



**HAL**  
open science

# Coupling Importance Sampling and Multilevel Monte Carlo using Sample Average Approximation

Ahmed Kebaier, Jérôme Lelong

► **To cite this version:**

Ahmed Kebaier, Jérôme Lelong. Coupling Importance Sampling and Multilevel Monte Carlo using Sample Average Approximation. 2017. hal-01214840v3

**HAL Id: hal-01214840**

**<https://hal.science/hal-01214840v3>**

Preprint submitted on 4 Jul 2017 (v3), last revised 7 Jul 2017 (v4)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Coupling Importance Sampling and Multilevel Monte Carlo using Sample Average Approximation

Ahmed Kebaier\* & Jérôme Lelong†

July 4, 2017

## Abstract

In this work, we propose a smart idea to couple importance sampling and Multilevel Monte Carlo (MLMC). We advocate a per level approach with as many importance sampling parameters as the number of levels, which enables us to compute the different levels independently. The search for parameters is carried out using sample average approximation, which basically consists in applying deterministic optimisation techniques to a Monte Carlo approximation rather than resorting to stochastic approximation. Our innovative estimator leads to a robust and efficient procedure reducing both the discretization error (the bias) and the variance for a given computational effort. In the setting of discretized diffusions, we prove that our estimator satisfies a strong law of large numbers and a central limit theorem with optimal limiting variance, in the sense that this is the variance achieved by the best importance sampling measure (among the class of changes we consider), which is however non tractable. Finally, we illustrate the efficiency of our method on several numerical challenges coming from quantitative finance and show that it outperforms the standard MLMC estimator.

**AMS 2000 Mathematics Subject Classification.** 60F05, 62F12, 65C05, 60H35.

**Key Words and Phrases.** Sample average approximation, Multilevel Monte Carlo, variance reduction, Uniform strong large law of numbers, Central limit theorem, Importance Sampling.

## 1 Introduction

Expectation involving a stochastic process are often computed using a Monte Carlo method combined with a discretization scheme. For instance, computing a hedging portfolio in finance uses these tools. Generally, the asset price follows a diffusion process  $(X_t)_{0 \leq t \leq T}$ , which a stochastic differential equation (SDE)

$$dX_t = b(X_t)dt + \sigma(X_t)dW_t, \quad X_0 = x \in \mathbb{R}^d \quad (1.1)$$

where  $b : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $\sigma : \mathbb{R}^d \rightarrow \mathcal{M}_{d \times q}$  and  $W$  is a Brownian motion with values in  $\mathbb{R}^q$  defined on some given probability space  $(\Omega, (\mathcal{F}_t)_{0 \leq t \leq T}, \mathbb{P})$  with finite time horizon  $T > 0$ . The process  $X$

---

\*Université Paris 13, Sorbonne Paris Cité, LAGA, CNRS (UMR 7539), kebaier@math.univ-paris13.fr. This research benefited from the support of the chair Risques Financiers, Fondation du Risque and the Laboratory of Excellence MME-DII (<http://labex-mme-dii.u-cergy.fr/>).

†Univ. Grenoble Alpes, Laboratoire Jean Kuntzmann, 51, rue des Mathématiques, BP 53, 38041 Grenoble, Cedex 09, France jerome.lelong@univ-grenoble-alpes.fr. This project was supported by the Finance for Energy Market Research Centre, [www.fime-lab.org](http://www.fime-lab.org).

hardly ever has an explicit solution, which implies the use of a discretization scheme in order to simulate it. Consider the continuous time Euler approximation  $X^n$  with time step  $\delta = T/n$  given by

$$dX_t^n = b(X_{\eta_n(t)}^n)dt + \sigma(X_{\eta_n(t)}^n)dW_t, \quad \eta_n(t) = \lfloor t/\delta \rfloor \delta.$$

This work aims at combining importance sampling with different discretization methods: first, we study the use of importance sampling for standard case of Euler Monte Carlo and then we apply it to multilevel Monte Carlo. Many different changes of measure can be seen to design importance sampling. When working with Lévy processes, it is common to use the Esscher transform to introduce a new family of measures. For Brownian driven SDEs, the Esscher transform actually corresponds to a Gaussian change of measure in the spirit of the Girsanov theorem. Following the ideas of Arouna [1], we consider a parametric family of stochastic processes  $(X_t(\theta))_{0 \leq t \leq T}$ , with  $\theta \in \mathbb{R}^q$ , driven by a Brownian motion with linear drift

$$dX_t(\theta) = (b(X_t(\theta)) + \sigma(X_t(\theta))\theta) dt + \sigma(X_t(\theta))dW_t.$$

We also define the continuous time Euler approximation  $X^n(\theta)$  of the process  $X(\theta)$ . From Girsanov's Theorem, the process  $(B_t^\theta \triangleq W_t + \theta t)_{t \leq T}$  is a Brownian motion under the new probability measure  $\mathbb{P}_\theta$  equivalent to  $\mathbb{P}$  and such that

$$\frac{d\mathbb{P}_\theta}{d\mathbb{P}} \Big|_{\mathcal{F}_t} = \exp \left( -\theta \cdot W_t - \frac{1}{2}|\theta|^2 t \right) \triangleq \mathcal{E}^-(W, \theta).$$

Therefore,

$$\mathbb{E}_\mathbb{P}[\psi(X_T)] = \mathbb{E}_{\mathbb{P}_\theta}[\psi(X_T(\theta))] = \mathbb{E}_\mathbb{P}[\psi(X_T(\theta))\mathcal{E}^-(W, \theta)]. \quad (1.2)$$

This equality still holds when replacing  $X$  (resp.  $X(\theta)$ ) by its Euler scheme  $X^n$  (resp.  $X^n(\theta)$ ). The l.h.s. and r.h.s. expectations are both computed under the same probability measure. In the following, we will always use the measure  $\mathbb{P}$  and therefore we will not write it anymore. The idea of importance sampling Monte Carlo is to use the r.h.s of (1.2) to build a Monte Carlo estimator using  $X^n(\theta)$  with  $\theta$  given by

$$\theta^* = \underset{\theta \in \mathbb{R}^q}{\operatorname{argmin}} \operatorname{Var}(\psi(X_T(\theta))\mathcal{E}^-(W, \theta)).$$

Importance sampling for Euler Monte Carlo is studied in Section 2: first, we investigate how to approximate  $\theta^*$  in practice and second we prove a Monte Carlo estimator using this approximation of  $\theta^*$  satisfies both a strong law of large numbers and a central limit theorem when both  $n$  and the number of samples go to infinity. This result extends the limit theorems obtained in [23], in which the authors investigated the case of a fixed number of discretization steps  $n$ . The error induced by using  $\mathbb{E}[\psi(X_T^n(\theta))]$  instead of  $\mathbb{E}[\psi(X_T(\theta))]$  is called the discretization error and is responsible for the bias of the Euler Monte Carlo estimator, while the Monte Carlo approximation only impacts the variance. The two errors are balanced when the number of samples  $N$  of the Monte Carlo method is chosen as  $N = n^2$ , which to overall complexity of  $n^3$ . In order to reduce the bias for a given computational effort, Kebaier [24] proposed to use the Statistical Romberg method, which combines discretization schemes on two nested time grids. This method was generalized by Giles [13] who proposed to use a multilevel Monte Carlo algorithm following the line of Heinrich's multilevel method for parametric integration [19].

Let  $m, L \in \mathbb{N}$  with  $m \geq 2$  and  $L > 0$ , the idea of the multilevel method is to write the expectation on the finest time grid as a telescopic sum involving all the other grids (referred to as levels)

$$\mathbb{E}[\psi(X_T^{m^L})] = \mathbb{E}[\psi(X_T^{m^0})] + \sum_{\ell=1}^L \mathbb{E}[\psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}})] \quad (1.3)$$

and then to approximate each expectation by a Monte Carlo method with a well chosen number of samples to balance the errors between the different terms. We refer the reader to the extensive literature linked to MLMC for more details, see e.g. [3, 9, 10, 11, 14, 16, 15, 18, 20, 27]. For a fixed computational budget, the use of multilevel techniques clearly reduces the bias error, but in many situations the high variance also brings in a significant inaccuracy, which naturally leads to trying to couple MLMC with variance reduction techniques.

In this work, we focus on coupling importance sampling with MLMC. In [5] and [17], the authors chose to apply MLMC to the right hand side of (1.2) coming up with

$$\begin{aligned} \mathbb{E}[\psi(X_T^{m^L})] &= \mathbb{E} \left[ \psi(X_T^{m^0}(\lambda)) \mathcal{E}^-(W, \lambda) \right] \\ &+ \sum_{\ell=1}^L \mathbb{E} \left[ (\psi(X_T^{m^\ell}(\lambda)) - \psi(X_T^{m^{\ell-1}}(\lambda))) \mathcal{E}^-(W, \lambda) \right]. \end{aligned} \quad (1.4)$$

This approach mixes all the levels through the optimization of the parameter  $\lambda$  and breaks the independence between the levels of the multilevel approach, which made it so popular and easy to implement.

Instead of using (1.4), we would rather apply importance sampling to each expectation in the telescopic sum of (1.3) to obtain for  $\lambda_1, \dots, \lambda_L \in \mathbb{R}^q$

$$\begin{aligned} \mathbb{E}[\psi(X_T^{m^L})] &= \mathbb{E} \left[ \psi(X_T^{m^0}(\lambda_0)) \mathcal{E}^-(W, \lambda_0) \right] \\ &+ \sum_{\ell=1}^L \mathbb{E} \left[ (\psi(X_T^{m^\ell}(\lambda_\ell)) - \psi(X_T^{m^{\ell-1}}(\lambda_\ell))) \mathcal{E}^-(W, \lambda_\ell) \right]. \end{aligned}$$

Our importance multilevel estimator is obtained by applying a Monte Carlo method to each of the levels  $\ell$  with  $N_\ell$  samples

$$\begin{aligned} Q_L(\lambda_0, \dots, \lambda_L) &= \frac{1}{N_0} \sum_{k=1}^{N_0} \psi(\tilde{X}_{T,0,k}^{m^0}(\lambda_0)) \mathcal{E}^-(\tilde{W}_{0,k}, \lambda_0) \\ &+ \sum_{\ell=1}^L \frac{1}{N_\ell} \sum_{k=1}^{N_\ell} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\lambda_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\lambda_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \lambda_\ell) \end{aligned} \quad (1.5)$$

The samples used in the different levels are independent and within each level they are i.i.d. For any  $\ell \geq 0$ , the variables  $\tilde{X}_{T,\ell,k}^{m^\ell}(\lambda_\ell)$  (resp.  $\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\lambda_\ell)$  when  $\ell > 0$ ) are the terminal values of the Euler schemes of  $X(\lambda_\ell)$  with  $m^\ell$  (resp.  $m^{\ell-1}$ ) time steps built using the same Brownian path  $\tilde{W}_{\ell,k}$ . The variance of the importance sampling MLMC estimator is given by

$$\text{Var}[Q_L] = N_0^{-1} \sigma_0(\lambda_0)^2 + \sum_{\ell=1}^L N_\ell^{-1} \frac{(m-1)T}{m^\ell} \sigma_\ell^2(\lambda_\ell)$$

where

$$\begin{aligned}\sigma_0^2(\lambda_0) &\triangleq \text{Var}[\psi(X_T^{m^0}(\lambda_0))\mathcal{E}^-(W, \lambda_0)] \\ \sigma_\ell^2(\lambda_\ell) &\triangleq \frac{m^\ell}{(m-1)T} \text{Var}\left[\left(\psi(X_T^{m^\ell}(\lambda_\ell)) - \psi(X_T^{m^{\ell-1}}(\lambda_\ell))\right)\mathcal{E}^-(W, \lambda_\ell)\right].\end{aligned}$$

By allowing for one importance sampling parameter  $\lambda_\ell$  per level, our approach has many advantages over [5, 17]. First, the computations within the different levels remain independent. Second, the variance of each level  $\ell$  only depends on  $\lambda_\ell$ , which reduces the global minimization problem to several smaller minimization problems. Third, we actually minimize the real variance of the estimator and not its asymptotic value and more importantly it can be implemented without knowing  $\nabla\psi$ , which however appears in the central limit theorem for MLMC. The new idea of using one importance sampling parameter per level was later taken up in [6] but coupled with stochastic approximation to build adaptive estimators.

Actually, minimizing  $\lambda \mapsto \sigma_\ell^2(\lambda)$  can be achieved by using the randomly truncated Robbins Monro algorithm proposed by Chen et al. [7, 8] and later investigated in the context of importance sampling by Lapeyre and Lelong [25] and Lelong [26]. The numerical stability of these stochastic algorithms strongly depends on the choice of the descent step — often referred to as the gain sequence — which proves to be highly sensitive in practice. To overcome this difficulty, Jourdain and Lelong [23] proposed to apply deterministic optimization techniques to sample average estimators to search for the optimal parameter. Following their methodology, we define  $\sigma_{\ell, N'_\ell}^2$  as the sample average approximation of  $\sigma_\ell^2$  with  $N'_\ell$  samples using the standard empirical Monte Carlo estimator of the variance. We assume that the samples used in  $\sigma_{\ell, N'_\ell}^2$  are independent of those used in  $Q_L$ . We refer to Section 3 for more details on the samples used in the different approximations. Now, we sketch the algorithm corresponding to our method.

```

1 for  $\ell = 0 : L$  do
2   Sample the random function  $\lambda \mapsto \sigma_{\ell, N'_\ell}(\lambda)$ . //  $\sigma_{\ell, N'_\ell}^2$  is the sample average
   approximation of  $\sigma_\ell^2$ , see Section 3.1
3   Compute  $\hat{\lambda}_\ell = \text{argmin} \sigma_{\ell, N'_\ell}^2(\lambda)$  using Newton–Raphson’s algorithm.
4   Independently of  $\sigma_{\ell, N'_\ell}^2$ , sample the level  $\ell$  of (1.5) using  $\hat{\lambda}_\ell$ .
5 end
6 Sum all the levels to obtain

```

$$\begin{aligned}Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L) &= \frac{1}{N_0} \sum_{k=1}^{N_0} \psi(\tilde{X}_{T,0,k}^{m^0}(\hat{\lambda}_0))\mathcal{E}^-(\tilde{W}_{0,k}, \hat{\lambda}_0) \\ &\quad + \sum_{\ell=1}^L \frac{1}{N_\ell} \sum_{k=1}^{N_\ell} \left(\psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\hat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\hat{\lambda}_\ell))\right)\mathcal{E}^-(\tilde{W}_{\ell,k}, \hat{\lambda}_\ell).\end{aligned}$$

**Algorithm 1.1:** Multilevel Importance Sampling (MLIS)

In Section 2, we investigate the standard Euler Monte Carlo method coupled with importance sampling. The importance sampling framework with MLMC is studied in Section 3. We

prove that  $Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L)$  satisfies a strong law of large numbers and a central limit theorem. Our MLIS estimator achieves the smallest possible variance within the family of MLMC estimators approximating  $\mathbb{E}[\psi(X_T)]$  using the class of processes  $(X(\lambda))_{\lambda \in \mathbb{R}^q}$ . Note that this is also the limiting variance obtained in [5] for the MLMC estimator built on (1.4) with the best possible parameter  $\lambda \in \mathbb{R}^q$ . The main difficulty in proving these results is the uniform control of the triangular arrays involved in the adaptive multilevel estimator. To overcome this issue, we prove in Section 4 new limit theorems for doubly indexed sequences of random variables in a general setting (see Propositions 4.1 and 4.3). In section 5, we illustrate the efficiency of MLIS on challenging problems coming from quantitative finance and show that it outperforms the standard MLMC estimator.

## 2 Importance sampling with Euler Monte Carlo

### 2.1 Notation and general assumptions

- For a vector  $x \in \mathbb{R}^q$ ,  $|x|$  denotes the Euclidean norm of  $x$ .
- The superscript  $*$  denotes the transpose operator.
- For a matrix  $A \in \mathcal{M}_{d \times q}$ ,  $|A|$  denotes the Frobenius norm of  $A$  defined by  $\sqrt{\text{Tr}(A^*A)}$ , which corresponds to the Euclidean norm on  $\mathbb{R}^{d \times q}$ .
- For  $q \in \mathbb{N}^*$ ,  $I_q$  denotes the identity matrix with size  $q \times q$ .
- For  $\alpha > 0$ , we define the set of functions

$$\mathcal{H}_\alpha = \left\{ \psi : \mathbb{R}^d \rightarrow \mathbb{R} \text{ s.t. } \exists c > 0, \beta \geq 1, \forall x \in \mathbb{R}^d, |\psi(x)| \leq c(1 + |x|^\beta) \right. \\ \left. \text{and } \forall x, y \in \mathbb{R}^d, |\psi(x) - \psi(y)| \leq c(1 + (|x|^\beta \wedge |y|^\beta))|x - y|^\alpha \right\} \quad (2.1)$$

- For a sequence of random variables  $(X_n)_n$ , “ $X_n \implies X$ ” means that  $(X_n)_n$  converges in distribution to  $X$ .

Here, we gather several standard assumptions required to ensure the convergence of the Euler scheme.

(H-1) *i.* The functions  $b$  and  $\sigma$  are Lipschitz

$$\forall x, y \in \mathbb{R}^d, |b(x) - b(y)| + |\sigma(x) - \sigma(y)| \leq C_{b,\sigma}|x - y|, \quad (\mathcal{H}_{b,\sigma})$$

for some real number  $C_{b,\sigma} > 0$ .

*ii.*  $\forall p \geq 1$ ,  $X, X^n \in L^p$  and there exists  $K_p(T) > 0$  s.t.

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} |X_t - X_t^n|^p \right] \leq \frac{K_p(T)}{n^{p/2}}.$$

*iii.* There exist  $\gamma \in [1/2, 1]$  and  $C_\psi(T, \gamma) > 0$  s.t.

$$n^\gamma (\mathbb{E}\psi(X_T^n) - \mathbb{E}\psi(X_T)) \rightarrow C_\psi(T, \gamma). \quad (\mathcal{H}_\gamma)$$

(H-2) The function  $\psi$  satisfies

$$\mathbb{P}(\psi(X_T) \neq 0) > 0 \quad \text{and} \quad \forall \theta \in \mathbb{R}^q, \mathbb{E} \left[ \psi(X_T)^2 e^{-\theta \cdot W_T} \right] < \infty. \quad (2.2)$$

## 2.2 General framework

In this section, we investigate the case of a Euler Monte Carlo. We consider the importance sampling representation of  $\mathbb{E}[\psi(X_T)]$  given by

$$\mathbb{E}[\psi(X_T(\theta))\mathcal{E}^-(W, \theta)].$$

The optimal value for  $\theta$  is given by

$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^q} \quad \text{with} \quad v(\theta) \triangleq \mathbb{E}[(\psi(X_T(\theta))\mathcal{E}^-(W, \theta))^2].$$

By using (1.2), we can rewrite  $v$  as

$$v(\theta) = \mathbb{E}[\psi(X_T)^2\mathcal{E}^+(W, \theta)] \quad \text{with} \quad \mathcal{E}^+(W, \theta) \triangleq e^{-W_T \cdot \theta + \frac{|\theta|^2}{2}}.$$

From a practical point of view, the quantity  $v(\theta)$  is not explicit so we use the Euler scheme to discretize  $X(\theta)$  and approximate  $\theta^*$  by

$$\theta_n \triangleq \operatorname{argmin}_{\theta \in \mathbb{R}^q} v_n(\theta) \quad \text{with} \quad v_n(\theta) \triangleq \mathbb{E}[\psi(X_T^n)^2\mathcal{E}(W, \theta)]. \quad (2.3)$$

Since the expectation is usually not tractable, we replace it by its sample average approximation and define

$$\theta_{n,N} \triangleq \operatorname{argmin}_{\theta \in \mathbb{R}^q} v_{n,N}(\theta) \quad \text{with} \quad v_{n,N}(\theta) \triangleq \frac{1}{N} \sum_{i=1}^N (\psi(X_{T,i}^n)^2\mathcal{E}(W_i, \theta)), \quad (2.4)$$

where  $(X_{T,i}^n, W_{T,i})_{1 \leq i \leq N}$  are i.i.d. samples with the law of  $(X_T^n, W_T)$ . The existence and uniqueness of  $\theta^*$ ,  $\theta_n$  and  $\theta_{n,N}$  are ensured by the following lemma whose proof can easily be adapted from [23, Lemma 1.1].

**Lemma 2.1.** *Under Condition  $(\mathcal{H}-2)$ , the functions  $v$ ,  $v_n$  and  $v_{n,N}$  are infinitely continuously differentiable for all  $n, N$  and the derivatives are obtained by exchanging expectation and differentiation. Moreover, the functions  $v$  and  $v_n$  are strongly convex and so is  $v_{n,N}$  for any  $N$  such that  $v_{n,N}$  is not identically zero.*

## 2.3 Convergence of the optimal importance sampling parameter

**Theorem 2.2.** *Suppose  $\sigma$  and  $b$  satisfy  $(\mathcal{H}_{b,\sigma})$ . Let  $\psi$  satisfy Condition  $(\mathcal{H}-2)$  and belong to  $\mathcal{H}_\alpha$  for some  $\alpha > 0$ . Then,  $\theta_n \rightarrow \theta^*$  a.s. when  $n \rightarrow +\infty$ .*

By Hölder's inequality, for any function  $\psi \in \mathcal{H}_\alpha$ ,  $(\mathcal{H}-2)$  implies that  $\sup_n \mathbb{E}[\psi(X_T^n)^2 e^{-\theta \cdot W_T}] < +\infty$ . Hence, the proof of the theorem ensues from [5, Theorem 2.2].

In the following, we let  $N$  depend on  $n$  so that  $N \triangleq N_n$  tends to infinity with  $n$ .

**Proposition 2.3.** *Assume that Assumption  $(\mathcal{H}_{b,\sigma})$  holds and that  $\psi \in \mathcal{H}_\alpha$  for some  $\alpha > 0$ . Then, for all  $K > 0$ , a.s. when  $n \rightarrow \infty$*

$$\sup_{|\theta| \leq K} |v_{n,N_n}(\theta) - v(\theta)| \rightarrow 0; \quad \sup_{|\theta| \leq K} |\nabla v_{n,N_n}(\theta) - \nabla v(\theta)| \rightarrow 0.$$

*Proof.* The proof of the two results are very similar, we omit the second one and concentrate on the uniform convergence for  $v_{n,N_n}$ . To do so, we will apply Proposition 4.3. Now, we check Assumptions  $(\mathcal{H}-4)$ ,  $(\mathcal{H}-5)$ ,  $(\mathcal{H}-6)$ . At first, note that under Assumption  $(\mathcal{H}_{b,\sigma})$ , we have the almost sure convergence of  $X_T^n$  towards  $X_T$ . As  $\psi \in \mathcal{H}_\alpha$ , it follows from Property  $(\mathcal{H}-1)$ -ii that for all  $a > 1$ ,  $\sup_{n \in \mathbb{N}} \mathbb{E} \left[ \left| \psi(X_T^n)^2 e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right|^a \right] < \infty$ . Note that for every fixed  $n$ , the sequence  $\left( \psi(X_{T,i}^n)^2 e^{-\theta \cdot W_{T,i} + \frac{1}{2} |\theta|^2 T} \right)_i$  is i.i.d. Then, we deduce that for all  $m \in \mathbb{N}^*$

$$\lim_{n \rightarrow \infty} \mathbb{E} [v_{n,m}(\theta)] = \mathbb{E} \left[ \psi(X_T)^2 e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right].$$

This yields  $(\mathcal{H}-4)$ . Let  $K > 0$ . As  $\psi \in \mathcal{H}_\alpha$  we obtain using the Cauchy Schwarz inequality and Property  $(\mathcal{H}-1)$ -ii that

$$\sup_n \sup_m m \text{Var} \left( \sup_{|\theta| \leq K} v_{n,m}(\theta) \right) \leq \sup_n \mathbb{E}^{1/2} [\psi(X_T^n)^8] \mathbb{E}^{1/2} \left[ \sup_{|\theta| \leq K} e^{-4\theta \cdot W_T + 2|\theta|^2 T} \right] < \infty.$$

Using the same arguments, we also get

$$\sup_n \sup_m \text{Var} \left( \psi(X_{T,m}^n)^2 \sup_{|\theta| \leq K} e^{-\theta \cdot W_{T,m} + \frac{1}{2} |\theta|^2 T} \right) < \infty.$$

This yields  $(\mathcal{H}-5)$ . Concerning the last assumption, if we fix  $\delta > 0$ ,  $\theta \in \mathbb{R}^d$  and set  $B(\theta, \delta)$  — the open ball with center  $\theta$  and radius  $\delta$  — then we have by Cauchy Schwarz inequality

$$\begin{aligned} \sup_n \mathbb{E} \left[ \psi(X_T^n)^2 \sup_{\theta' \in B(\theta, \delta)} \left| e^{-\theta' \cdot W_T + \frac{1}{2} |\theta'|^2 T} - e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right| \right]^2 &\leq \\ \sup_n \mathbb{E} [\psi(X_T^n)^4] \mathbb{E} \left[ \sup_{\theta' \in B(\theta, \delta)} \left| e^{-\theta' \cdot W_T + \frac{1}{2} |\theta'|^2 T} - e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right|^2 \right]. & \end{aligned}$$

Using the elementary algebraic inequality  $|e^x - e^y| \leq |x - y| (e^x + e^y)$ , we easily deduce that the quantity  $\mathbb{E} \left[ \sup_{\theta' \in B(\theta, \delta)} \left| e^{-\theta' \cdot W_T + \frac{1}{2} |\theta'|^2 T} - e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right|^2 \right]$  can be made arbitrarily small. Finally, Assumption  $(\mathcal{H}-6)$  is satisfied using Remark 4.4.  $\square$

**Theorem 2.4.** *Assume that Assumption  $(\mathcal{H}_{b,\sigma})$  holds and that  $\psi \in \mathcal{H}_\alpha$  for some  $\alpha > 0$ . Then,  $\theta_{n,N_n} \rightarrow \theta^*$  a.s. and  $\sqrt{N_n}(\theta_{n,N_n} - \theta^*) \Rightarrow N(0, \Gamma)$  when  $n \rightarrow \infty$  with*

$$\Gamma = [\nabla^2 v(\theta^*)]^{-1} \text{Var} \left[ (T\theta^* - W_T) \psi(X_T)^2 e^{-\theta^* \cdot W_T + \frac{1}{2} |\theta^*|^2 T} \right] [\nabla^2 v(\theta^*)]^{-1}.$$

*Proof.* We already know from Proposition 2.3 that a.s.  $v_{n,N_n}$  converges locally uniformly to  $v$ . Let  $\varepsilon > 0$ . By the strict convexity of  $v$ ,  $\delta \triangleq \inf_{|\theta - \theta^*| \geq \varepsilon} v(\theta) - v(\theta^*) > 0$ . The local uniform convergence of  $v_{n,N_n}$  to  $v$  ensures that

$$\exists n_\delta > 0, \forall n \geq n_\delta, \forall \theta \in \mathbb{R}^q \text{ s.t. } |\theta - \theta^*| \leq \varepsilon, |v_{n,N_n}(\theta) - v(\theta)| \leq \frac{\delta}{3}. \quad (2.5)$$



For  $n \geq n_\delta$  and  $\theta$  such that  $|\theta - \theta^*| \geq \varepsilon$ , we can deduce from the convexity of  $v_{n,N_n}$  that

$$\begin{aligned} v_{n,N_n}(\theta) - v_{n,N_n}(\theta^*) &\geq \frac{|\theta - \theta^*|}{\varepsilon} \left[ v_{n,N_n} \left( \theta^* + \varepsilon \frac{\theta - \theta^*}{|\theta - \theta^*|} \right) - v_{n,N_n}(\theta^*) \right] \\ &\geq \frac{|\theta - \theta^*|}{\varepsilon} \left[ v \left( \theta^* + \varepsilon \frac{\theta - \theta^*}{|\theta - \theta^*|} \right) - v(\theta^*) - \frac{2\delta}{3} \right] \geq \frac{\delta}{3} \end{aligned}$$

where the last two inequalities come from (2.5). If we apply this inequality for  $\theta = \theta_{n,N_n}$ , we obtain a contradiction since  $v_{n,N_n}(\theta_{n,N_n}) - v_{n,N_n}(\theta^*) \leq 0$ . Hence, we deduce that for all  $n \geq n_\delta$ ,  $|\theta_{n,N_n} - \theta^*| < \varepsilon$ . Therefore,  $\theta_{n,N_n}$  converges a.s. to  $\theta^*$ . If we combine this result with the local uniform convergence of  $v_{n,N_n}$  to the continuous function  $v$ , we deduce that  $v_{n,N_n}(\theta_{n,N_n})$  converges a.s. to  $v(\theta^*)$ .

Moreover, we get by Equation (3.9) that for all  $K > 0$

$$\begin{aligned} &\sup_{|\theta| \leq K} \left| \partial_{\theta^{(j)}} \psi(X_T)^2 e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right| \\ &\leq e^{K^2 T/2} \psi(X_T)^2 \left( K + (e^{KW_t^{(j)}} + e^{-KW_t^{(j)}}) \right) \prod_{i=1}^q (e^{KW_t^{(i)}} + e^{-KW_t^{(i)}}). \end{aligned}$$

The r.h.s is integrable by Condition  $(\mathcal{H}-2)$ . Hence,  $\mathbb{E} \left[ \sup_{|\theta| \leq K} \left| \nabla_{\theta} \psi(X_T)^2 e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right| \right] < +\infty$ . Similarly, one can prove that  $\mathbb{E} \left[ \sup_{|\theta| \leq K} \left| \nabla_{\theta}^2 \psi(X_T)^2 e^{-\theta \cdot W_T + \frac{1}{2} |\theta|^2 T} \right| \right] < +\infty$ . Then, to prove the central limit theorem governing the convergence of  $\theta_{n,N_n}$  to  $\theta^*$ , we reproduce the proof of [29, Theorem A2, pp. 74], which is mainly based on the a.s. local uniform convergence of  $\nabla v_{n,N_n}$  and on its asymptotic normality ensuing from Theorem A.1.  $\square$

## 2.4 A second stage Monte Carlo approach

In this section, we aim at building adaptive Monte Carlo estimators in the setting of discretized diffusion processes following the spirit of [23]. Our setting differs mainly because we want to let both the number of time steps and the number of samples go to infinity. Asymptotic results rely on a uniform controls of the triangular arrays involved in the adaptive importance sampling Monte Carlo estimator. The technical results from Section 4 will be tremendously useful to provide such controls.

Using the estimators of  $\theta^*$  studied in the previous section, we define a Monte Carlo estimator of  $\mathbb{E}[\psi(X_T)]$  based on Equation (1.2). We introduce the  $\sigma$ -algebra  $\mathcal{G}$  generated by the samples  $(W_i)_{i \geq 1}$  used to compute  $\theta_n$  and  $\theta_{n,N_n}$ .

Let  $(\tilde{W}_i)_i$  be i.i.d. samples according to the law of  $W$  but independent of  $\mathcal{G}$ . Conditionally on  $\mathcal{G}$ , we introduce i.i.d. samples  $(\tilde{X}_i(\theta_{n,N_n}))_i$  following the law of  $X(\theta_{n,N_n})$  such that for each  $i$ ,  $\tilde{X}_i(\theta_{n,N_n})$  is the solution of the SDE driven by  $\tilde{W}_i$ . We introduce  $(\tilde{\mathcal{G}}_k)_{k > 0}$  the filtration defined by  $\tilde{\mathcal{G}}_k = \sigma(\tilde{W}_i, 1 \leq i \leq k)$  and  $\mathcal{G}_k^\# = \mathcal{G} \vee \tilde{\mathcal{G}}_k$ . For each  $i > 0$ , we also consider  $\tilde{X}_i^n(\theta_{n,N_n})$  defined as the Euler discretization of  $\tilde{X}_i(\theta_{n,N_n})$ . Based on these new sets of samples, we define

$$M_{n,N_n} = \frac{1}{N_n} \sum_{i=1}^{N_n} g(\theta_{n,N_n}, \tilde{X}_{T,i}^n(\theta_{n,N_n}), \tilde{W}_{T,i}),$$

where the function  $g : \mathbb{R}^q \times \mathbb{R}^d \times \mathbb{R}^q \rightarrow \mathbb{R}$  is defined by

$$g(\theta, x, y) \triangleq \psi(x) e^{-\theta \cdot y - \frac{1}{2} |\theta|^2 T}. \quad (2.6)$$

For the clearness of the coming proofs, it is convenient to introduce the following notation

$$M_{n,N_n}(\theta) = \frac{1}{N_n} \sum_{i=1}^{N_n} g(\theta, \tilde{X}_{T,i}^n(\theta), \tilde{W}_{T,i}).$$

Note that  $M_{n,N_n} = M_{n,N_n}(\theta_{n,N_n})$ .

**Theorem 2.5.** *Assume that Assumption  $(\mathcal{H}_{b,\sigma})$  holds and that  $\psi \in \mathcal{H}_\alpha$  for some  $\alpha > 0$ . Then,  $M_{n,N_n} \rightarrow \mathbb{E}[\psi(X_T)]$  a.s. when  $n \rightarrow +\infty$ .*

*Proof.* Using the conditional independence of the samples  $(\tilde{X}_i^n(\theta_{n,N_n}), \tilde{W}_i)$ , we have

$$\mathbb{E}[g(\theta_{n,N_n}, \tilde{X}_{T,i}^n(\theta_{n,N_n}), \tilde{W}_{T,i}) | \mathcal{G}] = \mathbb{E}[\psi(X_T^n)] \triangleq e_n \quad \text{for all } i > 0.$$

Let  $\mathcal{V} \subset \mathbb{R}^q$  be a compact neighbourhood of  $\theta^*$ . We define the sequence

$$Y_{i,n} = \left( g(\theta_{n,N_n}, \tilde{X}_{T,i}^n(\theta_{n,N_n}), \tilde{W}_{T,i}) - e_n \right) 1_{\{\theta_{n,N_n} \in \mathcal{V}\}}$$

and its empirical average  $\bar{Y}_{m,n} = \frac{1}{m} \sum_{i=1}^m Y_{i,n}$  for all  $m > 0$ . It is obvious that  $\mathbb{E}[Y_{i,n}] = 0$  and using the conditional independence  $\mathbb{E}[|\bar{Y}_{m,n}|^2] = \frac{1}{m} \mathbb{E}[|Y_{1,n}|^2]$ .

$$\begin{aligned} \mathbb{E}[|Y_{1,n}|^2] &\leq \mathbb{E} \left[ \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_{T,i}^n(\theta_{n,N_n}), \tilde{W}_{T,i}) - e_n|^2 \middle| \mathcal{G} \right] 1_{\{\theta_{n,N_n} \in \mathcal{V}\}} \right] \\ &\leq \mathbb{E} \left[ v_n(\theta_{n,N_n}) 1_{\{\theta_{n,N_n} \in \mathcal{V}\}} \right] \leq \sup_{\theta \in \mathcal{V}} v_n(\theta). \end{aligned}$$

We know that  $v_n$  is convex and converges point-wise to  $v$ , which is also convex and continuous. Hence,  $v_n$  converges locally uniformly to  $v$ , which implies that for all compact sets  $K \subset \mathbb{R}^q$ ,  $\lim_{n \rightarrow +\infty} \sup_{\theta \in K} v_n(\theta) = \sup_{\theta \in K} v(\theta)$ . Hence,  $\sup_n \sup_{\theta \in \mathcal{V}} v_n(\theta) < +\infty$ . Applying Proposition 4.1 proves that  $\bar{Y}_{N_n,n} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$ . As  $\theta_{n,N_n}$  converges a.s. to  $\theta^* \in K$ , this also implies that  $\lim_{n \rightarrow +\infty} M_{n,N_n} = \mathbb{E}[\psi(X_T)]$  a.s.  $\square$

**Theorem 2.6.** *Under the assumptions of Theorem 2.5 and if Condition  $(\mathcal{H}_\gamma)$  holds, we have*

$$\sqrt{N_n}(M_{n,N_n} - \mathbb{E}[\psi(X_T)]) \Longrightarrow \mathcal{N}(C_\psi(T, \alpha), \sigma^2) \quad \text{when } n \rightarrow +\infty.$$

where  $\sigma^2 = \mathbb{E} \left[ \psi(X_T)^2 e^{-\theta^* \cdot W_T + \frac{1}{2} |\theta^*|^2 T} \right] - \mathbb{E}[\psi(X_T)]^2$ .

**Remark 2.7.** *Assume the number of time steps used in the Euler scheme is fixed to  $n = 1$  and consider the estimator  $M_{1,N}(\theta_{1,N})$ . Then, we know from [2, Theorem 3.4] that, when  $N \rightarrow \infty$ ,*

$$\begin{aligned} M_{1,N}(\theta_{1,N}) &\longrightarrow \mathbb{E}[g(\theta_1, X_T^1(\theta_1), W_T)] \quad a.s. \\ \sqrt{N}(M_{1,N}(\theta_{1,N}) - \mathbb{E}[g(\theta_1, X_T^1(\theta_1), W_T)]) &\Longrightarrow \mathcal{N}(0, \sigma_1^2) \end{aligned}$$

with  $\sigma_1^2 = \mathbb{E} \left[ \psi(X_T^1)^2 e^{-\theta_1 \cdot W_T + \frac{1}{2} |\theta_1|^2 T} \right] - \mathbb{E}[\psi(X_T^1)]^2$ .

*Proof.* We can write the left hand side of the convergence result by introducing  $M_{n,N_n}(\theta^*)$

$$\sqrt{N_n}(M_{n,N_n} - \mathbb{E}[\psi(X_T)]) = \sqrt{N_n}(M_{n,N_n}(\theta_{n,N_n}) - M_n(\theta^*)) + \sqrt{N_n}(M_{n,N_n}(\theta^*) - \mathbb{E}[\psi(X_T)])$$

The convergence of the last term on the r.h.s  $\sqrt{N_n}(M_{n,N_n}(\theta^*) - \mathbb{E}[\psi(X_T)])$  is governed by the central limit theorem for Euler Monte Carlo, which yields the announced limit (see [12]). It remains to prove that  $\sqrt{N_n}(M_{n,N_n}(\theta_{n,N_n}) - M_{n,N_n}(\theta^*))$  converges to zero in probability.

Let  $\varepsilon > 0$  and  $\alpha < \frac{1}{2}$ ,

$$\begin{aligned} & \mathbb{P} \left( \sqrt{N_n} |M_{n,N_n}(\theta_{n,N_n}) - M_{n,N_n}(\theta^*)| > \varepsilon \right) \\ &= \mathbb{P} \left( \sqrt{N_n} |M_{n,N_n}(\theta_{n,N_n}) - M_{n,N_n}(\theta^*)| > \varepsilon ; N_n^\alpha |\theta_{n,N_n} - \theta^*| > 1 \right) \\ & \quad + \mathbb{P} \left( \sqrt{N_n} |M_{n,N_n}(\theta_{n,N_n}) - M_{n,N_n}(\theta^*)| > \varepsilon ; N_n^\alpha |\theta_{n,N_n} - \theta^*| \leq 1 \right) \\ &= \mathbb{P} ( N_n^\alpha |\theta_{n,N_n} - \theta^*| > 1 ) \\ & \quad + \mathbb{P} \left( \sqrt{N_n} |M_{n,N_n}(\theta_{n,N_n}) - M_{n,N_n}(\theta^*)| 1_{\{N_n^\alpha |\theta_{n,N_n} - \theta^*| \leq 1\}} > \varepsilon \right). \end{aligned}$$

By Theorem 2.4,  $\mathbb{P} ( N_n^\alpha |\theta_{n,N_n} - \theta^*| > 1 )$  tends to zero when  $n$  goes to infinity. Let  $K > 0$  s.t. for all  $n$  large enough  $\{\theta \in \mathbb{R}^q : |\theta - \theta^*| \leq N_n^{-\alpha}\} \subset B(0, K)$ . We can bound the second term on the r.h.s. by using Markov's inequality

$$\begin{aligned} & \mathbb{P} \left( \sqrt{N_n} |M_{n,N_n}(\theta_{n,N_n}) - M_{n,N_n}(\theta^*)| 1_{\{N_n^\alpha |\theta_{n,N_n} - \theta^*| \leq 1\}} > \varepsilon \right) \\ & \leq \frac{N_n}{\varepsilon^2} \mathbb{E} \left[ |M_{n,N_n}(\theta_{n,N_n}) - M_{n,N_n}(\theta^*)|^2 1_{\{\theta_{n,N_n} \in B(0, K)\}} \right] \\ & \leq \frac{1}{\varepsilon^2} \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_T^n(\theta_{n,N_n}), \tilde{W}_T) - g(\theta^*, \tilde{X}_T^n(\theta^*), \tilde{W}_T)|^2 1_{\{\theta_{n,N_n} \in B(0, K)\}} \right] \\ & \leq \frac{1}{\varepsilon^2} \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_T^n(\theta_{n,N_n}), \tilde{W}_T) - g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T)|^2 1_{\{\theta_{n,N_n} \in B(0, K)\}} \right] \\ & \quad + \frac{1}{\varepsilon^2} \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T) - g(\theta^*, \tilde{X}_T^n(\theta^*), \tilde{W}_T)|^2 1_{\{\theta_{n,N_n} \in B(0, K)\}} \right]. \end{aligned}$$

We treat each of the two terms separately.

► **First term**

From the independence between  $\theta_{n,N_n}$  and  $\tilde{W}$ , we can write

$$\begin{aligned} & \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_T^n(\theta_{n,N_n}), \tilde{W}_T) - g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T)|^2 1_{\{\theta_{n,N_n} \in B(0, K)\}} \right] \\ &= \mathbb{E} \left[ |\psi(X_T^n) - \psi(X_T)|^2 \exp(-\theta_{n,N_n} \cdot \tilde{W}_T + \frac{1}{2} |\theta_{n,N_n}|^2 T) 1_{\{\theta_{n,N_n} \in B(0, K)\}} \right] \\ & \leq \mathbb{E} \left[ |\psi(X_T^n) - \psi(X_T)|^{2(1+\eta)} \right]^{\frac{1}{1+\eta}} e^{\frac{1+2\eta}{2\eta} K^2 T}, \quad \text{for some } \eta > 0. \end{aligned}$$

Relying on the uniform integrability ensured by property  $(\mathcal{H}-1)$ -ii and since  $\psi \in \mathcal{H}_\alpha$ , we can let  $n$  go to infinity inside the expectation to obtain that

$$\lim_{n \rightarrow +\infty} \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_T^n(\theta_{n,N_n}), \tilde{W}_T) - g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T)|^2 1_{\{\theta_{n,N_n} \in B(0, K)\}} \right] = 0.$$

► **Second term**

Since the function  $g$  is continuous w.r.t its first two parameters and  $X_T^\theta$  is continuous w.r.t the parameter  $\theta$ ,  $\lim_{n \rightarrow +\infty} g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T) - g(\theta^*, \tilde{X}_T^n(\theta^*), \tilde{W}_T) = 0$  a.s. To conclude the proof, we need to show that the family of r.v.

$$\left( |g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T) - g(\theta^*, \tilde{X}_T^n(\theta^*), \tilde{W}_T)|^2 \mathbf{1}_{\{\theta_{n,N_n} \in B(0,K)\}} \right)_n$$

is uniformly integrable.

First, for any  $\theta \in \mathbb{R}^q$  and  $2(1+\eta) > a > 2$

$$\begin{aligned} \mathbb{E} \left[ |g(\theta, \tilde{X}_T(\theta), \tilde{W}_T)|^a \right] &= \mathbb{E} \left[ |\psi(\tilde{X}_T)|^a e^{-(a-1)\theta \cdot \tilde{W}_T + \frac{(a-1)|\theta|^2 T}{2}} \right] \\ &\leq \mathbb{E} \left[ |\psi(\tilde{X}_T)|^{2(1+\eta)} \right]^{\frac{2(1+\eta)}{a}} e^{C|\theta|^2} \end{aligned} \quad (2.7)$$

where  $C$  is a constant only depending on  $a$  and  $T$ . This yields that for some  $\delta > 0$  and some constant  $C > 0$  independent of  $\theta$ ,  $\mathbb{E} \left[ |g(\theta, \tilde{X}_T(\theta), \tilde{W}_T)|^{2+\delta} \right] < C e^{C|\theta|^2}$ . Then, we get

$$\begin{aligned} &\sup_n \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T)|^{2+\delta} \mathbf{1}_{\{\theta_{n,N_n} \in B(0,K)\}} \right] \\ &= \sup_n \mathbb{E} \left[ \mathbb{E} \left[ |g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T)|^{2+\delta} | \theta_{n,N_n} \right] \mathbf{1}_{\{\theta_{n,N_n} \in B(0,K)\}} \right] \\ &\leq \sup_n C \mathbb{E} \left[ e^{C|\theta_{n,N_n}|^2} \mathbf{1}_{\{\theta_{n,N_n} \in B(0,K)\}} \right] \leq C e^{CK}. \end{aligned}$$

We can similarly prove that

$$\sup_n \mathbb{E} \left[ |g(\theta^*, \tilde{X}_T^n(\theta^*), \tilde{W}_T)|^{2+\delta} \right] \leq \sup_n \mathbb{E} \left[ |\psi(X_T^n)|^{2(1+\eta)} \right]^{\frac{2(1+\eta)}{2+\delta}} e^{C|\theta^*|^2}.$$

This prove that the family of r.v.

$$\left( |g(\theta_{n,N_n}, \tilde{X}_T(\theta_{n,N_n}), \tilde{W}_T) - g(\theta^*, \tilde{X}_T^n(\theta^*), \tilde{W}_T)|^2 \mathbf{1}_{\{\theta_{n,N_n} \in B(0,K)\}} \right)_n$$

is uniformly integrable, which ends the proof.  $\square$

### 3 Multilevel Importance sampling Monte Carlo

In the recent years, many works showed that MLMC supersedes Monte Carlo when combined with discretization schemes. Then, it has become natural to investigate how this new approach could be coupled with existing variance reduction techniques and in particular with importance sampling. In this section, we study the mathematical properties of our importance sampling MLMC estimator  $Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L)$ . First, we start by proving the existence and uniqueness of  $\hat{\lambda}_0, \dots, \hat{\lambda}_L$  in Section 3.2 and then we prove a strong law of large numbers and a central limit theorem for  $Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L)$  in Section 3.3.

### 3.1 General framework

Our multilevel importance sampling estimator writes

$$Q_L(\lambda_0, \dots, \lambda_L) = \frac{1}{N_0} \sum_{k=1}^{N_0} \psi(\tilde{X}_{T,0,k}^{m^0}(\lambda_0)) \mathcal{E}^-(\tilde{W}_{0,k}, \lambda_0) \\ + \sum_{\ell=1}^L \frac{1}{N_\ell} \sum_{k=1}^{N_\ell} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\lambda_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\lambda_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \lambda_\ell). \quad (3.1)$$

For any fixed  $\ell \in \{1, \dots, L\}$ , the random variables  $(\tilde{W}_{\ell,k})_{1 \leq k \leq N_\ell}$  are independent and are distributed according to the Brownian law. We assume that for  $\ell, \ell' \in \{1, \dots, L\}$ , with  $\ell \neq \ell'$ , the blocks  $(\tilde{W}_{\ell,k})_{1 \leq k \leq N_\ell}$  and  $(\tilde{W}_{\ell',k})_{1 \leq k \leq N_{\ell'}}$  are independent. For any fixed  $\ell \in \{1, \dots, L\}$  and  $k \in \{1, \dots, N_\ell\}$ , the variables  $\tilde{X}_{T,\ell,k}^{m^\ell}(\lambda_\ell)$  (resp.  $\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\lambda_\ell)$ ) are the terminal values of the Euler schemes of  $X(\lambda_\ell)$  with  $m^\ell$  (resp.  $m^{\ell-1}$ ) time steps built using the same Brownian path  $\tilde{W}_{\ell,k}$ . The key of the multilevel approach is to use the same Brownian path to compute  $\tilde{X}_{T,\ell,k}^{m^\ell}(\lambda_\ell)$  and  $\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\lambda_\ell)$ . The blocks of random variables used in two different levels are independent. From these assumptions, one can compute the variance of the multilevel estimator given by

$$\text{Var}[Q_L] = N_0^{-1} \sigma_0(\lambda_0)^2 + \sum_{\ell=1}^L N_\ell^{-1} \frac{(m-1)T}{m^\ell} \sigma_\ell^2(\lambda_\ell)$$

where

$$\sigma_0^2(\lambda_0) \triangleq \text{Var}[\psi(X_T^{m^0}(\lambda_0)) \mathcal{E}^-(W, \lambda_0)] \\ \sigma_\ell^2(\lambda_\ell) \triangleq \frac{m^\ell}{(m-1)T} \text{Var} \left[ \left\{ \psi(X_T^{m^\ell}(\lambda_\ell)) - \psi(X_T^{m^{\ell-1}}(\lambda_\ell)) \right\} \mathcal{E}^-(W, \lambda_\ell) \right].$$

By applying (1.2), the variances of each level  $\ell \geq 0$  can be written  $\sigma_\ell^2(\lambda_\ell) = v_\ell(\lambda_\ell) - \Xi_\ell^2$  with

$$v_0(\lambda_0) \triangleq \mathbb{E} \left[ \psi(X_T^{m^0})^2 \mathcal{E}^+(W, \lambda_0) \right], \quad \Xi_0 \triangleq \mathbb{E} \left[ \psi(X_T^{m^0}) \right] \quad (3.2)$$

$$v_\ell(\lambda_\ell) \triangleq \frac{m^\ell}{(m-1)T} \mathbb{E} \left[ \left| \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right|^2 \mathcal{E}^+(W, \lambda_\ell) \right], \quad (3.3)$$

$$\Xi_\ell \triangleq \sqrt{\frac{m^\ell}{(m-1)T}} \mathbb{E} \left[ \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right] \quad (3.4)$$

and  $\mathcal{E}^+(W, \lambda) \triangleq e^{-\lambda \cdot W_T + \frac{1}{2} |\lambda|^2 T}$ . Hence, the global variance is given by

$$\text{Var}[Q_L] = N_0^{-1} (v_0(\lambda_0) - \Xi_0^2) + \sum_{\ell=1}^L N_\ell^{-1} \frac{(m-1)T}{m^\ell} (v_\ell(\lambda_\ell) - \Xi_\ell^2).$$

To actually minimize the functions  $\lambda \mapsto v_\ell^2(\lambda)$ , we consider the sample average approximation of  $v_\ell$  with  $N'_\ell$  samples

$$\begin{aligned} v_{0,N'_0}(\lambda_0) &\triangleq \frac{1}{N'_0} \sum_{k=1}^{N'_0} \psi(X_{T,0,k}^{m^0})^2 \mathcal{E}^+(W_{0,k}, \lambda_0), \\ v_{\ell,N'_\ell}(\lambda_\ell) &\triangleq \frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 \mathcal{E}^+(W_{\ell,k}, \lambda_\ell). \end{aligned}$$

### 3.2 Convergence of the importance sampling parameters

From Lemma 2.1, we deduce that  $v_{\ell,N'_\ell}$  has a unique minimum

$$\widehat{\lambda}_\ell = \arg \min_{\lambda \in \mathbb{R}^q} v_{\ell,N'_\ell}(\lambda).$$

**Theorem 3.1.** *Assume  $b$  and  $\sigma$  are  $C^1$  with bounded derivatives,  $\psi \in \mathcal{H}_\alpha$  for some  $\alpha \geq 1$ ,  $\psi$  is  $C^1$  and  $\nabla \psi$  has polynomial growth. Then, the sequence of random functions  $(v_{\ell,N'_\ell} : \lambda \in \mathbb{R}^q \rightarrow v_{\ell,N'_\ell}(\lambda))_\ell$  converges a.s. locally uniformly to the strongly convex function  $v : \mathbb{R}^q \rightarrow \mathbb{R}$  defined by*

$$v(\lambda) \triangleq \mathbb{E} \left[ (\nabla \psi(X_T) \cdot U_T)^2 \mathcal{E}^+(W, \lambda) \right] \quad (3.5)$$

with

$$dU_t = \nabla b(X_t) U_t dt + \sum_{j=1}^q \nabla \sigma_j(X_t) U_t dW_t^j - \frac{1}{\sqrt{2}} \sum_{ij=1}^q \nabla \sigma_j(X_t) \sigma_i(X_t) d\check{W}_t^{i,j} \quad (3.6)$$

where  $\check{W}$  is a Brownian motion independent of  $W$  with values in  $\mathbb{R}^{q \times q}$ .

Moreover,  $\widehat{\lambda}_\ell$  converges a.s. to  $\lambda^* \triangleq \arg \min_{\lambda} v(\lambda)$ , when  $\ell \rightarrow +\infty$ .

*Proof.* Let us define the doubly indexed sequence

$$Y_{k,\ell}(\lambda) = \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,k}^{m^\ell}) - \psi(X_{T,k}^{m^{\ell-1}}) \right|^2 \mathcal{E}^+(W_k, \lambda).$$

For any fixed  $\ell$ , the sequence  $(Y_{k,\ell}(\lambda))_k$  is i.i.d. so that for any  $k$ ,  $\mathbb{E}[Y_{k,\ell}(\lambda)] = y_\ell(\lambda)$  with

$$y_\ell(\lambda) = \mathbb{E} \left[ \frac{m^\ell}{(m-1)T} \left| \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right|^2 \mathcal{E}^+(W, \lambda) \right].$$

We deduce from Proposition A.4 that the sequence  $(y_\ell)_\ell$  converges pointwise to the continuous function  $\mathbb{E} \left[ (\nabla \psi(X_T) \cdot U_T)^2 \mathcal{E}^+(W, \lambda) \right]$ , thus satisfying Assumption  $(\mathcal{H}\text{-}4)\text{-}i$ . The i.i.d. property of the sequence  $(Y_{k,\ell}(\lambda))_k$  also implies that

$$\mathbb{E} \left[ \sup_{|\lambda| \leq K} \frac{1}{N} \left( \sum_{k=1}^N Y_{k,\ell}(\lambda) \right)^2 \right] \leq \mathbb{E} \left[ \frac{1}{N} \sum_{k=1}^N \sup_{|\lambda| \leq K} Y_{k,\ell}(\lambda)^2 \right] \leq \frac{1}{N} \mathbb{E} \left[ \sup_{|\lambda| \leq K} Y_{1,\ell}(\lambda)^2 \right]. \quad (3.7)$$

$$\mathbb{E} \left[ \sup_{|\lambda| \leq K} Y_{1,\ell}(\lambda)^2 \right]^2 \leq \mathbb{E} \left[ \left( \frac{m^\ell}{(m-1)T} \left| \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right|^2 \right)^4 \right] \mathbb{E} \left[ \sup_{|\lambda| \leq K} \mathcal{E}^+(W, \lambda)^4 \right]. \quad (3.8)$$

Using the following upper bound

$$\sup_{|\lambda| \leq K} e^{-\lambda \cdot W_T + \frac{1}{2} |\lambda|^2 T} \leq e^{\frac{1}{2} K^2 T} \prod_{l=1}^q (e^{KW_T^{(l)}} + e^{-KW_T^{(l)}}), \quad (3.9)$$

$\mathbb{E} \left[ \sup_{|\lambda| \leq K} \mathcal{E}^+(W, \lambda)^4 \right] < +\infty$ . Let us have a closer look at the first term in (3.8). From Condition (2.1), we can write

$$\mathbb{E} \left[ \left( m^\ell \left| \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right|^2 \right)^4 \right] \leq C \mathbb{E} \left[ m^{4\ell} \left| X_T^{m^\ell} - X_T^{m^{\ell-1}} \right|^{8\alpha} \left( 1 + \left| X_T^{m^\ell} \right|^{8\beta} + \left| X_T^{m^{\ell-1}} \right|^{8\beta} \right) \right].$$

By using the strong rate of convergence of the Euler scheme, we notice that for any  $p > 1$ ,

$$\mathbb{E} \left[ m^{4\ell p} \left| X_T^{m^\ell} - X_T^{m^{\ell-1}} \right|^{8\alpha p} \right] \leq m^{4\ell p} C \left( m^{-4\alpha p \ell} + m^{-4\alpha p (\ell-1)} \right) \leq C m^{4\alpha p - 4\ell p (\alpha-1)}.$$

Hence, since  $\alpha \geq 1$ , by using the Cauchy Schwartz inequality we easily check that

$$\sup_{\ell} \mathbb{E} \left[ \left( \frac{m^\ell}{(m-1)T} \left| \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right|^2 \right)^4 \right] < +\infty.$$

By combining all these results into (3.8), we obtain that  $\sup_{\ell} \mathbb{E} \left[ \sup_{|\lambda| \leq K} Y_{1,\ell}^2(\lambda) \right] < +\infty$ . Then, we deduce along with (3.7) that the sequence  $(Y_{k,\ell})_{k,\ell}$  satisfies Assumption  $(\mathcal{H}-5)$  of Proposition 4.3.

Let  $\delta > 0$  and  $\lambda \in \mathbb{R}^d$ .

$$\begin{aligned} & \mathbb{E} \left[ \sup_{|\mu-\lambda| \leq \delta} |Y_{1,\ell}(\lambda) - Y_{1,\ell}(\mu)| \right]^2 \leq \\ & \mathbb{E} \left[ \left( \frac{m^\ell}{(m-1)T} \left| \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right|^2 \right)^2 \right] \mathbb{E} \left[ \sup_{|\mu-\lambda| \leq \delta} |\mathcal{E}^+(W, \lambda) - \mathcal{E}^+(W, \mu)|^2 \right]. \end{aligned}$$

We have just proved that the first expectation on the r.h.s is bounded uniformly in  $\ell$ . Since the exponential weights are a.s. continuous with respect to  $\lambda$ , it is clear that  $\lim_{\delta \rightarrow 0} \sup_{|\mu-\lambda| \leq \delta} |\mathcal{E}^+(W, \lambda) - \mathcal{E}^+(W, \mu)|^2 = 0$  a.s. Moreover, we can apply Lebesgue's theorem with the upper-bound given by (3.9) to deduce that

$$\lim_{\delta \rightarrow 0} \sup_{\ell} \mathbb{E} \left[ \sup_{|\mu-\lambda| \leq \delta} |Y_{1,\ell}(\lambda) - Y_{k,\ell}(\mu)| \right] = 0.$$

Thus, Assumption  $(\mathcal{H}-6)$  of Proposition 4.3 is satisfied. Finally, we can apply Proposition 4.3 to prove that the sequence  $\frac{1}{N_\ell'} \sum_{k=1}^{N_\ell'} Y_{k,\ell}$  converges a.s locally uniformly to 0. The convergence of  $\widehat{\lambda}_\ell$  to  $\lambda^*$  can be deduced by closely mimicking the proof of Theorem 2.4.  $\square$

### 3.3 Strong law of large numbers and central limit theorem

Let us introduce a sequence  $(a_\ell)_{\ell \in \mathbb{N}}$  of positive real numbers such that  $\lim_{L \rightarrow \infty} \sum_{\ell=1}^L a_\ell = \infty$ . We assume that the sample size  $N_\ell$  has the following form

$$N_{\ell,L}^\rho = \frac{\rho(L)}{m^\ell a_\ell} \sum_{k=1}^L a_k, \quad \ell \in \{0, \dots, L\} \quad (3.10)$$

for some increasing function  $\rho : \mathbb{N} \rightarrow \mathbb{R}$ .

We choose this form for  $N_\ell$  because it is a generic form allowing us a straightforward use of the Toeplitz Lemma, which is a key tool to prove the central limit theorem. Since  $\lim_{L \rightarrow \infty} \sum_{\ell=1}^L a_\ell = \infty$ , for any sequence  $(x_\ell)_{\ell \geq 1}$  converging to some limit  $x \in \mathbb{R}$ ,

$$\lim_{L \rightarrow +\infty} \frac{\sum_{\ell=1}^L a_\ell x_\ell}{\sum_{\ell=1}^L a_\ell} = x.$$

We define the  $\sigma$ -algebra  $\mathcal{G}$  generated by the samples  $(W_{\ell,k})_{\ell,k \geq 1}$  used to compute  $\hat{\lambda}_L$ . In the above framework, the variables  $(\tilde{W}_{\ell,k})_{\ell,k}$  are independent of  $\mathcal{G}$ . We also introduce the filtration  $(\tilde{\mathcal{G}}_\ell)_{\ell > 0}$  generated by  $(\tilde{W}_{\ell,k}, k \geq 1)_\ell$  and the filtration  $(\mathcal{G}_\ell^\#)_{\ell > 0}$  defined as  $\mathcal{G}_\ell^\# = \mathcal{G} \vee \tilde{\mathcal{G}}_\ell$ .

**Theorem 3.2.** *Assume that  $\sup_L \sup_\ell \frac{L^2 a_\ell}{\rho(L) \sum_{k=1}^L a_k} < +\infty$ . Then, under the assumptions of Theorem 3.1,  $Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L) \rightarrow \mathbb{E}[\psi(X_T)]$  a.s. when  $L \rightarrow +\infty$ .*

For the choice  $a_\ell = 1$  for all  $\ell$ , the condition on  $\rho$  reduces to  $\sup_L \frac{L}{\rho(L)} < +\infty$ .

*Proof.* As  $\mathbb{E}[\psi(X_T^L)]$  converges to  $\mathbb{E}[\psi(X_T)]$  as  $L$  goes to infinity, it is enough to show that  $Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L) - \mathbb{E}[\psi(X_T^L)]$  tends to 0.

$$\begin{aligned} Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L) - \mathbb{E}[\psi(X_T^L)] &= \frac{1}{N_{0,L}^\rho} \sum_{k=1}^{N_{0,L}^\rho} \psi(\tilde{X}_{T,0,k}^{m^0}(\hat{\lambda}_0)) \mathcal{E}^-(\tilde{W}_{0,k}, \hat{\lambda}_0) - \mathbb{E}[\psi(X_{T,0}^{m^0})] \\ &\quad + \sum_{\ell=1}^L \frac{1}{N_{\ell,L}^\rho} \left( \sum_{k=1}^{N_{\ell,L}^\rho} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\hat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\hat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \hat{\lambda}_\ell) \right. \\ &\quad \left. - \mathbb{E} \left[ \psi(\tilde{X}_{T,\ell}^{m^\ell}) - \psi(\tilde{X}_{T,\ell}^{m^{\ell-1}}) \right] \right). \end{aligned} \quad (3.11)$$

From Theorem 2.5 and Remark 2.7, we know that

$$\frac{1}{N_{0,L}^\rho} \sum_{k=1}^{N_{0,L}^\rho} \psi(\tilde{X}_{T,0,k}^{m^0}(\hat{\lambda}_0)) \mathcal{E}^-(\tilde{W}_{0,k}, \hat{\lambda}_0) - \mathbb{E}[\psi(X_{T,0}^{m^0})] \xrightarrow[L \rightarrow +\infty]{a.s.} 0.$$

Then, it suffices to prove that the remaining terms in (3.11) tend to 0 with  $L$ . Let  $\mathcal{V}$  be a



compact neighbourhood of  $\lambda^*$ .

$$\begin{aligned} & \sum_{\ell=1}^L \frac{1}{N_{\ell,L}^\rho} \left( \sum_{k=1}^{N_{\ell,L}^\rho} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\hat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\hat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \hat{\lambda}_\ell) - \mathbb{E} \left[ \psi(\tilde{X}_{T,\ell}^{m^\ell}) - \psi(\tilde{X}_{T,\ell}^{m^{\ell-1}}) \right] \right) = \\ & \sum_{\ell=1}^L \frac{1}{N_{\ell,L}^\rho} \left( \sum_{k=1}^{N_{\ell,L}^\rho} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\hat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\hat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \hat{\lambda}_\ell) - \mathbb{E} \left[ \psi(\tilde{X}_{T,\ell}^{m^\ell}) - \psi(\tilde{X}_{T,\ell}^{m^{\ell-1}}) \right] \right) 1_{\{\hat{\lambda}_\ell \in \mathcal{V}\}} \\ & + \sum_{\ell=1}^L \frac{1}{N_{\ell,L}^\rho} \left( \sum_{k=1}^{N_{\ell,L}^\rho} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\hat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\hat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \hat{\lambda}_\ell) - \mathbb{E} \left[ \psi(\tilde{X}_{T,\ell}^{m^\ell}) - \psi(\tilde{X}_{T,\ell}^{m^{\ell-1}}) \right] \right) 1_{\{\hat{\lambda}_\ell \notin \mathcal{V}\}} \end{aligned}$$

For  $\ell$  large enough (although random),  $1_{\{\hat{\lambda}_\ell \notin \mathcal{V}\}} = 0$ . Hence, the second term in the above equation tends to 0 a.s. when  $L$  goes to infinity. It remains to prove that the first term also converges to zero. To do so, we apply Proposition 4.1 to the sequence

$$\begin{aligned} Y_{\ell,q} = & q \frac{1}{N_{\ell,q}^\rho} \left( \sum_{k=1}^{N_{\ell,q}^\rho} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\hat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\hat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \hat{\lambda}_\ell) \right. \\ & \left. - \mathbb{E} \left[ \left( \psi(\tilde{X}_{T,\ell}^{m^\ell}) - \psi(\tilde{X}_{T,\ell}^{m^{\ell-1}}) \right) \right] \right) 1_{\{\hat{\lambda}_\ell \in \mathcal{V}\}} \end{aligned}$$

and set  $\bar{Y}_{L,q} = \frac{1}{L} \sum_{\ell=1}^L Y_{\ell,q}$ . Note that  $\mathbb{E}[Y_{\ell,q}] = 0$  for all  $\ell$  and  $q$ . Since the samples used in the different levels are independent and the  $\hat{\lambda}_\ell$ 's are independent of the filtration  $\tilde{\mathcal{G}}$ , we can write

$$\mathbb{E} \left[ |\bar{Y}_{L,q}|^2 \right] = \frac{1}{L^2} \mathbb{E} \left[ \mathbb{E} \left[ \left| \sum_{\ell=1}^L Y_{\ell,q} \right|^2 \middle| \mathcal{G} \right] \right] = \frac{1}{L^2} \sum_{\ell=1}^L \mathbb{E} \left[ |Y_{\ell,q}|^2 \right]. \quad (3.12)$$

Using the same kind of arguments, we obtain

$$\begin{aligned} \mathbb{E} \left[ |Y_{\ell,q}|^2 \right] & \leq q^2 \frac{1}{N_{\ell,q}^\rho} \mathbb{E} \left[ \left( \psi(\tilde{X}_{T,\ell}^{m^\ell}) - \psi(\tilde{X}_{T,\ell}^{m^{\ell-1}}) \right)^2 \mathcal{E}^+(\tilde{W}_\ell, \hat{\lambda}_\ell) 1_{\{\hat{\lambda}_\ell \in \mathcal{V}\}} \right] \\ & \leq \frac{q^2 a_\ell}{\rho(q) \sum_{k=1}^q a_k} \left\{ m^\ell \mathbb{E} \left[ \left( \psi(\tilde{X}_{T,\ell}^{m^\ell}) - \psi(\tilde{X}_{T,\ell}^{m^{\ell-1}}) \right)^2 \mathcal{E}^+(\tilde{W}_\ell, \hat{\lambda}_\ell) 1_{\{\hat{\lambda}_\ell \in \mathcal{V}\}} \right] \right\}. \end{aligned}$$

From Proposition A.4, the term into braces converges when  $\ell$  goes to infinity. Hence, using the assumptions on the function  $\rho$ , we get

$$\sup_q \sup_\ell \mathbb{E} \left[ |Y_{\ell,q}|^2 \right] < +\infty. \quad (3.13)$$

By combining Equations (3.12) and (3.13), we get that  $\sup_L \sup_q L \mathbb{E} \left[ |\bar{Y}_{L,q}|^2 \right] < +\infty$ . Hence, Proposition 4.1 yields that  $\bar{Y}_{L,L}$  vanishes when  $L$  goes to infinity and this ends the proof.  $\square$

**Theorem 3.3.** *Suppose that the assumptions of Theorem 3.1 hold and that Condition  $(\mathcal{H}_\gamma)$  is satisfied. If  $N_{\ell,L}^\rho$  is given by (3.10) with  $\rho(L) = m^{2\gamma L}(m-1)T$  and the sequence  $(a_\ell)_\ell$  satisfies*

$$\lim_{L \rightarrow \infty} \frac{1}{\left(\sum_{\ell=1}^L a_\ell\right)^{p/2}} \sum_{\ell=1}^L a_\ell^{p/2} = 0, \text{ for } p > 2, \quad (3.14)$$

then  $m^{\gamma L}(Q_L(\widehat{\lambda}_0, \dots, \widehat{\lambda}_L) - \mathbb{E}[\psi(X_T)]) \implies \mathcal{N}(C_\psi(T, \gamma), v(\lambda^*))$  when  $L \rightarrow \infty$ .

The convergence rate does not depend on the number of samples  $N'_\ell$  provided that they tend to infinity with  $\ell$ .

*Proof.* By assumption  $(\mathcal{H}_\gamma)$ , we have that  $\lim_{L \rightarrow +\infty} m^{\gamma L}(\mathbb{E}[\psi(X_T^{m^L}) - \psi(X_T)]) = C_\psi(T, \gamma)$ . The convergence of the level 0 is governed by Theorem 2.6 (see Remark 2.7) which yields that, when  $L \rightarrow \infty$ ,

$$\left( \frac{1}{\sqrt{N_{0,L}^\rho}} \sum_{k=1}^{N_{0,L}^\rho} \psi(\tilde{X}_{T,0,k}^{m^0}(\widehat{\lambda}_0)) \mathcal{E}^-(\tilde{W}_{0,k}, \widehat{\lambda}_0) - \mathbb{E}[\psi(X_T^{m^0})] \right) \implies \mathcal{N}(0, \sigma_0^2(\widehat{\lambda}_0)).$$

Then, we deduce from the choice of the function  $\rho$  that

$$m^{\gamma L} \left( \frac{1}{N_{0,L}^\rho} \sum_{k=1}^{N_{0,L}^\rho} \psi(\tilde{X}_{T,0,k}^{m^0}(\widehat{\lambda}_0)) \mathcal{E}^-(\tilde{W}_{0,k}, \widehat{\lambda}_0) - \mathbb{E}[\psi(X_T^{m^0})] \right) \xrightarrow[L \rightarrow +\infty]{\mathbb{P}} 0.$$

Since all the blocks are independent, it is sufficient to prove that

$$m^{\gamma L} \left( \sum_{\ell=1}^L \frac{1}{N_{\ell,L}^\rho} \sum_{k=1}^{N_{\ell,L}^\rho} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\widehat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\widehat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \widehat{\lambda}_\ell) - \mathbb{E}[\psi(X_T^n)] \right) \implies \mathcal{N}(0, v(\lambda^*)).$$

To do so, we introduce the  $(\mathcal{G}_l^\sharp)_{l \geq 1}$ -martingale array  $(Y_l^n)_{l \geq 1}$  defined by

$$Y_l^n \triangleq \sum_{\ell=1}^l \frac{m^{\gamma L}}{N_{\ell,L}^\rho} \sum_{i=1}^{N_{\ell,L}^\rho} \left[ \left( \psi(\tilde{X}_{T,\ell,i}^{m^\ell}(\widehat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,i}^{m^{\ell-1}}(\widehat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,i}, \widehat{\lambda}_\ell) - \mathbb{E} \left[ \psi(\tilde{X}_T^{m^\ell}) - \psi(\tilde{X}_T^{m^{\ell-1}}) \right] \right],$$

so  $\mathbb{E}[Y_l^n] = 0$  for all  $l, n$ . According to Theorem A.1, we need to study the asymptotic behaviors of the two quantities

$$\langle Y^n \rangle_L = \sum_{\ell=1}^L \mathbb{E} \left[ |Y_\ell^n - Y_{\ell-1}^n|^2 \middle| \mathcal{G}_{\ell-1}^\sharp \right] \text{ and } \sum_{\ell=1}^L \mathbb{E} \left[ |Y_\ell^n - Y_{\ell-1}^n|^p \middle| \mathcal{G}_{\ell-1}^\sharp \right], \text{ for } p > 2 \text{ as } n \rightarrow \infty.$$

Note that  $\widehat{\lambda}_\ell$  is  $\mathcal{G}_{\ell-1}^\sharp$ -measurable and for any  $\lambda \in \mathbb{R}^q$  the variables  $(\tilde{X}_{T,\ell,i}^{m^\ell}(\lambda), \tilde{X}_{T,\ell,i}^{m^{\ell-1}}(\lambda))_{1 \leq i \leq N_\ell}$  are independent of  $\mathcal{G}_{\ell-1}^\sharp$ , then using (3.10) with  $\rho(L) = m^{2\gamma L}(m-1)T$ , we rewrite the first quantity as follows

$$\langle Y^n \rangle_L = \frac{1}{\sum_{\ell=1}^L a_\ell} \sum_{\ell=1}^L a_\ell \left[ v_\ell(\widehat{\lambda}_\ell) - \Xi_\ell^2 \right]$$

with  $v_\ell$  defined by (3.3) and  $\Xi_\ell$  defined by (3.4). Let  $\mathcal{V}$  be a compact neighbourhood of  $\lambda^*$ . We can write

$$\langle Y^n \rangle_L = \frac{1}{\sum_{\ell=1}^L a_\ell} \sum_{\ell=1}^L a_\ell \left[ v_\ell(\widehat{\lambda}_\ell) - \Xi_\ell^2 \right] 1_{\{\widehat{\lambda}_\ell \in \mathcal{V}\}} + \frac{1}{\sum_{\ell=1}^L a_\ell} \sum_{\ell=1}^L a_\ell \left[ v_\ell(\widehat{\lambda}_\ell) - \Xi_\ell^2 \right] 1_{\{\widehat{\lambda}_\ell \notin \mathcal{V}\}}. \quad (3.15)$$

From Proposition A.4, we know that  $\Xi_\ell \rightarrow \mathbb{E}[\nabla\psi(X_T).U_T] = 0$ , where the last equality is a straightforward consequence of [24, Proposition 2.1]. From Proposition A.4, we know that the sequence of functions  $v_\ell$  converges pointwise to  $v$  defined by (3.5). Moreover, we can easily prove that this convergence is locally uniform. Hence, by the convergence of  $\widehat{\lambda}_\ell$  to  $\lambda^*$  (see Theorem 3.1), we deduce that  $v_\ell(\widehat{\lambda}_\ell)1_{\{\widehat{\lambda}_\ell \in \mathcal{V}\}}$  converges to  $v(\lambda^*)$  when  $\ell \rightarrow +\infty$ . Moreover, for  $\ell$  large enough (although random),  $1_{\{\widehat{\lambda}_\ell \notin \mathcal{V}\}} = 0$ .

Thus, we deduce from the Toeplitz lemma that  $\langle Y^n \rangle_L \rightarrow v(\lambda^*)$  a.s. Using Burkholder's inequality and Jensen's inequality together with the assumptions on  $\psi$  and Property (H-1)-ii, we obtain that for any  $p > 2$ , there exists  $C_p > 0$  such that

$$\sum_{\ell=1}^L \mathbb{E} \left[ |Y_\ell^n - Y_{\ell-1}^n|^p | \mathcal{G}_{\ell-1}^\# \right] \leq \frac{C_p}{\left( \sum_{\ell=1}^L a_\ell \right)^{p/2}} \sum_{\ell=1}^L a_\ell^{p/2} \xrightarrow{L \rightarrow \infty} 0$$

where the convergence to zero is ensured by (3.14). Consequently, we can apply Theorem A.1 to achieve the proof.  $\square$

**Remark 3.4.** *As usual, one can rescale  $m^{\gamma L}(Q_L(\widehat{\lambda}_0, \dots, \widehat{\lambda}_L) - \mathbb{E}[\psi(X_T)])$  by an estimator of  $v(\lambda^*)$  to obtain a central limit theorem with variance 1. Thanks to Theorem 3.1, we know that  $v_{\ell, N_\ell}(\widehat{\lambda}_\ell)$  is a convergent estimator of  $v(\lambda^*)$  and we can easily deduce from the proof of Theorem 3.3 that under its assumptions*

$$m^{2\gamma L} \left\{ \frac{1}{N_{0,L}^\rho} \left( \frac{1}{N_{0,L}^\rho} \sum_{k=1}^{N_{0,L}^\rho} (\psi(\tilde{X}_T^{m_0}) \mathcal{E}^+(\tilde{W}_{0,k}, \lambda_0))^2 - \left( \frac{1}{N_{0,L}^\rho} \sum_{k=1}^{N_{0,L}^\rho} \psi(\tilde{X}_T^{m_0}) \mathcal{E}^+(\tilde{W}_{0,k}, \lambda_0) \right)^2 \right) + \sum_{\ell=1}^L N_\ell^{-1} \frac{(m-1)T}{m^\ell} \left( \tilde{v}_{\ell, N_\ell}(\lambda_\ell) - \tilde{\Xi}_{\ell, N_\ell}^2 \right) \right\} \xrightarrow{L \rightarrow +\infty} v(\lambda^*).$$

Note the quantities  $\tilde{v}_{\ell, N_\ell}$  and  $\tilde{\Xi}_{\ell, N_\ell}$  are defined as in Equations (3.3) and (3.4) and but using the tilde sample paths  $(\tilde{X}_{\ell,k})$  and  $(\tilde{W}_{\ell,k})$ . The term into braces, which can be computed online during the multilevel Monte Carlo procedure, can be used to build confidence intervals. Any convergent estimator of  $v(\lambda^*)$  could of course be used, but this one has the advantage to correspond to the true variance of the multilevel Monte Carlo estimator for any finite number of levels  $L$  and not only asymptotically.

## 4 Strong law of large numbers for doubly indexed sequences

In this section, we prove two corner stone results used in the convergence of the multilevel approach. We tackle the convergence of empirical averages of doubly indexed random sequences when both indices tend to infinity together.

**Proposition 4.1.** Let  $(X_{n,m})_{n,m}$  be a doubly indexed sequence of vector valued random variables such that for all  $n$ ,  $\mathbb{E}[X_{n,m}] = x_m$  with  $\lim_{m \rightarrow +\infty} x_m = x$ . We define  $\bar{X}_{n,m} = \frac{1}{n} \sum_{i=1}^n X_{i,m}$ . Assume that the two following assumptions are satisfied

- ( $\mathcal{H}$ -3)    i.  $\sup_n \sup_m n \text{Var}(\bar{X}_{n,m}) < +\infty$ .  
               ii.  $\sup_n \sup_m \text{Var}(X_{n,m}) < +\infty$ .

Then, for all increasing functions  $\rho : \mathbb{N} \rightarrow \mathbb{N}$ ,  $\bar{X}_{n,\rho(n)} \rightarrow x$  a.s. and in  $\mathbb{L}^2$  when  $n \rightarrow \infty$ .

From this proposition, one can easily deduce the following corollary by extracting a bespoken subsequence

**Corollary 4.2.** Assume that  $(X_{n,m})_{n,m}$  be a doubly indexed sequence of vector valued random variables satisfying the assumptions of Proposition 4.1. Then, for any strictly increasing function  $\xi : \mathbb{N} \rightarrow \mathbb{N}$ ,  $\bar{X}_{\xi(n),n} \rightarrow x$  a.s. and in  $\mathbb{L}^2$  when  $n \rightarrow \infty$ .

*Proof of Proposition 4.1.* The proof of this result closely mimics the one of [28, Theorem IV.1.1]. We introduce the sequence  $(Y_{i,m})_{i,m}$  defined by  $Y_{i,m} = X_{i,m} - x_m$ , which satisfies  $\mathbb{E}[Y_{i,m}] = 0$ . As  $\lim_{m \rightarrow \infty} x_m = x$ , it is sufficient to prove that  $\bar{Y}_{n,\rho(n)} \rightarrow 0$  a.s.

Condition ( $\mathcal{H}$ -3)-i implies the  $\mathbb{L}^2$  convergence to 0. We introduce the sequence  $(Z_{n,m})_n$  defined by  $Z_{n,m} = \sup\{|Y_{k,m}| : n^2 \leq k < (n+1)^2\}$ . Let  $k$  be such that  $n^2 \leq k < (n+1)^2$ , then

$$|\bar{Y}_{k,m}| \leq n^{-2} \left( n^2 |\bar{Y}_{n^2,m}| + \sum_{i=n^2+1}^k |Y_{i,m}| \right),$$

$$Z_{n,m} \leq |\bar{Y}_{n^2,m}| + \frac{1}{n^2} \sum_{i=n^2+1}^{(n+1)^2} |Y_{i,m}|.$$

Then,

$$\mathbb{E}[Z_{n,m}^2] \leq \mathbb{E}[\bar{Y}_{n^2,m}^2] + \sum_{i=n^2+1}^{(n+1)^2} \left( \frac{\mathbb{E}[|Y_{i,m}|^2]}{n^4} + 2 \frac{\mathbb{E}[|\bar{Y}_{n^2,m}| |Y_{i,m}|]}{n^2} \right) + 2 \sum_{i,j=n^2+1; i \neq j}^{(n+1)^2} \frac{\mathbb{E}[|Y_{j,m}| |Y_{i,m}|]}{n^4}.$$

Let  $\kappa > 0$  denote the maximum of the upper bounds involved in Assumption ( $\mathcal{H}$ -3). Using the Cauchy Schwartz inequality, we get

$$\begin{aligned} \mathbb{E}[Z_{n,m}^2] &\leq \frac{\kappa}{n^2} + \frac{\kappa((n+1)^2 - n^2)}{n^4} + 2 \frac{\kappa^2((n+1)^2 - n^2)}{n^3} + 2 \frac{\kappa^2((n+1)^2 - n^2)^2}{n^4} \\ &\leq \frac{\kappa}{n^2} + \frac{\kappa(2n+1)}{n^4} + 2 \frac{\kappa^2(2n+1)}{n^3} + 2 \frac{\kappa^2(2n+1)^2}{n^4}. \end{aligned}$$

Hence, for any function  $\rho : \mathbb{N} \rightarrow \mathbb{N}$ ,  $\mathbb{E}[Z_{n,\rho(n)}^2] \leq Cn^{-2}$  where  $C > 0$  is a constant independent of  $\rho$ . Therefore, we have  $\mathbb{P}(Z_{n,\rho(n)} \geq n^{-1/4}) \leq Cn^{-3/2}$ . This inequality implies using the Borel Cantelli Lemma that, for  $n$  large enough  $Z_{n,\rho(n)} \leq n^{-1/4}$  a.s. which yields the a.s. convergence to 0.  $\square$

**Proposition 4.3.** Let  $(F_{n,m})_{n,m}$  be a doubly indexed sequence of random variables with values in the set of continuous functions, ie. for all  $n, m$ ,  $F_{n,m} : \Omega \rightarrow C^0(\mathbb{R}^d)$ . Moreover, we assume that there exists a sequence of deterministic functions  $(f_m)_m$  s.t. for all  $n$   $\mathbb{E}[F_{n,m}] = f_m$  for all  $m$ . We define  $\bar{F}_{n,m} = \frac{1}{n} \sum_{i=1}^n F_{i,m}$ . Assume that the two following assumptions are satisfied

( $\mathcal{H}$ -4) One of the following criteria holds

- i. The sequence  $(f_m)_m$  converges pointwise to some continuous function  $f$ .
- ii. The sequence  $(f_m)_m$  converges locally uniformly to some function  $f$ .

( $\mathcal{H}$ -5) For any compact set  $W \subset \mathbb{R}^d$ ,

- i.  $\sup_n \sup_m n \text{Var} \left( \sup_{x \in W} |\overline{F}_{n,m}(x)| \right) < +\infty$ .
- ii.  $\sup_n \sup_m \text{Var} \left( \sup_{x \in W} |F_{n,m}(x)| \right) < +\infty$ .

( $\mathcal{H}$ -6) For all  $y \in \mathbb{R}^d$ ,  $\lim_{\delta \rightarrow 0} \sup_n \sup_m \mathbb{E} \left[ \sup_{|x-y| \leq \delta} |F_{n,m}(x) - F_{n,m}(y)| \right] = 0$ .

Then, for all functions  $\rho : \mathbb{N} \rightarrow \mathbb{N}$ , the sequence of random functions  $\overline{F}_{n,\rho(n)}$  converges a.s. locally uniformly to the locally continuous function  $f$ .

**Remark 4.4.** • When for every fixed  $m$ , the sequence  $(F_{n,m})_n$  is independent and identically distributed, Assumption ( $\mathcal{H}$ -6) is ensured by

$$\forall y \in \mathbb{R}^d, \lim_{\delta \rightarrow 0} \limsup_m \mathbb{E} \left[ \sup_{|x-y| \leq \delta} |F_{1,m}(x) - F_{1,m}(y)| \right] = 0$$

and Assumption ( $\mathcal{H}$ -5)-ii implies ( $\mathcal{H}$ -5)-i.

- As in Corollary 4.2, for any strictly increasing function  $\xi : \mathbb{N} \rightarrow \mathbb{N}$ , the sequence  $\overline{F}_{\xi(n),n}$  converges a.s. locally uniformly to the locally continuous function  $f$ .

*Proof.* We can apply Proposition 4.1, to deduce that a.s.  $\overline{F}_{n,\rho(n)}$  converges pointwise to the function  $f$ . If we do not already know that  $f$  is continuous, then thanks to ( $\mathcal{H}$ -5)-ii, we can apply Lebesgue's theorem to deduce that the functions  $f_m$  are continuous. The uniform convergence of the sequence  $f_m$  to  $f$  (see ( $\mathcal{H}$ -4)-ii) proves that the function  $f$  is continuous.

Let  $W$  be a compact set of  $\mathbb{R}^d$ , we can cover  $W$  with a finite number  $K$  of open balls  $W_k$  with centers  $(x_k)_k$  and radiuses  $(r_k)_k$ , i.e.  $W_k = B(x_k, r_k)$  and  $W = \cup_{k=1}^K W_k$ . We want to prove that

$$\sup_{x \in W} |\overline{F}_{n,\rho(n)}(x) - f(x)| \xrightarrow[n \rightarrow +\infty]{a.s.} 0.$$

We write

$$\sup_{x \in W} |\overline{F}_{n,\rho(n)}(x) - f(x)| = \sum_{k=1}^K \sup_{x \in W_k} |\overline{F}_{n,\rho(n)}(x) - f(x)|. \quad (4.1)$$

We split each term

$$\begin{aligned} \sup_{x \in W_k} |\overline{F}_{n,\rho(n)}(x) - f(x)| &= \sup_{x \in W_k} |\overline{F}_{n,\rho(n)}(x) - \overline{F}_{n,\rho(n)}(x_k)| + \sup_{x \in W_k} |f(x) - f(x_k)| \\ &\quad + |\overline{F}_{n,\rho(n)}(x_k) - f(x_k)| \end{aligned} \quad (4.2)$$

Let  $\varepsilon > 0$ . The idea is to choose the radiuses  $r_k$  small enough to ensure that each term is controlled by a function of  $\varepsilon$ . Now, we make the idea precise. For all  $k = 1, \dots, K$ , the

last term can be made smaller than  $\varepsilon/K$  for  $n$  larger than some  $N_k$  using the pointwise convergence. For all  $n \geq \max_{k \leq K} N_k$ , and all  $1 \leq k \leq K$ ,  $|\bar{F}_{n,\rho(n)}(x_k) - f(x_k)| \leq \varepsilon/K$ . The function  $f$  being continuous, it is uniformly continuous on every  $W_k$ . If we choose the  $W_k$  such that their radii are small enough (we may need to increase  $K$ ), we can ensure that for all  $1 \leq k \leq K$   $\sup_{x \in W_k} |f(x) - f(x_k)| \leq \varepsilon/K$ . The first term on the r.h.s of (4.2) deserves more attention

$$\sup_{x \in W_k} |\bar{F}_{n,\rho(n)}(x) - \bar{F}_{n,\rho(n)}(x_k)| \leq \frac{1}{n} \sum_{i=1}^n \sup_{x \in W_k} |F_{i,\rho(n)}(x) - F_{i,\rho(n)}(x_k)|. \quad (4.3)$$

Now, for every  $1 \leq k \leq K$ , we want to apply Proposition 4.1 to the sequence of random variables  $(\sup_{x \in W_k} |F_{n,m}(x) - F_{n,m}(x_k)|)_{n,m}$ . Assumption  $(\mathcal{H}-3)$  is clearly satisfied using Minkowski's inequality.

Let us define the sequence  $(Y_{n,m})_{n,m}$  by

$$Y_{n,m} = \sup_{x \in W_k} |F_{n,m}(x) - F_{n,m}(x_k)| - \mathbb{E} \left[ \sup_{x \in W_k} |F_{n,m}(x) - F_{n,m}(x_k)| \right],$$

satisfying  $\mathbb{E}[Y_{n,m}] = 0$  and the assumptions of Proposition 4.1. Hence, it yields that

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{i=1}^n \sup_{x \in W_k} |F_{i,\rho(n)}(x) - F_{i,\rho(n)}(x_k)| - \mathbb{E} \left[ \sup_{x \in W_k} |F_{n,\rho(n)}(x) - F_{n,\rho(n)}(x_k)| \right] = 0. \quad (4.4)$$

From  $(\mathcal{H}-6)$ , we know that if the  $W_k$  are chosen small enough,

$$\sup_n \mathbb{E} \left[ \sup_{x \in W_k} |F_{n,\rho(n)}(x) - F_{n,\rho(n)}(x_k)| \right] \leq \varepsilon/K. \quad (4.5)$$

Then, combining (4.3), (4.4) and (4.5) yields that  $\sup_{x \in W_k} |\bar{F}_{n,\rho(n)}(x) - \bar{F}_{n,\rho(n)}(x_k)| \leq \varepsilon/K$ . We plus this inequality into (4.2) and deduce from (4.1), that for  $n$  large enough,

$$\sup_{x \in W} |\bar{F}_{n,\rho(n)}(x) - f(x)| \leq 3\varepsilon. \quad \square$$

## 5 Numerical experiments

### 5.1 Practical implementation

Our approach cleverly mixes the famous multilevel Monte Carlo technique with importance sampling to reduce the variance. A classical approach would have been to consider the multilevel approximation of  $\mathbb{E} \left[ \psi(X_T(\theta)) e^{-\theta \cdot W_T - \frac{1}{2} |\theta|^2 T} \right]$  while choosing the value of  $\theta$  which minimizes the variance of the central limit theorem for multilevel Monte Carlo (see [4]). This asymptotic variance involves both  $\nabla \psi$  and the process  $U$  given in (3.6). Hence, a classical approach to importance sampling for multilevel Monte Carlo would require extra knowledge than the function  $\psi$  and the underlying process  $X$ , thus precluding any kind of automation.

We have chosen a completely different approach allowing for one importance sampling parameter per level, which enables us to treat each level independently of the others. In each level, we use a sample average approximation as in [23] to compute the optimal importance

sampling parameter defined as the one minimizing the variance of the current level. From Theorem 3.3, we know that this approach is optimal in the sense that our multilevel estimator  $Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L)$  satisfies a central limit theorem with a limiting variance given by  $\inf v$  where  $v$  defined by (3.5) is the variance of the standard multilevel Monte Carlo estimator. We managed to provide an algorithm reaching the optimal limiting variance without computing  $\nabla\psi$  nor the process  $U$ , hence our approach can be made fully automatic.

**Computation of  $\hat{\lambda}_\ell$ .** The parameters  $\hat{\lambda}_\ell$  are defined as the solutions of strongly convex minimization problems. The minimization step is performed by the Newton–Raphson algorithm to  $\nabla v_{\ell, N'_\ell}$ . The samples required to compute  $\nabla v_{\ell, N'_\ell}$  and  $\nabla^2 v_{\ell, N'_\ell}$  are generated once and for all before starting the Newton–Raphson procedure such that the same samples are used through all the iterations of the gradient descent. This feature is specific to the optimisation step and may make the algorithm highly memory demanding as soon as the numbers  $N'_\ell$  become large. As the parameter  $\lambda$  is not involved in the function  $\psi$ , all the quantities  $\psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}})$  for  $k = 1, \dots, N_\ell$  can be precomputed before starting the minimization algorithm, which enables us to save a lot of computational time.

The efficiency of the Newton–Raphson algorithm very much depends on the convexity of the  $v_{\ell, N'_\ell}$  functions. As already pointed out in [23], the smallest eigenvalue of the Hessian matrix  $\nabla^2 v_{\ell, N'_\ell}$  is basically  $\frac{T}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 \mathcal{E}^+(W_{\ell,k}, \lambda)$ , which can become extremely small and then conflicts with the will to have the strongest possible convexity in order to speed up Newton–Raphson’s algorithm. This difficulty is circumvented by noticing the equality  $\nabla v_{\ell, N'_\ell}(\hat{\lambda}_\ell) = 0$  can be written as

$$\hat{\lambda}_\ell T - \frac{\frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} W_{k,\ell,T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\hat{\lambda}_\ell W_{T,\ell,k}}}{\frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\hat{\lambda}_\ell W_{T,\ell,k}}} = 0.$$

Hence,  $\hat{\lambda}_\ell$  can be interpreted as the root of  $\nabla u_{\ell, N'_\ell}$  with

$$u_{\ell, N'_\ell}(\lambda) = \frac{|\lambda|^2 T}{2} + \log \left( \frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\lambda W_{T,\ell,k}} \right).$$

The Hessian matrix of  $u_{\ell, N'_\ell}$  is given by

$$\begin{aligned} \nabla^2 u_{\ell, N'_\ell}(\lambda) = & T I_q + \frac{\frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} W_{k,\ell,T} (W_{k,\ell,T})^* \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\lambda W_{T,\ell,k}}}{\frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\lambda W_{T,\ell,k}}} \\ & - \frac{\left( \frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} W_{k,\ell,T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\lambda W_{T,\ell,k}} \right)}{\frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\lambda W_{T,\ell,k}}} \\ & - \frac{\left( \frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} W_{k,\ell,T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\lambda W_{T,\ell,k}} \right)^*}{\frac{1}{N'_\ell} \sum_{k=1}^{N'_\ell} \frac{m^\ell}{(m-1)T} \left| \psi(X_{T,\ell,k}^{m^\ell}) - \psi(X_{T,\ell,k}^{m^{\ell-1}}) \right|^2 e^{-\lambda W_{T,\ell,k}}}. \end{aligned} \quad (5.1)$$

From the Cauchy Schwartz inequality, it is clear that  $\nabla^2 u_{\ell, N'_\ell}(\lambda)$  is lower bounded by  $TI_q$ , where the inequality is to be understood in the sense of the order on symmetric matrices.

**Description of the algorithm.** Our algorithm splits in two steps: the minimization step to compute the optimal importance sampling measure and the MLMC step to actually provide an estimator of  $\mathbb{E}[\psi(X_T)]$ . The samples used in the two steps are independent. For the sake of clearness, we provide the pseudocode of our global method in in Algorithm 5.1.

<ol style="list-style-type: none"> <li>1 Generate <math>X_{T,0,1}^{m^0}, \dots, X_{T,0,N'_0}^{m^0}</math> i.i.d. samples following the law of <math>X_T^{m^0}</math> independently of the other blocks.</li> <li>2 Solve <math>\nabla u_{0, N'_0}(\hat{\lambda}_0) = 0</math> by using the Newton–Raphson algorithm.</li> <li>3 <b>for</b> <math>\ell = 1 : L</math> <b>do</b></li> <li>4     Generate <math>(X_{T,\ell,1}^{m^\ell}, X_{T,\ell,1}^{m^{\ell-1}}), \dots, (X_{T,\ell,N'_\ell}^{m^\ell}, X_{T,\ell,N'_\ell}^{m^{\ell-1}})</math> i.i.d. samples following the law of <math>(X_T^{m^\ell}, X_T^{m^{\ell-1}})</math> independently of the other blocks.</li> <li>5     Solve <math>\nabla u_{\ell, N'_\ell}(\hat{\lambda}_\ell) = 0</math> by using the Newton–Raphson algorithm.</li> <li>6 <b>end</b></li> <li>7 Conditionally on <math>\hat{\lambda}_0</math>, generate <math>\tilde{X}_{T,0,1}^{m^0}(\hat{\lambda}_0), \dots, \tilde{X}_{T,0,N_0}^{m^0}(\hat{\lambda}_0)</math> i.i.d. samples with the law of <math>X_T^{m^0}(\hat{\lambda}_0)</math> independently of the other blocks. The tilde and non tilde quantities are conditionally independent.</li> <li>8 <b>for</b> <math>\ell = 1 : L</math> <b>do</b></li> <li>9     Conditionally on <math>\hat{\lambda}_\ell</math>, generate <math>(\tilde{X}_{T,\ell,1}^{m^\ell}(\hat{\lambda}_\ell), \tilde{X}_{T,\ell,1}^{m^{\ell-1}}(\hat{\lambda}_\ell)), \dots, (\tilde{X}_{T,\ell,N_\ell}^{m^\ell}(\hat{\lambda}_\ell), \tilde{X}_{T,\ell,N_\ell}^{m^{\ell-1}}(\hat{\lambda}_\ell))</math> i.i.d. samples with the law of <math>(X_T^{m^\ell}(\hat{\lambda}_\ell), X_T^{m^{\ell-1}}(\hat{\lambda}_\ell))</math> independently of the other blocks. The tilde and non tilde quantities are conditionally independent.</li> <li>10 <b>end</b></li> <li>11 Compute the multilevel importance sampling estimator</li> </ol> $Q_L(\hat{\lambda}_0, \dots, \hat{\lambda}_L) = \frac{1}{N_0} \sum_{k=1}^{N_0} \psi(\tilde{X}_{T,0,k}^{m^0}(\hat{\lambda}_0)) \mathcal{E}^-(\tilde{W}_{0,k}, \hat{\lambda}_0) + \sum_{\ell=1}^L \frac{1}{N_\ell} \sum_{k=1}^{N_\ell} \left( \psi(\tilde{X}_{T,\ell,k}^{m^\ell}(\hat{\lambda}_\ell)) - \psi(\tilde{X}_{T,\ell,k}^{m^{\ell-1}}(\hat{\lambda}_\ell)) \right) \mathcal{E}^-(\tilde{W}_{\ell,k}, \hat{\lambda}_\ell).$
--

**Algorithm 5.1:** Multilevel Importance Sampling (MLIS)

**Complexity analysis.** In this paragraph, we focus on the impact of the number of levels  $L$  on the overall computational time of our algorithm. The computational cost of the standard multilevel estimator is proportional to

$$C_{ML} = \sum_{\ell=0}^L N_\ell m^\ell = m^{2L+1} L^2.$$



The global cost of our algorithm writes as the sum of the cost of the computation of the  $(\hat{\lambda}_\ell)_\ell$  and of the standard multilevel estimator

$$C_{MLIS} = \sum_{\ell=0}^L N'_\ell (m^\ell + 3K_\ell) + \sum_{\ell=0}^L N_\ell m^\ell$$

where  $K_\ell$  is the number of iterations of Newton–Raphson’s algorithm to approximate  $\hat{\lambda}_\ell$  and the factor 3 corresponds to the fact that building  $\nabla u_{\ell, N'_\ell}$  and  $\nabla^2 u_{\ell, N'_\ell}$  basically boils down to three Monte Carlo summations. In practice,  $K_\ell \leq 5$  as the problem is strongly convex. Because the same random variables are used at each iteration of the optimisation step, they must be stored, which makes the memory footprint of our algorithm proportional to  $N'_\ell$ .

So, if we choose  $N'_\ell = \frac{N_\ell m^\ell}{m^\ell + 15}$ , the total cost of our MLIS algorithm should be roughly twice the cost of the standard multilevel estimator. This choice of  $N'_\ell$  reduces the number of samples used to approximate the variance of the first levels compared to using directly  $N_\ell$ . However, when  $L$  increases,  $N'_\ell$  can become extremely large for small values of  $\ell$  which leads to an even larger memory footprint (see Section 5.1). Not to break the scalability of the algorithm, the values of  $N'_\ell$  have to be kept reasonable depending on the amount of memory available on the computer. For an instance, enforcing  $N'_\ell \leq 500000$  is reasonable on a computer with 8Gb of RAM. Anyway, it is crystal clear that a fairly good approximation of the variance  $v_\ell$  is enough and running for an ultimately accurate estimator would lead to a tremendous waste of computational time. Monitoring the convergence of  $v_{\ell, N'_\ell}$  would really help choosing sensible values for  $N'_\ell$ .

## 5.2 Comparison with existing algorithms

In Theorem 3.3, we obtain the same limiting variance as in [5], in which the authors apply MLMC to importance sampling (see (1.4)) and not vice-versa as we do. The way importance sampling and MLMC are coupled does not actually matter in terms of convergence rate but it does matter in practice. First, our approach preserves the independence of the different levels by solving one optimization problem per level instead of a global one. Hence, the contributions of the different levels are computed independently of each other as in the standard MLMC setting. Second, we use sample average approximation combined with Newton–Raphson’s algorithm to compute the best importance parameters, whereas in [5, 6], the authors rely on stochastic approximation, which is known to demand proper tuning to effectively converge in practice. Our approach inherits from the good convergence properties of Newton’s algorithm when applied to strongly convex problems with a tractable Hessian matrix. As already noted in [23], this approach is more stable and robust.

## 5.3 Experimental settings

We compare four methods in terms of their root mean squared error (RMSE): the crude Monte Carlo method (MC), the adaptive Monte Carlo method proposed in [23] (MC+IS), the Multilevel Monte Carlo method (ML) and our Importance Sampling Multilevel Monte Carlo estimator (ML+IS). We recall that the RMSE is defined by  $RMSE = \sqrt{\text{Bias}^2 + \text{Variance}}$ . In the computation of the bias, the true value is replaced by its multilevel Monte Carlo estimator with  $L = 9$  levels, which yields a very accurate approximation. Not to mention, the CPU times showed on the graphs take into account both the time to the search for the optimal parameter and the time for the second stage Monte Carlo, be it multilevel or not.

## 5.4 Multidimensional Dupire's framework

We consider a  $d$ -dimensional local volatility model, in which the dynamics, under the risk neutral measure, of each asset  $S^i$  is supposed to be given by

$$dS_t^i = S_t^i(r dt + \sigma(t, S_t^i)dW_t^i), \quad S_0 = (S_0^1, \dots, S_0^d)$$

where  $W = (W^1, \dots, W^d)$ , each component  $W^i$  being a standard Brownian motion with values in  $\mathbb{R}$ . For the numerical experiments, the covariance structure of  $W$  will be assumed to be given by  $\langle W^i, W^j \rangle_t = \rho t 1_{\{i \neq j\}} + t 1_{\{i=j\}}$ . We suppose that  $\rho \in (-\frac{1}{d-1}, 1)$ , which ensures that the matrix  $C = (\rho 1_{\{i \neq j\}} + 1_{\{i=j\}})_{1 \leq i, j \leq d}$  is positive definite. Let  $L$  denote the lower triangular matrix involved in the Cholesky decomposition  $C = LL^*$ . To simulate  $W$  on the time-grid  $0 < t_1 < t_2 < \dots < t_N$ , we need  $d \times N$  independent standard normal variables and set

$$\begin{pmatrix} W_{t_1} \\ W_{t_2} \\ \vdots \\ W_{t_{N-1}} \\ W_{t_N} \end{pmatrix} = \begin{pmatrix} \sqrt{t_1}L & 0 & 0 & \dots & 0 \\ \sqrt{t_1}L & \sqrt{t_2 - t_1}L & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \sqrt{t_{N-1} - t_{N-2}}L & 0 \\ \sqrt{t_1}L & \sqrt{t_2 - t_1}L & \dots & \sqrt{t_{N-1} - t_{N-2}}L & \sqrt{t_N - t_{N-1}}L \end{pmatrix} G,$$

where  $G$  is a normal random vector in  $\mathbb{R}^{d \times N}$ . The maturity time and the interest rate are respectively denoted by  $T > 0$  and  $r > 0$ . The local volatility function  $\sigma$  we have chosen is of the form

$$\sigma(t, x) = 0.6(1.2 - e^{-0.1t}e^{-0.001(xe^{rt}-s)^2})e^{-0.05\sqrt{t}}, \quad (5.2)$$

with  $s > 0$ . We know that there exists a duality between the variables  $(t, x)$  and  $(T, K)$  in Dupire's framework. Hence for formula (5.2) to make sense, one should choose  $s$  equal to the spot price of the underlying asset so that the bottom of the smile is located at the forward money. We refer to Figure 1 to have an overview of the smile.

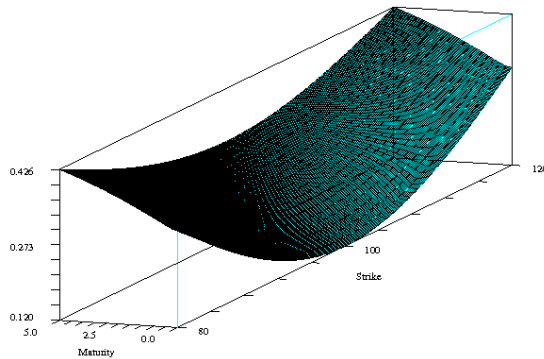


Figure 1: Local volatility function

**Basket option** We consider options with payoffs of the form  $(\sum_{i=1}^d \omega^i S_T^i - K)_+$  where  $(\omega^1, \dots, \omega^d)$  is a vector of algebraic weights. The strike value  $K$  can be taken negative to

deal with Put like options. With no surprise, we can see on Figure 2 that multilevel estimators always outperform their classical Monte Carlo counterpart. The comparison for very little accurate estimators may be meaningless as it is pretty difficult to reliably measure short execution times and the empirical variance of the estimator is in this case even less accurate than the estimator itself. Note that the points on the extreme right hand side are obtained for multilevel estimators with  $L = 2$ , respectively for Monte Carlo estimators with 256 samples. For RMSE between 0.1 and 0.005, our MLIS estimator is 10 times faster than the standard ML estimator. When a very high accuracy is required, namely when RMSE is smaller than 0.001, the MLIS estimator remains between 3 and 4 times faster than the standard multilevel estimator, which is already a great achievement since for this level of accuracy, the ML estimator may need several dozens of minutes to yield its result.

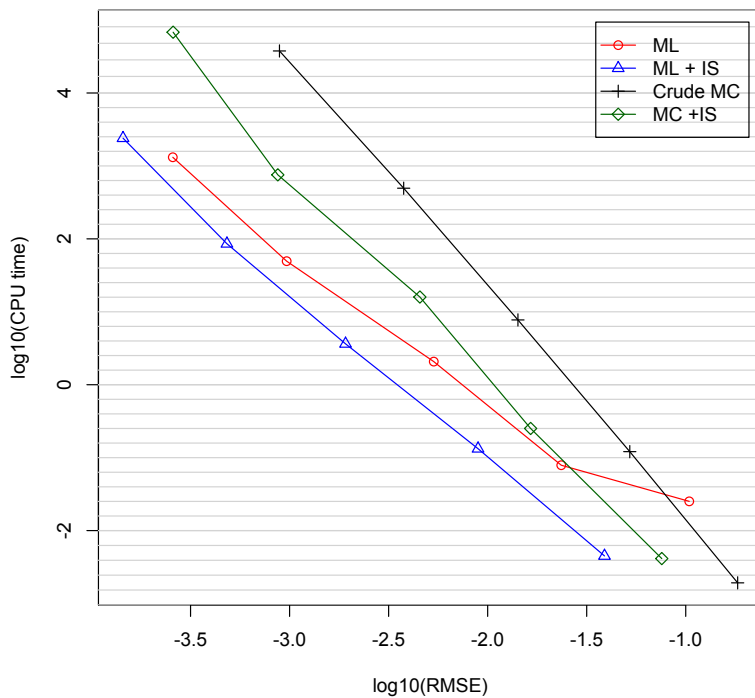


Figure 2:  $\sqrt{MSE}$  vs. CPU time for a basket option in the local volatility model with  $I = 5$ ,  $r = 0.05$ ,  $T = 1$ ,  $S_0 = 100$ ,  $K = 100$ ,  $m = 4$ .

## 5.5 Multidimensional Heston model

The multidimensional Heston model can be easily written by specifying on the one hand that each asset follows a 1-D Heston model and on the other hand the correlation structure between the involved Brownian motions. The asset price process  $S = (S^1, \dots, S^d)$  and the volatility

process  $\sigma = (\sigma^1, \dots, \sigma^d)$  solve

$$\begin{aligned} dS_t^i &= rS_t^i dt + \sqrt{\sigma_t^i} S_t^i dB_t^i \\ d\sigma_t^i &= \kappa^i (a^i - \sigma_t^i) dt + \nu_t^i \sqrt{\sigma_t^i} (\gamma^i dB_t^i + \sqrt{1 - (\gamma^i)^2} d\tilde{B}_t^i) \end{aligned}$$

where all the components of  $B = (B^1, \dots, B^d)$  and  $\tilde{B} = (\tilde{B}^1, \dots, \tilde{B}^d)$  are real valued Brownian motions. The vectors  $\kappa = (\kappa^1, \dots, \kappa^d)$  and  $a = (a^1, \dots, a^d)$  denote respectively the reversion rate and the mean level of each volatility process, while the vector  $\nu$  is the volatility of the volatility process. The vector  $\bar{\gamma} = (\gamma^1, \dots, \gamma^d)$  embodies the correlations between an asset and its volatility process, with  $\gamma^i \in ]-1, 1[$  for all  $1 \leq i \leq d$ . The vector valued processes  $B$  and  $\tilde{B}$  are independent and satisfy

$$d\langle B \rangle_t = \Gamma_S dt \quad \text{and} \quad d\langle \tilde{B} \rangle_t = I_d dt$$

where we assume for our experiments that the covariance matrix  $\Gamma_S$  has the structure

$$\Gamma_S = \begin{pmatrix} 1 & \rho & \dots & \rho \\ \rho & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho & \dots & \rho & 1 \end{pmatrix} \quad (5.3)$$

with  $\rho \in ]\frac{-1}{d-1}, 1[$ , such that the matrix  $\Gamma_S$  is positive definite. The processes  $B$  and  $\tilde{B}$  are Wiener processes with covariance matrices given by  $\Gamma_S$  and  $I_d$  respectively.

For the sake of simplicity, we decided not to add any extra correlation between the components of  $\tilde{B}$ , hence the choice  $d\langle \tilde{B} \rangle = I_d dt$  and we assume in the following that for all the  $\gamma^i$ 's are equal for  $1 \leq i \leq d$ ,  $\gamma^i = \gamma$ . The correlations between the volatilities are entirely specified by the correlations between the assets. Even though we do not aim at discussing the correlation structure of the multidimensional Heston model, we believe it is important to make precise the underlying correlation structure in the multidimensional model so that the experiments are easily reproducible.

The model can be equivalently written

$$\begin{aligned} dS_t^i &= rS_t^i dt + \sqrt{\sigma_t^i} S_t^i dB_t^i \\ d\sigma_t^i &= \kappa^i (a^i - \sigma_t^i) dt + \nu_t^i \sqrt{\sigma_t^i} dW_t^i \end{aligned}$$

where the processes  $W$  and  $B$  are Wiener processes satisfying

$$d\langle B \rangle_t = \Gamma_S dt; \quad d\langle B, W \rangle_t = \gamma \Gamma_S dt; \quad d\langle W \rangle_t = (\gamma^2 \Gamma_S + (1 - \gamma^2) I_d) dt.$$

The process  $(B, W)$  with values in  $\mathbb{R}^{2d}$  is a Wiener process with covariance matrix

$$\Gamma = \begin{pmatrix} \Gamma_S & \gamma \Gamma_S \\ \gamma \Gamma_S & \gamma^2 \Gamma_S + (1 - \gamma^2) I_d \end{pmatrix}.$$

Hence, the pair of processes  $(B, W)$  can be easily simulated by applying the Cholesky factorization of  $\Gamma$  to a standard Brownian motion with values in  $\mathbb{R}^{2d}$ .

**Basket Option** We consider a basket option as in the local volatility model. Figure 3 looks very much the same as in the case of the local volatility model (see Figure 2). The MLIS estimator always outperforms all the ML estimator by a factor of 3 to 4. Note that for small RMSE, the computational time can go beyond several hours, hence cutting it down by two or three times represents a real improvement.

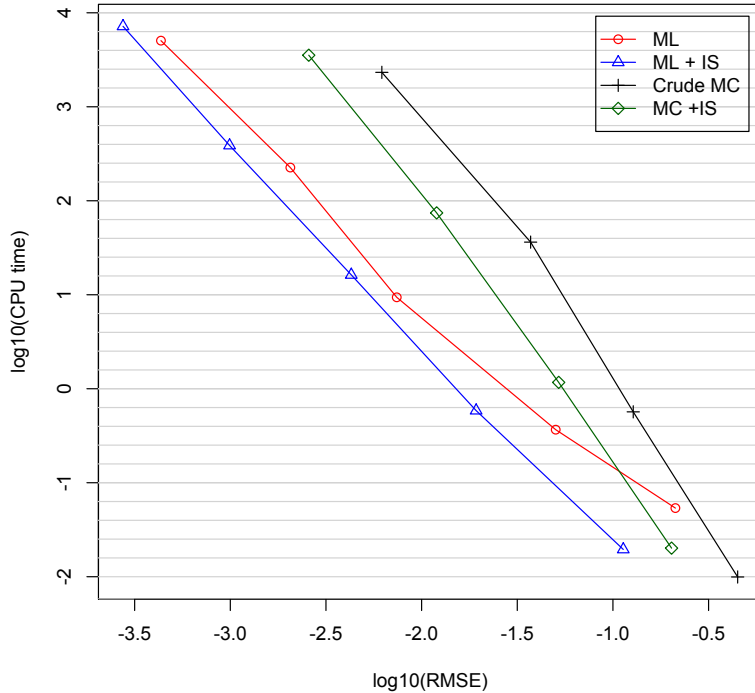


Figure 3:  $\sqrt{MSE}$  vs. CPU time for a best of option in the multidimensional Heston model with  $I = 10$ ,  $r = 0.03$ ,  $T = 1$ ,  $S_0 = 100$ ,  $K = 100$ ,  $\nu = 0.01$ ,  $\kappa = 2$ ,  $a = 0.04$ ,  $\gamma = -0.2$ ,  $\rho = 0.3$  and  $m = 4$ .

**Best of option** We consider options with payoffs of the form  $(\max_{1 \leq i \leq d} S_T^i - K)_+$ . The payoff of this option does obviously not satisfy the assumptions of Theorem 3.2 as the payoff of the “best of” options is not Hölder with  $\alpha \geq 1$ . Nonetheless, the multilevel approach beats the standard Monte Carlo technology by far (see Figure 4). Moreover, coupling importance sampling with the multilevel approach improves the accuracy. For a fixed RMSE, we can expect MLIS to be 3 faster than ML. This example shows the robustness of the method, which performs well whereas the theoretical assumptions are not satisfied.

## 6 Conclusion

We have presented a new estimator making the most of the recent works on multilevel Monte Carlo and on adaptive importance sampling. As expected, this new estimator outperforms the

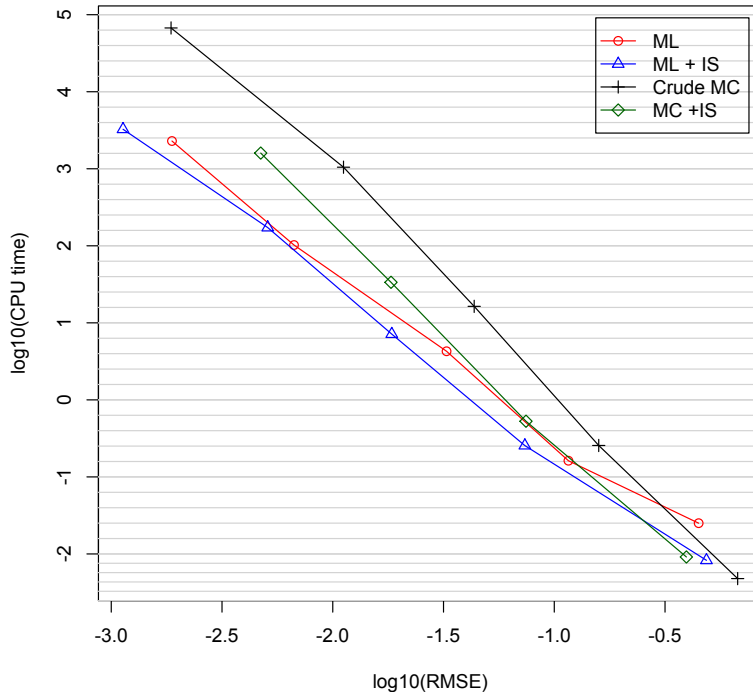


Figure 4:  $\sqrt{MSE}$  vs. CPU time for a best of option in the multidimensional Heston model with  $I = 5$ ,  $r = 0.03$ ,  $T = 1$ ,  $S_0 = 100$ ,  $K = 140$ ,  $\nu = 0.25$ ,  $\kappa = 2$ ,  $a = 0.04$ ,  $\gamma = 0.2$ ,  $\rho = 0.5$  and  $m = 4$ .

standard multilevel Monte Carlo estimator by a great deal. For a fixed accuracy measured in terms the mean squared error, the MLIS estimator is between 3 and 10 times faster than the standard multilevel Monte Carlo estimator. This efficiency of our MLIS approach could still be improved by monitoring the number of samples  $N'_\ell$  to be used to approximate the variance  $v_{\ell, N'_\ell}$  in each level. Actually, we believe that there is no need to compute a too accurate approximation of this variance as a slight decrease in the accuracy of  $\hat{\lambda}_\ell$  would not lead to a serious deterioration of the accuracy of the MLIS estimator but it could help to save a lot of computational time.

## Acknowledgment

We are grateful to the anonymous referees for their valuable comments and suggestions, which helped us greatly improve the paper.

## A Auxiliary lemmas

### A.1 Central limit theorems for martingale arrays

**Theorem A.1** (Central limit theorem for triangular array). *Suppose that  $(\Omega, \mathbb{F}, \mathbb{P})$  is a probability space and that for each  $n$ , we have a filtration  $\mathbb{F}_n = (\mathcal{F}_k^n)_{k \geq 0}$ , a sequence  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$  and a real vector martingale  $Y^n = (Y_k^n)_{k \geq 0}$  adapted to  $\mathbb{F}_n$ . We make the following two assumptions.*

(H-7) *i. There exists a deterministic symmetric positive semi-definite matrix  $\Gamma$ , such that*

$$\langle Y^n \rangle_{k_n} = \sum_{k=1}^{k_n} \mathbb{E} [|Y_k^n - Y_{k-1}^n|^2 | \mathcal{F}_{k-1}^n] \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \Gamma.$$

*ii. There exists a real number  $a > 1$ , such that*

$$\sum_{k=1}^{k_n} \mathbb{E} [|Y_k^n - Y_{k-1}^n|^{2a} | \mathcal{F}_{k-1}^n] \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0.$$

Then

$$Y_{k_n}^n \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Gamma) \quad \text{as } n \rightarrow \infty.$$

### A.2 Asymptotic behavior of the process $(X^{m^\ell} - X^{m^{\ell-1}})_{\ell \geq 0}$

In the following we recall some results around the stable convergence. Let  $Z_n$  be a sequence of random variables with values in a Polish space  $E$ , all defined on the same probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . Let  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$  be an extension of  $(\Omega, \mathcal{F}, \mathbb{P})$ , and let  $Z$  be an  $E$ -valued random variable on the extension. We say that  $(Z_n)$  converges in law to  $Z$  stably and write  $Z_n \xRightarrow{\text{stably}} Z$ , if

$$\mathbb{E}(Uh(Z_n)) \rightarrow \tilde{\mathbb{E}}(Uh(Z))$$

for all  $h : E \rightarrow \mathbb{R}$  bounded continuous and all bounded random variable  $U$  on  $(\Omega, \mathcal{F})$ . According to Section 2 of Jacod [21] and Lemma 2.1 of Jacod and Protter [22], we have the following result

**Lemma A.2.** *Let  $V_n$  and  $V$  be defined on  $(\Omega, \mathcal{F})$  with values in another metric space.*

$$\text{If } V_n \xrightarrow{\mathbb{P}} V, \quad Z_n \xRightarrow{\text{stably}} Z \quad \text{then } (V_n, Z_n) \xRightarrow{\text{stably}} (V, Z).$$

The following result proved by Ben Alaya and Kebaier [4, Theorem 3] is an improvement of Theorem 3.2 of Jacod and Protter [22], for the setting of Multilevel Euler scheme. More precisely, if  $(X_t^{m^\ell})_{t \geq 0}$  denotes the Euler scheme with time step  $m^\ell$ , with  $m, \ell \in \mathbb{N} \setminus \{0, 1\}$ , then we have the following weak convergence in the Skorohod topology.

**Theorem A.3.** *Assume that  $b$  and  $\sigma$  are  $\mathcal{C}^1$  with linear growth then the following result holds.*

$$\text{For all } m \in \mathbb{N} \setminus \{0, 1\}, \quad \sqrt{\frac{m^\ell}{(m-1)T}} (X^{m^\ell} - X^{m^{\ell-1}}) \xRightarrow{\text{stably}} U, \quad \text{as } \ell \rightarrow \infty,$$

with  $(U_t)_{0 \leq t \leq T}$  the  $d$ -dimensional diffusion process solution to (3.6)

**Proposition A.4.** Let  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$  be a  $C^1$  function such that  $\psi \in \mathcal{H}_\alpha$ , for some  $\alpha \geq 1$  and  $\nabla\psi$  has at most polynomial growth. For any real valued random variable  $Y$  defined on  $(\Omega, \mathcal{F})$  such that  $\mathbb{E}[|Y|^{1+\eta}] < \infty$ , for some  $\eta > 0$ , we have, for any  $\delta > 0$

$$\mathbb{E} \left[ \left( \frac{m^\ell}{(m-1)T} \right)^{\delta/2} \left( \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right)^\delta Y \right] \xrightarrow{\ell \rightarrow +\infty} \mathbb{E} \left[ (\nabla\psi(X_T) \cdot U_T)^\delta Y \right].$$

*Proof.* The Taylor expansion applied to the real valued function  $\psi$  yields

$$\begin{aligned} \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) &= \nabla\psi(X_T) \cdot (X_T^{m^\ell} - X_T^{m^{\ell-1}}) \\ &\quad + (X_T^{m^\ell} - X_T) \cdot \varepsilon(X_T^{m^\ell} - X_T) - (X_T^{m^{\ell-1}} - X_T) \cdot \varepsilon(X_T^{m^{\ell-1}} - X_T) \end{aligned}$$

with  $\varepsilon : \mathbb{R}^d \rightarrow \mathbb{R}^d$  satisfying  $\lim_{|x| \rightarrow 0} \varepsilon(x) = 0$ . From Property ((H-1)-ii), we easily get

$$\sqrt{\frac{m^\ell}{(m-1)T}} \left( (X_T^{m^\ell} - X_T) \cdot \varepsilon(X_T^{m^\ell} - X_T) - (X_T^{m^{\ell-1}} - X_T) \cdot \varepsilon(X_T^{m^{\ell-1}} - X_T) \right) \xrightarrow[\ell \rightarrow \infty]{\mathbb{P}} 0.$$

So, we conclude from Lemma A.2 and Theorem A.3 that

$$\sqrt{\frac{m^\ell}{(m-1)T}} \left( \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right) \xRightarrow{\text{stably}} \nabla\psi(X_T) \cdot U_T, \text{ as } \ell \rightarrow \infty.$$

Let  $\eta > \kappa > 0$ . From the assumptions on  $\psi$  together with Property (H-1)-ii, we get

$$\sup_{\ell \geq 0} \mathbb{E} \left[ \left| \left( \frac{m^\ell}{(m-1)T} \right)^{\delta/2} \left( \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right)^\delta Y \right|^{1+\kappa} \right] < \infty,$$

which yields the uniform integrability of the family  $\left( \left( \frac{m^\ell}{(m-1)T} \right)^{\delta/2} \left( \psi(X_T^{m^\ell}) - \psi(X_T^{m^{\ell-1}}) \right)^\delta Y \right)_\ell$ . The conclusion easily follows.  $\square$

## References

- [1] B. Arouna. Adaptative Monte Carlo method, a variance reduction technique. *Monte Carlo Methods Appl.*, 10(1):1–24, 2004. ISSN 0929-9629. doi: 10.1163/156939604323091180. URL <http://dx.doi.org/10.1163/156939604323091180>.
- [2] L. Badouraly Kassim, J. Lelong, and I. Loumrhari. Importance sampling for jump processes and applications to finance. *Journal of Computational Finance (to appear)*, 00(00), 2014. <http://hal.archives-ouvertes.fr/hal-00842362>.
- [3] M. Ben Alaya and A. Kebaier. Multilevel Monte Carlo for Asian options and limit theorems. *Monte Carlo Methods Appl.*, 20(3):181–194, 2014. ISSN 0929-9629. doi: 10.1515/mcma-2013-0025. URL <http://dx.doi.org/10.1515/mcma-2013-0025>.
- [4] M. Ben Alaya and A. Kebaier. Central limit theorem for the multilevel monte carlo euler method. *Ann. Appl. Probab.*, 25(1):211–234, 02 2015.



- [5] M. Ben Alaya, K. Hajji, and A. Kebaier. Importance sampling and statistical Romberg method. *Bernoulli*, 21(4):1947–1983, 2015. ISSN 1350-7265. doi: 10.3150/14-BEJ622. URL <http://dx.doi.org/10.3150/14-BEJ622>.
- [6] M. Ben Alaya, K. Hajji, and A. Kebaier. Improved adaptive Multilevel Monte Carlo and applications to finance. 2016. URL <http://adsabs.harvard.edu/abs/2016arXiv160302959B>.
- [7] H. F. Chen and Y. M. Zhu. Stochastic approximation procedures with randomly varying truncations. *Sci. Sinica Ser. A*, 29(9):914–926, 1986. ISSN 0253-5831.
- [8] H. F. Chen, G. Lei, and A. J. Gao. Convergence and robustness of the Robbins-Monro algorithm truncated at randomly varying bounds. *Stochastic Process. Appl.*, 27(2):217–231, 1988. ISSN 0304-4149. doi: 10.1016/0304-4149(87)90039-1. URL [http://dx.doi.org/10.1016/0304-4149\(87\)90039-1](http://dx.doi.org/10.1016/0304-4149(87)90039-1).
- [9] N. Collier, A.-L. Haji-Ali, F. Nobile, E. von Schwerin, and R. Tempone. A continuation multilevel Monte Carlo algorithm. *BIT*, 55(2):399–432, 2015. ISSN 0006-3835. doi: 10.1007/s10543-014-0511-3. URL <http://dx.doi.org/10.1007/s10543-014-0511-3>.
- [10] J. Creutzig, S. Dereich, T. Müller-Gronbach, and K. Ritter. Infinite-dimensional quadrature and approximation of distributions. *Found. Comput. Math.*, 9(4):391–429, 2009. ISSN 1615-3375. doi: 10.1007/s10208-008-9029-x. URL <http://dx.doi.org/10.1007/s10208-008-9029-x>.
- [11] S. Dereich. Multilevel Monte Carlo algorithms for Lévy-driven SDEs with Gaussian correction. *Ann. Appl. Probab.*, 21(1):283–311, 2011. ISSN 1050-5164. doi: 10.1214/10-AAP695. URL <http://dx.doi.org/10.1214/10-AAP695>.
- [12] D. Duffie and P. Glynn. Efficient Monte Carlo simulation of security prices. *Ann. Appl. Probab.*, 5(4):897–905, 1995.
- [13] M. B. Giles. Multilevel Monte Carlo path simulation. *Oper. Res.*, 56(3):607–617, 2008. ISSN 0030-364X. doi: 10.1287/opre.1070.0496. URL <http://dx.doi.org/10.1287/opre.1070.0496>.
- [14] M. B. Giles. Improved multilevel Monte Carlo convergence using the Milstein scheme. In *Monte Carlo and quasi-Monte Carlo methods 2006*, pages 343–358. Springer, Berlin, 2008. doi: 10.1007/978-3-540-74496-2\_20. URL [http://dx.doi.org/10.1007/978-3-540-74496-2\\_20](http://dx.doi.org/10.1007/978-3-540-74496-2_20).
- [15] M. B. Giles and L. Szpruch. Antithetic multilevel Monte Carlo estimation for multi-dimensional SDEs without Lévy area simulation. *Ann. Appl. Probab.*, 24(4):1585–1620, 2014. ISSN 1050-5164. doi: 10.1214/13-AAP957. URL <http://dx.doi.org/10.1214/13-AAP957>.
- [16] M. B. Giles, D. J. Higham, and X. Mao. Analysing multi-level Monte Carlo for options with non-globally Lipschitz payoff. *Finance Stoch.*, 13(3):403–413, 2009. ISSN 0949-2984. doi: 10.1007/s00780-009-0092-1. URL <http://dx.doi.org/10.1007/s00780-009-0092-1>.

- [17] K. Hajji. *Accélération de la méthode de Monte Carlo pour des processus de diffusions et applications en Finance*. PhD thesis, Université Paris 13, 2014.
- [18] S. Heinrich. Monte Carlo complexity of global solution of integral equations. *J. Complexity*, 14(2):151–175, 1998. ISSN 0885-064X. doi: 10.1006/jcom.1998.0471. URL <http://dx.doi.org/10.1006/jcom.1998.0471>.
- [19] S. Heinrich. Multilevel monte carlo methods. *Lecture Notes in Computer Science*, Springer-Verlag, 2179(1):58–67, 2001.
- [20] S. Heinrich and E. Sindambiwe. Monte carlo complexity of parametric integration. *J. Complexity*, 15(3):317–341, 1999. ISSN 0885-064X. doi: 10.1006/jcom.1999.0508. URL <http://dx.doi.org/10.1006/jcom.1999.0508>. Dagstuhl Seminar on Algorithms and Complexity for Continuous Problems (1998).
- [21] J. Jacod. On continuous conditional Gaussian martingales and stable convergence in law. In *Séminaire de Probabilités, XXXI*, volume 1655 of *Lecture Notes in Math.*, pages 232–246. Springer, Berlin, 1997.
- [22] J. Jacod and P. Protter. Asymptotic error distributions for the Euler method for stochastic differential equations. *Ann. Probab.*, 26(1):267–307, 1998. ISSN 0091-1798.
- [23] B. Jourdain and J. Lelong. Robust Adaptive Importance Sampling for Normal Random Vectors. *Ann. Appl. Probab.*, 19(5):1687–1718, 2009. doi: 10.1214/09-AAP595. URL <http://arxiv.org/pdf/0811.1496v2>.
- [24] A. Kebaier. Statistical Romberg extrapolation: a new variance reduction method and applications to option pricing. *Ann. Appl. Probab.*, 15(4):2681–2705, 2005. doi: 10.1214/105051605000000511.
- [25] B. Lapeyre and J. Lelong. A framework for adaptive Monte Carlo procedures. *Monte Carlo Methods Appl.*, 17(1):77–98, 2011. ISSN 0929-9629. doi: 10.1515/MCMA.2011.002. URL <http://dx.doi.org/10.1515/MCMA.2011.002>.
- [26] J. Lelong. Almost sure convergence for randomly truncated stochastic algorithms under verifiable conditions. *Statist. Probab. Lett.*, 78(16):2632–2636, 2008. ISSN 0167-7152. doi: 10.1016/j.spl.2008.02.034. URL <http://dx.doi.org/10.1016/j.spl.2008.02.034>.
- [27] V. Lemaire and G. Pagès. Multilevel Richardson-Romberg extrapolation. *Bernoulli*, 2016 (to appear). URL <http://arxiv.org/abs/1401.1177v3>.
- [28] D. Revuz. *Probabilités*. Hermann, 1997.
- [29] R. Y. Rubinstein and A. Shapiro. *Discrete event systems*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. John Wiley & Sons Ltd., Chichester, 1993. Sensitivity analysis and stochastic optimization by the score function method.