

Approximate Zero-Variance Simulation

Pierre L'Ecuyer

Université de Montréal, Canada

Bruno Tuffin

INRIA - Rennes Bretagne Atlantique, France

Winter Simulation Conference, Miami, December 2008

Outline

Introduction

Asymptotic framework for rare events

Approximate zero variance via importance sampling

A general Markov chain model

Approximate zero variance via control variates

Longer example: highly reliable Markovian system

Approximations and adaptive learning

Conclusions

For references and links to history: see the paper.

Monte Carlo integration: basics

Want to estimate $\mu = \mathbb{E}[X]$ for some random variable X .

Monte Carlo (simulation), in basic form:

- ▶ Generate n independent copies of X , say X_1, \dots, X_n ;
- ▶ estimate μ by $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$.

Almost sure convergence as $n \rightarrow \infty$ (strong law of large numbers).

For **confidence interval**, can use **central limit theorem**:

$$\mathbb{P} \left[\mu \in \left(\bar{X}_n - \frac{c_\alpha S_n}{\sqrt{n}}, \bar{X}_n + \frac{c_\alpha S_n}{\sqrt{n}} \right) \right] \approx 1 - \alpha$$

for a given confidence level $1 - \alpha$, where S_n^2 is a consistent estimator of $\sigma^2 = \text{Var}[X]$.

Other types of situations: estimate quantile, optimization, etc.

Efficiency and variance reduction

Accuracy of estimator X can be measured by half-width of confidence interval, $c_\alpha n^{-1/2}\sigma$, or by **relative** half-width, $c_\alpha n^{-1/2}\sigma/\mu$.

Variance reduction: want new estimator \tilde{X} with same mean μ but smaller variance σ^2 . Gives a narrower confidence interval for the same n , but convergence speed seems to remain $O(n^{-1/2})$, **at first sight**.

Efficiency and variance reduction

Accuracy of estimator X can be measured by half-width of confidence interval, $c_\alpha n^{-1/2}\sigma$, or by **relative** half-width, $c_\alpha n^{-1/2}\sigma/\mu$.

Variance reduction: want new estimator \tilde{X} with same mean μ but smaller variance σ^2 . Gives a narrower confidence interval for the same n , but convergence speed seems to remain $O(n^{-1/2})$, **at first sight**.

Does it?

Efficiency and variance reduction

Accuracy of estimator X can be measured by half-width of confidence interval, $c_\alpha n^{-1/2}\sigma$, or by **relative** half-width, $c_\alpha n^{-1/2}\sigma/\mu$.

Variance reduction: want new estimator \tilde{X} with same mean μ but smaller variance σ^2 . Gives a narrower confidence interval for the same n , but convergence speed seems to remain $O(n^{-1/2})$, **at first sight**.

Does it? In fact, there are situations where the **rate can be improved**, say to $O(n^{-k})$ for $k > 1/2$, or even to $O(e^{-kn})$ for $k > 0$.

In other cases, μ and σ are parameterized by some parameter ε , and a well-crafted estimator may give us $\sigma(\varepsilon)/\mu(\varepsilon) \rightarrow 0$ when $\varepsilon \rightarrow 0$, whereas the original estimator has $\sigma(\varepsilon)/\mu(\varepsilon) \rightarrow \infty$.

Efficiency and variance reduction

Accuracy of estimator X can be measured by half-width of confidence interval, $c_\alpha n^{-1/2}\sigma$, or by **relative** half-width, $c_\alpha n^{-1/2}\sigma/\mu$.

Variance reduction: want new estimator \tilde{X} with same mean μ but smaller variance σ^2 . Gives a narrower confidence interval for the same n , but convergence speed seems to remain $O(n^{-1/2})$, **at first sight**.

Does it? In fact, there are situations where the **rate can be improved**, say to $O(n^{-k})$ for $k > 1/2$, or even to $O(e^{-kn})$ for $k > 0$.

In other cases, μ and σ are parameterized by some parameter ε , and a well-crafted estimator may give us $\sigma(\varepsilon)/\mu(\varepsilon) \rightarrow 0$ when $\varepsilon \rightarrow 0$, whereas the original estimator has $\sigma(\varepsilon)/\mu(\varepsilon) \rightarrow \infty$.

Bias (if any) and computational cost are also important factors for the quality of an estimator. Could look at **work-normalized mean-square error** (absolute or relative). In this talk, we focus on the variance only.

Ultimate dream: a zero-variance estimator

Can be achieved **in theory** via **importance sampling (IS)** or via **control variates (CV)**, as we will see later in the talk.

Unfortunately, an exact implementation is impractical: It requires the knowledge of μ (and usually much more) in the first place!

Ultimate dream: a zero-variance estimator

Can be achieved **in theory** via **importance sampling (IS)** or via **control variates (CV)**, as we will see later in the talk.

Unfortunately, an exact implementation is impractical: It requires the knowledge of μ (and usually much more) in the first place!

On the other hand, by plugging crude approximations of these unknown quantities in place of the exact ones in the zero-variance sampling strategies, we may reduce the variance tremendously, and sometimes the convergence rate as well.

We will discuss various (efficient) approximation schemes in the paper, with an emphasis on **rare-event** situations.

This is what we call **approximate zero-variance simulation**.

Ultimate dream: a zero-variance estimator

Can be achieved **in theory** via **importance sampling (IS)** or via **control variates (CV)**, as we will see later in the talk.

Unfortunately, an exact implementation is impractical: It requires the knowledge of μ (and usually much more) in the first place!

On the other hand, by plugging crude approximations of these unknown quantities in place of the exact ones in the zero-variance sampling strategies, we may reduce the variance tremendously, and sometimes the convergence rate as well.

We will discuss various (efficient) approximation schemes in the paper, with an emphasis on **rare-event** situations.

This is what we call **approximate zero-variance simulation**.

Has been studied for IS by **Booth (1985, 1987)**, **Kollman et al. (1999)**, **Baggerly et al. (2000)**, and for CV by **Henderson and Glynn 2002**, **Gobet and Maire (2006)**, and **Kim and Henderson (2006, 2007)**, among others.

A rare-event setting

We estimate a small quantity $\mu_0 = \mu_0(\varepsilon) > 0$, where $\mu_0(\varepsilon) \rightarrow 0$ when the rarity parameter $\varepsilon \rightarrow 0$, by an unbiased estimator $X = X(\varepsilon) \geq 0$.

A rare-event setting

We estimate a small quantity $\mu_0 = \mu_0(\varepsilon) > 0$, where $\mu_0(\varepsilon) \rightarrow 0$ when the **rarity parameter** $\varepsilon \rightarrow 0$, by an unbiased estimator $X = X(\varepsilon) \geq 0$.

In a **queueing** system with buffer size B and s servers, we can take $\varepsilon = 1/B$ if we are interested in very large values of B , and $\varepsilon = 1/s$ if we are interested in what happens when there is a large number of servers.

In a **reliability** model, the failure rates may be taken as polynomial functions of ε .

A rare-event setting

We estimate a small quantity $\mu_0 = \mu_0(\varepsilon) > 0$, where $\mu_0(\varepsilon) \rightarrow 0$ when the **rarity parameter** $\varepsilon \rightarrow 0$, by an unbiased estimator $X = X(\varepsilon) \geq 0$.

In a **queueing** system with buffer size B and s servers, we can take $\varepsilon = 1/B$ if we are interested in very large values of B , and $\varepsilon = 1/s$ if we are interested in what happens when there is a large number of servers.

In a **reliability** model, the failure rates may be taken as polynomial functions of ε .

The asymptotic behavior when $\varepsilon \rightarrow 0$ should be a good indicator of what happens when ε is very small.

Inaccuracy of standard Monte Carlo for rare events

With standard Monte Carlo, $\mu_0(\varepsilon)$ becomes harder to estimate as $\varepsilon \rightarrow 0$.

Example. Suppose $X(\varepsilon)$ is an **indicator function**, so $\mu_0(\varepsilon) = \mathbb{P}[X(\varepsilon) = 1]$. Then the **relative variance** (squared relative error) blows up:

$$\frac{\text{Var}[X]}{\mu_0^2(\varepsilon)} = \frac{1 - \mu_0(\varepsilon)}{\mu_0(\varepsilon)} \approx \frac{1}{\mu_0(\varepsilon)} \rightarrow \infty \text{ when } \varepsilon \rightarrow 0.$$

Standard Monte Carlo estimates $\mu_0(\varepsilon)$ by \bar{X}_n , and needs $n = \underline{O}(1/\mu_0(\varepsilon))$ for a meaningful estimate.

If $\mu_0(\varepsilon) = 10^{-10}$, for example, we need $n = 10^{14}$ for 1% relative error.

We then need more clever estimators.

Examples of situations where this happens

- ▶ Expected amount of radiation that crosses a given protection shield.
- ▶ Probability of a large loss from an investment portfolio.
- ▶ Value-at-risk (quantile estimation).
- ▶ Ruin probability for an insurance firm.
- ▶ Probability that the completion time of a large project exceeds a given threshold.
- ▶ Probability of buffer overflow, or mean time to overflow, in a queueing system.
- ▶ Proportion of packets lost in a communication system.
- ▶ Air traffic control.
- ▶ Mean time to failure or other reliability or availability measure for a highly reliable system (e.g., fault-tolerant computers, safety systems).

Classical robustness properties in this context

Commonly-used characterizations of $X(\varepsilon)$ in rare-event setting:

- ▶ It has **bounded relative error (BRE)** (bounded relative variance) if

$$\lim_{\varepsilon \rightarrow 0} \frac{\text{Var}[X(\varepsilon)]}{\mu_0^2(\varepsilon)} < \infty.$$

Classical robustness properties in this context

Commonly-used characterizations of $X(\varepsilon)$ in rare-event setting:

- ▶ It has **bounded relative error (BRE)** (bounded relative variance) if

$$\lim_{\varepsilon \rightarrow 0} \frac{\text{Var}[X(\varepsilon)]}{\mu_0^2(\varepsilon)} < \infty.$$

- ▶ It is **logarithmically efficient (LE)** or **asymptotically optimal** if

$$\lim_{\varepsilon \rightarrow 0} \frac{\ln \mathbb{E}[X^2(\varepsilon)]}{2 \ln \mu_0(\varepsilon)} = 1.$$

Means (roughly) that if $\mu_0(\varepsilon) \rightarrow 0$ at an exponential rate, then the standard deviation converges at least at the same exponential rate.

- ▶ BRE is stronger than LE, and can be more difficult to reach.

Generalization: BRM- k and LE- k

L., Blanchet, Glynn, Tuffin (2009)

An estimator $X(\varepsilon)$ with mean $\mu_0(\varepsilon)$ has **bounded relative moment of order k (BRM- k)** if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[X^k(\varepsilon)]}{\mu_0^k(\varepsilon)} < \infty.$$

Generalization: BRM- k and LE- k

L., Blanchet, Glynn, Tuffin (2009)

An estimator $X(\varepsilon)$ with mean $\mu_0(\varepsilon)$ has **bounded relative moment of order k (BRM- k)** if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[X^k(\varepsilon)]}{\mu_0^k(\varepsilon)} < \infty.$$

It has **logarithmic efficiency of order k (LE- k)** if

$$\lim_{\varepsilon \rightarrow 0} \frac{\ln \mathbb{E}[X^k(\varepsilon)]}{k \ln \mu_0(\varepsilon)} = 1.$$

Generalization: BRM- k and LE- k

L., Blanchet, Glynn, Tuffin (2009)

An estimator $X(\varepsilon)$ with mean $\mu_0(\varepsilon)$ has **bounded relative moment of order k (BRM- k)** if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[X^k(\varepsilon)]}{\mu_0^k(\varepsilon)} < \infty.$$

It has **logarithmic efficiency of order k (LE- k)** if

$$\lim_{\varepsilon \rightarrow 0} \frac{\ln \mathbb{E}[X^k(\varepsilon)]}{k \ln \mu_0(\varepsilon)} = 1.$$

Interesting and relevant for situations where we need estimators of the variance or of other moments higher than the mean.

Relevant for the validity of Berry-Esseen bound, for example.

Vanishing Relative Moment of order k (VRCM- k)

$X(\varepsilon)$ has vanishing relative centered moment of order k (VRCM- k) if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[|X(\varepsilon) - \mu_0(\varepsilon)|^k]}{\mu_0^k(\varepsilon)} = 0.$$

Vanishing Relative Moment of order k (VRCM- k)

$X(\varepsilon)$ has vanishing relative centered moment of order k (VRCM- k) if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[|X(\varepsilon) - \mu_0(\varepsilon)|^k]}{\mu_0^k(\varepsilon)} = 0.$$

Theorem: True if and only if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[X^k(\varepsilon)]}{\mu_0^k(\varepsilon)} = 1.$$

Vanishing Relative Moment of order k (VRCM- k)

$X(\varepsilon)$ has vanishing relative centered moment of order k (VRCM- k) if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[|X(\varepsilon) - \mu_0(\varepsilon)|^k]}{\mu_0^k(\varepsilon)} = 0.$$

Theorem: True if and only if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[X^k(\varepsilon)]}{\mu_0^k(\varepsilon)} = 1.$$

It has vanishing relative variance or relative error (VRE), if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\sigma(\varepsilon)}{\mu_0(\varepsilon)} = 0.$$

Vanishing Relative Moment of order k (VRCM- k)

$X(\varepsilon)$ has vanishing relative centered moment of order k (VRCM- k) if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[|X(\varepsilon) - \mu_0(\varepsilon)|^k]}{\mu_0^k(\varepsilon)} = 0.$$

Theorem: True if and only if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\mathbb{E}[X^k(\varepsilon)]}{\mu_0^k(\varepsilon)} = 1.$$

It has vanishing relative variance or relative error (VRE), if

$$\limsup_{\varepsilon \rightarrow 0} \frac{\sigma(\varepsilon)}{\mu_0(\varepsilon)} = 0.$$

When VRCM occurs, the rare event difficulty is reversed! May seem strange and perhaps unachievable at first sight, but does happen, in cases where zero-variance approximation improves when $\varepsilon \rightarrow 0$.

Goal in rare-event simul.: build estimators with these properties.

Importance Sampling (IS)

Suppose $X = h(Y)$ where Y is a random vector with density f .

Instead of generating Y from f , we can generate it from another density \tilde{f} , such that $\tilde{f}(y) > 0$ whenever $h(y)f(y) \neq 0$. We have

$$\mathbb{E}[X] = \int h(y)f(y)dy = \int \left[\frac{h(y)f(y)}{\tilde{f}(y)} \right] \tilde{f}(y)dy = \tilde{\mathbb{E}} \left[\frac{h(Y)f(Y)}{\tilde{f}(Y)} \right],$$

where $\tilde{\mathbb{E}}$ is the expectation under the new density.

Importance Sampling (IS)

Suppose $X = h(Y)$ where Y is a random vector with density f .

Instead of generating Y from f , we can generate it from another density \tilde{f} , such that $\tilde{f}(y) > 0$ whenever $h(y)f(y) \neq 0$. We have

$$\mathbb{E}[X] = \int h(y)f(y)dy = \int \left[\frac{h(y)f(y)}{\tilde{f}(y)} \right] \tilde{f}(y)dy = \tilde{\mathbb{E}} \left[\frac{h(Y)f(Y)}{\tilde{f}(Y)} \right],$$

where $\tilde{\mathbb{E}}$ is the expectation under the new density.

More generally, in a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let $X = h(\omega)$ where $h : \Omega \rightarrow \mathbb{R}$ and ω obeys \mathbb{P} . Instead of sampling ω from \mathbb{P} , we sample it from $\tilde{\mathbb{P}}$ such that $d\tilde{\mathbb{P}}(y) \neq 0$ when $h(y)d\mathbb{P}(y) \neq 0$. We have

$$\mathbb{E}[X] = \int_{\Omega} h(\omega)d\mathbb{P}(\omega) = \int_{\Omega} h(\omega) \frac{d\mathbb{P}(\omega)}{d\tilde{\mathbb{P}}(\omega)} d\tilde{\mathbb{P}}(\omega) = \tilde{\mathbb{E}} [h(\omega)L(\omega)]$$

where $L = d\mathbb{P}/d\tilde{\mathbb{P}}$.

For an unbiased estimator of $\mu_0 = \mathbb{E}[X]$ with IS, we generate $\omega_1, \dots, \omega_n$ i.i.d. under $\tilde{\mathbb{P}}$ and compute

$$\bar{X}_{\text{is},n} = \frac{1}{n} \sum_{i=1}^n h(\omega_i)L(\omega_i)$$

For an unbiased estimator of $\mu_0 = \mathbb{E}[X]$ with IS, we generate $\omega_1, \dots, \omega_n$ i.i.d. under $\tilde{\mathbb{P}}$ and compute

$$\bar{X}_{\text{is},n} = \frac{1}{n} \sum_{i=1}^n h(\omega_i)L(\omega_i)$$

We want to select $\tilde{\mathbb{P}}$ so that $X_{\text{is}} = h(\omega)L(\omega)$ (under $\tilde{\mathbb{P}}$) has smallest possible variance, and certainly smaller than that of X .

For an unbiased estimator of $\mu_0 = \mathbb{E}[X]$ with IS, we generate $\omega_1, \dots, \omega_n$ i.i.d. under $\tilde{\mathbb{P}}$ and compute

$$\bar{X}_{\text{is},n} = \frac{1}{n} \sum_{i=1}^n h(\omega_i)L(\omega_i)$$

We want to select $\tilde{\mathbb{P}}$ so that $X_{\text{is}} = h(\omega)L(\omega)$ (under $\tilde{\mathbb{P}}$) has smallest possible variance, and certainly smaller than that of X .

If $h \geq 0$, then taking $d\tilde{\mathbb{P}}(\omega)$ proportional to $h(\omega)d\mathbb{P}(\omega)$ gives **zero variance**:

$$X_{\text{is}} = h(\omega)L(\omega) = h(\omega) \frac{d\mathbb{P}(\omega)}{d\tilde{\mathbb{P}}(\omega)}$$

is the proportionality **constant**. This constant must be equal to μ_0 .

For an unbiased estimator of $\mu_0 = \mathbb{E}[X]$ with IS, we generate $\omega_1, \dots, \omega_n$ i.i.d. under $\tilde{\mathbb{P}}$ and compute

$$\bar{X}_{\text{is},n} = \frac{1}{n} \sum_{i=1}^n h(\omega_i)L(\omega_i)$$

We want to select $\tilde{\mathbb{P}}$ so that $X_{\text{is}} = h(\omega)L(\omega)$ (under $\tilde{\mathbb{P}}$) has smallest possible variance, and certainly smaller than that of X .

If $h \geq 0$, then taking $d\tilde{\mathbb{P}}(\omega)$ proportional to $h(\omega)d\mathbb{P}(\omega)$ gives **zero variance**:

$$X_{\text{is}} = h(\omega)L(\omega) = h(\omega) \frac{d\mathbb{P}(\omega)}{d\tilde{\mathbb{P}}(\omega)}$$

is the proportionality **constant**. This constant must be equal to μ_0 .

How can we implement this? Or approximate it?

VRCM implies convergence to a zero-variance sampling measure

Suppose

$$\mu_0(\varepsilon) = \mathbb{E}_{\mathbb{P}_\varepsilon}[X(\varepsilon)] = \int_{\Omega} X(\varepsilon, \omega) d\mathbb{P}_\varepsilon(\omega) \rightarrow 0$$

when $\varepsilon \rightarrow 0$. The zero-variance measure \mathbb{P}_ε^* here satisfies

$$\frac{d\mathbb{P}_\varepsilon^*(\omega)}{d\mathbb{P}_\varepsilon(\omega)} = \frac{X(\varepsilon, \omega)}{\mu_0(\varepsilon)}.$$

Proposition (L., Blanchet, Tuffin, Glynn 2009).

If $X(\varepsilon)$ is VRCM- $(1 + \delta)$ for some $\delta > 0$, then

$$\lim_{\varepsilon \rightarrow 0} \sup_{A \in \mathcal{F}} |\mathbb{P}_\varepsilon(A) - P_\varepsilon^*(A)| = 0.$$

That is, the sampling distribution **must** converge in total variation to the zero-variance measure associated with $X(\varepsilon)$, regardless of what sampling strategy we use (IS or not).

Proof:

$$\begin{aligned}
 \sup_{A \in \mathcal{F}} |\mathbb{P}_\varepsilon^*(A) - \mathbb{P}_\varepsilon(A)| &\leq \sup_{A \in \mathcal{F}} |\mathbb{E}_{\mathbb{P}_\varepsilon} [(d\mathbb{P}_\varepsilon^*/d\mathbb{P}_\varepsilon) \mathbb{I}(A)] - \mathbb{E}_{\mathbb{P}_\varepsilon} [\mathbb{I}(A)]| \\
 &\leq \mathbb{E}_{\mathbb{P}_\varepsilon} |d\mathbb{P}_\varepsilon^*/d\mathbb{P}_\varepsilon - 1| \\
 &\leq \mathbb{E}_{\mathbb{P}_\varepsilon}^{1/(1+\delta)} \left[|d\mathbb{P}_\varepsilon^*/d\mathbb{P}_\varepsilon - 1|^{(1+\delta)} \right] \\
 &\leq \mathbb{E}_{\mathbb{P}_\varepsilon}^{1/(1+\delta)} \left[|X(\varepsilon)/\mu_0(\varepsilon) - 1|^{(1+\delta)} \right] \\
 &= [c_{1+\delta}(\varepsilon)]^{1/(1+\delta)} \\
 &\xrightarrow{\varepsilon \rightarrow 0} 0.
 \end{aligned}$$

A discrete-time Markov chains framework

Simulation model represented by DTMC $\{Y_j, j \geq 0\}$ with (large) state space \mathcal{Y} , and a set of absorbing states $\Delta \subset \mathcal{Y}$.

Transition kernel: $P(B | y) = \mathbb{P}[Y_j \in B | Y_{j-1} = y]$.

Stopping time: $\tau = \inf\{j : Y_j \in \Delta\}$.

One-step cost $c(y, y')$ for each transition $y \rightarrow y'$.

Total cost: $X = \sum_{j=1}^{\tau} c(Y_{j-1}, Y_j)$.

Expected cost-to-go from state y : $\mu(y) = \mathbb{E}[X | Y_0 = y]$.

We assume that $\mathbb{E}[\tau | Y_0 = y] < \infty$ and $\mu(y) < \infty$ for all $y \in \mathcal{Y}$.

Want to estimate $\mu_0 = \mu(y_0)$ for some initial state y_0 .

This covers a wide range of situations, including a finite time horizon.

Recurrence equation for μ

The function $\mu : \mathcal{Y} \rightarrow \mathbb{R}$ satisfies the recurrence (Poisson) equation

$$\mu(y) = \mathbb{E}_y[c(y, Y_1) + \mu(Y_1)] = \int_{\mathcal{Y}} [c(y, z) + \mu(z)] dP[dz | y]$$

for $y \notin \Delta$, and $\mu(y) = 0$ for $y \in \Delta$

Recurrence equation for μ

The function $\mu : \mathcal{Y} \rightarrow \mathbb{R}$ satisfies the recurrence (Poisson) equation

$$\mu(y) = \mathbb{E}_y[c(y, Y_1) + \mu(Y_1)] = \int_{\mathcal{Y}} [c(y, z) + \mu(z)] dP[dz | y]$$

for $y \notin \Delta$, and $\mu(y) = 0$ for $y \in \Delta$

Solving this for μ amounts to solve a huge linear system if state space \mathcal{Y} is large but finite. If \mathcal{Y} is continuous, can think of approximating μ by a linear combination of basis functions, or more generally by tuning the parameters of a parameterized function.

These techniques are used in machine learning and approximate dynamic programming (including least-squares Monte Carlo).

Recurrence equation for μ

The function $\mu : \mathcal{Y} \rightarrow \mathbb{R}$ satisfies the recurrence (Poisson) equation

$$\mu(y) = \mathbb{E}_y[c(y, Y_1) + \mu(Y_1)] = \int_{\mathcal{Y}} [c(y, z) + \mu(z)] dP[dz | y]$$

for $y \notin \Delta$, and $\mu(y) = 0$ for $y \in \Delta$

Solving this for μ amounts to solve a huge linear system if state space \mathcal{Y} is large but finite. If \mathcal{Y} is continuous, can think of approximating μ by a linear combination of basis functions, or more generally by tuning the parameters of a parameterized function.

These techniques are used in machine learning and approximate dynamic programming (including least-squares Monte Carlo).

Limitation: often, either the error is large and difficult to estimate, or the approximation is too costly to compute.

Interpretation as a Markov decision process (MDP)¹⁸

At each step of the Markov chain, we can change the transition kernel P for another kernel \tilde{P} . That is, $\tilde{P}(B | y) = \tilde{\mathbb{P}}[Y_j \in B | Y_{j-1} = y]$.

Want to select the transition kernels **dynamically** in a way that minimizes the variance. The **decision** (selection) at each step may depend on past and current history.

An optimal (selection) policy gives zero variance, as we will see soon.

Interpretation as a Markov decision process (MDP)¹⁸

At each step of the Markov chain, we can change the transition kernel P for another kernel \tilde{P} . That is, $\tilde{P}(B | y) = \tilde{\mathbb{P}}[Y_j \in B | Y_{j-1} = y]$.

Want to select the transition kernels **dynamically** in a way that minimizes the variance. The **decision** (selection) at each step may depend on past and current history.

An optimal (selection) policy gives zero variance, as we will see soon.

This **optimal policy** takes $d\tilde{P}(y_1 | y)$ proportional to $dP(y_1 | y)[c(y, y_1) + \mu(y_1)]$, with proportionality constant $1/\mu(y)$.

We can implement an **approximation** of it using an approximation of μ . Often, a crude approximation of μ can be computed cheaply.

IS for a discrete-time Markov chain

We change P to \tilde{P} such that $\tilde{\mathbb{E}}[\tau] < \infty$ and $\tilde{P}(B | y) > 0$ whenever $\int_B [c(y, y_1) + \mu(y_1)] dP(y_1 | y) > 0$.

The estimator X is replaced by

$$X_{\text{is}} = \sum_{j=1}^{\tau} c(Y_{j-1}, Y_j) \prod_{i=1}^j L(Y_{i-1}, Y_i),$$

where $L(Y_{i-1}, Y_i) = (dP/d\tilde{P})(Y_i | Y_{i-1})$.

IS for a discrete-time Markov chain

We change P to \tilde{P} such that $\tilde{\mathbb{E}}[\tau] < \infty$ and $\tilde{P}(B | y) > 0$ whenever $\int_B [c(y, y_1) + \mu(y_1)] dP(y_1 | y) > 0$.

The estimator X is replaced by

$$X_{\text{is}} = \sum_{j=1}^{\tau} c(Y_{j-1}, Y_j) \prod_{i=1}^j L(Y_{i-1}, Y_i),$$

where $L(Y_{i-1}, Y_i) = (dP/d\tilde{P})(Y_i | Y_{i-1})$.

Theorem. If we choose \tilde{P} so that

$$\frac{d\tilde{P}(y_1 | y)}{dP(y_1 | y)} = \begin{cases} \frac{c(y, y_1) + \mu(y_1)}{\mu(y)} & \text{if } \mu(y) > 0, \\ 1 & \text{if } \mu(y) = 0 \end{cases}$$

(this density integrates to 1), then X_{is} has **zero variance**.

IS for a discrete-time Markov chain (cont.)

Proof. The proof can be done by backward induction on the step number, starting from step τ , using the fact that $\mu(Y_\tau) = 0$.

IS for a discrete-time Markov chain (cont.)

Proof. The proof can be done by backward induction on the step number, starting from step τ , using the fact that $\mu(Y_\tau) = 0$.

Unique Markov chain implementation of the zero-variance estimator.

IS for a discrete-time Markov chain (cont.)

Proof. The proof can be done by backward induction on the step number, starting from step τ , using the fact that $\mu(Y_\tau) = 0$.

Unique Markov chain implementation of the zero-variance estimator.

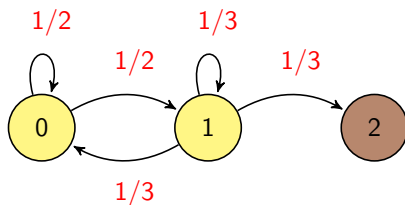
Simple special case: finite state space \mathcal{Y} .

The DTMC has transition probabilities $p(y_1 | y) = \mathbb{P}[Y_1 = y_1 | Y_0 = y]$, which are replaced by $\tilde{p}(y_1 | y) = \tilde{\mathbb{P}}[Y_1 = y_1 | Y_0 = y]$.

We have $L(y, y_1) = p(y_1 | y) / \tilde{p}(y_1 | y)$. For the zero variance:

$$\tilde{p}(y_1 | y) = \begin{cases} p(y_1 | y) \frac{c(y, y_1) + \mu(y_1)}{\mu(y)} & \text{if } \mu(y) > 0, \\ p(y_1 | y) & \text{if } \mu(y) = 0. \end{cases}$$

An example where zero-variance gives $\tilde{\mathbb{E}}[\tau] = \infty$

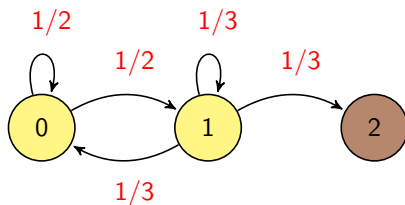


Here, $\mu(y)$ is the expected number of transitions before reaching state 2, given that we are in state y . We have $\mu(2) = 0$.

Zero-variance IS gives $\tilde{p}(y, 2) = 0$ for $y = 0, 1$, so the chain will never reach the stopping time τ under these new probabilities!

We have zero variance but infinite computing cost.

An example where zero-variance gives $\tilde{\mathbb{E}}[\tau] = \infty$



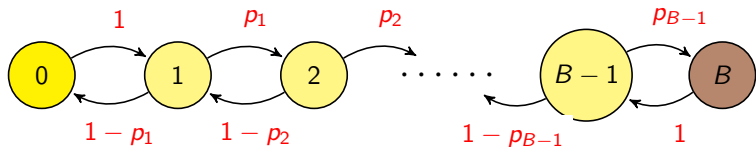
Here, $\mu(y)$ is the expected number of transitions before reaching state 2, given that we are in state y . We have $\mu(2) = 0$.

Zero-variance IS gives $\tilde{p}(y, 2) = 0$ for $y = 0, 1$, so the chain will never reach the stopping time τ under these new probabilities!

We have zero variance but infinite computing cost.

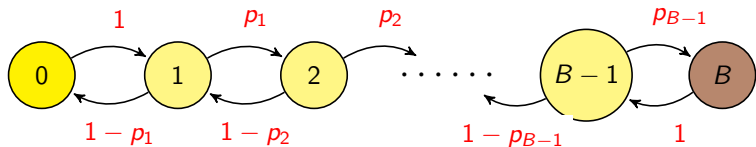
Trick to resolve this: add a cost $\delta > 0$ to any transition that enters Δ . Afterwards, subtract δ to the final (zero-variance) estimator.

Example 2: A birth-and-death process



Let $\tau = \inf\{j > 0 : Y_j \in \{0, B\}\}$ and define $\mu(y) = \mathbb{P}[Y_\tau = B \mid Y_0 = y]$. We want to estimate $\mu_0 = \mu(1)$, the probability of reaching B before coming back to 0 .

Example 2: A birth-and-death process



Let $\tau = \inf\{j > 0 : Y_j \in \{0, B\}\}$ and define $\mu(y) = \mathbb{P}[Y_\tau = B \mid Y_0 = y]$. We want to estimate $\mu_0 = \mu(1)$, the probability of reaching B before coming back to 0. We have the recurrence:

$$\mu(y) = p_y \mu(y+1) + (1 - p_y) \mu(y-1)$$

for $y = 1, \dots, B-1$, with $\mu(0) = 0$ and $\mu(B) = 1$. Zero-variance change of measure gives

$$\tilde{p}_y = p_y \mu(y+1) / \mu(y) \quad \text{for } y \geq 1.$$

Because $\mu(0) = 0$, we also see that $\tilde{p}_1 = 1$ and that no sample path will ever return to 0 under zero-variance IS.

Example 2 with $p_y = p$ for $1 \leq y \leq B - 1$

For $\rho = p/(1-p) \neq 1/2$, it is known that $\mu(y) = (1 - \rho^{-y})/(1 - \rho^{-B})$, so

$$\tilde{p}_y = \frac{1 - \rho^{-y-1}}{1 - \rho^{-y}} p = \frac{(1 - \rho^{y+1})}{(1 - \rho^y)} \frac{1}{\rho} p.$$

Example 2 with $p_y = p$ for $1 \leq y \leq B - 1$

For $\rho = p/(1-p) \neq 1/2$, it is known that $\mu(y) = (1 - \rho^{-y})/(1 - \rho^{-B})$, so

$$\tilde{p}_y = \frac{1 - \rho^{-y-1}}{1 - \rho^{-y}} p = \frac{(1 - \rho^{y+1})}{(1 - \rho^y)} \frac{1}{\rho} p.$$

Those probabilities **do not depend on B** .

We have $(p_{y-1}/\tilde{p}_{y-1})(1 - p_y)/(1 - \tilde{p}_y) = 1$, i.e., the cycles do not contribute to the likelihood ratio.

Example 2 with $p_y = p$ for $1 \leq y \leq B - 1$

For $\rho = p/(1-p) \neq 1/2$, it is known that $\mu(y) = (1 - \rho^{-y})/(1 - \rho^{-B})$, so

$$\tilde{p}_y = \frac{1 - \rho^{-y-1}}{1 - \rho^{-y}} p = \frac{(1 - \rho^{y+1})}{(1 - \rho^y)} \frac{1}{\rho} p.$$

Those probabilities **do not depend on B** .

We have $(p_{y-1}/\tilde{p}_{y-1})(1 - p_y)/(1 - \tilde{p}_y) = 1$, i.e., the cycles do not contribute to the likelihood ratio.

For large B , $\mu(y) = (\rho^{B-y} - \rho^B)/(1 - \rho^B) \approx \rho^{B-y}$. This approximation is the probability of the sample path that goes directly from y to B . When ρ is small, it is the **dominating** path.

Example 2 with $p_y = p$ for $1 \leq y \leq B - 1$

For $\rho = p/(1-p) \neq 1/2$, it is known that $\mu(y) = (1 - \rho^{-y})/(1 - \rho^{-B})$, so

$$\tilde{p}_y = \frac{1 - \rho^{-y-1}}{1 - \rho^{-y}} p = \frac{(1 - \rho^{y+1})}{(1 - \rho^y)} \frac{1}{\rho} p.$$

Those probabilities **do not depend on B** .

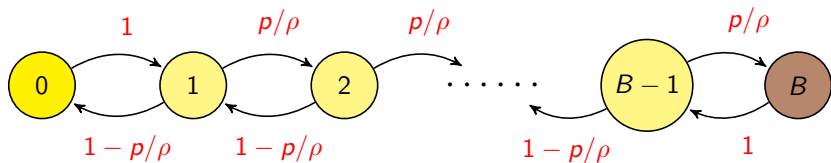
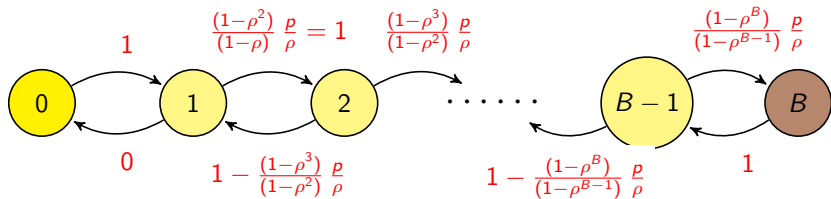
We have $(p_{y-1}/\tilde{p}_{y-1})(1 - p_y)/(1 - \tilde{p}_y) = 1$, i.e., the cycles do not contribute to the likelihood ratio.

For large B , $\mu(y) = (\rho^{B-y} - \rho^B)/(1 - \rho^B) \approx \rho^{B-y}$. This approximation is the probability of the sample path that goes directly from y to B . When ρ is small, it is the **dominating** path.

Using this approximation leads to **VRCM- k** if $p \rightarrow 0$ and fixed B , because then $\rho^{B-y}/\mu(y) \rightarrow 1$.

If $B \rightarrow \infty$ for fixed $p < 1/2$, it gives (only) **BRM- k** .

$$\rho = p/(1 - p).$$



Zero-Variance via control variates (CV)

Same DTMC model. Still want to estimate $\mu_0 = \mu(y_0) = \mathbb{E}[X]$ where $X = \sum_{j=1}^{\tau} c(Y_{j-1}, Y_j)$. We can write

$$\mu(y_0) = X - M_{\tau}$$

where

$$\begin{aligned} M_{\tau} &= \sum_{j=1}^{\tau} [c(Y_{j-1}, Y_j) + \mu(Y_j) - \mu(Y_{j-1})] \\ &= \sum_{j=1}^{\tau} [c(Y_{j-1}, Y_j) + \mu(Y_j) - \mathbb{E}[c(Y_{j-1}, Y_j) + \mu(Y_j) \mid Y_{j-1}]]. \end{aligned}$$

Zero-Variance via control variates (CV)

Same DTMC model. Still want to estimate $\mu_0 = \mu(y_0) = \mathbb{E}[X]$ where $X = \sum_{j=1}^{\tau} c(Y_{j-1}, Y_j)$. We can write

$$\mu(y_0) = X - M_{\tau}$$

where

$$\begin{aligned} M_{\tau} &= \sum_{j=1}^{\tau} [c(Y_{j-1}, Y_j) + \mu(Y_j) - \mu(Y_{j-1})] \\ &= \sum_{j=1}^{\tau} [c(Y_{j-1}, Y_j) + \mu(Y_j) - \mathbb{E}[c(Y_{j-1}, Y_j) + \mu(Y_j) \mid Y_{j-1}]]. \end{aligned}$$

So if we could compute and subtract M_{τ} (as a CV) we would have **zero-variance**. But of course, μ is unknown.

Approximate zero variance via control variates

Replace μ in M_τ by an approximation v such that $v(y) = 0$ for $y \in \Delta$:

$$M_\tau = \sum_{j=1}^{\tau} [c(Y_{j-1}, Y_j) + v(Y_j) - \mathbb{E}[c(Y_{j-1}, Y_j) + v(Y_j) \mid Y_{j-1}]],$$

and define the CV estimator $X_{cv} = X - M_\tau$.

We have $\mathbb{E}[M_\tau] = 0$, and thus $\mathbb{E}[X_{cv}] = \mathbb{E}[X]$ (unbiased) regardless of v .

Variance can be reduced significantly if v is a good approximation of μ .

Approximate zero variance via control variates

Replace μ in M_τ by an approximation v such that $v(y) = 0$ for $y \in \Delta$:

$$M_\tau = \sum_{j=1}^{\tau} [c(Y_{j-1}, Y_j) + v(Y_j) - \mathbb{E}[c(Y_{j-1}, Y_j) + v(Y_j) \mid Y_{j-1}]],$$

and define the CV estimator $X_{cv} = X - M_\tau$.

We have $\mathbb{E}[M_\tau] = 0$, and thus $\mathbb{E}[X_{cv}] = \mathbb{E}[X]$ (unbiased) regardless of v .

Variance can be reduced significantly if v is a good approximation of μ .

Warning: Not the right tool for rare-event simulation, because it does not make the rare events more frequent.

Approximate zero variance via control variates

Replace μ in M_τ by an approximation v such that $v(y) = 0$ for $y \in \Delta$:

$$M_\tau = \sum_{j=1}^{\tau} [c(Y_{j-1}, Y_j) + v(Y_j) - \mathbb{E}[c(Y_{j-1}, Y_j) + v(Y_j) \mid Y_{j-1}]],$$

and define the CV estimator $X_{cv} = X - M_\tau$.

We have $\mathbb{E}[M_\tau] = 0$, and thus $\mathbb{E}[X_{cv}] = \mathbb{E}[X]$ (unbiased) regardless of v .

Variance can be reduced significantly if v is a good approximation of μ .

Warning: Not the right tool for **rare-event simulation**, because it does not make the rare events more frequent.

Extensions to regenerative simulation (Kim & Henderson 07), and to infinite-horizon models with discounting and stochastic differential equations (Henderson & Glynn 02).

Model of Highly Reliable Markovian System (HRMS)

c component types, n_i components of type i .

Markov chain step: failure or repair of one component.

$Y_j = (Y_j^{(1)}, \dots, Y_j^{(c)}) = \text{num. failed compon. of each type at step } j$.

$\{Y_j, j \geq 0\}$ DTMC with $p(y, y') = \mathbb{P}[Y_j = y' \mid Y_{j-1} = y]$.

Suppose that failure probabilities are much smaller than repair probabilities. This is typical of highly reliable systems.

The state space \mathcal{Y} is partitioned in: (1) a (decreasing) set of up states \mathcal{U} and (2) the set of failure states \mathcal{F} .

For any set A , let $\tau_A = \text{first hitting time of } A$, and

$$\mu(y) = \mathbb{P}[\tau_{\mathcal{F}} < \tau_{\mathbf{0}} \mid Y_0 = y],$$

the prob. of visiting \mathcal{F} before returning to $\mathbf{0}$.

Goal: estimate $\mu_0 = \mu(\mathbf{0})$. This can be difficult when $\mu(\mathbf{0})$ is very small.

Some previous IS heuristics:

Balanced failure biasing (BFB) (Shahabuddin 1994) changes p to \tilde{p} as follows, for $x \notin B$:

$$\tilde{p}(x, y) = \begin{cases} \frac{1}{|F(x)|} & \text{if } y \in F(x) \text{ and } p_R(x) = 0; \\ \rho \frac{1}{|F(x)|} & \text{if } y \in F(x) \text{ and } p_R(x) > 0; \\ (1 - \rho) \frac{p(x, y)}{p_R(x)} & \text{if } y \in R(x); \\ 0 & \text{otherwise.} \end{cases}$$

Simple failure biasing (SFB) (Shahabuddin 1988): Replace $1/|F(x)|$ above by $p(x, y) / \sum_{y \in F(x)} p(x, y)$.

SBLR (Alexopoulos and Shultes 2001) changes the probabilities in a way that over any cycle in the visited states during the simulation, the cumulated likelihood ratio remains bounded

These methods do not attempt to mimic zero-variance sampling.

Proposed approximation (ZVA)

Approximate μ by some easily computable function v , and plug into zero-variance formula.

For any state $y \in \mathcal{U}$, let $\Gamma(y)$ be the set of all paths

$\pi = (y = y_0 \rightarrow y_1 \rightarrow \dots \rightarrow y_k)$ where $y_1, \dots, y_{k-1} \notin \mathcal{F} \cup \{\mathbf{0}\}$, $y_k \in \mathcal{F}$, and having positive probability

$$p(\pi) = \prod_{j=1}^k p(y_{j-1}, y_j) > 0.$$

Because these paths represent disjoint events, we have

$$\mu(y) = \sum_{\pi \in \mu_0(y)} p(\pi).$$

This last sum may contain a huge (perhaps ∞) number of terms.

A very crude approximation is to just take the path with largest probability, i.e., approximate

$$\mu(y) = \sum_{\pi \in \Gamma(y)} p(\pi)$$

by its lower bound

$$v_0(y) = \max_{\pi \in \Gamma(y)} p(\pi).$$

Computing $v_0(y)$ amounts to computing a shortest path from y to \mathcal{F} , where the length of a link $y' \rightarrow y''$ is $-\log p(y', y'')$. Easy.

This would work fine if a single path dominates the sum (this may happen when failure transitions have very small probabilities), but this v_0 will often underestimate the bound significantly.

Refinements

Typically, the farther we are from \mathcal{F} , the more v_0 underestimates μ . Close to \mathcal{F} , things are fine, but not close to $\mathbf{0}$.

First simple correction:

1. Estimate $\mu(\mathbf{0})$ in preliminary runs with crude IS strategy;
2. Find constant $\alpha \leq 1$ such that $(v_0(\mathbf{0}))^\alpha$ equals this estimate;
3. Use $v_1(y) = (v_0(y))^\alpha$ for all $y \in \mathcal{U}$ as approx. of $\mu(y)$.

This v_1 matches μ for $y \in \mathcal{F}$ and matches its estimate at $y = \mathbf{0}$.

Second refinement: Replace α by a state-dependent exponent

$$\alpha(y) = 1 + [\alpha(\mathbf{0}) - 1] \frac{\log v_0(y)}{\log v_0(\mathbf{0})},$$

where $\alpha(\mathbf{0}) = \alpha$ as above. This $\alpha(y)$ changes progressively from 1 near \mathcal{F} to $\alpha(\mathbf{0}) < 1$ in state $\mathbf{0}$. The correction here is milder than in the previous case when we are close to \mathcal{F} .

Let $v_2(y) = (v_0(y))^{\alpha(y)}$ be the resulting approximation.

Example: Three types of components

$c = 3$ and $n_1 = n_2 = n_3$. We have $(n_1 + 1)(n_2 + 1)(n_3 + 1)$ states.

Expon. repair times with mean 1.

Failure rate λ_i for component type i ,
with $\lambda_1 = \varepsilon$, $\lambda_2 = 1.5\varepsilon$, and $\lambda_3 = 2\varepsilon^2$, for some small real number ε .

We will try different values of (n_i, ε) .

\mathcal{F} = states where at least one component type has fewer than 2 operational units.

To define $v_0(y)$, we consider all three paths to \mathcal{F} that result from failures of a single component type, and sum their probabilities.

The table contains results with $n = 2^{20}$ runs.

Best estimate of $\mu(\mathbf{0})$: obtained from a large number of runs with our best IS strategies.

Mean

n_i	ε	$\mu(\mathbf{0})$	$v_0(\mathbf{0})$	BFB	SBLR
3	0.001	2.6×10^{-3}	1.3×10^{-3}	2.7×10^{-3}	2.6×10^{-3}
6	0.01	1.8×10^{-7}	3.4×10^{-8}	1.9×10^{-7}	(9.9×10^{-11})
6	0.001	1.7×10^{-11}	3.4×10^{-12}	1.8×10^{-11}	(1.8×10^{-16})
12	0.1	6.0×10^{-8}	3.2×10^{-9}	4.8×10^{-8}	1.3×10^{-8}
12	0.001	3.9×10^{-28}	3.5×10^{-29}	(1.8×10^{-40})	(2.9×10^{-45})

Variance

n_i	ε	BFB	SBLR
3	0.001	1.8×10^{-2}	8.0×10^{-3}
6	0.01	6.3×10^{-11}	(4.5×10^{-16})
6	0.001	8.8×10^{-19}	(2.0×10^{-26})
12	0.1	8.1×10^{-10}	1.7×10^{-10}
12	0.001	(3.2×10^{-74})	(3.5×10^{-84})

Mean

n_i	ε	$\mu(\mathbf{0})$	ZVA(v_0)	ZVA(v_1)	ZVA(v_2)
3	0.001	2.6×10^{-3}	2.6×10^{-3}	2.6×10^{-3}	2.6×10^{-3}
6	0.01	1.8×10^{-7}	1.8×10^{-7}	1.8×10^{-7}	1.8×10^{-7}
6	0.001	1.7×10^{-11}	1.7×10^{-11}	1.7×10^{-11}	1.7×10^{-11}
12	0.1	6.0×10^{-8}	6.0×10^{-8}	6.2×10^{-8}	6.7×10^{-8}
12	0.001	3.9×10^{-28}	3.9×10^{-28}	3.9×10^{-28}	3.9×10^{-28}

Variance

n_i	ε	α	ZVA(v_0)	ZVA(v_1)	ZVA(v_2)	RE(v_2)
3	0.001	0.906	6.5×10^{-4}	2.7×10^{-3}	9.3×10^{-9}	0.04
6	0.01	0.903	2.0×10^{-14}	1.2×10^{-14}	7.7×10^{-15}	0.48
6	0.001	0.939	1.2×10^{-23}	1.1×10^{-23}	7.6×10^{-24}	0.16
12	0.1	0.851	1.6×10^{-10}	2.9×10^{-10}	1.5×10^{-11}	64.50
12	0.001	0.963	1.4×10^{-55}	9.3×10^{-56}	9.4×10^{-56}	0.78

We have $\alpha \rightarrow 1$ when $\varepsilon \rightarrow 0$ or when $n_i \nearrow$

Approxim. and adaptive learning (for CV and IS)

- ▶ Zero-variance schemes require the knowledge or a good approximation of the function μ .
- ▶ More practical approach: approximate μ in a parametric class of functions $\mathcal{V} = \{v(\cdot; \theta) : \mathcal{Y} \rightarrow \mathbb{R}, \theta \in \Theta\}$, where $\Theta \subseteq \mathbb{R}^m$, and $\theta = (\theta_1, \dots, \theta_m)$ is a vector of parameters that we try to optimize so that $v = v(\cdot; \theta)$ is close to μ in some sense.
- ▶ Convenient way (in terms of computations): fixed set of functions v_1, \dots, v_m , independent over \mathbb{R} , and define \mathcal{V} as the space of all linear combinations of these functions:

$$\mathcal{V} = \left\{ v = v(\cdot; \theta) = \sum_{i=1}^m \theta_i v_i(\cdot) \right\}$$

where $\theta = (\theta_1, \dots, \theta_m) \in \mathbb{R}^m$.

- ▶ Can also use **nonlinear** parameterization.

Approximate zero variance and VRE

$$\theta^* = \arg \min_{\theta \in \Theta} \sum_{\ell=1}^r |v(y_\ell; \theta) - \hat{\mu}(y_\ell)|^2$$

found by [least-squares regression](#) or standard temporal differences method ([stochastic approximation](#)):

$$\theta^{(n+1)} = \theta^{(n)} + a_n \left(c(y_n, y_{n+1}) + v(y_{n+1}, \theta^{(n)}) - v(y_n, \theta^{(n)}) \right) \nabla_{\theta} v(y_n, \theta^{(n+1)}).$$

- ▶ The method converges to the zero-variance sampling only if the parameterized class \mathcal{V} contains the function μ .
- ▶ No better convergence than the $O(n^{-1/2})$ (in the probabilistic sense) otherwise.
- ▶ Issue: selecting a good parameterized space \mathcal{V} . Basis function adaptation possible during the learning phase.
- ▶ For IS, instead of parameterizing v , we can directly parameterize the IS distribution, and then minimize the variance of X_{is} with respect to these IS parameters. [In any case, needs to be sampled easily!](#)

Other adaptive techniques in the literature

Adaptive Monte Carlo (AMC) for IS: (Booth, 1985)

Learns iteratively function $\mu(\cdot)$, and still plug the approximation into the zero-variance change of measure formula instead of $\mu(\cdot)$.

- ▶ Considers several steps and n_i independent simulation replications at step i .
- ▶ At step i , replaces $\mu(x)$ by a guess $\mu^{(i)}(x)$
- ▶ use probabilities

$$\tilde{P}_{y,z}^{(i)} = \frac{P_{y,z}(c_{y,z} + \mu^{(i)}(z))}{\sum_w P_{y,w}(c_{y,w} + \mu^{(i)}(w))}.$$

- ▶ Gives a new estimation $\mu^{(i+1)}(y)$ of $\mu(y)$, from which a new transition matrix $\tilde{P}^{(i+1)}$ is defined.

Adaptive stochastic approximation (ASA)

- ▶ (Ahamed et al, 06) just uses a single sample path (y_0, \dots, y_n) .
- ▶ Initial distribution for y_0 , matrix $\tilde{P}^{(0)}$ and guess $\mu^{(0)}(\cdot)$.
- ▶ At step j of the path, if $y_j \notin \Delta$,
 - ▶ matrix $\tilde{P}^{(j)}$ used to generate y_{j+1} .
 - ▶ From y_{j+1} , update the estimate of $\mu(y_j)$ by temporal difference

$$\begin{aligned} \mu^{(j+1)}(y_j) &= (1 - a_j(y_j))\mu^{(j)}(y_j) \\ &+ a_j(y_j) \left[c(y_j, y_{j+1}) + \mu^{(j)}(y_{j+1}) \right] \frac{P(y_j, y_{j+1})}{\tilde{P}^{(j)}(y_j, y_{j+1})}, \end{aligned}$$

where $\{a_j(y), j \geq 0\}$, sequence of **step sizes**

- ▶ For $\delta > 0$ constant,

$$\tilde{P}^{(j+1)}(y_j, y_{j+1}) = \max \left(P(y_j, y_{j+1}) \frac{[c(y_j, y_{j+1}) + \mu^{(j+1)}(y_{j+1})]}{\mu^{(j+1)}(y_j)}, \delta \right).$$
- ▶ Otherwise $\mu^{(j+1)}(y) = \mu^{(j)}(y)$, $\tilde{P}^{(j+1)}(y, z) = P^{(j)}(y, z)$.
- ▶ Normalize: $P^{(j+1)}(y_j, y) = \frac{\tilde{P}^{(j+1)}(y_j, y)}{\sum_z \tilde{P}^{(j+1)}(y_j, z)}$.
- ▶ If $y_j \in \Delta$, y_{j+1} generated from initial distribution, but estimations of $P(\cdot, \cdot)$ and $\mu(\cdot)$ kept.
- ▶ Batching techniques used to get a confidence interval.

Drawbacks of these learning techniques

- ▶ Must store the vectors $\mu^{(n)}(\cdot)$.
- ▶ State-space typically very large when we use simulation... This limits the practical effectiveness of the methods.
- ▶ We can use adaptive learning to optimize the parameter vector θ of an approximation instead of learning $\mu(y)$ directly.
- ▶ Example (Bolia et al. 04):
 - ▶ To approximate the value function in the context of pricing an American option by IS;
 - ▶ least-squares regression to approx. the value function and the continuation value function (the value of the option conditional on not exercising it at the current step)
 - ▶ Basis functions $v_i(y) = \exp[a_i \log^2 y + b_i \log y]$, where the a_i and b_i are constants: IS then mixtures of lognormal, easy to sample!
 - ▶ In a first-stage, they learn good vectors of weights θ and use them to approximate both the zero-variance IS and the optimal stopping rule.
 - ▶ In a second stage, they use these approximations to estimate the value function more accurately.

Iterated CVs for computing $\int f(y)dy$

(Gobet and Maire 2006)

- ▶ Uses an orthonormal basis of functions $e_k(x)$, $1 \leq k \leq p$.
- ▶ From n sample values Y_i , computes $b_i^{(1)} = \frac{1}{n} \sum_{i=1}^n f(Y_i)e_i(Y_i)$ and first approx

$$f^{(1)}(y) = \frac{1}{n} \sum_{i=1}^n b_i^{(1)} e_i(x).$$

- ▶ At step $m \geq 2$, the same technique for n independent copies is iteratively applied to function again applied to $f - f^{(m-1)}$.
- ▶ Eventually, at step M , the approximation is

$$f^{(M)}(y) = \sum_{k=1}^p a_k e_k(y)$$

with $\hat{a}_k = \sum_{m=1}^M b_k^{(m)}$.

Iterated CV (2)

- ▶ Can use s -dimensional basis functions (Legendre or Tchebichev polynomials, Korobov space) used such that the real coefficient a_m (now multidimensional) decreases fast when m increases:

$$|a_m| \leq \frac{C_1}{(m_1 \cdots m_s)^L}.$$

- ▶ The variance of the \hat{a}_m decreases fast too then.
- ▶ One-dimensional case: optimal choice: $M = \frac{\ln(n)}{(2L-1)\ln(2)}$ and, then using $n \ln(n)$ sample values,

$$\sigma^2(\hat{a}_m) \leq C_2 \frac{1}{n^{L-1/2-\epsilon}}.$$

Conclusions

- ▶ Both IS and CV can achieve zero-variance [in theory](#).
- ▶ But zero-variance sampling can only be approximated, usually by approximating the value function μ .
- ▶ Fits the framework of approximate dynamic programming.
- ▶ Approximation by a linear combination of a fixed set of basis functions; natural in the case of continuous or very large state spaces.
- ▶ Difficulty: choice of those basis functions.
- ▶ Another important hurdle for IS: approximation must be constructed to allow efficient random variate generation.
- ▶ In practice, approximate zero-variance sampling can provide very significant gains in efficiency.
- ▶ More applications could benefit from this technology.
- ▶ Other methods can improve the $O(n^{-1/2})$ convergence rate: [adaptive stratification](#), [generalized antithetic variates](#), [randomized quasi-Monte Carlo](#), for example.