights  $\tilde{w}^i = \tilde{\pi}(\mathbf{x}^i)/q(\mathbf{x}^i)$ .
3. Normalize the weights  $w^i = \tilde{w}^i / \sum_{j=1}^N \tilde{w}^j$ .

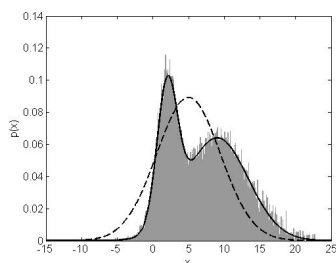
Each step is carried out for  $i = 1, \dots, N$ .

The convergence of the resulting approximation  $\hat{\pi}^N(\mathbf{x}) = \sum_{i=1}^N w^i \delta_{\mathbf{x}^i}(\mathbf{x})$  is since long well established.

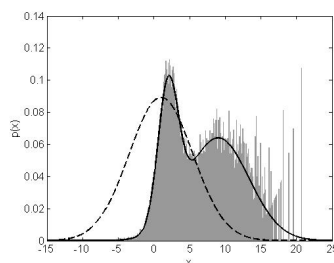
The fact that we are sampling from a user-chosen proposal distribution  $q(\mathbf{x})$  is corrected for by the weights, which **accounts for the discrepancy** between the proposal  $q(\mathbf{x})$  and the target  $\pi(\mathbf{x})$ .

11/16

## The importance of a good proposal density



$q_1(\mathbf{x}) = \mathcal{N}(\mathbf{x} | 5, 20)$  (dashed)



$q_2(\mathbf{x}) = \mathcal{N}(\mathbf{x} | 1, 20)$  (dashed)

50 000 samples were used in both simulations.

**Lesson learned:** It is important to be careful in selecting the proposal distribution.

12/16

## Ex) Importance sampling of the joint filtering PDF

**Problem statement:** Use importance sampling to compute the joint filtering PDF  $p(\mathbf{x}_{1:t} | \mathbf{y}_{1:t})$  for  $(\mathbf{x} = \mathbf{x}_{1:t}, \pi(\mathbf{x}) = p(\mathbf{x}_{1:t} | \mathbf{y}_{1:t}))$

$$\begin{aligned} \mathbf{X}_{t+1} | (\mathbf{X}_t = \mathbf{x}_t) &\sim p(\mathbf{x}_{t+1} | \mathbf{x}_t), & \mathbf{X}_{t+1} &= f(\mathbf{X}_t) + \mathbf{V}_t, \\ \mathbf{Y}_t | (\mathbf{X}_t = \mathbf{x}_t) &\sim p(\mathbf{y}_t | \mathbf{x}_t), & \mathbf{Y}_t &= g(\mathbf{X}_t) + \mathbf{E}_t, \\ \mathbf{X}_0 &\sim p(\mathbf{x}_0), & \mathbf{X}_0 &\sim p(\mathbf{x}_0). \end{aligned}$$

**Key challenge:** Nontrivial to design proposal distributions for high-dimensional problems. Here the dimension of the space  $\mathcal{X}^t$  grows with  $t$ ! ( $\mathbf{x}_t \in \mathcal{X}$ ).

13/16

## Ex) Importance sampling of the joint filtering pdf

**Idea:** Reuse computations over time by exploiting the sequential structure of the SSM via a proposal distribution that factorizes as

$$q(\mathbf{x}_{0:t} | y_{1:t}) = q(\mathbf{x}_0) \prod_{s=1}^t q(\mathbf{x}_s | \mathbf{x}_{0:s-1}, y_{1:s}) = q(\mathbf{x}_0) \prod_{s=1}^t q(\mathbf{x}_s | \mathbf{x}_{s-1}, y_s)$$

Next we derive the weight function

$$\begin{aligned} \omega_t(\mathbf{x}_{0:t}) &= \frac{\tilde{\pi}(\mathbf{x}_{0:t})}{q(\mathbf{x}_{0:t})} = \frac{p(\mathbf{x}_{0:t}, y_{1:t})}{q(\mathbf{x}_{0:t} | y_{1:t})} = \dots \\ &= \frac{p(y_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{x}_{t-1})}{q(\mathbf{x}_t | \mathbf{x}_{t-1}, y_t)} \underbrace{\frac{p(\mathbf{x}_{0:t-1}, y_{1:t-1})}{q(\mathbf{x}_{0:t-1} | y_{1:t-1})}}_{\omega_{t-1}(\mathbf{x}_{0:t-1})} \end{aligned}$$

Hence, the weights can also be computed sequentially

$$\tilde{w}_t = \frac{p(y_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{x}_{t-1})}{q(\mathbf{x}_t | \mathbf{x}_{t-1}, y_t)} \tilde{w}_{t-1}$$

14/16

## Ex) Importance sampling of the joint filtering pdf

**Sequential importance sampling:** New samples are proposed sequentially and weights are computed sequentially.

**Show stopper:** It can be shown that the variance of the weights will grow unboundedly (**weight degeneracy**).

Practical consequence of weight degeneracy: after some time there will only be one weight with non-zero value (more in lecture 5).

---

Next lecture we will derive a working importance sampler by directly target the (marginal) filtering density  $p(\mathbf{x}_t | y_{1:t})$ .

Note that the dimension of  $\mathbf{x}_t$  is fixed, whereas the dimension of  $p(\mathbf{x}_{0:t} | y_{1:t})$  grows with  $t$ .

15/16

## A few concepts to summarize lecture 3

**Monte Carlo method:** Computational method making use of random sampling to obtain numerical solutions.

**Target density:** The probability density function that we are interested in.

**Empirical approximation:** An approximation of a distribution made up of weighted samples.

**Importance sampling:** A general technique for estimating properties of some target distribution when we only have access to samples from a distribution that is different from the target distribution.

**Proposal distribution:** A user-chosen distribution that it should be simple to sample from.

**Sequential importance sampling:** An importance sampler where the proposal distribution is defined sequentially and where the weights can be evaluated sequentially.

16/16