# Adaptive Monte Carlo variance reduction with two-time-scale stochastic approximation

REIICHIRO KAWAI

**Abstract**

Combined control variates and importance sampling variance reduction and its two-fold optimality are investigated. Two-time-scale stochastic approximation algorithm is applied in parameter search for the combination and almost sure convergence of the algorithm to the unique optimum is proved. The parameter search procedure is further incorporated into adaptive Monte Carlo simulation, and its law of large numbers and central limit theorem are proved to hold. An numerical example is provided to illustrate the effectiveness of the method.

*Keywords:* Control variates, Girsanov theorem, importance sampling, Monte Carlo methods, stochastic approximation, two time scales, variance reduction.
*2000 Mathematics Subject Classification:* 65C05, 62L20, 93E35, 60G15.

## 1 Introduction

In this paper, we develop an adaptive Monte Carlo variance reduction procedure with three steps: (i) we combine two variance reduction techniques, the control variates (CV) and the importance sampling (IS), and investigate the two-fold optimality, (ii) we apply the two-time-scale stochastic approximation algorithm in parameter search for the combination and prove almost sure convergence of the algorithm to the unique optimum, and (iii) we incorporate the parameter search procedure into adaptive Monte Carlo simulation.

There is a vast literature on optimal use of a single variance reduction technique, and some combinations have been investigated, for example, IS and the stratification in Glasserman *et al.*[6], while IS and the temporal difference control variates in Randhawa and Juneja [8]. To our knowledge, it seems that the combination of CV and IS has not been studied elsewhere. It however turns out to be unclear whether or not there exists global optimum, and thus by "simultaneous optimality" in this paper, we mean either (I) "optimality in IS parameter with CV always held optimal with

respect to the IS parameter" or (II) "optimality in CV parameter with IS always held optimal with respect to the CV parameter."

The optimal parameter search is a challenging problem even in the IS only. In financial engineering settings, for example, a systematic way to derive nearly optimal parameter is proposed in [6], while the application of the Robbins-Monro algorithm, which is a *single-time-scale* stochastic approximation algorithm, in a fairly general formulation is studied in Su and Fu [9] and Arouna [1], and a pure-jump Lévy process framework in Kawai [7]. In order to search two-fold parameters of CV and of IS at the same time, we apply the *two-time-scale* stochastic approximation algorithm, which is a stochastic recursive algorithm in which some of the components are updated using step-sizes that are very small compared to those of the remaining components and whose almost sure convergence is first rigorously proved in Borkar [3]. In our case, depending on which step-sizes dominate the others, the two-time-scale algorithm are "equilibrated CV and quasi-static IS" or "equilibrated IS and quasi-static CV", corresponding respectively to (I) and (II), appeared in the last paragraph, and one of main contributions of this paper lies in proving the almost sure convergence to each optimum.

It is often not quite practical to perform the algorithm first to search optimal parameters and then re-run the Monte Carlo simulation in cooperation with the optimized variance reduction techniques, simply because this double implementation is evidently more time-consuming than the plain Monte Carlo simulation first performed. The idea of incorporating the parameter search with a stochastic approximation into adaptive Monte Carlo simulation has been suggested in Arouna [2] with full theoretical support and is also applicable to our method, which becomes more attractive when used as an adaptive Monte Carlo simulation.

The rest of this paper is organized as follows. Section 2 recalls in brief the principle of the control variates and the importance sampling. Section 3 formulates the combined control variates and importance sampling and discusses optimality when either one component is fixed. Section 4 proves the almost sure convergence of two-time-scale algorithms to the unique limiting point for the combination and shows resulting variance reduction effect. Then, in Section 5, the parameter search using the two-time-scale algorithm is applied in adaptive Monte Carlo simulation, and the strong law of large numbers for empirical mean and the central limit theorem for empirical variance are proved to hold. Section 6 presents and discusses results of a numerical experiment in a financial engineering example, and Section 7 concludes this study.

## 2  Preliminaries

Let us begin with some general notations which will be used throughout the text. $\mathbb{R}^d$ is the $d$-dimensional Euclidean space with the norm $\|\cdot\|$ and the inner product $\langle\cdot,\cdot\rangle$, $\mathbb{R}_0^d := \mathbb{R}^d \setminus \{0\}$. $(\Omega,\mathscr{F},\mathbb{P})$ is our underlying probability space. $\mathbb{C}([0,\infty);\mathbb{R}^d)$ is the space of continuous functions from $[0,\infty)$ into $\mathbb{R}^d$. $\mathbb{P}|_{\mathscr{F}_t}$ is the restriction of a probability measure $\mathbb{P}$ to $\sigma$-field $\mathscr{F}_t$, while $\overset{\mathscr{L}}{\to}$ denotes convergence in distribution. As usual, $I_d$ denotes the $d$-dimensional identity matrix, $\nabla_x$ the

gradient, and $\text{Hess}_x[\cdot]$ the Hessian matrix, with respect to the variable $x$. By $\Pi_H(x)$, we indicate the projection of $x$ onto the set $H$, that is, the closest point in $H$ to $x$. We denote by $\|\cdot\|_o$ the operator norm of a linear transformation, so if $A \in \mathbb{R}^{d \times d}$, then $\|A\|_o = \sup_{\|x\| \leq 1} \|Ax\|$. The characteristic function of the marginal distributions of the Brownian motion in $\mathbb{R}^d$ is uniquely given by

$$\mathbb{E}_{\mathbb{P}}\left[e^{i\langle y, W_t\rangle}\right] = \exp\left[t\left(i\langle y, \gamma\rangle - \frac{1}{2}\langle y, Ay\rangle\right)\right],$$

where $\gamma \in \mathbb{R}^d$ and where $A$ is a symmetric nonnegative-definite $d \times d$ matrix. We will say that a Brownian motion satisfying the above characteristic function is generated by $(\gamma, A)$, and as usual we denote by $\mathcal{N}(\gamma, A)$ the marginal distribution of the Brownian motion at unit time, or equivalently, the normal distribution with mean $\gamma$ and with variance-covariance matrix $A$. *In this paper, we restrict ourselves to the standard Brownian motion, or to the standard normal random vector, that is, we set $A \equiv I_d$ throughout.* Clearly, this simplification loses no generality since any $d$-dimensional normal random vector can easily be generated from the standard normal random vector. Also, we define by $(\mathscr{F}_t)_{t \geq 0}$ the natural filtration of $\{W_t : t \geq 0\}$.

Let $F : \mathbb{C}([0, T]; \mathbb{R}^d) \mapsto \mathbb{R}$ be such that for some $c > 1$,

$$F(W) := F(\{W_t : t \in [0, T]\}) \in L^{4c}(\Omega, \mathscr{F}_T, \mathbb{P}), \tag{2.1}$$

and $\mathbb{P}(F(W) \neq 0) > 0$. For ease in notation, we will write $F(W - \lambda) := F(\{W_t - t\lambda : t \in [0, T]\})$.

Throughout this paper, we are interested in the variance reduction in evaluating

$$C := \mathbb{E}_{\mathbb{P}}[F(W)]$$

by Monte Carlo simulation.

## 2.1 Control variates variance reduction

The control variates method we will consider in this paper is of a linear type based upon the following equality,

$$C = \mathbb{E}_{\mathbb{P}}[F(W)] = \mathbb{E}_{\mathbb{P}}[F(W) - \langle \theta, W_T\rangle]. \tag{2.2}$$

We assume that $F(W)$ and $W_T$ are correlated, that is,

$$\mathbb{E}_{\mathbb{P}}[(F(W) - C)W_T] \neq 0. \tag{2.3}$$

The variance of the component inside the right hand side expectation in (2.2) is well defined due to (2.1) and is given by

$$V_1(\theta) := \mathbb{E}_{\mathbb{P}}\left[(F(W) - \langle \theta, W_T\rangle)^2\right] - C^2.$$

Thanks to the quadratic form of $V_1(\theta)$ in $\theta$, it can be shown that $V_1$ is minimized uniquely at

$$\theta^* := \mathbb{E}_{\mathbb{P}}\left[W_T W_T'\right]^{-1} \mathbb{E}_{\mathbb{P}}[F(W)W_T] = T^{-1}\mathbb{E}_{\mathbb{P}}[F(W)W_T], \tag{2.4}$$

and the resulting minimal variance is given by

$$
\begin{aligned}
V_1(\theta^*) &= V_1(0) - \|\mathbb{E}_{\mathbb{P}}[F(W)W_T]\|^2 \\
&= V_1(0) - T\|\theta^*\|^2.
\end{aligned}
\tag{2.5}
$$

The condition (2.3) ensures that $V_1(\theta^*) < V_1(0)$ and implies that the method is particularly effective when $F(W)$ and $W_T$ are highly correlated.

## 2.2 Importance sampling variance reduction via Esscher transform

The importance sampling method is aimed at reducing the variance of iid Monte Carlo summands by appropriately transforming the underlying probability measure, from which interested random variables or stochastic processes are generated, so as to put more weight on important events and less on undesirable ones.

Let $\{W_t : t \geq 0\}$ be a Brownian motion in $\mathbb{R}^d$ generated by $(0, I_d)$. Under the probability measure $\mathbb{Q}_\lambda$, where $\lambda \in \mathbb{R}^d$ and which is defined via the Radon-Nikodym derivative, $\mathbb{P}$-*a.s.*,

$$
\frac{d\mathbb{Q}_\lambda}{d\mathbb{P}}|_{\mathscr{F}_t} = \frac{e^{\langle \lambda, W_t \rangle}}{\mathbb{E}_{\mathbb{P}}\left[e^{\langle \lambda, W_t \rangle}\right]} = e^{\langle \lambda, W_t \rangle - \frac{1}{2}t\|\lambda\|^2},
$$

the stochastic process $\{W_t : t \geq 0\}$ is again a Brownian motion generated by $(\lambda, I_d)$. Then, the probability measures $\mathbb{P}$ and $\mathbb{Q}_\lambda$ are mutually absolutely continuous, and we also get $\mathbb{E}_{\mathbb{Q}_\lambda}[e^{-\langle \lambda, W_1 \rangle}] < +\infty$, and $\mathbb{Q}_\lambda$-*a.s.*,

$$
\frac{d\mathbb{P}}{d\mathbb{Q}_\lambda}|_{\mathscr{F}_t} = \left(\frac{d\mathbb{Q}_\lambda}{d\mathbb{P}}|_{\mathscr{F}_t}\right)^{-1} = e^{-\langle \lambda, W_t \rangle + \frac{1}{2}t\|\lambda\|^2}.
$$

The importance sampling method we will consider in this paper is based upon the equality

$$
\begin{aligned}
\mathbb{E}_{\mathbb{P}}[F(W)] &= \mathbb{E}_{\mathbb{Q}_\lambda}\left[\left(\frac{d\mathbb{Q}_\lambda}{d\mathbb{P}}|_{\mathscr{F}_T}\right)^{-1} F(W)\right] \\
&= \mathbb{E}_{\mathbb{Q}_\lambda}\left[e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} F(W)\right].
\end{aligned}
$$

The variance under the probability measure $\mathbb{Q}_\lambda$ is well defined due to (2.1) and is given by

$$
\begin{aligned}
V_2(\lambda) &:= \mathbb{E}_{\mathbb{Q}_\lambda}\left[\left(\frac{d\mathbb{P}}{d\mathbb{Q}_\lambda}|_{\mathscr{F}_T}\right)^2 F(W)^2\right] - C^2 \\
&= \mathbb{E}_{\mathbb{P}}\left[e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} F(W)^2\right] - C^2.
\end{aligned}
$$

The gradient and the Hessian matrix of $V_2$ are well defined due to (2.1) and are given respectively by

$$
\nabla_\lambda V_2(\lambda) = \mathbb{E}_{\mathbb{P}}\left[(T\lambda - W_T)e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} F(W)^2\right],
$$

and

$$
\text{Hess}_\lambda[V_2(\lambda)] = \mathbb{E}_{\mathbb{P}}\left[(TI_d + (T\lambda - W_T)(T\lambda - W_T)')e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} F(W)^2\right].
$$

4

In the above setting, [1] formulates a stochastic approximation algorithm as follows. (See also [9].) Let $\{W_{n,t} : t \in [0,T]\}_{n \in \mathbb{N}}$ be a sequence of iid replications of the stochastic process $\{W_t : t \in [0,T]\}$. For ease of notation, we will write $W_n := W_{n,T}$ for $n \in \mathbb{N}$, and $F(W)_n := F(\{W_{n,t} : t \in [0,T]\})$. Let $\{H_n\}_{n \in \{0\} \cup \mathbb{N}}$ be an increasing sequence of compact sets in $\mathbb{R}^d$ such that $\{0\} \in H_0$ and $\cup_{n=0}^{\infty} H_n = \mathbb{R}^d$, and define a sequence $\{Y_n\}_{n \in \mathbb{N}}$ of random vectors in $\mathbb{R}^d$ by

$$Y_{n+1} = (T\lambda_n - W_{n+1}) e^{-\langle \lambda_n, W_{n+1} \rangle + \frac{1}{2} T \|\lambda_n\|^2} F(W)_{n+1}^2,$$

where $\lambda_0 \in H_0$, $\{\lambda_n\}_{n \in \mathbb{N}}$ is a sequence of random vectors in $\mathbb{R}^d$ iteratively generated by

$$\lambda_{n+1} = \Pi_{H_{\sigma(n)}} [\lambda_n - \varepsilon_n Y_{n+1}], \tag{2.6}$$

where $\sigma(n)$ is the number counter of projections up to the $n$-th step, and where $\{\varepsilon_n\}_{n \in \{0\} \cup \mathbb{N}}$ is a sequence of non-increasing positive constants satisfying

$$\sum_{n=0}^{+\infty} \varepsilon_n = +\infty, \quad \sum_{n=0}^{+\infty} \varepsilon_n^2 < +\infty.$$

The following is a summary of the main results of [1].

**Theorem 2.1.** *The function $V_2$ is strictly convex on $\mathbb{R}^d$. Moreover, the sequence $\{\lambda_n\}_{n \in \mathbb{N}}$ in (2.6) converges $\mathbb{P}$-a.s. to $\lambda^*$ such that $\nabla_\lambda V_2(\lambda^*) = 0$, with $\lim_{n \uparrow +\infty} \sigma(n) < +\infty$, $\mathbb{P}$-a.s.*

# 3    Combination of control variates and importance sampling

The aim of this paper is at the simultaneous effectiveness of the control variates and the importance sampling in view of the following equation, for $\theta \in \mathbb{R}^d$ and $\lambda \in \mathbb{R}^d$,

$$
\begin{aligned}
\mathbb{E}_{\mathbb{P}}[F(W)] &= \mathbb{E}_{\mathbb{P}}[F(W) - \langle \theta, W_T \rangle] \\
&= \mathbb{E}_{\mathbb{Q}_\lambda} \left[ \frac{d\mathbb{P}}{d\mathbb{Q}_\lambda} |_{\mathscr{F}_T} (F(W) - \langle \theta, W_T \rangle) \right] \\
&= \mathbb{E}_{\mathbb{Q}_\lambda} \left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2} T \|\lambda\|^2} (F(W) - \langle \theta, W_T \rangle) \right].
\end{aligned}
\tag{3.1}
$$

For $\theta \in \mathbb{R}^d$ and $\lambda \in \mathbb{R}^d$, the variance of the component inside the expectation (3.1) is well defined and is given by

$$V_3(\theta, \lambda) := \mathbb{E}_{\mathbb{P}} \left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2} T \|\lambda\|^2} (F(W) - \langle \theta, W_T \rangle)^2 \right] - C^2.$$

Observe that $V_3(0,0) = V_1(0) = V_2(0)$, $V_3(\theta, 0) = V_1(\theta)$, and $V_3(0, \lambda) = V_2(\lambda)$.

## 3.1 Importance sampling with fixed control variates

We first consider a stochastic approximation algorithm to find the minimum of the function $V_3(\theta, \lambda)$ with $\theta$ held fixed, that is, the control variates component is fixed. Also, define $G_1 : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{C}([0,T];\mathbb{R}^d) \mapsto \mathbb{R}^d$ by

$$G_1(\theta, \lambda, X) := e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} (F(W) - \langle \theta, W_T \rangle)^2 (T\lambda - W_T),$$

so that

$$\nabla_\lambda V_3(\theta, \lambda) = \mathbb{E}_\mathbb{P}[G_1(\theta, \lambda, W)].$$

Then, for each $\theta \in \mathbb{R}^d$, consider a stochastic iteration

$$\lambda_{n+1} = \Pi_{H_{\sigma(n)}} [\lambda_n - \varepsilon_n G_1(\theta, \lambda_n, W_{n+1})], \tag{3.2}$$

where $\{H_n\}_{n \in \{0\} \cup \mathbb{N}}$ and $\sigma(n)$ are as the ones in the unconstrained algorithm (2.6), and where $\lambda_0 \in H_0$.

**Proposition 3.1.** *Fix $\theta \in \mathbb{R}^d$. The function $V_3(\theta, \lambda)$ is strictly convex in $\lambda$ on $\mathbb{R}^d$. Moreover, the sequence $\{\lambda_n\}_{n \in \mathbb{N}}$ in (3.2) converges $\mathbb{P}$-a.s. to $g(\theta) \in \mathbb{R}^d$, such that $[\nabla_\lambda V_3(\theta, \lambda)]_{\lambda = g(\theta)} = 0$, with $\lim_{n \uparrow +\infty} \sigma(n) < +\infty$, $\mathbb{P}$-a.s.*

*Proof.* The claims can be proved essentially in a similar manner to those of Theorem 2.1 with the help of the assumption (2.1). (See [1] for details.) □

Clearly, the primal interest out of Proposition 3.1 is the inequality $V_3(\theta^*, g(\theta^*)) < V_3(\theta^*, 0)$, or more illustratively, the density transform

$$\mathbb{E}_\mathbb{P}[F(W) - \langle \theta^*, W_T \rangle] \Rightarrow \mathbb{E}_{\mathbb{Q}_{g(\theta^*)}} \left[ \frac{d\mathbb{P}}{d\mathbb{Q}_{g(\theta^*)}} |_{\mathscr{F}_T} (F(W) - \langle \theta^*, W_T \rangle) \right].$$

## 3.2 Control variates with fixed importance sampling

Next, we consider the variance $V_3(\theta, \lambda)$ with $\lambda$ held fixed in $\mathbb{R}^d$, that is, the importance sampling component is fixed. For convenience, define $G_2 : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{C}([0,T];\mathbb{R}^d) \mapsto \mathbb{R}^d$ by

$$G_2(\theta, \lambda, W) := -2e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} (F(W) - \langle \theta, W_T \rangle) W_T,$$

so that

$$\nabla_\theta V_3(\theta, \lambda) = \mathbb{E}_\mathbb{P}[G_2(\theta, \lambda, W)].$$

Then, for each $\lambda \in \mathbb{R}^d$, consider a stochastic iteration

$$\theta_{n+1} = \Pi_{H_{\sigma(n)}} [\theta_n - \varepsilon_n G_2(\theta_n, \lambda, X_{n+1})], \tag{3.3}$$

where $\theta_0 \in H_0$. Then, a similar result to Proposition 3.1 also holds true for the algorithm (3.3).

**Proposition 3.2.** *Fix $\lambda \in \mathbb{R}^d$. The function $V_3(\theta, \lambda)$ is strictly convex in $\theta$ on $\mathbb{R}^d$. Moreover, the sequence $\{\theta_n\}_{n \in \mathbb{N}}$ in (3.3) converges $\mathbb{P}$-a.s. to $h(\lambda) \in \mathbb{R}^d$, such that $[\nabla_\theta V_3(\theta, \lambda)]|_{\theta = h(\lambda)} = 0$, with $\lim_{n \uparrow +\infty} \sigma(n) < +\infty$, $\mathbb{P}$-a.s.*

*Proof.* Since $\lambda$ is fixed in $\mathbb{R}^d$, the gradient and the Hessian matrix of $V_3$ with respect to $\theta$ are well defined and are given by

$$\nabla_\theta V_3(\theta, \lambda) = -2\mathbb{E}_\mathbb{P}\left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} (F(W) - \langle \theta, W_T \rangle) W_T \right],$$

and

$$\text{Hess}_\theta[V_3(\theta, \lambda)] = 2\mathbb{E}_\mathbb{P}\left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} W_T W_T' \right].$$

The Hessian matrix is positive-definite since for any $y \in \mathbb{R}_0^d$,

$$y'\text{Hess}_\theta[V_3(\theta, \lambda)]y = 2\mathbb{E}_\mathbb{P}\left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} \langle y, W_T \rangle^2 \right] > 0,$$

and this proves the first assertion.

Next, it follows from the quadratic form of $V_3(\theta, \lambda)$ in $\theta$ that $h(\lambda)$ uniquely exists and

$$h(\lambda) = \mathbb{E}_\mathbb{P}\left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} W_T W_T' \right]^{-1} \mathbb{E}_\mathbb{P}\left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} F(W) W_T \right], \tag{3.4}$$

where the matrix inverse above is well defined since the matrix is positive-definite as well. The resulting minimal variance is given by

$$V_3(h(\lambda), \lambda) = V_3(0, \lambda) - \left\langle h(\lambda), \mathbb{E}_\mathbb{P}\left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2}T\|\lambda\|^2} F(W) W_T \right] \right\rangle,$$

which is strictly smaller than $V_3(0, \lambda)$. For the almost sure convergence to $h(\lambda)$, following the results of Chen and Zhu [4] and Delyon [5], we will show that for each $m \in \mathbb{N}$,

$$\mathbb{E}_\mathbb{P}\left[ \|G_1(\theta, \lambda, W)\|^2 \mathbb{1}(\|\theta\| \le m) \right] < +\infty.$$

To prove the above, it is sufficient to show that

$$\mathbb{E}_\mathbb{P}\left[ e^{-2\langle \lambda, W_T \rangle} F(W)^2 \|W_T\|^2 \right] < +\infty,$$

and

$$\mathbb{E}_\mathbb{P}\left[ e^{-2\langle \lambda, W_T \rangle} \langle \theta, W_T \rangle^2 \|W_T\|^2 \mathbb{1}(\|\theta\| \le m) \right] < +\infty.$$

The first one is straightforward by the Cauchy-Schwartz inequality and the Girsanov theorem,

$$\mathbb{E}_\mathbb{P}\left[ e^{-2\langle \lambda, W_T \rangle} F(W)^2 \|W_T\|^2 \right]^2 \le e^{8T\|\lambda\|^2} \mathbb{E}_\mathbb{P}\left[ F(W)^4 \right] \mathbb{E}_\mathbb{P}\left[ \|W_T - 4T\lambda\|^4 \right] < +\infty,$$

while for the second,

$$\begin{aligned}
\mathbb{E}_\mathbb{P}\left[ e^{-2\langle \lambda, W_T \rangle} \langle \theta, W_T \rangle^2 \|W_T\|^2 \mathbb{1}(\|\theta\| \le m) \right] &\le m^2 \mathbb{E}_\mathbb{P}\left[ e^{-2\langle \lambda, W_T \rangle} \|W_T\|^4 \right] \\
&\le m^2 e^{2T\|\lambda\|^2} \mathbb{E}_\mathbb{P}\left[ \|W_T - 2T\lambda\|^4 \right] < +\infty.
\end{aligned}$$

The proof is complete. $\qquad \square$

The following lemma is used in Proposition 4.1 and is very useful for the computation of $h(\lambda)$.

**Lemma 3.3.** *For $\lambda \in \mathbb{R}^d$,*

$$\begin{aligned}
h(\lambda) &= T^{-1}\left(I_d - \frac{T\lambda\lambda'}{1+T\|\lambda\|^2}\right)\mathbb{E}_\mathbb{P}\left[e^{-\langle\lambda,W_T\rangle-\frac{1}{2}T\|\lambda\|^2}F(W)W_T\right] \\
&= T^{-1}\left(I_d - \frac{T\lambda\lambda'}{1+T\|\lambda\|^2}\right)\mathbb{E}_\mathbb{P}\left[F(W-\lambda)(W_T-T\lambda)\right].
\end{aligned}$$

*Proof.* With the help of the Girsanov theorem, we get

$$\begin{aligned}
\mathbb{E}_\mathbb{P}\left[e^{-\langle\lambda,W_T\rangle+\frac{1}{2}T\|\lambda\|^2}W_T W_T'\right] &= e^{T\|\lambda\|^2}\mathbb{E}_\mathbb{P}\left[e^{-\langle\lambda,W_T\rangle-\frac{1}{2}T\|\lambda\|^2}W_T W_T'\right] \\
&= e^{T\|\lambda\|^2}\mathbb{E}_{\mathbb{Q}_{-\lambda}}\left[W_T W_T'\right] \\
&= e^{T\|\lambda\|^2}\mathbb{E}_\mathbb{P}\left[(W_T-T\lambda)(W_T-T\lambda)'\right] \\
&= e^{T\|\lambda\|^2}T\left(I_d+T\lambda\lambda'\right).
\end{aligned}$$

The Sherman-Morrison-Woodbury formula yields

$$\left(I_d+T\lambda\lambda'\right)^{-1} = I_d - \frac{T\lambda\lambda'}{1+T\|\lambda\|^2},$$

and the rest is straightforward. □

With the results of this subsection, what we aim at is the inequality $V_3(h(\lambda^*),\lambda^*) < V_3(0,\lambda^*)$, or more illustratively, the transform

$$\mathbb{E}_{\mathbb{Q}_{\lambda^*}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}_{\lambda^*}}|_{\mathscr{F}_T}F(W)\right] \Longrightarrow \mathbb{E}_{\mathbb{Q}_{\lambda^*}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}_{\lambda^*}}|_{\mathscr{F}_T}(F(W)-\langle h(\lambda^*),W_T\rangle)\right].$$

Let us close this section with a summarizing corollary.

**Corollary 3.4.** *Let $\theta^*$ and $\lambda^*$ be constants in $\mathbb{R}^d$ defined, respectively, by (2.4) and in Theorem 2.1, and let $g$ and $h$ be defined respectively in Proposition 3.1 and (3.4), or Lemma 3.3. Then,*

$$\begin{aligned}
V_3(\theta^*,g(\theta^*)) &\le V_3(\theta^*,0) < V_3(0,0), \\
V_3(h(\lambda^*),\lambda^*) &< V_3(0,\lambda^*) \le V_3(0,0).
\end{aligned}$$

# 4 Parameter search with two-time-scale stochastic approximation

The above corollary justifies the effectiveness of the combination of the two variance reduction techniques. It should be noted that there might exist $\theta$ such that $V_3(\theta,g(\theta)) < V_3(\theta^*,g(\theta^*))$, or

$\lambda$ such that $V_3(h(\lambda),\lambda) < V_3(h(\lambda^*),\lambda^*)$. The natural interest now directs to the parameter search for $(\theta,\lambda)$ such that

$$V_3(\theta,\lambda) < V_3(\theta^*,g(\theta^*)) \wedge V_3(h(\lambda^*),\lambda^*).$$

To this end, we will employ the so-called two-time-scale stochastic approximation algorithm,

$$\lambda_{n+1} = \Pi_{H^\lambda_{\sigma_1(n)}} \left[ \lambda_n - \varepsilon_n G_1(\theta_n,\lambda_n,W_{n+1}) \right],$$
$$\theta_{n+1} = \Pi_{H^\theta_{\sigma_2(n)}} \left[ \theta_n - \delta_n G_2(\theta_n,\lambda_n,W_{n+1}) \right],$$

where $\{H^\lambda_n\}_{n\in\{0\}\cup\mathbb{N}}$ and $\{H^\theta_n\}_{n\in\{0\}\cup\mathbb{N}}$ are suitable sequences of constraint sets and $(\theta_0,\lambda_0) \in H^\theta_0 \times H^\lambda_0$, and where $\{\delta_n\}_{n\in\{0\}\cup\mathbb{N}}$ is a sequences of non-increasing positive constants satisfying

$$\sum_{n=0}^{+\infty} \delta_n = +\infty, \quad \sum_{n=0}^{+\infty} \delta_n^2 < +\infty,$$

while its decay rate is *not identical* to that of the sequence $\{\varepsilon_n\}_{n\in\{0\}\cup\mathbb{N}}$.

*In what follows, $\{\varepsilon_n\}_{n\in\{0\}\cup\mathbb{N}}$ and $\{\delta_n\}_{n\in\{0\}\cup\mathbb{N}}$ are as the sequences given just above, while the constants $\theta^*$, $\lambda^*$, and the functions g, and h are as the ones appeared in Corollary 3.4.*

## 4.1  Equilibrated control variates and quasi-static importance sampling

Define the function $V_4(\lambda) := V_3(h(\lambda),\lambda)$. Before proceeding to the two-time-scale algorithm, we still need consider a stochastic iteration

$$\lambda_{n+1} = \Pi_{H_{\sigma(n)}} \left[ \lambda_n - \varepsilon_n G_1(h(\lambda_n),\lambda_n,W_{n+1}) \right]. \tag{4.1}$$

Then, we have the following.

**Proposition 4.1.** *The function $V_4$ is strictly convex on $\mathbb{R}^d$. Moreover, the sequence $\{\lambda_n\}_{n\in\mathbb{N}}$ in (4.1) converges $\mathbb{P}$-a.s. to $\lambda^\dagger$ such that $\nabla_\lambda V_4(\lambda^\dagger) = 0$, with $\lim_{n\uparrow+\infty}\sigma(n) < +\infty$, $\mathbb{P}$-a.s.*

*Proof.* Using the chain rule of the gradient, we get

$$
\begin{aligned}
\nabla_\lambda V_4(\lambda) = \nabla_\lambda V_3(h(\lambda),\lambda) &= \nabla_\lambda \left[ h(\lambda)' \right] \left[ \nabla_\theta V_3(\theta,\lambda) \right]|_{\theta=h(\lambda)} + \left[ \nabla_\lambda V_3(\theta,\lambda) \right]|_{\theta=h(\lambda)} \\
&= \left[ \nabla_\lambda V_3(\theta,\lambda) \right]|_{\theta=h(\lambda)},
\end{aligned}
$$

where the last equality follows from the definition of $h(\lambda)$. Moreover,

$$
\begin{aligned}
\operatorname{Hess}_\lambda \left[ V_4(\lambda) \right] &= \nabla'_\lambda \left( \left[ \nabla_\lambda V_3(\theta,\lambda) \right]|_{\theta=h(\lambda)} \right) \\
&= \left[ \operatorname{Hess}_\lambda V_3(\theta,\lambda) + \nabla'_\theta \nabla_\lambda V_3(\theta,\lambda) \right]|_{\theta=h(\lambda)} \\
&= \left[ \operatorname{Hess}_\lambda V_3(\theta,\lambda) \right]|_{\theta=h(\lambda)},
\end{aligned}
$$

where the last equality follows from $\nabla'_\theta \nabla_\lambda (V_3(\theta,\lambda)) = (\nabla'_\lambda \nabla_\theta (V_3(\theta,\lambda)))'$ and again from the definition of $h(\lambda)$. The positive-definiteness of $[\operatorname{Hess}_\lambda V_3(\theta,\lambda)]|_{\theta=h(\lambda)}$ proves the first assertion.

9

Next, with the help of the definition of $h(\lambda)$, we get for $y \in \mathbb{R}_0^d$,

$$
\begin{aligned}
y' \mathrm{Hess}_\lambda \left[ V_4(\lambda) \right] y &= \mathbb{E}_\mathbb{P} \left[ \left( T + \langle y, T\lambda - W_T \rangle^2 \right) e^{-\langle \lambda, W_T \rangle + \frac{1}{2} T \|\lambda\|^2} \left( F(W) - \langle h(\lambda), W_T \rangle \right)^2 \right] \\
&\geq T \mathbb{E}_\mathbb{P} \left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2} T \|\lambda\|^2} \left( F(W) - \langle h(\lambda), W_T \rangle \right)^2 \right] \\
&= T \mathbb{E}_\mathbb{P} \left[ e^{-\langle \lambda, W_T \rangle + \frac{1}{2} T \|\lambda\|^2} F(W)^2 \right],
\end{aligned}
$$

which tends to $+\infty$ as $\|\lambda\| \uparrow +\infty$. (The proof for this explosion can be found in the proof of Proposition 1 [1].) Together with $V_4(\lambda) > 0$ for $\lambda \in \mathbb{R}^d$, we get $\lim_{\|\lambda\| \uparrow +\infty} V_4(\lambda) = +\infty$. From the strict convexity of $V_4$, it then follows that there exists $\lambda^\dagger$ satisfying $\nabla_\lambda V_4(\lambda^\dagger) = 0$.

Note that the first assertion also ensures that the iteration (4.1) is a gradient-based stochastic approximation algorithm for a strictly convex function. To prove the almost sure convergence of the algorithm to $\lambda^\dagger$, it suffices to show that for each $m \in \mathbb{N}$,

$$
\mathbb{E}_\mathbb{P} \left[ \| G_1(h(\lambda), \lambda, W) \|^2 \mathbb{1} \left( \|\lambda\| \leq m \right) \right] < +\infty.
$$

To this end, observe first that for $\lambda$ such that $\|\lambda\| \leq m$, the Cauchy-Schwartz inequality and the Girsanov theorem yields

$$
\begin{aligned}
\| h(\lambda) \|^2 &\leq T^{-2} \left\| I_d - \frac{T \lambda \lambda'}{1 + T \|\lambda\|^2} \right\|_o^2 \mathbb{E}_\mathbb{P} \left[ F(W)^2 \right] \mathbb{E}_\mathbb{P} \left[ e^{-2\langle \lambda, W_T \rangle - T \|\lambda\|^2} \|W_T\|^2 \right] \\
&\leq \mathbb{E}_\mathbb{P} \left[ F(W)^2 \right] e^{T m^2} \left( d + 4 m^2 \right) < +\infty.
\end{aligned}
$$

Here, we use the fact that the squared matrix norm $\|A\|_o^2$ equals the maximum of eigenvalues of the positive-semidefinite matrix $A'A$ and in the above case there exists only two distinct eigenvalues of $1$ and $(1 + T \|\lambda\|^2)^{-2} \leq 1$. Hence,

$$
\begin{aligned}
\mathbb{E}_\mathbb{P} & \left[ \| G_1(h(\lambda), \lambda, W) \|^2 \mathbb{1} \left( \|\lambda\| \leq m \right) \right] \\
&\leq e^{\frac{3c-1}{c-1} T m^2} \mathbb{E}_\mathbb{P} \left[ |F(W) - \langle h(\lambda), W_T \rangle|^{4c} \mathbb{1} \left( \|\lambda\| \leq m \right) \right]^{\frac{1}{c}} \\
&\qquad \times \mathbb{E}_\mathbb{P} \left[ \left\| W_T - \frac{3c-1}{c-1} T \lambda \right\|^{\frac{2c}{c-1}} \mathbb{1} \left( \|\lambda\| \leq m \right) \right]^{1 - \frac{1}{c}},
\end{aligned}
$$

which is finite due to (2.1) and where the constant $c$ is given also in (2.1). $\qquad \square$

At first glance, it looks ideal to directly perform the algorithm (4.1), while not quite so from a practical point of view since the computation of $h(\lambda_n)$ at each step is not a smart choice (although $h(\lambda)$ is simplified to some extent in Lemma 3.3).

Consider then the following two-time-scale stochastic approximation algorithm,

$$
\lambda_{n+1} = \Pi_{H^\lambda_{\sigma_1(n)}} \left[ \lambda_n - \varepsilon_n G_1(\theta_n, \lambda_n, W_{n+1}) \right], \tag{4.2}
$$

$$
\theta_{n+1} = \Pi_{H^\theta_{\sigma_2(n)}} \left[ \theta_n - \delta_n G_2(\theta_n, \lambda_n, W_{n+1}) \right], \tag{4.3}
$$

with

$$\lim_{n\uparrow+\infty} \frac{\varepsilon_n}{\delta_n} = 0, \tag{4.4}$$

that is, the iteration (4.2) moves much slower than the other (4.3). The following is the main result of this subsection.

**Theorem 4.2.** *The sequence* $\{(\theta_n, \lambda_n)\}_{n\in\mathbb{N}}$ *defined in (4.2)-(4.3) with (4.4) converges* $\mathbb{P}$*-a.s. to* $(h(\lambda^\dagger), \lambda^\dagger)$*, with* $\lim_{n\uparrow+\infty} \sigma_1(n) < +\infty$ *and* $\lim_{n\uparrow+\infty} \sigma_2(n) < +\infty$, $\mathbb{P}$*-a.s.*

*Proof.* By Proposition 3.2, for each $\lambda \in \mathbb{R}^d$, the sequence $\{\theta_n\}_{n\in\mathbb{N}}$ in (3.3) converges $\mathbb{P}$*-a.s.* to $h(\lambda)$. Moreover, by Proposition 4.1, the sequence $\{\lambda_n\}_{n\in\mathbb{N}}$ in (4.1) converges $\mathbb{P}$*-a.s.* to $\lambda^\dagger$. The claim thus holds by Theorem 1.1 of [3] with the condition (4.4). $\square$

An interpretation of the two-time-scale stochastic approximation algorithm (4.3)-(4.2) with the condition (4.4) is as follows. The fast $\theta$ sees $\lambda$ as "quasi-static", and thus the algorithm (4.3) has an effect similar to fixing $\lambda$ and running the (single-time-scale) algorithm (3.3) for a long time. Meanwhile, in the algorithm (4.3), $\theta_n$ can be seen as a close approximation of $h(\lambda_n)$, that is, "equilibrated", and thus almost same as running the (single-time-scale) algorithm (4.1).

## 4.2 Equilibrated importance sampling and quasi-static control variates

The opposite scheme, that is, with equilibrated importance sampling and quasi-static control variates, also works. Define a function $V_5(\theta) := V_3(\theta, g(\theta))$, let $\{\varepsilon_n\}_{n\in\{0\}\cup\mathbb{N}}$ be defined for (2.6), and set a stochastic iteration

$$\theta_{n+1} = \Pi_{H^\theta_{\sigma_2(n)}} \left[ \theta_n - \delta_n G_2(\theta_n, g(\theta_n), W_{n+1}) \right]. \tag{4.5}$$

Also, set a two-time-scale stochastic approximation algorithm

$$\lambda_{n+1} = \Pi_{H^\lambda_{\sigma_1(n)}} \left[ \lambda_n - \varepsilon_n G_1(\theta_n, \lambda_n, W_{n+1}) \right],$$

$$\theta_{n+1} = \Pi_{H^\theta_{\sigma_2(n)}} \left[ \theta_n - \delta_n G_2(\theta_n, \lambda_n, W_{n+1}) \right],$$

with

$$\lim_{n\uparrow+\infty} \frac{\delta_n}{\varepsilon_n} = 0, \tag{4.6}$$

that is, the iteration for $\lambda$ moves much faster than the one for $\theta$. Then, the following holds true.

**Theorem 4.3.** *The function* $V_5$ *is strictly convex on* $\mathbb{R}^d$, *and the sequence* $\{\theta_n\}_{n\in\mathbb{N}}$ *in (4.5) converges* $\mathbb{P}$*-a.s. to* $\theta^\dagger$ *such that* $\nabla_\theta V_5(\theta^\dagger) = 0$. *Moreover, the sequence* $\{(\theta_n, \lambda_n)\}_{n\in\mathbb{N}}$ *in (4.2)-(4.3) with (4.6) converges* $\mathbb{P}$*-a.s. to* $(\theta^\dagger, g(\theta^\dagger))$*, with* $\lim_{n\uparrow+\infty} \sigma_1(n) < +\infty$ *and* $\lim_{n\uparrow+\infty} \sigma_2(n) < +\infty$, $\mathbb{P}$*-a.s.*

11

*Proof.* For each $\theta \in \mathbb{R}^d$, the chain rule of the gradient yields

$$
\begin{aligned}
\nabla_\theta V_5(\theta) &= [\nabla_\theta V_3(\theta, \lambda)]|_{\lambda=g(\theta)} + [\nabla'_\theta g(\theta)] [\nabla_\lambda V_3(\theta, \lambda)]|_{\lambda=g(\theta)} \\
&= [\nabla_\theta V_3(\theta, \lambda)]|_{\lambda=g(\theta)},
\end{aligned}
$$

and moreover,

$$
\begin{aligned}
\text{Hess}_\theta [V_5(\theta)] &= \nabla'_\theta \left( [\nabla_\theta V_3(\theta, \lambda)]|_{\lambda=g(\theta)} \right) \\
&= \left[ \text{Hess}_\theta V_3(\theta, \lambda) + \nabla'_\lambda \nabla_\theta V_3(\theta, \lambda) \right]|_{\lambda=g(\theta)} \\
&= [\text{Hess}_\theta V_3(\theta, \lambda)]|_{\lambda=g(\theta)},
\end{aligned}
$$

whose positive-definiteness proves the first assertion.

Next, observe that

$$
\begin{aligned}
\text{Hess}_\theta [V_5(\theta)] &= 2\mathbb{E}_\mathbb{P} \left[ e^{-\langle g(\theta), W_T \rangle + \frac{1}{2}T\|g(\theta)\|^2} W_T W'_T \right] \\
&= 2 e^{T\|g(\theta)\|^2} T \left( I_d + T g(\theta) g(\theta)' \right).
\end{aligned}
$$

Hence, for $y \in \mathbb{R}^d$ such that $\|y\| = 1$, we get

$$
y' \text{Hess}_\theta [V_5(\theta)] y = 2 e^{T\|g(\theta)\|^2} T \left( 1 + T \langle y, g(\theta) \rangle^2 \right) \geq 2T,
$$

uniformly in $\theta \in \mathbb{R}^d$. Together with $V_5(\theta) > 0$ for $\theta \in \mathbb{R}^d$, we get $\lim_{\|\theta\|\uparrow+\infty} V_5(\theta) = +\infty$. From the strict convexity of $V_5$, it follows that there exists $\theta^\dagger$ satisfying $\nabla_\theta V_5(\theta^\dagger) = 0$.

Finally, recall that we have shown in Proposition 3.1 that for each $\theta \in \mathbb{R}^d$, $\{\lambda_n\}_{n\in\mathbb{N}}$ in the stochastic iteration (3.2) converges $\mathbb{P}$-*a.s.* to $g(\theta)$. The last assertion thus holds by Theorem 1.1 of [3] with the second assertion and with the condition (4.6). $\qquad\square$

# 5 Adaptive Monte Carlo variance reduction

We have proved so far that the two-time-scale stochastic approximation algorithms converge almost surely to the root of the gradient of the interested variance. The aim of this section is to prove the asymptotic normality of the empirical mean and empirical variance of adaptive Monte Carlo simulations, in our "two-time-scale" framework. To this end, define $Y : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{C}([0,T]; \mathbb{R}^d) \mapsto \mathbb{R}$ by

$$
Y(\theta, \lambda, W) := e^{-\langle \lambda, W_T \rangle - \frac{1}{2}T\|\lambda\|^2} \left( F(W + \lambda) - \langle \theta, W_T + \lambda T \rangle \right),
$$

with its empirical mean and empirical variance defined respectively by

$$
EM_n := \frac{1}{n} \sum_{k=1}^n Y(\theta_{k-1}, \lambda_{k-1}, W_k),
$$

and

$$
EV_n^2 := \frac{1}{n} \sum_{k=1}^n Y^2(\theta_{k-1}, \lambda_{k-1}, W_k) - EM_n^2.
$$

The following is the main result of this section.

**Theorem 5.1.** *Let the sequences $\{\theta_n\}_{n\in\mathbb{N}}$ and $\{\lambda_n\}_{n\in\mathbb{N}}$ are generated by the iterations (4.2)-(4.3), either with (4.4) or with (4.6). Then, it holds that*

$$\mathbb{P}\left[\lim_{n\uparrow+\infty} EM_n = C\right] = 1,$$

*and as $n\uparrow+\infty$,*

$$\sqrt{n}\frac{EM_n - C}{EV_n} \xrightarrow{\mathscr{L}} \mathscr{N}(0,1),$$

*under the probability measure $\mathbb{P}$. Moreover,*

*(i) If the iterations (4.2)-(4.3) are performed with (4.4), then it holds under the probability measure $\mathbb{P}$ that*

$$\mathbb{P}\left[\lim_{n\uparrow+\infty} EV_n^2 = V_4(\lambda^\dagger)\right] = 1, \text{ and } \sqrt{n}(EM_n - C) \xrightarrow{\mathscr{L}} \mathscr{N}(0, V_4(\lambda^\dagger)), \text{ as } n\uparrow+\infty.$$

*(ii) If the iterations (4.2)-(4.3) are performed with (4.6), then it holds under the probability measure $\mathbb{P}$ that*

$$\mathbb{P}\left[\lim_{n\uparrow+\infty} EV_n^2 = V_5(\theta^\dagger)\right] = 1, \text{ and } \sqrt{n}(EM_n - C) \xrightarrow{\mathscr{L}} \mathscr{N}(0, V_5(\theta^\dagger)), \text{ as } n\uparrow+\infty.$$

*Proof.* Following [2] (Theorem 1-2 and Corollary 1), it suffices to show that

$$\sup_{n\in\mathbb{N}} \mathbb{E}_{\mathbb{P}}\left[Y^4(\theta_{n-1}, \lambda_{n-1}, W_n)\right] < +\infty,$$

and that the expectations $\mathbb{E}_{\mathbb{P}}[Y^2(\theta, \lambda, W)]$ and $\mathbb{E}_{\mathbb{P}}[Y^4(\theta, \lambda, W)]$ are continuous at $(\theta, \lambda) = (h(\lambda^\dagger), \lambda^\dagger)$ for (i), while at $(\theta, \lambda) = (\theta^\dagger, g(\theta^\dagger))$ for (ii). First, with the help of the Girsanov theorem and the Hölder's inequality, we get

$$\mathbb{E}_{\mathbb{P}}\left[Y^4(\theta, \lambda, W)\right] = \mathbb{E}_{\mathbb{P}}\left[e^{-3\langle\lambda, W_T\rangle + \frac{3}{2}T\|\lambda\|^2}(F(W) - \langle\theta, W_T\rangle)^4\right] \qquad (5.1)$$

$$\leq e^{\frac{3(4c-1)}{2(c-1)}T\|\lambda\|^2}\mathbb{E}_{\mathbb{P}}\left[|F(W) - \langle\theta, W_T\rangle|^{4c}\right]^{\frac{1}{c}},$$

where the constant $c$ is the one given in (2.1). Due to the definition of the constraint sequences $\{H_n^\theta\}_{n\in\mathbb{N}}$ and $\{H_n^\lambda\}_{n\in\mathbb{N}}$ in the iteration (4.2)-(4.3), and $\lim_{n\uparrow+\infty}\sigma_1(n) < +\infty$ and $\lim_{n\uparrow+\infty}\sigma_2(n) < +\infty$, $\mathbb{P}$-*a.s.*, shown in Theorem 4.2-4.3, it is then clear that $\sup_{n\in\mathbb{N}}\mathbb{E}_{\mathbb{P}}[Y^4(\theta_{n-1}, \lambda_{n-1}, W_n)] < +\infty$.

Next, we have $\mathbb{E}_{\mathbb{P}}[Y^2(\theta, \lambda, W)] = V_3(\theta, \lambda) + C^2$, and thus its continuity in $(\theta, \lambda)$ is clear, while the continuity of the fourth moment follows from its finiteness uniformly on a neighborhood of either $(h(\lambda^\dagger), \lambda^\dagger)$ or of $(\theta^\dagger, g(\theta^\dagger))$, and from that the ($\omega$-pointwise) continuity of the expression inside the right hand side expectation of (5.1). $\qquad\square$

# 6 Numerical illustration

Let $W := (W_1, \ldots, W_d)$ be a standard normal random vector in $\mathbb{R}^d$ (with independent components) under the probability measure $\mathbb{P}$. We consider a discrete approximation of the so-called Asian payoff, one of the most standard financial structures, in the well known Black-Scholes framework,

$$F(W) = e^{-rT} \max \left[ \frac{1}{d} \sum_{n=1}^{d} S_0 \exp \left[ \sum_{k=1}^{n} \left( \left( r - \frac{1}{2}\sigma^2 \right) \frac{T}{d} + \sqrt{\sigma^2 \frac{T}{d}} W_k \right) \right] - K, 0 \right],$$

for a strike level $K > 0$. In this experiments, we fix $S_0 = 50$, $r = 0.05$, $\sigma = 0.10$ or $0.30$, $T = 1$, and $d = 16$ and generate $N = 5e+4$ of iid Monte Carlo summands. Based upon the iid Monte Carlo summands, we run the two-time-scale stochastic approximation algorithm (4.3)-(4.2) with

$$\delta_n = \frac{\alpha_\delta}{(n+1)^{\beta_\delta}}, \quad \varepsilon_n = \frac{\alpha_\varepsilon}{(n+1)^{\beta_\varepsilon}},$$

where $\alpha_\delta$ and $\alpha_\varepsilon$ are non-negative and $\beta_\delta$ and $\beta_\varepsilon$ are constants in $(1/2, 1)$. Recall that the case $\beta_\delta < \beta_\varepsilon$ accelerates the iteration for the control variates component compared to the importance sampling one, that is, "equilibrated control variates and quasi-static importance sampling" (AD2), while the case $\beta_\delta > \beta_\varepsilon$ corresponds to "equilibrated importance sampling and quasi-static control variates" (AD3). Evidently, we can perform the adaptive method with the importance sampling component only (AD1) by setting $\alpha_\delta = 0$, which reduces the two-time-scale algorithm to the single-time-scale algorithm (2.6), and this reduces to the framework of [2]. We set the increasing compact constraints as $H_n^\lambda = H_n^\theta = \{x \in \mathbb{R}^d : \|x\| \leq 100\ln(100(n+1))\}$, and start each iteration with the origin, that is, $\theta_0 = \lambda_0 = (0, \ldots, 0)'$. We will evaluate the variance reduction efficiency via the ratio of variances (vratio), defined by

$$\text{(vratio)} := \frac{V_3(0,0)}{EV_N}.$$

Results for our experiment is reported in Table 1. Moreover, for illustration purpose, we provide in Figure 1 various results along progress of the adaptive Monte Carlo variance reduction procedure in the case $(\sigma, K) = (0.30, 55)$. The upper two and the lower left are of (AD2), that is, the case $(\alpha_\delta, \beta_\delta, \alpha_\varepsilon, \beta_\varepsilon) = (1e\text{-}2, 7/11, 1.1e\text{-}2, 10/11)$. The lower right figure compares convergences of Monte Carlo average in plain Monte Carlo simulation (MC), and (AD1) and (AD2), while three dotted lines indicate $2.21 \pm 0.5\%$.

In this specific example, the implementation of (AD1) takes approximately twice as long per replication as (MC), mainly due to the addition of the stochastic approximation algorithm. This implies that the effective improvements of (AD1) are about half of their variance ratios. We have observed that the adaptive component of Monte Carlo simulations itself requires very little additional computation burden. The implementation of either (AD2) or (AD3) takes about three times as (MC). Thus, the actual effective improvements are about a third of the variance ratios, and in this sense, there are some cases where it is clever to stop at (AD1) and not to go further into either
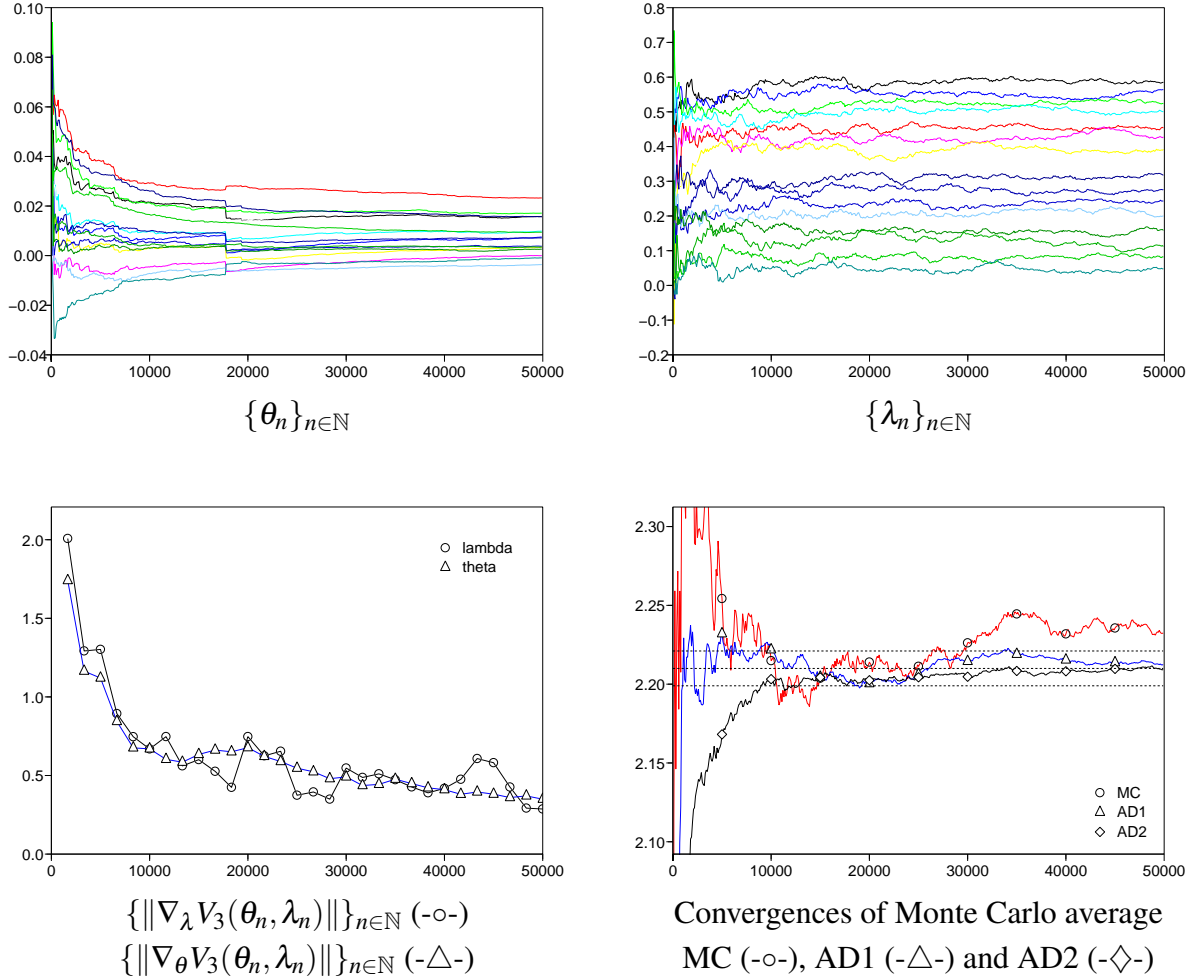
$$\{\theta_n\}_{n\in\mathbb{N}}$$

$$\{\lambda_n\}_{n\in\mathbb{N}}$$

$$\{\|\nabla_\lambda V_3(\theta_n,\lambda_n)\|\}_{n\in\mathbb{N}} \text{ (-○-)}$$
$$\{\|\nabla_\theta V_3(\theta_n,\lambda_n)\|\}_{n\in\mathbb{N}} \text{ (-△-)}$$

Convergences of Monte Carlo average
MC (-○-), AD1 (-△-) and AD2 (-◇-)

Figure 1: Numerical results of the case $(\sigma, K) = (0.3, 55)$ with $(\alpha_\delta, \beta_\delta, \alpha_\varepsilon, \beta_\varepsilon) = (1\text{e-}2, 7/11, 1.1\text{e-}2, 10/11)$ for (AD2).

(AD2) or (AD3). In practice, it is however rare to consider variance reduction for Monte Carlo simulation when the structure of the function $F$ is simple, just as for this experiment, and clearly the additional computational cost of the stochastic approximation algorithm becomes negligible as $F$ are more complicated. Finally, to reduce computational cost (but with some variance ratio given up), it certainly a great scheme to stop the stochastic approximation iteration at some point $N_s$, and to only perform Monte Carlo simulation from that point on in cooperation with variance reduction with the obtained parameter $(\theta_{N_s}, \lambda_{N_s})$. For example, in the case given in Figure 1, the absolute gradients drop very fast until, say $N_s = 1.5\text{e+}3$, at which may be a good point to stop the stochastic approximation algorithms. In Table 2, we report variance ratios based on 5e+4 replications for various $N_s$ in the case $\sigma = 0.30$. As expected, stopping during very early stage can achieve most part of variance reduction effect. The choice of the stopping point $N_s$ is a difficult trade-off problem between variance reduction effect and computation cost, which can be different for different

problems, that is, the structure of the function $F$, and also depends highly on the choice of $(\theta_0, \lambda_0)$.

# 7 Concluding remarks

In this paper, we have developed and analyzed an application of the two-time-scale stochastic approximation algorithm in the optimal parameter search for the combined control variates and importance sampling framework. The algorithm converges almost surely to the unique root of the gradients of the variance either in the sense of "equilibrated control variates and quasi-static importance sampling," or of "equilibrated importance sampling and quasi-static control variates." In each case, the almost sure convergence guarantees the incorporation of the algorithms into the adaptive Monte Carlo variance reduction procedure. It is a great advantage that our method is built on a quite general formulation of the control variates and the importance sampling and is thus directly applicable to various Monte Carlo simulations, unlike many problem-specific variance reduction methods in the literature. We have illustrated the effectiveness of our method in a financial engineering example through numerical results. Finally, it would also be interesting to analyze our method in a pure-jump Lévy process framework, which is left as a future research.

# References

[1] Arouna, B. (2004) Robbins-Monro algorithms and variance reduction in finance, *Journal of Computational Finance*, **7**, 35-61.

[2] Arouna, B (2004) Adaptive Monte Carlo method, a variance reduction technique, *Monte Carlo Methods and Applications*, **10**, 1-24.

[3] Borkar, V.S. (1997) Stochastic approximation with two time scales, *Systems & Control Letters*, **29**, 291-294.

[4] Chen, K.-F., Zhu, Y. (1986) *Stochastic Approximation Procedure with randomly varying truncations,* Scientia Sinica Series.

[5] Delyon, B. (1996) General results on the convergence of stochastic algorithms, *IEEE Transactions on Automatic Control*, **41**, 1245-1255.

[6] Glasserman, P., Heidelberger, P., Shahabuddin, P. (1999) Asymptotic optimal importance sampling and stratification for pricing path-dependent options, *Mathematical Finance*, **9**, 117-152.

[7] Kawai, R. (2007) Optimal importance sampling via stochastic approximation for Lévy processes with exponential tilting, *Preprint*.

[8] Randhawa, R.S., Juneja, S. (2004) Combining importance sampling and temporal difference control variates to simulate Markov Chains, *ACM Transactions on Modeling and Computer Simulation*, **14**, 1-30.

[9] Su, Y., Fu. M.-C. (1998) Optimal importance sampling in securities pricing, *Journal of Computational Finance,* **5**, 27-50.

| | | $\alpha_\delta$ | $\beta_\delta$ | $\alpha_\varepsilon$ | $\beta_\varepsilon$ | price | vratio |
|---|---|---|---|---|---|---|---|
| $\sigma = 0.10$ | | | | | | | |
| $K$ | | $\alpha_\delta$ | $\beta_\delta$ | $\alpha_\varepsilon$ | $\beta_\varepsilon$ | price | vratio |
| 45 | (MC) | - | - | - | - | 6.058 | (8.76) |
| | (AD1) | - | - | 8e-3 | 10/11 | 6.051 | 10.84 |
| | (AD2) | 5e-3 | 7/11 | 6e-3 | 10/11 | 6.053 | 33.15 |
| | (AD3) | 0.1 | 10/11 | 1e-3 | 7/11 | 6.049 | 40.78 |
| 50 | (MC) | - | - | - | - | 1.931 | (4.96) |
| | (AD1) | - | - | 3e-2 | 10/11 | 1.916 | 6.90 |
| | (AD2) | 2e-2 | 7/11 | 2e-2 | 10/11 | 1.929 | 10.23 |
| | (AD3) | 2e-2 | 10/11 | 2e-2 | 7/11 | 1.918 | 9.95 |
| 55 | (MC) | - | - | - | - | 0.207 | (0.56) |
| | (AD1) | - | - | 1.0 | 10/11 | 0.199 | 19.82 |
| | (AD2) | 1e-4 | 7/11 | 1.6 | 10/11 | 0.203 | 21.82 |
| | (AD3) | 1e-4 | 10/11 | 0.5 | 7/11 | 0.201 | 22.09 |
| $\sigma = 0.30$ | | | | | | | |
| $K$ | | $\alpha_\delta$ | $\beta_\delta$ | $\alpha_\varepsilon$ | $\beta_\varepsilon$ | price | vratio |
| 45 | (MC) | - | - | - | - | 7.187 | (59.72) |
| | (AD1) | - | - | 4e-3 | 10/11 | 7.149 | 8.55 |
| | (AD2) | 1e-2 | 7/11 | 3e-3 | 10/11 | 7.152 | 14.11 |
| | (AD3) | 6e-3 | 10/11 | 1e-3 | 7/11 | 7.145 | 12.94 |
| 50 | (MC) | - | - | - | - | 4.202 | (40.35) |
| | (AD1) | - | - | 8e-3 | 10/11 | 4.165 | 8.99 |
| | (AD2) | 5e-3 | 7/11 | 6e-3 | 10/11 | 4.174 | 11.42 |
| | (AD3) | 6e-3 | 10/11 | 6e-3 | 7/11 | 4.162 | 11.62 |
| 55 | (MC) | - | - | - | - | 2.231 | (23.19) |
| | (AD1) | - | - | 1e-2 | 10/11 | 2.212 | 7.41 |
| | (AD2) | 1e-2 | 7/11 | 1.1e-2 | 10/11 | 2.209 | 12.58 |
| | (AD3) | 1e-2 | 10/11 | 1e-2 | 7/11 | 2.213 | 14.24 |

Table 1: Numerical results (values in the parenthesis for (MC) are raw variance, not variance ratio)

| | | $\sigma = 0.30$ | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $K$ | $N_s$ | 5e+3 | 1e+4 | 1.5e+4 | 2e+4 | 2.5e+4 | 3e+4 | 5e+4 |
| 45 | (AD1) | 8.15 | 8.39 | 8.47 | 8.50 | 8.52 | 8.54 | 8.55 |
| | (AD2) | 13.25 | 13.71 | 13.80 | 13.94 | 14.03 | 14.06 | 14.11 |
| | (AD3) | 12.40 | 12.79 | 12.85 | 12.81 | 12.89 | 12.90 | 12.94 |
| 50 | (AD1) | 8.27 | 8.66 | 8.82 | 8.88 | 8.93 | 8.95 | 8.99 |
| | (AD2) | 10.47 | 10.86 | 10.96 | 11.29 | 11.36 | 11.38 | 11.42 |
| | (AD3) | 10.94 | 11.44 | 11.34 | 11.46 | 11.62 | 11.58 | 11.62 |
| 55 | (AD1) | 6.33 | 6.83 | 7.06 | 7.20 | 7.28 | 7.34 | 7.41 |
| | (AD2) | 9.68 | 10.47 | 10.79 | 11.20 | 11.41 | 11.48 | 12.58 |
| | (AD3) | 13.27 | 13.93 | 13.88 | 14.05 | 14.16 | 14.25 | 14.24 |

Table 2: Variance ratios based on 5e+4 replications when stochastic approximation algorithm is stopped at $N_s$.