



**HAL**  
open science

# Large deviations principle for the Adaptive Multilevel Splitting Algorithm in an idealized setting

Charles-Edouard Bréhier

► **To cite this version:**

Charles-Edouard Bréhier. Large deviations principle for the Adaptive Multilevel Splitting Algorithm in an idealized setting. *ALEA : Latin American Journal of Probability and Mathematical Statistics*, 2015. hal-01118745

**HAL Id: hal-01118745**

**<https://hal.science/hal-01118745v1>**

Submitted on 20 Feb 2015

**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.

# Large deviations principle for the Adaptive Multilevel Splitting Algorithm in an idealized setting

Charles-Edouard Bréhier \*

## Abstract

The Adaptive Multilevel Splitting (AMS) algorithm is a powerful and versatile method for the simulation of rare events. It is based on an interacting (via a mutation-selection procedure) system of replicas, and depends on two integer parameters:  $n \in \mathbb{N}^*$  the size of the system and the number  $k \in \{1, \dots, n-1\}$  of the replicas that are eliminated and resampled at each iteration.

In an idealized setting, we analyze the performance of this algorithm in terms of a Large Deviations Principle when  $n$  goes to infinity, for the estimation of the (small) probability  $\mathbb{P}(X > a)$  where  $a$  is a given threshold and  $X$  is real-valued random variable. The proof uses the technique introduced in [BLR15]: in order to study the log-Laplace transform, we rely on an auxiliary functional equation.

Such Large Deviations Principle results are potentially useful to study the algorithm beyond the idealized setting, in particular to compute rare transitions probabilities for complex high-dimensional stochastic processes.

*Keywords:* Monte-Carlo simulation, rare events, multilevel splitting, large deviations  
*MSC:* 65C05; 65C35; 62G30; 60F10

## 1 Introduction

In many problems from engineering, biology, chemistry, physics or finance, rare events are often critical and have a huge impact on the phenomena which are studied. From a general mathematical perspective, we may consider the following situation: let  $(X_t)_{t \in \mathbb{T}}$ , where  $\mathbb{T} = \mathbb{N}$  or  $\mathbb{R}$ , be a (discrete or continuous in time) stochastic process, taking values in  $\mathbb{R}^d$ . Assume that  $A, B \subset \mathbb{R}^d$  are two *metastable regions*: starting from a neighborhood of  $A$  (resp. of  $B$ ), the probability that the process reaches  $B$  (resp.  $A$ ) before hitting  $A$  (resp.  $B$ ) is very small (typically, less than  $10^{-10}$ ). As a consequence, a direct numerical Monte-Carlo with an ensemble of size  $N$  does not provide significant results when  $N$  is reasonably large (typically, less than  $10^{10}$ ) in real-life applications.

Even if theoretical asymptotic expansions on quantities of interest are available - such as the Kramers-Arrhenius law given for instance by the Freidlin-Wentzell Large Deviations Theory or Potential Theory for the exit problem of a diffusion process in the small noise regime - in practice their explicit computation is not possible (for instance when the dimension is large) and numerical simulations are unavoidable.

---

\*Institut de Mathématiques, Université de Neuchâtel, Rue Emile Argand 11, CH-2000 Neuchâtel. e-mail: charles-edouard.brehier@unine.ch

It is thus essential to propose efficient and general methods, and to rigorously study their consistency and efficiency properties. Two main families of methods have been introduced in the 1950's and studied extensively since then, in order to improve the Monte-Carlo simulation algorithms, in particular for rare events: importance sampling and importance splitting (see for instance [AG07], [RT09] for general reviews of these methods and [KH51] for the historical introduction of importance splitting). The main difference between these two methods is the following: the first one is *intrusive*, meaning that the dynamics of the stochastic process (more generally, the distribution of the random variable of interest) is modified so that the probability that the event of interest increases and in a Monte-Carlo simulation it is realized more often, while the second is not intrusive and can thus be used more directly for complex problems. Instead, for importance splitting strategies, the state space is decomposed as a nested sequence of regions which are visited sequentially and more easily by an interacting system of replicas.

In this paper, we focus on an importance splitting strategy which is known as the Multilevel Splitting approach and describe it in the following setting. Let  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a given function and assume we want to estimate the probability  $p = \mathbb{P}(X > a)$  that a real-valued random variable  $X = h(Y)$  (where  $Y$  is a  $\mathbb{R}^d$ -valued random variable) belongs to  $(a, +\infty)$  for a given threshold  $a \in \mathbb{R}$ . This situation is not restrictive for many applications; indeed, we may take  $X = \mathbb{1}_{\tau_B < \tau_A}$  and any  $a \in (0, 1)$  in the situation described above, where  $\tau_A$  and  $\tau_B$  are the hitting times of  $A$  and  $B$  by the process  $X$ . A key assumption on the distribution of  $X$  is the following: we assume that the cumulative distribution function  $F$  of  $X$  - *i.e.*  $F(x) = \mathbb{P}(X \leq x)$  for any  $x \in \mathbb{R}$  - is continuous; for convenience, we also assume that  $F(0) = 0$  - *i.e.*  $X > 0$  almost surely.

The multilevel splitting approach (see [KH51], [GHSZ99], [CDMFG12] for instance) is based on the following decomposition of  $p$  as a telescopic product of conditional probabilities:

$$p = \mathbb{P}(X > a) = \prod_{i=1}^N \mathbb{P}(X > a_i | X > a_{i-1}), \quad (1)$$

where  $a_0 = 0 < a_1 < \dots < a_N = a$  is a sequence of non-decreasing *i.e.* levels. In other words, the realization of the event  $\{X > a\}$  is split into the realizations of the  $N$  events  $\{X > a_i\}$  conditional on  $\{X > a_{i-1}\}$ ; each event has a larger probability than the initial one and is thus much easier to realize. Then each of the conditional probabilities is estimated separately, for instance with independent Monte-Carlo simulations, or using a Sequential Monte-Carlo technique with a splitting of successful trajectories. This approach have been studied with different viewpoints and variants under different names in the literature - nested sampling [Ski06], [Ski07], subset simulation [AB01], RESTART, [VAVA91], [VAVA94].

For future reference, we introduce the following (unbiased) estimator of  $p$  given by the multilevel splitting approach with  $N$  levels and  $n$  replicas:

$$\hat{p}_n^N = \prod_{i=1}^N \frac{1}{n} \sum_{m=1}^n \mathbb{1}_{X_m^{(i)} > a_i}, \quad (2)$$

where the random variables  $(X_m^{(i)})_{1 \leq m \leq N, 1 \leq i \leq N}$  are independent and the distribution of  $X_m^{(i)}$  is  $\mathcal{L}(X | X > a_{i-1})$ . Thus  $\hat{p}_n^N$  is a product of  $N$  independent Monte-Carlo estimators of the conditional probabilities in (1).

The efficiency of the algorithm depends crucially on the choice of the sequence of levels  $(a_i)_{1 \leq i \leq N}$ : for a fixed size  $N$ , the variance of the estimator is minimized when the conditional probabilities are equal (to  $p^{1/N}$ ); moreover the associated variance converges (to  $-p^2 \log(p)/n$ ) when  $N$  goes to infinity - see for instance [CDMFG12] for more details.

To get a more flexible algorithms, a possible approach is to compute levels adaptively, as proposed in [CG07], and studied extensively in the last years, see for instance [BLR15], [BGT14], [CG14], [GHML11], [Sim14], [Wal14]. It is essential to check that these adaptive versions still give reliable results, and to prove they do it efficiently.

More precisely, we consider the Algorithm 2.2 defined below, which depends on two parameters  $n$  and  $k$ , with the condition  $1 \leq k \leq n - 1$ . We let evolve a system of  $n$  interacting replicas, and at each iteration a selection-mutation procedure leads to resample the system as follows: we compute the  $k$ -th order statistic  $Z$  - which corresponds to the so-called level at the given iteration - of the system and eliminate the  $k$  replicas with values less than  $Z$ ; they are then resampled using the conditional distribution  $\mathcal{L}(X|X > Z)$  of  $X$  conditional on  $\{X > Z\}$ . The algorithm stops when  $Z \geq a$ , and we define an estimator  $\hat{p}^{n,k}$  depending on the number of iterations and of the terminal configuration of the system of replicas, see (5). In practice, we require to be able to sample according to the conditional distribution  $\mathcal{L}(X|X > z)$  for any value of  $z$ : this is part of the idealized setting assumption; even if it is rarely satisfied in real-life applications, the study of the algorithm in that setting is already challenging and yields very interesting results, that can usually be generalized beyond this simplified case at the price of a much more intricate analysis.

Let us recall a few fundamental results. In [GHML11] (see also [Sim14], [Wal14]), it was proved that for any value of  $n \geq 2$  then  $\hat{p}^{n,1}$  is an unbiased estimator of  $p$  - meaning that  $\mathbb{E}[\hat{p}^{n,1}] = p$ . This result was extend to general  $1 \leq k \leq n - 1$  in [BLR15]. Efficiency properties have been studied with the proof of Central Limit Theorems in two different kinds of regimes: either  $k$  is fixed and  $n \rightarrow +\infty$  (see [BGT14] as well as [GHML11] and [Sim14] when  $k = 1$ ), or both  $k$  and  $n$  go to infinity, in such a way that  $k/n$  converges to  $\alpha \in (0, 1)$  - which gives a fixed proportion of resampled replicas at each iteration, see [CG07] and the more recent work [CG14] in a very general framework.

The efficiency is ensured by the observation that the asymptotic variance is the same for both the adaptive and the non-adaptive versions. Moreover, it is much smaller than when using a crude Monte-Carlo estimator, *i.e.* the empirical average

$$\bar{p}_n = \frac{1}{n} \sum_{m=1}^n \mathbb{1}_{X_m > a}, \quad (3)$$

where the random variables  $(X_m)_{1 \leq m \leq n}$  are independent and identically distributed, with distribution  $\mathcal{L}(X)$ .

In this paper, we prove a similar result with a different criterion, which seems to be original compared with existing literature: we prove a Large Deviations Principle principle for the distribution of the estimator  $\hat{p}^{n,k}$  given by the adaptive algorithm when  $k$  is fixed and  $n \rightarrow +\infty$ . Our main result is Theorem 3.1, which in particular yields for any given  $\epsilon > 0$

$$\frac{1}{n} \log \left( \mathbb{P}(|\hat{p}^{n,k} - p| \geq \epsilon) \right) \xrightarrow{n \rightarrow +\infty} -\min(I(p + \epsilon), I(p - \epsilon)) < 0.$$

The rate function  $I$  - see (7) - obtained in Theorem 3.1 does not depend on  $k$ . We then compare this rate function with  $\mathcal{I}$  - see (23) - the rate function obtained for a crude Monte-Carlo estimator  $\bar{p}_n$  given by (3) (thanks to Cramer Theorem, see [DZ10]) and show that for

any  $y \in (0, 1) \setminus p$  we have  $I(y) > \mathcal{I}(y)$  - we have  $\mathcal{I}(p) = I(p) = 0$ , and  $\mathcal{I}(y) = I(y) = +\infty$  if  $y \notin (0, 1)$  - and thus

$$\frac{\mathbb{P}(\hat{p}^{n,k} - p > \epsilon)}{\mathbb{P}(\bar{p}_n - p > \epsilon)} \xrightarrow{n \rightarrow +\infty} 0.$$

In other words, for large  $n$ , the probability that  $\hat{p}^{n,k}$  deviates from  $p$  from above (and similarly from below) with threshold  $\epsilon > 0$  decreases exponentially fast, at a faster rate than for  $\bar{p}_n$ .

Moreover, we prove that the non-adaptive, fixed-levels estimator  $\hat{p}_n^N$  satisfies a Large Deviations Principle when  $n \rightarrow +\infty$  with rate function  $\mathcal{I}_N$  for a fixed number of levels  $N$  and when the levels are chosen in an optimal way, namely such that  $\mathbb{P}(X > a_i | X > a_{i-1}) = p^{1/N}$  does not depend on  $i$ . We then show that  $\lim_{N \rightarrow +\infty} \mathcal{I}_N(y) \leq I(y)$  for any  $y \in \mathbb{R}$ : this inequality is sufficient to prove that asymptotically the adaptive algorithm performs (at least) as well as the non-adaptive version in this setting, in terms of Large Deviations.

The proof of Theorem 3.1 relies on the technique introduced in [BLR15]. First, we restrict the study of the properties of the algorithm to the case when  $X$  is exponentially distributed with parameter 1 (this key remark was introduced first in [GHML11] and used also in [Sim14], [Wal14]). Instead of working on  $\hat{p}^{n,k}$  directly, we focus on its logarithm  $\log(\hat{p}^{n,k})$ , and prove that when considering the algorithm as depending on an initial condition  $x$ , the Laplace transform of the latter is solution of a functional (integral) equation (with respect to the  $x$  variable) - thanks to a decomposition of the realizations of the algorithm according to the value of the first level. To study the equation in the asymptotic regime considered in this paper, we then derive a linear ordinary differential equation of order  $k$  and perform an asymptotic expansion. Note that we do not give all details for the derivation of the differential equations and the basic properties of its coefficients; for some points we refer the reader to [BLR15] where all the arguments are proved with details and here we mainly focus on the proof of the new asymptotic results as well as on the interpretation of the Large Deviations Principle for our purpose.

It seems that studying the performance of multilevel splitting algorithms via Large Deviations Principle is an original approach, which can complement the more classical studies which are all based on Central Limit Theorems. In this paper, we proved a result in a specific regime ( $k$  is fixed,  $n \rightarrow +\infty$ ) in the idealized setting. To go further, it would be interesting to look at other regimes ( $k, n \rightarrow +\infty$  with  $k/n \rightarrow \alpha \in (0, 1)$ ) and to go beyond the idealized setting. This will be the subject of future investigation.

The paper is organized as follows. In Section 2, we introduce our main assumptions (Section 2.1), describe the Adaptive Multilevel Splitting algorithm (Section 2.2) and recall several of its fundamental properties used in the sequel of the article (Section 2.3). The main result of this paper is given in Section 3: it is the Large Deviations Principle for the estimator of the probability given by the AMS estimator, see Theorem 3.1. An important auxiliary result is stated in Section 4, and proofs are carried over in Section 5 - some technical estimates being proved in Section 7. We compare the performance in terms of the Large Deviations Principle of the AMS algorithm with two other methods in Section 6: a crude Monte-Carlo method and a fixed-level splitting method. Finally, we give some concluding remarks and perspectives in Section 8.

## 2 Description of the Adaptive Multilevel Splitting algorithm

### 2.1 Assumptions

Let  $X$  be some real random variable. For simplicity, we assume that  $X > 0$  almost surely.

We want to estimate the probability  $p = \mathbb{P}(X > a)$ , where  $a > 0$  is some threshold. When  $a$  goes to  $+\infty$ ,  $p$  goes to 0 and we have to estimate the probability of a rare event.

We make a fundamental assumption on the distribution of  $X$ .

**Assumption 2.1.** *Let  $F$  denote the cumulative distribution function of  $X$ : we assume that  $F$  is continuous.*

More generally, for both theoretical and practical purpose, we introduce for  $0 \leq x \leq a$  the conditional probability

$$P(x) = \mathbb{P}(X > a | X > x); \quad (4)$$

we also denote by  $\mathcal{L}(X|X > x)$  the associated conditional distribution, and  $F(\cdot; x)$  its cumulative distribution function: for any  $y > x$  we have  $F(y; x) = \frac{F(y) - F(x)}{1 - F(x)}$  whenever  $F(x) < 1$ .

We notice two important equalities:  $P(a) = 1$ , and the estimated probability is  $p = P(0)$ ; in fact, the distribution of  $X$  is equal to  $\mathcal{L}(X|X > 0)$ .

The idealized setting refers to the following assumptions:

- Assumption 2.1 is satisfied (*theoretical condition*);
- it is possible to sample according to the conditional distribution  $\mathcal{L}(X|X > x)$  for any  $x \in [0, a)$  (*practical condition*).

In view of a practical implementation of the algorithm, the second condition is probably the most restrictive. One may rely on some approximation of the conditional distribution  $\mathcal{L}(X|X > x)$  thanks to a Metropolis-Hastings algorithm: in that case (see [CG14] for instance), the analysis we develop here does not apply, but gives an interesting insight for the behavior in the case of a large number of steps in the Metropolis-Hastings auxiliary scheme (rigorously, we treat the case of an infinite number of steps).

### 2.2 The algorithm

We now present the Adaptive Multilevel Splitting algorithm, under the assumptions of Section 2.1 above.

The algorithm depends on two parameters:

- the number of replicas  $n$ ;
- the number  $k \in \{1, \dots, n - 1\}$  of replicas that are resampled at each iteration.

The other necessary parameters are the initial condition  $x$  and the stopping threshold  $a$ : the aim is to estimate the conditional probability  $P(x)$  introduced in (4). For future reference, we denote by  $\text{AMS}(n, k; a, x)$  the algorithm.

The dependence with respect to  $x$  allows us below to state fundamental functional equations on useful observables of the estimator computed at the end of the iterations of the algorithm, as a function of  $x$ . In practice, we are interested in the case  $x = 0$ ; in this situation, the algorithm is denoted by  $\text{AMS}(n, k; a)$ .

Before we detail the algorithm, we introduce important notation. First, when we consider a random variable  $X_i^j$ , the subscript  $i$  denotes the index in  $\{1, \dots, n\}$  of a replica, while the superscript  $j$  denotes the iteration of the algorithm.

Moreover, we use the following notation for order statistics. Let  $Y = (Y_1, \dots, Y_n)$  be independent and identically distributed (i.i.d.) real valued random variables with continuous cumulative distribution function; then there exists almost surely a unique (random) permutation  $\sigma$  of  $\{1, \dots, n\}$  such that  $Y_{\sigma(1)} < \dots < Y_{\sigma(n)}$ . For any  $k \in \{1, \dots, n\}$ , we then denote by  $Y_{(k)} = Y_{\sigma(k)}$  the so-called  $k$ -th order statistic of the sample  $Y$ . Sometimes we need to specify the size of the sample of which we consider the order statistics: we then use the notation  $Y_{(k,n)}$ .

We are now in position to write the AMS( $n, k; a, x$ ) algorithm.

**Algorithm 2.2** (Adaptive Multilevel Splitting, AMS( $n, k; a, x$ )).

**Initialization:** Set the initial level  $Z^0 = x$ .

Sample  $n$  i.i.d. realizations  $X_1^0, \dots, X_n^0$ , with distribution  $\mathcal{L}(X|X > x)$ .

Define  $Z^1 = X_{(k)}^0$ , the  $k$ -th order statistics of the sample  $X^0 = (X_1^0, \dots, X_n^0)$ , and  $\sigma^1$  the (a.s.) unique associated permutation:  $X_{\sigma^1(1)}^0 < \dots < X_{\sigma^1(n)}^0$ .

Set  $j = 1$ .

**Iterations (on  $j \geq 1$ ):** While  $Z^j < a$ :

- Conditional on  $Z^j$ , sample  $k$  new independent random variables  $(Y_1^j, \dots, Y_k^j)$ , according to the law  $\mathcal{L}(X|X > Z^j)$ .
- Set

$$X_i^j = \begin{cases} Y_{(\sigma^j)^{-1}(i)}^j & \text{if } (\sigma^j)^{-1}(i) \leq k \\ X_i^{j-1} & \text{if } (\sigma^j)^{-1}(i) > k. \end{cases}$$

In other words, we resample exactly  $k$  out of the  $n$  replicas, namely those with index  $i$  such that  $X_i^{j-1} \leq Z^j$ , i.e. such that  $i \in \{\sigma^j(1), \dots, \sigma^j(k)\}$  (which is equivalent to  $(\sigma^j)^{-1}(i) \leq k$ ). They are resampled according to the conditional distribution  $\mathcal{L}(X|X > Z^j)$ . The other replicas are not modified.

- Define  $Z^{j+1} = X_{(k)}^j$ , the  $k$ -th order statistics of the sample  $X^j = (X_1^j, \dots, X_n^j)$ , and  $\sigma^{j+1}$  the (a.s.) unique associated permutation:  $X_{\sigma^{j+1}(1)}^j < \dots < X_{\sigma^{j+1}(n)}^j$ .
- Finally increment  $j \leftarrow j + 1$ .

**End of the algorithm:** Define  $J^{n,k}(x) = j - 1$  as the (random) number of iterations. Notice that  $J^{n,k}(x)$  is such that  $Z^{J^{n,k}(x)} < a$  and  $Z^{J^{n,k}(x)+1} \geq a$ .

Notice for instance that  $J^{n,k}(x) = 0$  if and only if  $Z^1 > a$ : we mean that in this case the algorithm has required 0 iteration, since the stopping condition at the beginning of the loop (on  $j$ ) is satisfied without entering into the loop.

The estimator of the probability  $P(x)$  is defined by

$$\hat{p}^{n,k}(x) = C^{n,k}(x) \left(1 - \frac{k}{n}\right)^{J^{n,k}(x)}, \quad (5)$$

with

$$C^{n,k}(x) = \frac{1}{n} \text{Card} \left\{ i; X_i^{J^{n,k}(x)} \geq a \right\}. \quad (6)$$

The interpretation of the factor  $C^{n,k}(x)$  is the following: it is the proportion of the replicas  $X_i^{J^{n,k}(x)}$  which satisfy  $X_i^j \geq a$ : since  $X_{(k)}^{J^{n,k}(x)} = Z^{J^{n,k}(x)+1} \geq a$ , we have  $C^{n,k}(x) \geq \frac{n-k+1}{n}$ . Notice that  $C^{n,1}(x) = 1$ .

When  $x = 0$ , to simplify notations we set  $\hat{p}^{n,k} = \hat{p}^{n,k}(0)$ .

## 2.3 Properties of the AMS Algorithm 2.2

### Well-posedness

We first recall some important results on the well-posedness of the algorithm. For more detailed statements and complete proofs, see Section 3.2 in [BLR15], in particular Proposition 3.2 there.

First, at each iteration  $j$  of the algorithm, conditional on the level  $Z^j$ , the resampling produces a family of  $n$  random variables  $(X_i^j)_{1 \leq i \leq n}$  which are independent and identically distributed, with distribution  $\mathcal{L}(X|X > Z^j)$ . By Assumption 2.1, conditional on  $Z^j$  the latter conditional distribution also admits a continuous cumulative distribution function  $F(\cdot; Z^j)$ ; as a consequence, almost surely the permutation  $\sigma^{j+1}$  is unique, and the level  $Z^{j+1}$  is well-defined.

Moreover, if we assume that  $P(x) > 0$ , almost surely the algorithm stops after a finite number of steps, for any values of  $k$  and  $n$  such that  $1 \leq k \leq n - 1$ : the random variable  $J^{n,k}(x)$  almost surely takes values in  $\mathbb{N}$ , and the estimator  $\hat{p}^{n,k}(x)$  is well-defined and takes values in  $(0, 1]$ .

### Reduction to the exponential case

We now state properties that are essential for our theoretical study of the algorithm below.

One of the main tools in [BLR15] and [BGT14], which was also used in [GHML11] in the case  $k = 1$ , is the restriction to the case where the random variables are exponentially distributed. More precisely, assume that  $P(x) > 0$ , and denote by  $\mathcal{E}(1)$  the exponential distribution with mean 1. Then in distribution the algorithm  $\text{AMS}(n, k; a)$  is equal to the algorithm  $\text{AMS}_{\text{expo}}(n, k; -\log(p))$  in which we assume that the distribution is  $\mathcal{E}(1)$ ; a similar result holds for  $\text{AMS}(n, k; a, x)$  when  $x \in [0, a)$ . In particular, the associated estimators are equally distributed. The main argument is the well-known equality of distribution  $F(X) = U$  where  $U$  is uniformly distributed on  $(0, 1)$ .

In the sequel, we state in Section 3 our results in the general setting - *i.e.* for  $\text{AMS}(n, k; a)$ , with the probability  $p$  and the estimator  $\hat{p}^{n,k}$  - but in the remaining of the paper we give proofs in the exponential case, namely for  $\text{AMS}_{\text{expo}}(n, k; a_{\text{expo}}, x)$  with  $a_{\text{expo}} = -\log(p)$ , and we omit the reference to the exponential case to simplify the notation. Whether we consider the general or the exponential case will be clear from the context.

## 3 The Large Deviations Principle result for the AMS algorithm

The main result of this article is the following Theorem 3.1, which states a Large Deviations Principle (in the sense of [DZ10]) for the distribution  $\mu^{n,k} = \mathcal{L}(\hat{p}^{n,k})$  of  $\hat{p}^{n,k}$  for fixed probability



$p > 0$  and  $k \in \mathbb{N}^*$ , in the limit  $n \rightarrow +\infty$ .

**Theorem 3.1.** *Assume that  $p \in (0, 1)$  and  $k \in \mathbb{N}^*$  are fixed. Then the sequence  $(\mu^{n,k})_{n \in \mathbb{N}, n > k}$  of distributions of the estimator  $\hat{p}^{n,k}$  of  $p$  obtained by the AMS( $n, k; a$ ) algorithm satisfies a Large Deviations Principle with the rate function  $I$  defined by*

$$I(y) = \begin{cases} +\infty & \text{if } y \notin (0, 1) \\ \log(y) \log\left(\frac{\log(p)}{\log(y)}\right) + \log\left(\frac{y}{p}\right) & \text{if } y \in (0, 1). \end{cases} \quad (7)$$

We observe that the rate function does not depend on  $k$ .

Notice that the statement above is restricted to  $p \in (0, 1)$ . Indeed, when  $p = 1$ , we have almost surely  $\hat{p}^{n,k} = 1$  (the algorithm stops after 0 iteration). Moreover, we always estimate the probability of events which have a positive probability (otherwise the algorithm does not stop after a finite number of iterations).

The following Proposition describes some properties of the rate function  $I$ .

**Proposition 3.2.** *The rate function  $I$  is of class  $C^\infty$  on its domain  $(0, 1)$ .*

*Moreover,  $p$  is the unique minimizer of  $I$ : we have  $I(p) = I'(p) = 0$ ,  $I''(p) = \frac{1}{-p^2 \log(p)} > 0$ .*

*Finally, for any  $y \in (0, 1) \setminus \{p\}$  we have  $I(y) > 0$ ;  $I$  is decreasing on  $(0, p)$  and is increasing on  $(p, 1)$ .*

*Proof.* Straightforward computations yield that for  $y \in (0, 1)$  we have

$$\begin{aligned} \frac{dI(y)}{dy} &= \frac{\log(\log(p)) - \log(\log(y))}{y}, \\ \frac{d^2I(y)}{dy^2} &= -\frac{\log(\log(p)) - \log(\log(y))}{y^2} - \frac{1}{y^2 \log(y)}. \end{aligned}$$

□

Let  $\epsilon \in (0, \max(p, 1 - p))$ ; then from Theorem 3.1 we have when  $n \rightarrow +\infty$

$$\frac{1}{n} \log\left(\mathbb{P}\left(|\hat{p}^{n,k} - p| \geq \epsilon\right)\right) \xrightarrow{n \rightarrow +\infty} -\min(I(p + \epsilon), I(p - \epsilon)) < 0. \quad (8)$$

Applying the Borel-Cantelli Lemma, we get the almost sure convergence  $\hat{p}^{n,k} \rightarrow p$ .

**Remark 3.3.** *The almost sure limit is consistent with the unbiasedness result ( $\mathbb{E}[\hat{p}^{n,k}] = p$ ) from [BLR15]. There we were only able to prove the convergence in probability of  $\hat{p}^{n,k}$  to  $p$ .*

*Notice also that in [BGT14] we proved a Central Limit Theorem:*

$$\sqrt{n}(\hat{p}^{n,k} - p) \rightarrow \mathcal{N}(0, -p^2 \log(p)).$$

*The asymptotic variance is given by  $I''(p)$ .*

We conclude this section with a result showing that the choice of the regime  $p$  (and  $k$ ) fixed and  $n \rightarrow +\infty$  is crucial to get Theorem 3.1. Indeed, set  $k = 1$ , and for a given  $\sigma > 0$  assume that  $n$  and  $p$  are related through the following formula:  $-\log(p) = \sigma^2 n$ . Then  $\frac{\hat{p}^{n,k}}{p}$  converges (in law) to a log-normal distribution, as stated in the following proposition.

**Proposition 3.4.** *If  $-\log(p) = \sigma^2 n$ , we have the convergence in distribution*

$$\lim_{n \rightarrow \infty} \frac{\hat{p}^{n,1}}{p} = \exp(\sigma Z - \sigma^2/2),$$

where  $Z \sim \mathcal{N}(0, 1)$ .

The proof is postponed to Section 5.1, since it uses the same arguments as the proof of Theorem 3.1 in the case  $k = 1$ .

Let  $\epsilon > 0$ . Then (compare with (8) with  $\epsilon p$  instead of  $\epsilon$ )

$$\mathbb{P}\left(\left|\frac{\hat{p}^{n,1}}{p} - 1\right| \geq \epsilon\right) \xrightarrow[n = -\frac{\log(p)}{\sigma^2} \rightarrow +\infty]{} \mathbb{P}_{Z \sim \mathcal{N}(0,1)}(|\exp(\sigma Z - \sigma^2/2) - 1| \geq \epsilon) > 0,$$

where the limit is positive, while owing to (8) when  $p$  fixed,  $\mathbb{P}(|\frac{\hat{p}^{n,1}}{p} - 1| \geq \epsilon)$  converges to 0 exponentially fast when  $n \rightarrow +\infty$ .

## 4 Strategy of the proof

To prove Theorem 3.1, we in fact first prove a Large Deviations Principle for  $\tilde{\mu}^{n,k} = \mathcal{L}(\log(\hat{p}^{n,k}))$ , with rate function  $J$  given below.

**Proposition 4.1.** *Assume that  $p \in (0, 1)$  and  $k \in \mathbb{N}^*$  are fixed. Then the sequence  $(\tilde{\mu}^{n,k})_{n \in \mathbb{N}, n > k}$  of distributions of  $\log(\hat{p}^{n,k})$  obtained by the AMS( $n, k; a$ ) algorithm satisfies a Large Deviations Principle with the rate function  $J$  defined by*

$$J(z) = \begin{cases} +\infty & \text{if } z \geq 0 \\ z - \log(p) - z \log(\frac{z}{\log(p)}) & \text{if } z < 0. \end{cases} \quad (9)$$

Then Theorem 3.1 immediately follows from Proposition 4.1 and the application of the contraction principle (see [DZ10], Theorem 4.2.1): we have  $\hat{p}^{n,k} = \exp(\log(\hat{p}^{n,k}))$ , and we obtain the rate function with the identity  $I(y) = J(\log(y))$ .

The proof of Proposition 4.1 relies on the use of the Gärtner-Ellis Theorem (see Theorem 2.3.6 in [DZ10]) and the asymptotic analysis when  $n \rightarrow +\infty$  of the log-Laplace transform of  $\tilde{\mu}^{n,k}$ .

**Proposition 4.2.** *Set for any  $1 \leq k \leq n - 1$  and any  $\lambda \in \mathbb{R}$*

$$\Lambda_{n,k}(\lambda) = \log\left(\mathbb{E}\left[\exp(\lambda \log(\hat{p}^{n,k}))\right]\right). \quad (10)$$

*Then for any fixed  $k \in \mathbb{N}^*$  and any  $\lambda \in \mathbb{R}$  we have the convergence*

$$\frac{1}{n} \Lambda_{n,k}(n\lambda) \rightarrow \Lambda(\lambda) = -\log(p)(\exp(-\lambda) - 1). \quad (11)$$

*The Fenchel-Legendre transform  $\Lambda^*$  of  $\Lambda$  satisfies:*

$$\begin{aligned} \Lambda^*(z) &= \sup_{\lambda \in \mathbb{R}} (\lambda z - \Lambda(\lambda)) \\ &= \begin{cases} +\infty & \text{if } z \geq 0 \\ z - \log(p) - z \log(\frac{z}{\log(p)}) & \text{if } z < 0. \end{cases} \end{aligned} \quad (12)$$

Then for any  $k \in \mathbb{N}^*$ , the sequence of distributions  $(\tilde{\mu}^{n,k})_{n \in \mathbb{N}, n > k}$  satisfies a Large Deviations Principle, with the rate function  $J = \Lambda^*$ .

The proof of (11) is the main task of this paper. In Section 5.1, we give a first easy proof in the case  $k = 1$ , relying on the knowledge of the distribution of  $J^{n,1}$ : it is a Poisson distribution with mean  $-n \log(p)$ . We can then compute explicitly  $\Lambda_{n,1}(\lambda)$  and prove (11). In Section 5.2, we study the general case  $k \geq 1$  with the method introduced in [BLR15], in the exponential case: for the algorithm  $\text{AMS}_{\text{expo}}(n, k; a, x)$ , we derive a functional equation on the Laplace transform  $\exp(\Lambda_{n,k}(\lambda))$  as a function of the initial condition  $x$ , for fixed parameter  $\lambda$ .

For completeness, we close this Section with the computation of the Fenchel-Legendre transform  $J = \Lambda^*$  of  $\Lambda$  in Proposition 4.2.

*Proof.* First, assume that  $z \geq 0$ . Then  $\lambda z - \Lambda(\lambda) \rightarrow +\infty$  when  $\lambda \rightarrow +\infty$ : thus  $\Lambda^*(z) = +\infty$ . Notice that this result is not surprising, since  $\log(\hat{p}^{n,k}) < 0$  almost surely.

If  $z < 0$ , the map  $\lambda \in \mathbb{R} \mapsto \lambda z - \Lambda(\lambda)$  admits the limit  $-\infty$  for  $z \rightarrow \pm\infty$ , and attains its maximum at the unique solution  $\lambda_z$  of the equation  $z - \frac{d\Lambda(\lambda)}{d\lambda}(\lambda_z) = 0$ , which is given by  $\lambda_z = -\log(\frac{z}{\log(p)})$ . Then  $\Lambda^*(z) = \lambda_z z - \Lambda(\lambda_z)$ , which gives (12).  $\square$

## 5 Proof of Proposition 4.2

### 5.1 The case $k = 1$

We start with a proof of Theorem 3.1 when  $k = 1$ : in this case, we have  $C^{n,1} = 1$  almost surely, and the number of iterations  $J^{n,1}$  follows a Poisson distribution  $\mathcal{P}(-n \log(p))$  (see for instance [BLR15], [GHML11]).

As a consequence, it is very easy to prove Proposition 4.1. Let  $\lambda \in \mathbb{R}$ . Then

$$\begin{aligned} \Gamma_{n,1}(\lambda) &= \exp(\Lambda_{n,1}(\lambda)) \\ &= \mathbb{E}[\exp(\lambda \log(\hat{p}^{n,1}))] \\ &= \mathbb{E}[\exp(\lambda \log(1 - 1/n) J^{n,1})] \\ &= \exp\left(-n \log(p) (\exp(\lambda \log(1 - 1/n)) - 1)\right). \end{aligned}$$

It is now easy to conclude: when  $n \rightarrow +\infty$

$$\begin{aligned} \frac{1}{n} \log(\Lambda(n\lambda)) &= -\log(p) (\exp(n\lambda \log(1 - 1/n)) - 1) \\ &\xrightarrow{n \rightarrow +\infty} -\log(p) (\exp(-\lambda) - 1). \end{aligned}$$

We have performed explicit calculations, using the knowledge of the distribution of  $J^{n,1}$ . However for  $k > 1$ , we cannot rely on such simple arguments and we need other tools.

We would like to use the connexion with the Poisson distribution in order to give an interpretation of the rate functions  $I$  and  $J$ . More precisely,  $I$  is obtained from  $J$  by the contraction principle ( $I(y) = J(\log(y))$ ), and  $J$  is the rate function obtained in the Cramer theorem where the distribution  $R$  is such that  $-R \sim \mathcal{P}(-\log(p))$ . Indeed, let  $(R_m)_{m \in \mathbb{N}^*}$  be independent, with the same distribution as  $X$ ; if we denote by  $\bar{R}_n = \frac{1}{n} \sum_{m=1}^n R_m$  the empirical

average, we compute for any  $\lambda \in \mathbb{R}$

$$\begin{aligned}\mathbb{E}[\exp(n\lambda\bar{R}_n)] &= \left(\mathbb{E}[\exp(\lambda R)]\right)^n \\ &= \left(\exp(-\log(p)(\exp(-\lambda) - 1))\right)^n.\end{aligned}$$

To conclude this section on the case  $k = 1$ , we prove Proposition 3.4. We use again the explicit knowledge of the distribution of  $J^{n,1}$  and use a Central Limit Theorem on exponential distributions to conclude.

*Proof of Proposition 3.4.* We write (with  $a = -\log(p) = \sigma^2 n$ )

$$\begin{aligned}\frac{\hat{p}^{n,1}}{p} &= \exp(J^{n,1} \log(1 - 1/n) + a) \\ &= \exp\left(\frac{J^{n,1} - na}{\sqrt{na}} \sqrt{na} \log(1 - 1/n) + a + na \ln(1 - 1/n)\right).\end{aligned}$$

By the Central Limit Theorem on the Poisson distribution, one gets, in the limit  $n \rightarrow +\infty$ , the following convergence in distribution

$$\frac{J^{n,1} - na}{\sqrt{na}} \rightarrow \mathcal{N}(0, 1).$$

Moreover, when  $n \rightarrow +\infty$ , we have  $\sqrt{na} \log(1 - 1/n) = n\sigma \log(1 - 1/n) \rightarrow -\sigma$  and  $a + na \log(1 - 1/n) = \sigma^2 (n + n^2 \ln(1 - 1/n))$  tends to  $\frac{-\sigma^2}{2}$ . This concludes the proof.  $\square$

## 5.2 The general case

In this section, we give the main arguments used to prove Proposition 4.2 in the general case  $k \in \mathbb{N}^*$ . In particular, we want to show that the rate function we obtain does not depend on  $k$ . The proof of some important but technical results is postponed to Section 7.

Even if in Section 5.1 above we have proved Proposition 4.2 in the case  $k = 1$ , we include this case in our general framework, and obtain an alternative proof.

To this aim, we make use of the strategy introduced in [BLR15] to study the properties of the  $\text{AMS}(n, k; a)$  algorithm. First, as explained in Section 2.3, we are allowed to restrict the study to the case where  $X$  is exponentially distributed: it is enough to study the  $\text{AMS}_{\text{expo}}(n, k; a_{\text{expo}})$  algorithm, where  $a_{\text{expo}} = -\log(p)$ .

Moreover, one of the main ideas is to consider the initial condition of the algorithm as an extra variable: for  $x \in [0, a)$ , we study the  $\text{AMS}_{\text{expo}}(n, k; a_{\text{expo}}, x)$  algorithm. From now on, in this Section, and in Section 7, we only consider the exponential case and we omit the dependence.

**Definition 5.1.** We use the following notation: for any  $(x, y) \in \mathbb{R}^2$

$$\begin{aligned} f(y) &= \exp(-y)\mathbb{1}_{y \geq 0} \quad , \quad F(y) = (1 - \exp(-y))\mathbb{1}_{y \geq 0} = \int_{-\infty}^y f(z)dz; \\ f(y; x) &= \frac{f(y)}{1 - F(x)}\mathbb{1}_{y \geq x} \quad , \quad F(y; x) = \frac{F(y) - F(x)}{1 - F(x)}\mathbb{1}_{y \geq x} = \int_{-\infty}^y f(z; x)dz; \\ f_{n,k}(y; x) &= k \binom{n}{k} F(y; x)^{k-1} f(y; x) (1 - F(y; x))^{n-k}, \\ F_{n,k}(y; x) &= \int_x^y f_{n,k}(z; x)dz. \end{aligned}$$

Let  $X$  be exponentially distributed with parameter 1. Then  $f$  (resp.  $F$ ) is the density (resp. the c.d.f.) of  $\mathcal{L}(X)$ . For  $x \geq 0$ ,  $f(\cdot; x)$  (resp.  $F(\cdot; x)$ ) is the density (resp. the c.d.f.) of the conditional distribution  $\mathcal{L}(X|X > x)$ .

Finally, let  $(X_1, \dots, X_n)$  be i.i.d. with the distribution of  $\mathcal{L}(X)$ , with the associated order statistics  $X_{(1)} < \dots < X_{(n)}$ . Then  $f_{n,k}(\cdot; x)$  (resp.  $F_{n,k}(\cdot; x)$ ) is the density (resp. the c.d.f.) of the  $k$ -th order statistic  $X_{(k)}$ .

The main object we need to study is the following function  $\Gamma_{n,k}$  of  $\lambda \in \mathbb{R}$  (considered as a fixed parameter) and the initial condition  $x \in [0, a]$

$$\begin{aligned} \Gamma_{n,k}(\lambda; x) &= \mathbb{E} \left[ \exp(n\lambda \log(\hat{p}^{n,k}(x))) \right] \\ &= \exp(\Lambda_{n,k}(n\lambda; x)). \end{aligned} \tag{13}$$

Notice that we include  $x = a$  in the domain of definition of the functions  $\Gamma_{n,k}$  and  $\Lambda_{n,k}$  (defined by (10)). It is also important to remark that we evaluate the latter at  $(n\lambda; x)$ .

We state several fundamental results which together yield Proposition 4.2 in the  $x$ -dependent case; to get (11) it is then enough to take  $x = 0$ .

First, Proposition 5.2 gives a functional equation satisfied by  $\Gamma_{n,k}(\lambda; \cdot)$  on  $[0, a]$ , for any value of the parameters  $1 \leq k < n$  and  $\lambda \in \mathbb{R}$ .

We use the following auxiliary function:

$$\Theta_{n,k}(\lambda; x) = \sum_{\ell=0}^{k-1} \exp\left(n\lambda \log\left(1 - \frac{\ell}{n}\right)\right) \left(F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x)\right), \tag{14}$$

with the convention  $F_{n,0}(y; x) = \mathbb{1}_{y \geq x}$ .

**Proposition 5.2.** For any  $n \in \mathbb{N}^*$ ,  $k \in \{1, \dots, n-1\}$ , and  $\lambda \in \mathbb{R}$ , the function  $\Gamma_{n,k}(\lambda; \cdot)$  is solution on the interval  $[0, a]$  of the functional equation (with the unknown  $\Gamma$ ):

$$\Gamma(x) = \int_x^a \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \Gamma(y) f_{n,k}(y; x) dy + \Theta_{n,k}(\lambda; x). \tag{15}$$

Notice that for the moment, it is not clear that  $\Gamma_{n,k}$  is the unique solution of the functional equation (15). We will prove this property below.

For completeness, we include a proof of this result, even if follows the same lines as Proposition 4.2 in [BLR15].

*Proof of Proposition 5.2.* We decompose the expectation according to the value of the (random) number of iterations  $J^{n,k}(x)$  in the algorithm starting from  $x$ :

$$\begin{aligned}\Gamma_{n,k}(\lambda; x) &= \mathbb{E}\left[\exp(n\lambda \log(\hat{p}^{n,k}(x)))\right] \\ &= \mathbb{E}\left[\exp(n\lambda \log(\hat{p}^{n,k}(x)))\mathbb{1}_{J^{n,k}(x)=0}\right] + \mathbb{E}\left[\exp(n\lambda \log(\hat{p}^{n,k}(x)))\mathbb{1}_{J^{n,k}(x)\geq 1}\right].\end{aligned}$$

First, since  $\{J^{n,k}(x) = 0\} = \{Z^1 \geq a\} = \bigcup_{\ell=0}^{k-1} \{X_{(\ell+1)} \geq a > X_{(\ell)}\}$ , we have

$$\begin{aligned}\mathbb{E}\left[\exp(n\lambda \log(\hat{p}^{n,k}(x)))\mathbb{1}_{J^{n,k}(x)=0}\right] &= \mathbb{E}\left[\exp(n\lambda \log(C^{n,k}(x)))\mathbb{1}_{J^{n,k}(x)=0}\right] \\ &= \sum_{\ell=0}^{k-1} \exp\left(n\lambda \log\left(1 - \frac{\ell}{n}\right)\right) \left(F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x)\right) \\ &= \Theta_{n,k}(\lambda; x).\end{aligned}$$

Second, we use  $\{J^{n,k}(x) \geq 1\} = \{Z^1 \leq a\}$  and condition with respect to  $Z^1$ :

$$\begin{aligned}\mathbb{E}\left[\exp(n\lambda \log(\hat{p}^{n,k}(x)))\mathbb{1}_{J^{n,k}(x)\geq 1}\right] &= \mathbb{E}\left[\mathbb{E}\left[\exp(n\lambda \log(\hat{p}^{n,k}(x)))\middle|Z^1\right]\mathbb{1}_{Z^1 < a}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\exp(n\lambda \log((1 - k/n)^{J^{n,k}(x)-1} C^{n,k}(x)) + n\lambda \log(1 - k/n))\middle|Z^1\right]\mathbb{1}_{Z^1 < a}\right] \\ &= \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \mathbb{E}\left[\mathbb{E}\left[\exp(n\lambda \log((1 - k/n)^{J^{n,k}(Z^1)} C^{n,k}(Z^1)))\middle|Z^1\right]\mathbb{1}_{Z^1 < a}\right] \\ &= \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \mathbb{E}\left[\Gamma_{n,k}(Z^1; x)\mathbb{1}_{Z^1 < a}\right] \\ &= \int_x^a \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \Gamma_{n,k}(\lambda; y) f_{n,k}(y; x) dy.\end{aligned}$$

We have used a kind of Markov property for the algorithm: up to taking into account for one more iteration, the algorithm behaves the same starting from  $x$  or from  $Z^1 \in (x, a]$ .  $\square$

Notice that the functional equation (15) involves a simple factor depending only on  $\lambda$ ,  $n$  and  $k$  in the integral, and that on both the left and the right-hand sides the function  $\Gamma$  is evaluated at the same value of the parameter  $\lambda$ . These observations are consequences of the choice to prove a Large Deviations Principle for  $\log(\hat{p}^{n,k})$  (instead of  $\hat{p}^{n,k}$ ) thanks to the Gärtner-Ellis Theorem, and to conclude with the use of the contraction principle; the same trick was used in [BGT14] to prove the Central Limit Theorem, thanks to the delta-method and the use of Levy Theorem. If one replaces  $\log(\hat{p}^{n,k}(x))$  with  $\hat{p}^{n,k}(x)$  in (13), then one obtains a more complicated functional equation where the observations above do not hold, and which is not easily exploitable. In particular, one does not obtain a nice counterpart of the fundamental result, Proposition 5.3 below.

We now state in Proposition 5.3 that solutions  $\Gamma$  of the functional equation (15) are in fact solutions of a linear Ordinary Differential Equation (ODE) of order  $k$ , with constant coefficients.

**Proposition 5.3.** *For any  $n \in \mathbb{N}^*$ ,  $k \in \{1, \dots, n-1\}$ , and  $\lambda \in \mathbb{R}$ , let  $\Gamma$  be a solution of the functional equation (15). Then it is solution of the following linear ODE of order  $k$ :*

$$\frac{d^k}{dx^k} \Gamma_{n,k}(\lambda; x) = \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \mu^{n,k} \Gamma_{n,k}(\lambda; x) + \sum_{m=0}^{k-1} r_m^{n,k} \frac{d^m}{dx^m} \Gamma_{n,k}(\lambda; x). \quad (16)$$

The coefficients  $\mu^{n,k}$  and  $(r_m^{n,k})_{0 \leq m \leq k-1}$  satisfy the following properties:

$$\begin{aligned} \mu^{n,k} &= (-1)^k n \dots (n - k + 1) \\ \nu^k - \sum_{m=0}^{k-1} r_m^{n,k} \nu^m &= (\nu - n) \dots (\nu - n + k - 1) \quad \text{for all } \nu \in \mathbb{R}. \end{aligned} \quad (17)$$

A sketch of proof of this result is postponed to Section 7. It uses the same arguments as to prove the corresponding functional equation in [BLR15]. For the proof of (17) in particular, we refer to that article.

To conclude on uniqueness of the solution of (15), and then prove asymptotic expansions on  $\Gamma_{n,k}$ , we prove the following Lemma.

**Lemma 5.4.** *For any fixed  $k \in \{1, \dots, \}$  and any  $\lambda \in \mathbb{R}$ , we have for any  $m \in \{0, \dots, k-1\}$*

$$\begin{aligned} \left. \frac{d^m}{dx^m} \Gamma_{n,k}(\lambda; x) \right|_{x=a} &= \left. \frac{d^m}{dx^m} \Theta_{n,k}(\lambda; x) \right|_{x=a} \\ &\underset{n \rightarrow \infty}{\sim} n^m (1 - \exp(-\lambda))^m. \end{aligned} \quad (18)$$

By Cauchy-Lipschitz theory, the linear ODE (16) with the conditions (18) at  $x = a$  admits a unique solution; therefore it is clear that  $\Gamma_{n,k}$  is the unique solution of (15).

**Remark 5.5.** *To prove the Central Limit Theorem in [BGT14], we used a similar result although in a weaker form: we only needed to prove  $\left. \frac{d^m}{dx^m} \Theta_{n,k}(\lambda; x) \right|_{x=a} = O(n^m)$ . Here we require a more precise asymptotic result in order to prove that the coefficient  $\gamma_{n,k}^1(\lambda)$  defined in Proposition 5.6 below converges to 1 (in fact, we only need that it is bounded from below by a positive constant).*

We finally explain how to obtain asymptotic knowledge on  $\Gamma_{n,k}(\lambda; x)$  and  $\Lambda_{n,k}(n\lambda, x)$  when  $n \rightarrow +\infty$ . First, the  $k$  roots  $(\nu_{n,k}^\ell(\lambda))_{1 \leq \ell \leq k}$  of the polynomial equation associated with the linear ODE (16) are pairwise distinct for  $n$  large enough (the other parameters  $\lambda$  and  $k$  being fixed), and more precisely they satisfy (20). As a consequence, the solution  $\Gamma_{n,k}$  can be written (see (19)) as a linear combination of exponential functions  $x \mapsto \exp(\nu_{n,k}^\ell(\lambda)(x - a))$ . Finally, using the asymptotic expression for the derivatives of order  $0, \dots, k-1$  at  $x = a$ , we obtain a linear system of equations, solve it using the Cramer's formulae and obtain the asymptotic expression (21). The proof is postponed to Section 7.

**Proposition 5.6.** *Let  $k \in \{1, \dots, \}$  and  $\lambda \in \mathbb{R}$  be fixed. Then for  $n$  large enough, we have for any  $x \in [0, a]$*

$$\Gamma_{n,k}(\lambda, x) = \sum_{\ell=1}^k \gamma_{n,k}^\ell(\lambda) \exp\left(\nu_{n,k}^\ell(\lambda)(x - a)\right), \quad (19)$$

where

$$\nu_{n,k}^\ell(\lambda) \sim n \left(1 - e^{-\lambda} e^{i2\pi \frac{(\ell-1)}{k}}\right) \quad (20)$$

and

$$\gamma_{n,k}^\ell(\lambda) \rightarrow \mathbb{1}_{\ell=1}. \quad (21)$$

We now conclude and prove Proposition 4.1, namely the Large Deviations Principle for  $(\mathcal{L}(\log(\hat{p}^{n,k})))_{n>k}$ .

We start with the case  $k > 1$ . Then for any  $\ell \in \{2, \dots, k\}$  we have for any  $\lambda \in \mathbb{R}$

$$\operatorname{Re}\left(1 - e^{-\lambda} e^{i2\pi(\ell-1)/k}\right) > \operatorname{Re}\left(1 - e^{-\lambda}\right).$$

As a consequence, for  $x < a$  we have when  $n \rightarrow +\infty$

$$e^{\nu_{n,k}^\ell(\lambda)(x-a)} = o\left(e^{1-\exp(-\lambda)}(x-a)\right),$$

and thus

$$\frac{1}{n} \Lambda_{n,k}(n\lambda; x) = \frac{1}{n} \log(\Gamma_{n,k}(\lambda; x)) \underset{n \rightarrow +\infty}{\sim} \nu_{n,k}^1(\lambda)(x-a) \underset{n \rightarrow +\infty}{\rightarrow} (1 - e^{-\lambda})(x-a).$$

When  $k = 1$ , the linear ODE (16) is of order 1, and it is easy to check that

$$\Gamma_{n,1}(\lambda; x) = \exp\left(\nu_{n,k}^1(\lambda)(x-a)\right),$$

so that the same asymptotic result as above holds.

It remains to take  $x = a$ , and to recall that  $a = -\log(p)$  if  $p = \mathbb{P}(X > a)$  and  $X$  is exponentially distributed with parameter 1.

This concludes the proof of Proposition 4.1.

## 6 Comparison with other algorithms

We propose a comparison (in terms of large deviations) of the Adaptive Multilevel Splitting algorithm with the two other methods described in the Introduction: a direct, naive Monte-Carlo method, based on a non-interacting system of replicas with the same size (see the estimator (3)), and a non-adaptive version of multilevel splitting (see the estimator (2)).

In the first case, we obtain that large deviations are much less likely for the AMS algorithm than for the crude Monte-Carlo method. In the second case, we show that the AMS estimator is more efficient than the non-adaptive one taken in the limit of a large number  $N$  of fixed levels.

These results are consistent with the cost analysis and the comparison based on the central limit theorem, see [BLR15], [BGT14], [CDMFG12], [CG14].

### 6.1 Crude Monte-Carlo

We compare the performance of the AMS algorithm with the use of a Crude Monte-Carlo estimation in the large  $n$  limit.

Let  $(X_m)_{m \in \mathbb{N}^*}$  a sequence of independent and identically distributed random variables, each one being equal in law with  $X$ .

Then for any  $n \in \mathbb{N}^*$

$$\bar{p}_n = \frac{1}{n} \sum_{m=1}^n \mathbb{1}_{X_m > a} \tag{22}$$

is an unbiased estimator of  $p$ .



It is a classical result (Theorem 2.2.3 in [DZ10]) due to Cramer that the sequence  $(\mathcal{L}(\bar{p}_n))_{n \in \mathbb{N}^*}$  satisfies a Large Deviations Principle with the rate function (case of Bernoulli random variables, see Exercice 2.2.23 in [DZ10]):

$$\mathcal{I}(y) = \begin{cases} +\infty & \text{if } y \notin (0, 1) \\ y \log\left(\frac{y}{p}\right) + (1-y) \log\left(\frac{1-y}{1-p}\right) & \text{if } y \in (0, 1). \end{cases} \quad (23)$$

The comparison between the algorithms is based on the following result:

**Proposition 6.1.** *For any  $p \in (0, 1)$  and any  $y \in (0, 1)$ , we have*

$$\begin{aligned} I(y) &\geq \mathcal{I}(y), \\ I(y) &= \mathcal{I}(y) \quad \text{if and only if } y = p. \end{aligned}$$

*Proof.* We explicitly mention the dependence of  $I$  and of  $\mathcal{I}$  with respect to  $p$ , and we define

$$D(y, p) = I(y, p) - \mathcal{I}(y, p).$$

It is clear that  $D(p, p) = 0$  for any  $p \in (0, 1)$ . We compute that

$$\frac{\partial D(y, p)}{\partial p} = \frac{1-y}{p \log(p)} \left( \frac{\log(y)}{1-y} - \frac{\log(p)}{1-p} \right);$$

since the function  $t \mapsto \frac{\log(t)}{1-t}$  is strictly decreasing on  $(0, 1)$  (as can be seen by computing its first and second order derivatives), we see that for any  $y, p \in (0, 1)^2$  we have the inequalities

$$\frac{\partial D(y, p)}{\partial p} > 0 \quad \text{if } y > p \quad \text{and} \quad \frac{\partial D(y, p)}{\partial p} < 0 \quad \text{if } y < p.$$

Using  $D(p, p) = 0$ , it is easy to conclude. □

Now let  $\epsilon \in (0, \max(p, 1-p))$ ; then for  $n$  large we have

$$\frac{\mathbb{P}(\hat{p}^{n,k} - p > \epsilon)}{\mathbb{P}(\bar{p}_n - p > \epsilon)} = \exp(n\Delta(\epsilon, n)) \rightarrow 0,$$

exponentially fast, since we have by the Large Deviations Principles  $\Delta(\epsilon, n) \rightarrow \mathcal{I}(p + \epsilon) - I(p + \epsilon) < 0$  when  $n \rightarrow +\infty$  (notice that both  $\mathcal{I}$  and  $I$  are increasing on  $(p, 1)$ ).

The same arguments apply to get

$$\frac{\mathbb{P}(\hat{p}^{n,k} - p < -\epsilon)}{\mathbb{P}(\bar{p}_n - p < -\epsilon)} \rightarrow 0.$$

As a consequence, the probability of observing large deviations from the mean  $p$  is much smaller for the AMS algorithm than when using a crude Monte-Carlo estimator, in the large  $n$  limit. This statement is a new way of expressing the efficiency of the AMS algorithm.

Notice that in the discussion above we have not assumed that we are estimating a probability in a rare event regime: the conclusion holds for any  $p \in (0, 1)$ . Now it is also instructive

to compare  $I((1 + \epsilon)p)$  and  $\mathcal{I}((1 + \epsilon)p)$  for a given  $\epsilon \in (0, 1)$  and when  $p \rightarrow 0$ : it amounts at looking at deviations of the relative error, and we have

$$\begin{aligned} \lim_{n \rightarrow +\infty} \frac{1}{n} \log \left( \mathbb{P} \left( \frac{\hat{p}^{n,k} - p}{p} > \epsilon \right) \right) &= -I(p(1 + \epsilon)) \sim_{p \rightarrow 0} -\frac{(\log(1 + \epsilon))^2}{-2 \log(p)} \\ \lim_{n \rightarrow +\infty} \frac{1}{n} \log \left( \mathbb{P} \left( \frac{\bar{p}_n - p}{p} > \epsilon \right) \right) &= -\mathcal{I}(p(1 + \epsilon)) \sim_{p \rightarrow 0} -p((1 + \epsilon) \log(1 + \epsilon) - \epsilon). \end{aligned}$$

Given  $\delta > 0$ , in order to have a probability lower than  $\delta$  that the relative error is larger than  $\epsilon$ , in the small  $p$  limit, one thus needs a number of replicas  $n$  which scales like  $1/p$  when using the crude Monte-Carlo method, while it scales like  $-\log(p)$  (which is much smaller) when using the AMS algorithm. Moreover, since the expected workload is of size  $n$  when using the Monte-Carlo method and of size  $-n \log(p)$  when using the AMS algorithm, it is clear that in terms of large deviations from the mean the AMS algorithm is more efficient than the crude Monte-Carlo method.

Notice that this discussion is consistent with the conclusions coming from the Central Limit Theorem, where in the regime  $p \rightarrow 0$  the asymptotic variance is equivalent to  $p$  when using the crude Monte-Carlo method and  $-p^2 \log(p)$  when using the AMS algorithm: to obtain reliable confidence intervals on the relative error, the number of replicas  $n$  scales in the same way.

## 6.2 Non-adaptive Multilevel Splitting

We now compare the rate function  $I$  obtained for the Large Deviations Principle on the AMS algorithm, with the one we obtain when using a deterministic (non-adaptive) sequence of levels.

Namely, using Assumption 2.1, we decompose the probability as a telescoping product of  $N \in \mathbb{N}^*$  conditional probabilities

$$p = \mathbb{P}(X > a) = \prod_{i=1}^N \mathbb{P}(X > a_i | X > a_{i-1}), \quad (24)$$

associated with a given non-decreasing sequence of levels  $a_0 = 0 < a_1 < \dots < a_N = a$ . We denote by  $p^{(i)} = \mathbb{P}(X > a_i | X > a_{i-1})$  the  $i$ -th conditional probability. The sequence is of size  $N$  and we study the asymptotic regime  $N \rightarrow +\infty$ .

We can define an unbiased estimator of  $p$  as follows: let  $n \in \mathbb{N}^*$  and set

$$\hat{p}_n^N = \prod_{i=1}^N \bar{p}_n^{(i)}, \quad (25)$$

where  $(\bar{p}_n^{(i)})_{1 \leq i \leq N}$  is a family of independent random variables, where each  $\bar{p}_n^{(i)}$  is a Crude Monte-Carlo estimator (as defined in the section above) for the probability  $p^{(i)}$  with  $n$  realizations. More precisely, let  $(X_m^{(i)})_{1 \leq m \leq n, 1 \leq i \leq N}$  be independent random variables, such that  $\mathcal{L}(X_m^{(i)}) = \mathcal{L}(X | X > a_{i-1})$ , and set

$$\bar{p}_n^{(i)} = \frac{1}{n} \sum_{m=1}^n \mathbb{1}_{X_m^{(i)} > a_i}. \quad (26)$$

From a practical point of view, notice that the computation of these estimators requires the sampling of random variables according to the conditional distribution  $\mathcal{L}(X|X > a_{i-1})$  for each  $i \in \{1, \dots, N\}$ , just like for the adaptive version.

Here  $n$  thus denotes the number of replicas used for the estimation of the probabilities in both the adaptive and the non-adaptive versions. We needed the extra parameter  $N$  to denote the number of iterations (*i.e.* the length of the sequence of levels) of the algorithm, while we know that the average number of iterations is of the order  $-\frac{n \log(p)}{k}$  in the adaptive case. Therefore, to study the non-adaptive version, we first let  $n \rightarrow +\infty$ , and then analyze the behavior of the asymptotic quantities with respect to  $N$  (in the limit  $N \rightarrow +\infty$ ), while for the adaptive version we need to pass to the limit only once, namely  $n \rightarrow +\infty$ .

Clearly, by the independence properties of the random variables introduced here we have

$$\mathbb{E}[\hat{p}_n^N] = p.$$

Moreover, it is well-known that, for a given value of  $N$  (the length of the sequence of levels) the asymptotic variance (when  $n$  goes to  $+\infty$ ) is minimized when  $p^{(i)} = p^{1/N}$  for any  $i \in \{1, \dots, N\}$  (*i.e.* the conditional probabilities in (24) are equal); moreover the asymptotic variance is a decreasing function of  $N$ , which converges to  $\frac{-p^2 \log(p)}{n}$  when  $N \rightarrow +\infty$ . From a practical point of view, the computation of the associated sequence of levels  $a_1, \dots, a_{N-1}$  is *a priori* difficult: the adaptive version overcomes this issue, and in the regime  $N \rightarrow +\infty$  both the non-adaptive and the adaptive version have the same statistical properties.

As a consequence, from now on we assume that  $p^{(i)} = p^{1/N}$  for any  $i \in \{1, \dots, N\}$ .

For any  $i \in \{1, \dots, N\}$ ,  $(\mathcal{L}(\bar{p}_n^{(i)}))_{n \in \mathbb{N}^*}$  satisfies a Large Deviations Principle with the rate function (see (23))

$$\mathcal{I}_N(y) = \begin{cases} +\infty & \text{if } y \notin (0, 1) \\ y \log\left(\frac{y}{p^{1/N}}\right) + (1-y) \log\left(\frac{1-y}{1-p^{1/N}}\right) & \text{if } y \in (0, 1). \end{cases} \quad (27)$$

Since for any  $n \in \mathbb{N}^*$  the random variables  $(\bar{p}_n^{(i)})_{1 \leq i \leq N}$  are independent, it is easy to generalize this statement as follows. The sequence  $(\mathcal{L}(\bar{p}_n^{(1)}, \dots, \bar{p}_n^{(N)}))_{n \in \mathbb{N}^*}$  satisfies a Large Deviations Principle in  $\mathbb{R}^N$  with the rate function (with abuse of notation  $\mathcal{I}_N$  refers both to the function depending on a 1-dimensional or a  $N$ -dimensional variable)

$$\mathcal{I}_N(y_1, \dots, y_N) = \sum_{i=1}^N \mathcal{I}_N(y_i). \quad (28)$$

Now as a consequence of the contraction principle, since  $\hat{p}_n^N = \prod_{i=1}^N \bar{p}_n^{(i)}$ , the sequence  $(\hat{p}_n^N)_{n \in \mathbb{N}^*}$  also satisfies a Large Deviations Principle with the rate function

$$I_N(y) = \inf \left\{ \mathcal{I}_N(y_N, \dots, y_1) ; y = \prod_{i=1}^N y_i \right\}. \quad (29)$$

On the one hand, it is clear that if  $y \notin (0, 1)$ , then  $I_N(y) = +\infty$ . Indeed, for any  $(y_1, \dots, y_N)$  satisfying the constraint  $y = \prod_{i=1}^N y_i \notin (0, 1)$ , at least one of the  $y_i$ 's satisfies  $y_i \notin (0, 1)$ , which yields  $\mathcal{I}_N(y_i) = \mathcal{I}_N(y_1, \dots, y_N) = +\infty$ .

On the other hand, by definition of  $I_N$ , we have for any  $y \in (0, 1)$

$$\begin{aligned} I_N(y) &\leq \mathcal{I}_N(y^{1/N}, \dots, y^{1/N}) = N\mathcal{I}_N(y^{1/N}) \\ &= Ny^{1/N} \log\left(\frac{y^{1/N}}{p^{1/N}}\right) + N(1 - y^{1/N}) \log\left(\frac{1 - y^{1/N}}{1 - p^{1/N}}\right) \\ &\xrightarrow{N \rightarrow \infty} \log(y) - \log(p) - \log(y) \log\left(\frac{\log(y)}{\log(p)}\right) = I(y). \end{aligned}$$

For our purpose, this inequality is sufficient.

We now interpret the previous inequality in terms of asymptotic estimates for deviations of  $\hat{p}_n^N$  and of  $\hat{p}^{n,k}$  with respect to their expected value  $p$ . Let  $\epsilon > 0$ , then we have by definition of the Large Deviations Principle with rate function  $I_N$

$$\begin{aligned} \liminf_{n \rightarrow +\infty} \frac{1}{n} \log\left(\mathbb{P}(|\hat{p}_n^N - p| > \epsilon)\right) &\geq -\inf\{I_N(y) ; |y - p| \geq \epsilon\} \\ &\geq -\inf\left\{N\mathcal{I}_N(y^{1/N}) ; |y - p| \geq \epsilon\right\} \\ &\geq -\min\left\{N\mathcal{I}_N((p + \epsilon)^{1/N}), N\mathcal{I}_N((p - \epsilon)^{1/N})\right\}, \end{aligned}$$

using that  $\mathcal{I}_N$  is non-increasing on  $(-\infty, p^{1/N})$  and non-decreasing on  $(p^{1/N}, +\infty)$ .

To conclude, notice that

$$\begin{aligned} \lim_{N \rightarrow +\infty} -\min\left\{N\mathcal{I}_N((p + \epsilon)^{1/N}), N\mathcal{I}_N((p - \epsilon)^{1/N})\right\} &= -\min\{I(p + \epsilon), I(p - \epsilon)\} \\ &= \lim_{n \rightarrow +\infty} \frac{1}{n} \log\left(\mathbb{P}(|\hat{p}^{n,k} - p| > \epsilon)\right). \end{aligned}$$

We can thus assess that the Adaptive Multilevel Splitting algorithm is more efficient (in a large sense) than the non-adaptive version in terms of large deviations when the number of replicas  $n$  goes to  $+\infty$  and in the limit of large number  $N$  of levels.

## 7 Proof of the technical estimates

In this section, we give detailed proofs for the technical auxiliary results used in Section 5.2.

*Proof of Proposition 5.3.* We proceed by recursion, like in the proof of Proposition 6.4 in [BLR15] and Lemma 2 in [BGT14]. We fix the values of  $1 \leq k < n$  and of  $\lambda \in \mathbb{R}$ .

Differentiating recursively with respect to  $x$ , for any  $0 \leq l \leq k - 1$  and for any  $0 \leq x \leq a$  we have (for a family of coefficients described by (32) below)

$$\begin{aligned} \frac{d^l}{dx^l} (\Gamma_{n,k}(\lambda; x) - \Theta_{n,k}(\lambda; x)) &= \mu_l^{n,k} \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \int_x^a \Gamma_{n,k}(\lambda; y) f_{n,k-l}(y; x) dy \\ &\quad + \sum_{m=0}^{l-1} r_{m,l}^{n,k} \frac{d^m}{dx^m} (\Gamma_{n,k}(\lambda; x) - \Theta_{n,k}(\lambda; x)), \end{aligned} \quad (30)$$

and that differentiating once more we get

$$\begin{aligned} \frac{d^k}{dx^k} (\Gamma_{n,k}(\lambda; x) - \Theta_{n,k}(\lambda; x)) &= \mu^{n,k} \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \Gamma_{n,k}(\lambda; x) \\ &\quad + \sum_{m=0}^{k-1} r_m^{n,k} \frac{d^m}{dx^m} (\Gamma_{n,k}(\lambda; x) - \Theta_{n,k}(\lambda; x)), \end{aligned} \quad (31)$$

with  $\mu^{n,k} := \mu_k^{n,k}$  and  $r_m^{n,k} := r_{m,k}^{n,k}$ .

The coefficients satisfy

$$\begin{cases} \mu_0^{n,k} = 1, \mu_{l+1}^{n,k} = -(n-k+l+1)\mu_l^{n,k}; \\ r_{0,l+1}^{n,k} = -(n-k+l+1)r_{0,l}^{n,k}, \quad \text{if } l > 0, \\ r_{m,l+1}^{n,k} = r_{m-1,l}^{n,k} - (n-k+l+1)r_{m,l}^{n,k}, \quad 1 \leq m \leq l, \\ r_{l,l}^{n,k} = -1. \end{cases} \quad (32)$$

Notice that these coefficients do not depend on  $\lambda$ , and are the same as in [BLR15] and [BGT14]. Properties (17) are proved in [BLR15].

Thanks to (17), for all  $j \in \{0, \dots, k-1\}$  and any  $x \in [0, a]$  we have

$$\frac{d^k}{dx^k} \exp((n-k+j+1)(x-a)) = \sum_{m=0}^{k-1} r_m^{n,k} \frac{d^m}{dx^m} \exp((n-k+j+1)(x-a)).$$

Using the expression of  $F_{n,k}$ , straightforward computations show that  $\Theta_{n,k}(\lambda; \cdot)$  is a linear combination of the exponential functions  $z \mapsto \exp(-nz), \dots, \exp(-(n-k+1)z)$ ; therefore

$$\frac{d^k}{dx^k} \Theta_{n,k}(t, x) = \sum_{m=0}^{k-1} r_m^{n,k} \frac{d^m}{dx^m} \Theta_{n,k}(t, x),$$

and thus (31) gives (16). □

*Proof of Lemma 5.4.* From (30), the equality in (18) is clear.

We claim that for any  $0 \leq m \leq k-1$  and any  $0 \leq \ell \leq k-1$

$$\frac{d^m}{dx^m} \left( F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x) \right) \Big|_{x=a} \underset{n \rightarrow \infty}{\sim} n^m \binom{m}{\ell} (-1)^\ell. \quad (33)$$

In particular,  $\frac{d^m}{dx^m} \left( F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x) \right) \Big|_{x=a} = 0 = \binom{m}{\ell}$  for  $n$  large enough as soon as  $\ell > m$ . Conclusion is then straightforward: using the definition (14) of  $\Theta_{n,k}$ , we get

$$\begin{aligned} \frac{1}{n^m} \frac{d^m}{dx^m} \Theta_{n,k}(\lambda; x) \Big|_{x=a} &= \frac{1}{n^m} \sum_{\ell=0}^{k-1} \frac{d^m}{dx^m} \exp(n\lambda \log(1 - \frac{\ell}{n})) \left( F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x) \right) \Big|_{x=a} \\ &\xrightarrow{n \rightarrow \infty} \sum_{\ell=0}^m \binom{m}{\ell} (-1)^\ell \exp(-\ell\lambda) \\ &= \left( 1 - \exp(-\lambda) \right)^m. \end{aligned}$$

We now prove (33) by induction on  $m$ .

We first consider  $m = 0$ . Then for any  $\ell \in \mathbb{N}^*$  we have  $F_{n,\ell}(a; a) = 0$  and  $F_{n,0}(a; a) = 1$  (by the convention  $F_{n,0}(y; x) = \mathbf{1}_{y \geq x}$ ), and (33) holds.

Let us also consider  $m = 1$ , when  $k \geq 2$ . Then  $\frac{d}{dx} F_{n,0}(a; x) \Big|_{x=a} = 0$ , while for any  $x \leq a$

$$\frac{d}{dx} F_{n,\ell}(a; x) = \frac{d}{dx} F_{n,\ell}(a-x; 0) = -f_{n,\ell}(a-x; 0) = -f_{n,\ell}(a; x)$$

as a consequence of the absence of memory property of the exponential distribution.

Now since  $f_{n,\ell}(a, a) = n\mathbb{1}_{\ell=1}$ , we get (33) for  $m = 1$ .

The induction is based on the following relations (deduced from elementary computations; for a proof see [BLR15], Section 6.3)

$$\left\{ \begin{array}{l} \frac{d}{dx} f_{n,1}(y; x) = n f_{n,1}(y; x). \\ \text{for } \ell \in \{2, \dots, n-1\}, \frac{d}{dx} f_{n,\ell}(y; x) = (n - \ell + 1)(f_{n,\ell}(y; x) - f_{n,\ell-1}(y; x)). \end{array} \right. \quad (34)$$

Thanks to the first formula in (34), we easily get (33) for  $\ell = 0$  by induction on  $m$ .

If now  $\ell \in \{1, \dots, k-1\}$ , we have the recursive formula for  $m \geq 1$

$$\begin{aligned} \frac{d^{m+1}}{dx^{m+1}} \left( F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x) \right) \Big|_{x=a} &= \frac{d^m}{dx^m} \left( f_{n,\ell+1}(a; x) - f_{n,\ell}(a; x) \right) \Big|_{x=a} \\ &= (n - \ell) \frac{d^{m-1}}{dx^{m-1}} \left( f_{n,\ell+1}(a; x) - f_{n,\ell}(a; x) \right) \Big|_{x=a} \\ &\quad - (n - \ell + 1) \frac{d^{m-1}}{dx^{m-1}} \left( f_{n,\ell}(a; x) - f_{n,\ell-1}(a; x) \right) \Big|_{x=a} \\ &= (n - \ell) \frac{d^m}{dx^m} \left( F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x) \right) \Big|_{x=a} \\ &\quad - (n - \ell + 1) \frac{d^m}{dx^m} \left( F_{n,\ell-1}(a; x) - F_{n,\ell}(a; x) \right) \Big|_{x=a} \end{aligned}$$

Finally using the induction hypothesis and obtain

$$\begin{aligned} \frac{1}{n^{m+1}} \frac{d^{m+1}}{dx^{m+1}} \left( F_{n,\ell}(a; x) - F_{n,\ell+1}(a; x) \right) \Big|_{x=a} &\xrightarrow{n \rightarrow +\infty} (-1)^\ell \binom{m}{\ell} - (-1)^{\ell-1} \binom{m}{\ell-1} \\ &= (-1)^\ell \binom{m+1}{\ell}. \end{aligned}$$

This concludes the proof of Lemma 5.4. □

*Proof of Proposition 5.6.* The  $\nu_{n,k}^\ell(\lambda)$  are the roots of the characteristic equation associated with the linear ODE (16):

$$\frac{(n - \nu) \dots (n - k + 1 - \nu)}{n \dots (n - k + 1)} - \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) = 0,$$

which can be rewritten as a polynomial equation of degree  $k$  with respect to the variable  $\bar{\nu}_n = \frac{\nu}{n}$ :

$$\frac{(1 - \bar{\nu}_n) \dots \left(1 - \frac{k-1}{n} - \bar{\nu}_n\right)}{1 \dots \left(1 - \frac{k-1}{n}\right)} - \exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) = 0,$$

where  $\exp\left(n\lambda \log\left(1 - \frac{k}{n}\right)\right) \xrightarrow{n \rightarrow +\infty} \exp(-k\lambda)$ .

By continuity of the roots of polynomials of degree  $k$  with respect to the coefficients, we get that for all  $\ell \in \{1, \dots, k\}$  (with an appropriate ordering of the roots)

$$\frac{\nu_{n,k}^\ell(\lambda)}{n} \rightarrow \bar{\nu}_{\infty,k}^\ell(\lambda)$$

where  $(1 - \bar{\nu}_{\infty,k}^\ell(\lambda))^k = e^{-k\lambda}$ . This identity immediately yields (20).

As a consequence, for  $n$  large enough the roots  $\nu_{n,k}^\ell(\lambda)$  are pairwise distinct. Then (19) holds for some complex numbers  $\gamma_{n,k}^\ell(\lambda)$ , where  $\ell \in \{1, \dots, k\}$ . Thanks to (19) evaluated at  $x = a$ , these coefficients are solution of the following linear system of equations:

$$\begin{cases} \gamma_{n,k}^1(\lambda) + \dots + \gamma_{n,k}^k(\lambda) = \Gamma_{n,k}(\lambda; x)|_{x=a}, \\ \gamma_{n,k}^1(\lambda)\nu_{n,k}^1(\lambda) + \dots + \gamma_{n,k}^k(\lambda)\nu_{n,k}^k(\lambda) = \frac{d}{dx}\Gamma_{n,k}(\lambda; x)|_{x=a}, \\ \vdots \\ \gamma_{n,k}^1(\lambda)\left(\nu_{n,k}^1(\lambda)\right)^{k-1} + \dots + \gamma_{n,k}^k(\lambda)\left(\nu_{n,k}^k(\lambda)\right)^{k-1} = \frac{d^{k-1}}{dx^{k-1}}\Gamma_{n,k}(\lambda; x)|_{x=a}. \end{cases} \quad (35)$$

This system is equivalent with

$$\begin{cases} \gamma_{n,k}^1(\lambda) + \dots + \gamma_{n,k}^k(\lambda) = \Gamma_{n,k}(\lambda; x)|_{x=a} \xrightarrow{n \rightarrow +\infty} 1, \\ \gamma_{n,k}^1(\lambda)\bar{\nu}_{n,k}^1(\lambda) + \dots + \gamma_{n,k}^k(\lambda)\bar{\nu}_{n,k}^k(\lambda) = \frac{1}{n} \frac{d}{dx}\Gamma_{n,k}(\lambda; x)|_{x=a} \xrightarrow{n \rightarrow +\infty} \bar{\nu}_{\infty,k}^1(\lambda), \\ \vdots \\ \gamma_{n,k}^1(\lambda)\bar{\nu}_{n,k}^1(\lambda)^{k-1} + \dots + \gamma_{n,k}^k(\lambda)\bar{\nu}_{n,k}^k(\lambda)^{k-1} = \frac{1}{n^{k-1}} \frac{d^{k-1}}{dx^{k-1}}\Gamma_{n,k}(\lambda; x)|_{x=a} \xrightarrow{n \rightarrow +\infty} \bar{\nu}_{\infty,k}^1(\lambda)^{k-1}, \end{cases} \quad (36)$$

thanks to (18) and (20), where  $\bar{\nu}_{n,k}^\ell(\lambda) = \frac{\nu_{n,k}^\ell(\lambda)}{n} \xrightarrow{n \rightarrow +\infty} \bar{\nu}_{\infty,k}^\ell(\lambda)$ .

It is now easy to get (21), which concludes the proof of Proposition 5.6.  $\square$

## 8 Conclusion and perspectives

We have established (Theorem 3.1) a Large Deviations Principle result for the Adaptive Multilevel Splitting AMS( $n, k$ ) Algorithm in an idealized setting, when the number of replicas  $n$  goes to infinity while the parameter  $k$  and the threshold  $a$  remain fixed. The rate function does not depend on  $k$ : when  $k = 1$ , the proof is very simple and uses an interpretation of the algorithm with a Poisson process (the number of iterations follows a Poisson distribution). When  $k > 1$ , we rely on a functional equation technique which was already used to prove unbiasedness and asymptotic normality of the estimator in the previous works [BLR15] and [BGT14].

We were able to relate the efficiency of the algorithm with this Large Deviations result, with a comparison with two algorithms (see Section 6): a crude Monte-Carlo method and a non-adaptive version. More generally, in other situations Large Deviations could be a powerful tool to compare adaptive or non-adaptive multilevel splitting algorithms, instead of resorting only on comparison of asymptotic variances associated with central limit theorems.

Let us mention a few open directions for future works. First, it should be interesting to look at the regime where  $k$  also goes to infinity, with  $k/n$  converging to a proportion  $\alpha \in (0, 1)$ . We expect to prove that the optimal rate function is obtained for  $\alpha$  decreasing to 0: indeed, the asymptotic variance is minimized in this regime. A comparison with a non-adaptive version of the algorithm is expected to show that the adaptive algorithm behaves (in terms of large deviations) like the non-adaptive version when the number of replicas and of levels goes to infinity, like in the regime we have studied in this paper.

A severe restriction is given by the so-called idealized setting: we need to know how to sample according to the conditional distribution  $\mathcal{L}(X|X > x)$ . In practice, and especially when computing crossing probabilities for high dimensional metastable stochastic processes, it is not satisfied and the multilevel splitting algorithm needs to use an importance function to define appropriate levels, and at each step the computation of the new sample uses the one at the previous iteration (thanks to a branching procedure of the successful trajectories). A natural question is whether one can prove a Large Deviations Principle in such a framework, and study quantitatively how the rate function depends on the importance function.

In fact, when using both non-adaptive (see [GKvO02], [GHSZ98]) and adaptive ([BGG<sup>+</sup>], in preparation) multilevel splitting algorithms, one may observe a very large difference between the value of the estimator (averaged over a number  $M$  of independent realizations) and the true result, or between the results obtained for different choices of the importance function. Even if the estimator of the probability is unbiased, in such situations one observes an *apparent bias* toward smaller values if  $M$  is not sufficiently large. This phenomenon is explained by specificity of the models: there are several channels to reach the region  $B$  from  $A$  (in the case of the estimation of crossing probabilities between metastable states of a Markov process), which may be sampled very differently when the importance function changes. It should be interesting to investigate the relation between this phenomenon and the Large Deviations Principle for the associated estimator.

## Acknowledgments

The author would like to thank B. Bercu and A. Richou for suggesting this work, and F. Cérou, A. Guyader and M. Rousset for helpful discussions and advice.

## References

- [AB01] S. K. Au and J. L. Beck. Estimation of small failure probabilities in high dimensions by subset simulation. *Journal of Probabilistic Engineering Mechanics*, 16:263–277, 2001.
- [AG07] S. Asmussen and P. W. Glynn. *Stochastic simulation: algorithms and analysis*, volume 57 of *Stochastic Modelling and Applied Probability*. Springer, New York, 2007.
- [BGG<sup>+</sup>] C. E. Bréhier, M. Gazeau, L. Goudenège, T. Lelièvre, and M. Rousset. Unbiasedness for some generalized adaptive multilevel splitting algorithms. *in preparation*.
- [BGT14] C. E. Bréhier, L. Goudenège, and L. Tudela. Central limit theorem for adaptive multilevel splitting estimators in an idealized setting. *preprint*, 2014.
- [BLR15] C. E. Bréhier, T. Lelièvre, and M. Rousset. Analysis of adaptive multilevel splitting algorithms in an idealized setting. *ESAIM Probability and Statistics*, to appear, 2015.
- [CDMFG12] F. Cérou, P. Del Moral, T. Furon, and A. Guyader. Sequential Monte Carlo for rare event estimation. *Stat. Comput.*, 22(3):795–808, 2012.



- [CG07] F. Cérou and A. Guyader. Adaptive multilevel splitting for rare event analysis. *Stoch. Anal. Appl.*, 25(2):417–443, 2007.
- [CG14] F. Cérou and A. Guyader. Fluctuations of adaptive multilevel splitting. *preprint*, 2014.
- [DZ10] A. Dembo and O. Zeitouni. *Large deviations techniques and applications. 2nd ed., corrected 2nd printing*. Berlin: Springer, 2nd ed., corrected 2nd printing edition, 2010.
- [GHML11] A. Guyader, N. Hengartner, and E. Matzner-Løber. Simulation and estimation of extreme quantiles and extreme probabilities. *Appl. Math. Optim.*, 64(2):171–196, 2011.
- [GHSZ98] P. Glasserman, P. Heidelberger, P. Shahabuddin, and T. Zajic. A large deviations perspective on the efficiency of multilevel splitting. *IEEE Trans. Automat. Control*, 43(12):1666–1679, 1998.
- [GHSZ99] P. Glasserman, P. Heidelberger, P. Shahabuddin, and T. Zajic. Multilevel splitting for estimating rare event probabilities. *Oper. Res.*, 47(4):585–600, 1999.
- [GKvO02] M. J. J. Garvels, D. P. Kroese, and J. C. W. van Ommeren. On the importance function in splitting simulation. *European Transactions on Telecommunications*, 13(4):363–371, 2002.
- [KH51] H. Kahn and T. E. Harris. Estimation of particle transmission by random sampling. *National Bureau of Standards*, 12:27–30, 1951.
- [RT09] G. Rubino and B. Tuffin. Introduction to rare event simulation. In *Rare event simulation using Monte Carlo methods*, pages 1–13. Wiley, Chichester, 2009.
- [Sim14] E. Simonnet. Combinatorial analysis of the adaptive last particle method. *Statistics and Computing*, pages 1–20, 2014.
- [Ski06] J. Skilling. Nested sampling for general Bayesian computation. *Bayesian Anal.*, 1(4):833–859 (electronic), 2006.
- [Ski07] J. Skilling. Nested sampling for Bayesian computations. In *Bayesian statistics 8*, Oxford Sci. Publ., pages 491–524. Oxford Univ. Press, Oxford, 2007.
- [VAVA91] M. Villén-Altamirano and J. Villén-Altamirano. Restart: A method for accelerating rare events simulations. In *Proceeding of the thirteenth International Teletraffic Congress*, volume Copenhagen, Denmark, June 19-26 of *Queueing, performance and control in ATM: ITC-13 workshops*, pages 71–76. North-Holland, Amsterdam-New York, 1991.
- [VAVA94] M. Villén-Altamirano and J. Villén-Altamirano. Restart: a straightforward method for fast simulation of rare events. In *Proceedings of the 1994 Winter Simulation Conference*, volume Orlando 1994, December 1994, pages 282–289. 1994.

- [Wal14] C. Walter. Moving particles: a parallel optimal multilevel splitting method with applications in quantiles estimation and meta-model-based algorithms. *preprint*, 2014.