

Advanced Simulation - Lecture 4

Patrick Rebeschini

January 24th, 2018

Generic Sampling Methods

- For Monte Carlo methods, you need samples from distributions.
- Seen: inversion, transformation, composition, rejection.
- Today: importance sampling.

Importance Sampling

- We want to compute

$$I = \mathbb{E}_{\pi}(\varphi(X)) = \int_{\mathbb{X}} \varphi(x) \pi(x) dx.$$

- We do not know how to sample from the target π but have access to a proposal distribution of density q .
- We only require that

$$\pi(x) > 0 \Rightarrow q(x) > 0;$$

i.e. the support of q includes the support of π .

- q is called the proposal, or importance, distribution.

Importance Sampling

- We have the following identity

$$I = \mathbb{E}_\pi(\varphi(X)) = \mathbb{E}_q(\varphi(X)w(X)),$$

where $w : \mathbb{X} \rightarrow \mathbb{R}^+$ is the importance weight function

$$w(x) = \frac{\pi(x)}{q(x)}.$$

- Hence for $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} q$,

$$\hat{I}_n^{\text{IS}} = \frac{1}{n} \sum_{i=1}^n \varphi(X_i)w(X_i).$$

- It can be interpreted as performing the following approximation of π

$$\hat{\pi}_n^{\text{IS}}(dx) = \frac{1}{n} \sum_{i=1}^n w(X_i)\delta_{X_i}(dx).$$

Importance Sampling Properties

■ **SLLN:** If $\mathbb{E}_q(|\varphi(X)|w(X)) < \infty$ then $\lim_{n \rightarrow \infty} \widehat{I}_n^{\text{IS}} = I$.

■ **Unbiased:** $\mathbb{E}_q(\widehat{I}_n^{\text{IS}}) = I$.

■ **Variance & CLT:** $\mathbb{V}_q(\widehat{I}_n^{\text{IS}}) = \sigma_{\text{IS}}^2/n$ where

$$\sigma_{\text{IS}}^2 := \mathbb{V}_q(\varphi(X)w(X))$$

and

$$\lim_{n \rightarrow \infty} \sqrt{n}(\widehat{I}_n^{\text{IS}} - I) \xrightarrow{D} \mathcal{N}(0, \sigma_{\text{IS}}^2).$$

Importance Sampling: Practical Advices

- Consistency does not require $\sigma_{\text{IS}}^2 < \infty$ but highly recommended in practice (!).
- **Sufficient condition:** If $\mathbb{E}_{\pi}(\varphi^2(X)) < \infty$ and $w(x) \leq M$ for all x for some $M < \infty$, then $\sigma_{\text{IS}}^2 < \infty$.
- In practice ensure $w(x) \leq M$ although it is neither necessary nor sufficient, as seen in the following example.

Importance Sampling: Example

- $\pi(x) = \mathcal{N}(x; 0, 1)$, $q(x) = \mathcal{N}(x; 0, \sigma^2)$.
- For $\sigma^2 \geq 1$, $w(x) \leq M$ for all x ,
and for $\sigma^2 < 1$, $w(x) \rightarrow \infty$ as $|x| \rightarrow \infty$.
- For $\varphi(x) = x^2$, we have $\sigma_{\text{IS}}^2 < \infty$ for all $\sigma^2 > 1/2$.
- For $\varphi(x) = \exp\left(\frac{\beta}{2}x^2\right)$, we have $I < \infty$ for $\beta < 1$
but $\sigma_{\text{IS}}^2 = \infty$ for $\beta > 1 - \frac{1}{2\sigma^2}$.

Optimal Importance Distribution

- **Proposition:** The optimal proposal minimising $\mathbb{V}_q \left(\widehat{I}_n^{\text{IS}} \right)$ is given by

$$q_{\text{opt}}(x) = \frac{|\varphi(x)| \pi(x)}{\int_{\mathbb{X}} |\varphi(x)| \pi(x) dx}.$$

- **Proof.** We have indeed

$$\mathbb{V}_q(\varphi(X)w(X)) = \mathbb{E}_q(\varphi^2(X)w^2(X)) - I^2.$$

For $q = q_{\text{opt}}$, we have

$$\begin{aligned} \mathbb{E}_{q_{\text{opt}}}(\varphi^2(X)w^2(X)) &= \int_{\mathbb{X}} \frac{\varphi^2(x)\pi^2(x)}{|\varphi(x)|\pi(x)} dx \cdot \int_{\mathbb{X}} |\varphi(x)|\pi(x) dx \\ &= \left(\int_{\mathbb{X}} |\varphi(x)|\pi(x) dx \right)^2 \end{aligned}$$

We also have by Jensen's inequality for any q

$$\mathbb{E}_q(\varphi^2(X)w^2(X)) \geq \mathbb{E}_q^2(|\varphi(X)|w(X)) = \left(\int_{\mathbb{X}} |\varphi(x)|\pi(x) dx \right)^2.$$

Optimal Importance Distribution

- $q_{\text{opt}}(x)$ can never be used in practice!
- For $\varphi(x) > 0$ we have $q_{\text{opt}}(x) = \varphi(x)\pi(x)/I$ and $\mathbb{V}_{q_{\text{opt}}}(\widehat{I}_n^{\text{IS}}) = 0$ but this is because

$$\varphi(x) w(x) = \varphi(x) \frac{\pi(x)}{q_{\text{opt}}(x)} = I,$$

it requires knowing I !

- This can be used as a guideline to select q ; i.e. select $q(x)$ such that $q(x) \approx q_{\text{opt}}(x)$.
- Particularly interesting in rare event simulation, not quite in statistics.

Normalised Importance Sampling

- Standard IS has limited applications in statistics as it requires knowing $\pi(x)$ and $q(x)$ exactly.
- Assume $\pi(x) = C_\pi \times \pi_u(x)$ and $q(x) = C_q \times q_u(x)$, $\pi(x) > 0 \Rightarrow q(x) > 0$ and define

$$w_u(x) = \frac{\pi_u(x)}{q_u(x)}.$$

- An alternative identity is

$$I = \mathbb{E}_\pi(\varphi(X)) = \frac{\int_{\mathbb{X}} \varphi(x) w_u(x) q(x) dx}{\int_{\mathbb{X}} w_u(x) q(x) dx}.$$

- Let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} q$ then

$$\widehat{I}_n^{\text{NIS}} = \frac{\sum_{i=1}^n \varphi(X_i) w_u(X_i)}{\sum_{i=1}^n w_u(X_i)}$$

is strongly consistent through the SLLN as long as $\mathbb{E}_q(|\varphi(X)| w(X)) < \infty$.

- Variance of IS:

$$\mathbb{V}(\widehat{I}_n^{\text{IS}}) = \frac{1}{n} \int \frac{(\varphi(x)\pi(x) - Iq(x))^2}{q(x)} dx$$

while variance of NIS (using the Delta method):

$$\mathbb{V}(\widehat{I}_n^{\text{NIS}}) = \frac{1}{n} \int \frac{\pi(x)^2 (\varphi(x) - I)^2}{q(x)} dx.$$

Details on the Delta method

If

$$\sqrt{n} \left(\begin{pmatrix} X_n \\ Y_n \end{pmatrix} - \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix} \right) \xrightarrow{D} \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \right)$$

and $g : (x \ y)^T \mapsto g(x, y)$ then

$$\sqrt{n} \left(g \begin{pmatrix} X_n \\ Y_n \end{pmatrix} - g \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix} \right) \xrightarrow{D} \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \nabla g|_{\mu}^T \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \nabla g|_{\mu} \right).$$

With $g : (x \ y)^T \mapsto x/y$ we have

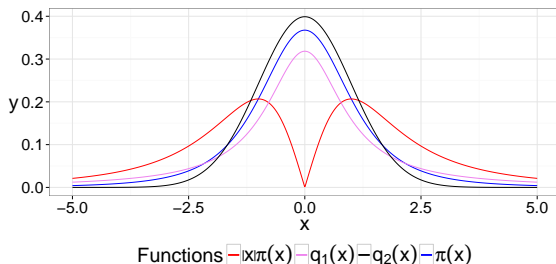
$$\nabla g : (x \ y)^T \mapsto \left(\frac{1}{y} \quad -\frac{x}{y^2} \right)^T$$

and thus

$$\nabla g|_{\mu} = \left(\frac{1}{\mu_y} \quad -\frac{\mu_x}{\mu_y^2} \right)^T$$

Toy Example: t-distribution

- We want to compute $I = \mathbb{E}_\pi(|X|)$ where $\pi(x) \propto (1 + x^2/3)^{-2}$ (t_3 -distribution).
- 1 Directly sample from π .
- 2 Use $q_1(x) = g_{t_1}(x) \propto (1 + x^2)^{-1}$ (t_1 -distribution).
- 3 Use $q_2(x) \propto \exp(-x^2/2)$ (normal).



Toy Example: t-distribution

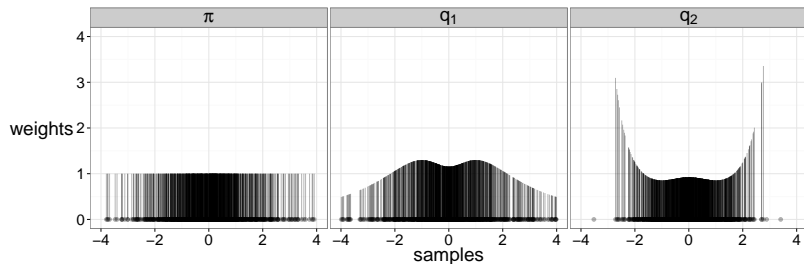


Figure: Sample weights obtained for 1000 realisations of X_i , from the different proposal distributions.

Toy Example: t-distribution

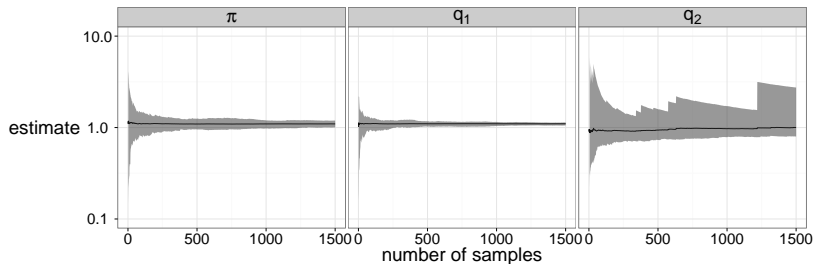


Figure: Estimates \hat{I}_n of I obtained after 1 to 1500 samples. The grey shaded areas correspond to the range of 100 independent replications.

Variance of importance sampling estimators

- Standard Importance Sampling: $X_1, \dots, X_n \stackrel{iid}{\sim} q$,

$$\widehat{I}_n^{\text{IS}} = \frac{1}{n} \sum_{i=1}^n \varphi(X_i) w(X_i).$$

- Asymptotic Variance:

$$\begin{aligned} \mathbb{V}_{as} \left(\widehat{I}_n^{\text{IS}} \right) &= \mathbb{E}_q \left[\left(\varphi(X) w(X) - \mathbb{E}_q \left(\varphi(X) w(X) \right) \right)^2 \right] \\ &\approx \frac{1}{n} \sum_{i=1}^n \left(\varphi(X_i) w(X_i) - \widehat{I}_n^{\text{IS}} \right)^2. \end{aligned}$$

- Thus the asymptotic variance can be estimated consistently with

$$\frac{1}{n} \sum_{i=1}^n \left(\varphi(X_i) w(X_i) - \widehat{I}_n^{\text{IS}} \right)^2.$$

Variance of importance sampling estimators

- Normalised Importance Sampling: $X_1, \dots, X_n \stackrel{iid}{\sim} q$,

$$\hat{I}_n^{\text{NIS}} = \frac{\sum_{i=1}^n \varphi(X_i) w_u(X_i)}{\sum_{i=1}^n w_u(X_i)}.$$

- Asymptotic Variance:

$$\mathbb{V}_{as} \left(\hat{I}_n^{\text{NIS}} \right) = \frac{\mathbb{E}_q \left[(\varphi(X)w(X) - I \times w(X))^2 \right]}{\mathbb{E}_q [w(X)]^2}.$$

- Thus the asymptotic variance can be estimated consistently with

$$\frac{\frac{1}{n} \sum_{i=1}^N w_u(X_i)^2 \left(\varphi(X_i) - \hat{I}_n^{\text{NIS}} \right)^2}{\left(\frac{1}{n} \sum_{i=1}^N w_u(X_i) \right)^2}.$$

- If only one weight, say $w_u(X_j)$, is significant compared to the others, then

$$\widehat{I}_n^{\text{NIS}} = \frac{\sum_{i=1}^n \varphi(X_i) w_u(X_i)}{\sum_{i=1}^n w_u(X_i)} \approx \varphi(X_j).$$

The “effective sample size” is one.

- To how many unweighted samples correspond our weighted samples of size n ? Solve for n_e in

$$\frac{1}{n} \mathbb{V}_{as} \left(\widehat{I}_n^{\text{NIS}} \right) = \frac{\sigma^2}{n_e},$$

where σ^2/n_e corresponds to the variance of an unweighted sample of size n_e .

- We solve by matching $\varphi(X_i) - \hat{I}^{\text{NIS}}$ with $\varphi(X_i) - I \approx \sigma$ as if they were i.i.d samples:

$$\frac{1}{n} \frac{\frac{1}{n} \sum_{i=1}^N w_u(X_i)^2 \left(\varphi(X_i) - \hat{I}_n^{\text{NIS}} \right)^2}{\left(\frac{1}{n} \sum_{i=1}^N w_u(X_i) \right)^2} \approx \frac{\sigma^2}{n_e}$$

i.e.

$$\frac{1}{n} \frac{\frac{1}{n} \sum_{i=1}^N w_u(X_i)^2}{\left(\frac{1}{n} \sum_{i=1}^N w_u(X_i) \right)^2} = \frac{1}{n_e}.$$

- The solution is

$$n_e = \frac{\left(\sum_{i=1}^n w_u(X_i) \right)^2}{\sum_{i=1}^n w_u(X_i)^2},$$

and is called the effective sample size.

Rejection and Importance Sampling in High Dimensions

- **Toy example:** Let $\mathbb{X} = \mathbb{R}^d$ and

$$\pi(x) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\sum_{i=1}^d x_i^2}{2}\right)$$

and

$$q(x) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{\sum_{i=1}^d x_i^2}{2\sigma^2}\right).$$

- How do Rejection sampling and Importance sampling scale in this context?

Performance of Rejection Sampling

- We have

$$w(x) = \frac{\pi(x)}{q(x)} = \sigma^d \exp\left(-\frac{\sum_{i=1}^d x_i^2}{2} \left(1 - \frac{1}{\sigma^2}\right)\right) \leq \sigma^d$$

for $\sigma > 1$.

- Acceptance probability is

$$\mathbb{P}(X \text{ accepted}) = \frac{1}{\sigma^d} \rightarrow 0 \text{ as } d \rightarrow \infty,$$

i.e. exponential degradation of performance.

- For $d = 100$, $\sigma = 1.2$, we have

$$\mathbb{P}(X \text{ accepted}) \approx 1.2 \times 10^{-8}.$$

Performance of Importance Sampling

- We have

$$w(x) = \sigma^d \exp\left(-\frac{\sum_{i=1}^d x_i^2}{2} \left(1 - \frac{1}{\sigma^2}\right)\right).$$

- Variance of the weights:

$$\mathbb{V}_q[w(X)] = \left(\frac{\sigma^4}{2\sigma^2 - 1}\right)^{d/2} - 1$$

where $\sigma^4 / (2\sigma^2 - 1) > 1$ for any $\sigma^2 > 1/2$.

- For $d = 100$, $\sigma = 1.2$, we have

$$\mathbb{V}_q[w(X)] \approx 1.8 \times 10^4.$$

Wait a minute...

Lecture 1:

- Simpson's rule for approximating integrals: error in $\mathcal{O}(n^{-1/d})$.

Lecture 2:

- Monte Carlo for approximating integrals: error in $\mathcal{O}(n^{-1/2})$ with rate independent of d .

And now:

- Importance Sampling standard deviation in the Gaussian example in $\exp(d)n^{-1/2}$.

The rate is indeed independent of d but the “constant” (in n) explodes exponentially (in d).

Markov chain Monte Carlo

- Revolutionary idea introduced by Metropolis et al., J. Chemical Physics, 1953.
- **Key idea:** Given a target distribution π , build a Markov chain $(X_t)_{t \geq 1}$ such that, as $t \rightarrow \infty$, $X_t \sim \pi$ and

$$\frac{1}{n} \sum_{t=1}^n \varphi(X_t) \rightarrow \int \varphi(x) \pi(x) dx$$

when $n \rightarrow \infty$ e.g. almost surely.

- Central limit theorems with a rate in $1/\sqrt{n}$.
- In some cases the constant (in n) does not explode exponentially with the dimension d , but polynomially.

Side Dish: Control Variates

- Variance reduction techniques, not always applicable but useful in some cases.
- Suppose that we want to compute

$$I = \int \varphi(x)\pi(x)dx$$

and that we know exactly

$$J = \int \psi(x)\pi(x)dx.$$

- Sample X_1, \dots, X_n from π and compute

$$\hat{I}_n = \frac{1}{n} \sum_{i=1}^n (\varphi(X_i) - \lambda(\psi(X_i) - J)).$$

- What is the benefit of \hat{I}_n over the standard Monte Carlo estimator?