

m-ary search trees when $m \geq 27$: a strong asymptotics for the space requirements

(updated version, september 03)

BRIGITTE CHAUVIN & NICOLAS POUYANNE

Département de mathématiques
Université de Versailles - Saint-Quentin
45, avenue des Etats-Unis
78035 Versailles cedex

chauvin@math.uvsq.fr
pouyanne@math.uvsq.fr

Abstract

It is known that the joint distribution of the number of nodes of each type of an m -ary search tree is asymptotically multivariate normal when $m \leq 26$.

When $m \geq 27$, we show the following strong asymptotics of the random vector $X_n = (X_n^{(1)}, \dots, X_n^{(m-1)})$, where $X_n^{(i)}$ denotes the number of nodes containing $i - 1$ keys after having introduced $n - 1$ keys in the tree: there exist (nonrandom) vectors X , C and S and random variables ρ and φ such that

$$\frac{X_n - nX}{n^{\sigma_2}} - \rho(C \cos(\tau_2 \log n + \varphi) + S \sin(\tau_2 \log n + \varphi)) \xrightarrow[n \rightarrow \infty]{} 0$$

almost surely and in L^2 ; σ_2 and τ_2 denote the real and imaginary parts of one of the eigenvalues of the transition matrix, having the second greatest real part.

1 Introduction

An m -ary search tree is a data structure that grows by the progressive insertion of keys into a tree with branch factor m (first sentence in Lew's and Mahmoud's paper [8]). Each node of such a tree contains $0, 1, \dots$ or $m - 1$ keys and gives rise to m branches (see section 2 for the detailed definition of an m -ary search tree).

Our purpose is to make precise the asymptotic behaviour of the vector X_n whose coordinates are (with our notations) the number of nodes containing $0, \dots, m - 2$ keys in a random m -ary search tree holding $n - 1$ keys, as n tends to infinity.

Several cost measures on random m -ary search trees have been studied in the literature, one of them being frequently the total number of nodes S_n also

called space requirement; in the vectorial frame, just notice that S_n is an affine function of X_n . This cost is classically studied using generating functions and the method of moments. Mahmoud and Pittel ([11]) describe the asymptotics of the mean and the variance of S_n and derive a normal limit distribution for $m \leq 15$. Lew and Mahmoud ([8]) extend this range to $m \leq 26$. Smythe ([16]) and Mahmoud and Smythe ([12]) conjecture that the limit distribution is not normal for $m > 26$. In the related frame of branching processes, the change of normal limit laws to non-normal ones depends on the second eigenvalue of the transition matrix (which corresponds for m -ary search trees to the transition at $m = 26$) and already appears for instance in Athreya and Ney's book ([1]). It has been often noticed by the previous authors dealing with m -ary search trees (see for instance exemple 3.1 in Smythe [16]).

The state of the art can be found in Chern and Hwang's paper ([3]): a phase transition occurs between $m = 26$ and $m = 27$. It is quite readable on the variance of the space requirement, the asymptotics of which having two types of behaviour depending on the values of m : for small m ($m \leq 26$), the variance is of order n and the rescaled space requirement is asymptotically normal, but for $m \geq 27$, the variance is of order $n^{2\sigma}$ for some (known) real number σ , $\sigma > 1/2$ and a periodic phenomenon appears.

In the range $m \geq 27$, the challenge comes from the questions asked by Chern and Hwang: they prove (in [3], Corollary 2) that the distribution of S_n , even conveniently renormalized, does not approach any fixed distribution function but fluctuates via some periodic function. They ask for more intuitive explanations of the phase transition than pure analytic reasons.

The asymptotic normality for $m \leq 26$ can also be found by contraction method (see for instance Neininger and Rueschendorf [14]). Interestingly, the same phase transition for the variance is noticed by physicists in the close context of random fragmentation problem (for instance in Dean and Majumdar [4]).

The literature on the subject, including limit distribution results by contraction method, mostly takes advantage of the "divide-and-conquer" recursivity (sometimes called "backward" method). Another point of view on these processes is based on the dynamical recursivity (sometimes called "forward" method), already used in Smythe's ([16]) and Mahmoud's and Smythe's ([12]) papers.

We consider $(X_n)_{n \geq 1}$ as a Markov process, and we notice that X_n is a kind of Pólya urn model, or a random walk or a multitype branching process, depending on one's favourite background.

In section 2 and 3, we see how $(X_n)_{n \geq 1}$ can be viewed as a Markovian process with values in \mathbb{R}^{m-1} and that its evolution is driven by a transition-type matrix A in the following remarkable (since *linear*) way:

$$E^{\mathcal{F}_n}(X_{n+1}) = \left(\text{Id} + \frac{A}{n} \right) X_n, \quad (1)$$

where \mathcal{F}_n is the past before time n and Id is the identity matrix. Our method is based on exploiting the linearity of this evolution.

Thus \mathbb{C}^{m-1} is decomposed along the eigenspaces of A and if $\text{Sp}(A)$ denotes the set of eigenvalues of A (all its eigenvalues are simple), we have

$$\mathbb{C}^{m-1} = \bigoplus_{\lambda \in \text{Sp}(A)} \ker(A - \lambda \text{Id}) \quad (2)$$

$$\text{Id} = \sum_{\lambda \in \text{Sp}(A)} \pi_\lambda$$

$$A = \sum_{\lambda \in \text{Sp}(A)} \lambda \pi_\lambda$$

where π_λ denotes the projection on the eigenspace $\ker(A - \lambda \text{Id})$ relatively to the decomposition (2). Moreover, 1 is an eigenvalue, the other ones having a real part strictly less than 1. If λ_2 and $\overline{\lambda_2}$ are the eigenvalues having the greatest real part, say $\sigma_2, \sigma_2 < 1$, we write the following fundamental decomposition of vector X_n (\overline{x} denotes the conjugate of a complex number x):

$$X_n = \pi_1 X_n + \pi_{\lambda_2} X_n + \pi_{\overline{\lambda_2}} X_n + \sum_{\lambda \neq 1, \lambda_2, \overline{\lambda_2}} \pi_\lambda X_n. \quad (3)$$

This spectral decomposition of X_n coincides, as $m \geq 27$, with the almost sure asymptotic expansion of X_n for the first three terms; it is a key phenomenon. For this purpose, the analysis of each projection $\pi_\lambda X_n$ is performed by rescaling it in order to get a martingale. Notice that appearance of martingale methods is not surprising, considering the evolution given by formula (1). The result then comes from the spectral decomposition (3) and from the lemmas in section 4 explaining successively that the first projection is of order n , the projections $\pi_\lambda X_n$ for $\Re(\lambda) > 1/2$ are of order n^λ by a L^2 -convergence theorem of martingales and the remaining projections $\pi_\lambda X_n$ for $\Re(\lambda) \leq 1/2$ are asymptotically almost surely negligible. One can find the complete theorem with its proof in section 5. Simulations in section 6 help to visualize the phenomena.

Notice that our approach is somehow complementary to Mahmoud's one in his recent paper [10] where the frame (Pólya schemes) is quite large, including m -ary search trees, and focuses on the leading term of X_n ; our study goes further in the expansion of X_n but is restricted here to m -ary search trees.

Using similar arguments, we hope that the asymptotics of the “profile” (meaning the number of nodes level by level in the tree) of an m -ary search tree is tractable: a natural generalization of the binary search tree case ([6]) to higher dimensions would consist in considering the number of nodes of each type level by level, and introducing some “level polynomial” vectors. It will be the subject of a forthcoming paper.

2 Definition and Markovianity of the process

One throws a sequence of numbers in $[0, 1]$, named the *keys*, uniformly in $[0, 1]^{\mathbb{N}^*}$. The keys are placed one after another in an m -ary tree (one node-root, from

each node grow m branches). The following recursive rule describes the way a key named k is inserted in the tree.

i) If the root contains strictly less than $m - 1$ keys, then k is inserted in the root. One draws usually keys in a root from left to right in increasing order.

ii) If the root is already saturated, *i.e.* if it contains $m - 1$ keys named k_1, \dots, k_{m-1} , ordered such that $k_i < k_{i+1}$, then corresponds to each interval $I_1 =]-\infty, k_1[$, $I_{j+1} =]k_j, k_{j+1}[$ ($1 \leq j \leq m - 2$), $I_m =]k_{m-1}, +\infty[$ a subtree being itself an m -ary search tree. One draws usually the branches corresponding to I_1, \dots, I_m from left to right. In this situation, k is inserted in the subtree that corresponds to the interval I_j such that $k \in I_j$.¹

Figure 1 is an example of 4-ary tree obtained by insertion of the numbers 0.3, 0.1, 0.4, 0.15, 0.9, 0.2, 0.6, 0.5, 0.35, 0.8, 0.97, 0.93, 0.23, 0.84, 0.62, 0.64, 0.33, 0.83 in this order.

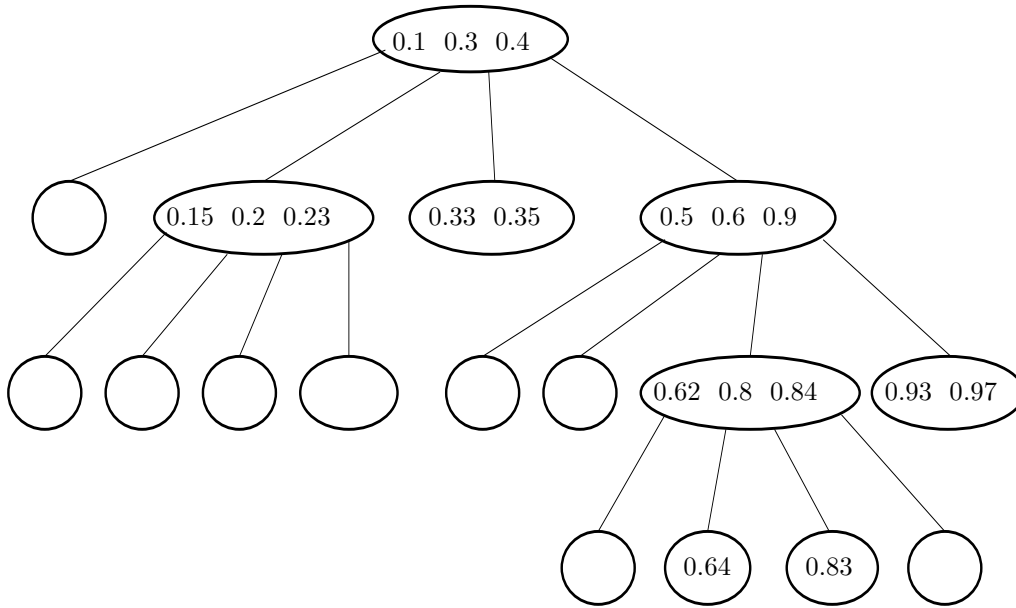


Figure 1: insertion of the keys 0.3, 0.1, 0.4, 0.15, 0.9, 0.2, 0.6, 0.5, 0.35, 0.8, 0.97, 0.93, 0.23, 0.84, 0.62, 0.64, 0.33, 0.83 in a 4-ary tree

Although it is not explicitly used later on, let us mention (see Mahmoud's book [9] for details) that such a sequence $(T_n)_{n \in \mathbb{N}}$ of trees has the same distribution as the one obtained by construction of T_n from a random permutation

¹In this paper, our convention is that empty nodes (corresponding to the m above intervals) appear when the concerned internal node has just been saturated by the insertion of a $m - 1$ -st key. Other conventions are possible, for instance, empty nodes could appear once the first key is stored in the concerned internal node. Anyway, the choice of any convention has no impact on the results.

of n integers, with a uniform distribution on the set of permutations. It is the so-called random permutation model.

In the sequel, $(T_n)_{n \in \mathbb{N}}$ and other parameters of interest are random variables on the space Ω of infinite m -ary trees ². The space is endowed with the natural filtration $(\mathcal{F}_n)_{n \in \mathbb{N}}$, \mathcal{F}_n being the σ -field generated up to time n .

For each $i = \{1, \dots, m\}$ and $n \geq 1$, we define the number $X_n^{(i)}$ as the number of nodes which contain $i-1$ keys after insertion of the $(n-1)$ -st key; such nodes are named nodes of type i . The question consists in describing the asymptotic behaviour of the $X_n^{(i)}$ as n tends to infinity.

Counting the total number of keys in nodes of each type in the tree holding $(n-1)$ keys leads to the formula

$$n-1 = \sum_{i=1}^m (i-1)X_n^{(i)}. \quad (4)$$

This formula which binds the $X_n^{(i)}$ allows to limit the study to the $m-1$ first indices i . It should not be confused with the relation (8) later on which counts the number of free places (or gaps) in the tree.

In the whole paper, we call

$$V = \mathbb{R}^{m-1}$$

(or more exactly the real vector space of matrices having one column and $m-1$ rows). The random vector $X_n \in V$ is defined for all $n \geq 1$ as

$$X_n = \begin{pmatrix} X_n^{(1)} \\ \vdots \\ X_n^{(m-1)} \end{pmatrix},$$

and evolves as follows. The first $m+1$ vectors X_n are nonrandom:

$$\begin{cases} X_1 = {}^t(1, 0, \dots, 0) \\ X_2 = {}^t(0, 1, \dots, 0) \\ \vdots \\ X_{m-1} = {}^t(0, \dots, 0, 1) \\ X_m = {}^t(m, 0, \dots, 0) \\ X_{m+1} = {}^t(m-1, 1, 0, \dots). \end{cases}$$

The following ones are random. For instance, $X_{m+2} = {}^t(m-2, 2, 0, \dots)$ with probability $\frac{m-1}{m+1}$, and $X_{m+2} = {}^t(m-1, 0, 1, 0, \dots)$ with probability $\frac{2}{m+1}$. These probabilities are computed with the rules of the random permutation model:

²It is necessary to define random variables on this (big) probability space in order to give a meaning to almost sure and L^2 -convergences later on. For tree probability spaces, see for instance Neveu ([15]).

when $n - 1$ keys are inserted, the probability that the n -th one falls between two of them is $1/n$ (the probability that it falls on the left-hand side of the smallest one or on the right-hand side of the greatest one is $1/n$ too). Consequently, only the relative order of the keys is taken into account (not their values).

More generally, the transition rules between the states at time n and $n + 1$ are the following: for each i between 1 and $m - 1$, if the n -th key falls on a node of type i , then

$$X_{n+1} = X_n + \Delta_i,$$

where

$$\left\{ \begin{array}{l} \Delta_1 = {}^t(-1, 1, 0, 0, \dots) \\ \Delta_2 = {}^t(0, -1, 1, 0, \dots) \\ \vdots \\ \Delta_{m-2} = {}^t(0, \dots, 0, -1, 1) \\ \Delta_{m-1} = {}^t(m, 0, \dots, 0, -1) \end{array} \right. ,$$

and this event takes place with probability $\frac{i}{n} X_n^{(i)}$ because each node of type i contains i free places.

Let us emphasize here that this last probability, containing the randomness of the evolution of the process is *linear* in X_n . For this reason, for each $i \in \{1, \dots, m - 1\}$, let l_i be the *linear* form of V defined as

$$l_i = i dx_i,$$

where dx_i is the i -th coordinate form of $V = \mathbb{R}^{m-1}$. The process $(X_n)_n$ in V is now defined by the first vector X_1 and the transition condition for each $n \geq 1$:

$$\left\{ \begin{array}{l} X_{n+1} = X_n + \Delta_1, \text{ with probability } \frac{1}{n} l_1(X_n), \\ \vdots \\ X_{n+1} = X_n + \Delta_{m-1}, \text{ with probability } \frac{1}{n} l_{m-1}(X_n). \end{array} \right. \quad (5)$$

In other words, the process is a random walk in V defined by X_1 and a random increment $\Delta(n + 1)$ between times n and $n + 1$:

$$X_{n+1} = X_n + \Delta(n + 1), \quad (6)$$

with the transition probabilities

$$P(\Delta(n + 1) = \Delta_i | X_n) = \frac{1}{n} l_i(X_n), \quad 1 \leq i \leq m - 1. \quad (7)$$

Note that the process (X_n) satisfies the relation

$$\sum_{i=1}^{m-1} i X_n^{(i)} = \sum_{i=1}^{m-1} l_i(X_n) = n \quad (8)$$

available for each $n \geq 1$, meaning that the numbers $l_i(X_n)/n$ are probabilities of disjoint events whose union is the total probability space. The interpretation of

iii) *Note on this kind of process.* The above random walk of an m -ary search tree belongs to a larger family of vector processes $(Z_n)_n$ in \mathbb{R}^s (for any integer $s \geq 1$). Such a process can be defined as a random walk starting from some $Z_1 \in \mathbb{R}^s$, with random increments which take their values in a finite set of vectors $\{\Delta_1, \dots, \Delta_s\}$:

$$\forall n \geq 1, Z_{n+1} = Z_n + \Delta(n+1),$$

with the transition probabilities

$$\forall n \geq 1, P(\Delta(n+1) = \Delta_i | Z_n) = \frac{1}{n} l_i(Z_n), \quad 1 \leq i \leq s$$

where the l_i 's are linear forms on \mathbb{R}^s . The process is Markovian and the transition probabilities between time n and time $n+1$ depend linearly on the state at time n .

In order to guarantee that such a process is well defined, that is to say that the numbers $l_i(Z_n)/n$ are almost surely nonnegative and that their sum equals 1 for all n , one needs further assumptions on the parameters, namely on Z_1 , the l_i and the Δ_i (all these assumptions are satisfied by m -ary search trees). First, hypotheses that allow Z_2 to be well defined:

$$\sum_{i=1}^s l_i(Z_1) = 1 \quad \text{and} \quad \forall j \in \{1, \dots, s\}, l_j(Z_1) \geq 0;$$

then the hypotheses on the increments (an elementary induction shows that they are enough to make sure that the process is well defined): for all $j, k \in \{1, \dots, s\}$,

$$\begin{cases} \sum_{i=1}^s l_i(\Delta_j) = 1, \\ j \neq k \implies l_j(\Delta_k) \geq 0, \\ l_j(\Delta_j) = 0 \quad \text{or} \quad l_j(Z_1)Z + \sum_{i=1}^s l_j(\Delta_i)Z = l_j(\Delta_j)Z. \end{cases}$$

Only the diagonal terms $l_j(\Delta_j)$ are allowed to be negative. The last arithmetical condition just indicates that if $l_j(\Delta_j)$ is nonzero for some j , it divides (as real number) $l_j(Z_1)$ and all the $l_j(\Delta_i)$.

The conditions defining such a model remain stable after an invertible linear change of coordinates. Keeping in mind remark i), it means that these conditions are sufficient to guarantee that the corresponding generalized Pólya urn is tenable. The choice of a good basis of V is the key point of what follows.

3 Evolution of the process and average-case analysis

Both are based on the computation of the conditional expectation:

$$E^{\mathcal{F}_n}(X_{n+1}) = \sum_{i=1}^{m-1} \frac{1}{n} l_i(X_n)(X_n + \Delta_i).$$

If one denotes by A the endomorphism of V defined by

$$\forall Z \in V, AZ = \sum_{i=1}^{m-1} l_i(Z)\Delta_i,$$

one gets the following formula which makes precise that the above conditional expectation is a linear function of the state at time n :

$$E^{\mathcal{F}_n}(X_{n+1}) = \left(\text{Id}_V + \frac{A}{n} \right) X_n \quad (10)$$

where Id_V is the identity map of V . An immediate consequence of this fact is the computation of the expectation of the random vector X_n : define $\Gamma_1 = \text{Id}_V$ and

$$\Gamma_n = \prod_{k=1}^{n-1} \left(\text{Id}_V + \frac{A}{k} \right)$$

for all $n \geq 2$, so that one gets the expression

$$E(X_n) = \Gamma_n X_1.$$

In the canonical basis of V , the matrix of A is

$$A = \begin{pmatrix} -1 & & & & & & & & m(m-1) \\ 1 & -2 & & & & & & & \\ & 2 & -3 & & & & & & \\ & & \ddots & \ddots & & & & & \\ & & & \ddots & -(m-2) & & & & \\ & & & & m-2 & -(m-1) & & & \end{pmatrix}, \quad (11)$$

where an empty entry means a zero entry. This matrix A is the *transition matrix* (or *endomorphism*) of the process. The characteristic polynomial of A is

$$\chi_A(z) = \prod_{k=1}^{m-1} (z+k) - m!. \quad (12)$$

The matrix A has only simple (complex) eigenvalues and 1 is the eigenvalue having the greatest real part. Furthermore, when m is even, 1 is the only real eigenvalue; when m is odd, the only other real eigenvalue is $-m-1$. Figure 2, made with the help of Maple, shows the complex eigenvalues of A when m equals 50. The plot of all roots of A in the complex plane seems to have always the same shape: regularly spaced points on the algebraic curve defined by equation $\prod_{1 \leq k \leq m-1} |z+k|^2 = (m!)^2$. An important fact for the sequel is that **all the eigenvalues different from 1 have a real part less than 1/2 if and only if $m \leq 26$.**

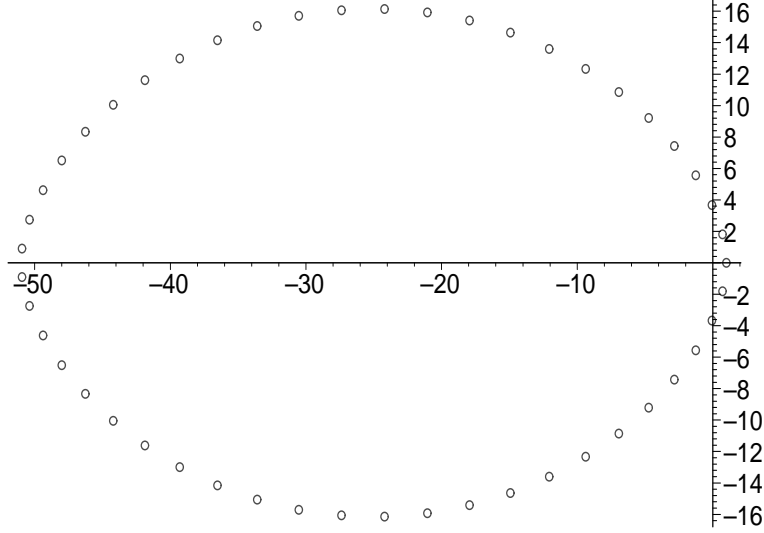


Figure 2: roots of χ_A when $m = 50$

We denote by $\text{Sp}(A)$ the set of (complex) eigenvalues of A , and

$$V_{\mathbb{C}} = \mathbb{C}^{m-1},$$

or more exactly $V_{\mathbb{C}} = V \otimes_{\mathbb{R}} \mathbb{C}$. For every $\lambda \in \text{Sp}(A)$, we denote by π_{λ} the projection of $V_{\mathbb{C}}$ on the eigenspace $\ker(A - \lambda \text{Id}_{V_{\mathbb{C}}})$ relatively to the decomposition

$$V_{\mathbb{C}} = \bigoplus_{\mu \in \text{Sp}(A)} \ker(A - \mu \text{Id}_{V_{\mathbb{C}}}).$$

Then, define $\gamma_1(\lambda) = 1$ and

$$\gamma_n(\lambda) = \prod_{k=1}^{n-1} \left(1 + \frac{\lambda}{k}\right)$$

for all $n \geq 2$, so that the endomorphism Γ_n splits into the sum $\Gamma_n = \sum_{\lambda} \gamma_n(\lambda) \pi_{\lambda}$, and the expectation of X_n equals $E(X_n) = \sum_{\lambda} \gamma_n(\lambda) \pi_{\lambda} X_1$. Since 1 is the eigenvalue of A having the greatest real part, and since by Stirling formula

$$\gamma_n(\lambda) = \frac{\Gamma(n + \lambda)}{\Gamma(\lambda + 1)\Gamma(n)} = \frac{n^{\lambda}}{\Gamma(\lambda + 1)} + O(n^{\lambda-1}) \quad (13)$$

as n tends to infinity, the first term in the above expansion of $E(X_n)$ is $\gamma_n(1)\pi_1 X_1$ and one gets

$$\lim_{n \rightarrow \infty} \frac{E(X_n)}{n} = \pi_1 X_1. \quad (14)$$

Note that this limit is nonzero, since otherwise $E(X_n) = o(n)$ and taking expectation in formula (8) provides a contradiction. The coordinates of the vector

$\pi_1 X_1$ are explicitly given in section 5 (see (17)) proving once again that $\pi_1 X_1$ is nonzero.

4 Local study along principal directions

Keeping in mind the spectral decomposition of the process

$$X_n = \pi_1 X_n + \pi_{\lambda_2} X_n + \pi_{\overline{\lambda_2}} X_n + \sum_{\lambda \neq 1, \lambda_2, \overline{\lambda_2}} \pi_\lambda X_n ,$$

we achieve in this section the local study of the projections $\pi_\lambda X_n$ for every eigenvalue λ of the transition matrix A , via three lemmas.

The first lemma describes explicitly the projection π_1 on the fixed points of the matrix A . It is available for every process as the one defined at the end of section 2 (see remark iii)) as soon as 1 is a *simple* eigenvalue of A . It is applied for the final result to $Y = X_n$ and $\pi_1 X_1$ is explicitly computed in section 5, thus giving the first term of the expansion of X_n .

Lemma 1 (First projection lemma)

$$\forall Y \in V, \quad \pi_1 Y = \left(\sum_{j=1}^{m-1} l_j(Y) \right) \pi_1 X_1 .$$

Proof of lemma 1. Let L be the endomorphism of V defined for all Y in V by

$$LY = \left(\sum_{j=1}^{m-1} l_j(Y) \right) \pi_1 X_1 .$$

Note that L is nonzero because $\pi_1 X_1 \neq 0$ (recall the end of the previous section). Because of the relation $\sum_j l_j(\Delta_i) = 1$ for all i (see (9)), the value of L at each Δ_i is $\pi_1 X_1$. Thus for all Y in V ,

$$LA(Y) = L \left(\sum_{i=1}^{m-1} l_i(Y) \Delta_i \right) = \sum_{i=1}^{m-1} l_i(Y) \pi_1 X_1 = L(Y) ,$$

hence $LA = L$. But since $\pi_1 X_1$ is a fixed point of A , one has $AL = L$ too. Then A and L commute, and this product is L .

Since π_1 is a polynomial in A , the endomorphisms π_1 and L commute. Because $\sum_i l_i(X_1) = 1$ (relation (8) for $n = 1$), the endomorphisms π_1 and $L\pi_1$ (and $\pi_1 L$ too) have the same value at X_1 . Since they are zero on the hyperplane spanned by the eigenvectors associated to eigenvalues different from 1, they are equal. Then, $\pi_1 = L\pi_1 = \pi_1 L$. But $\pi_1 L = L$, obviously. Thus $\pi_1 = L$. \square

Lemma 2 gives the asymptotics of the moments of all projections $\pi_\lambda X_n$. There appear the different behaviours of these moments, depending on the position of the real part of λ with respect to $1/2$. This result is to be compared

with similar ones for the second moments in the literature and contains the technical reason of the *phase transition* mentioned by these authors. It is available for every process defined at the end of section 2. Notice that the L^2 -convergence in the final result only requires the second moment asymptotics, but the almost sure convergence comes from the higher moments asymptotics.

Lemma 2 (Moments lemma) *Let λ be an eigenvalue of A , σ its real part, $(\cdot|\cdot)$ any positive definite Hermitian form on $V_{\mathbb{C}}$, and $Z \in V_{\mathbb{C}}$. Then, for every nonnegative integer p , if $\sigma \neq 1/2$,*

$$\begin{cases} E(|(Z|\pi_{\lambda}X_n)|^{2p}) = O(n^p + n^{2p\sigma}) \\ E(|(Z|\pi_{\lambda}X_n)|^{2p+1}) = O(n^{p+1} + n^{2p\sigma+1}) \end{cases}$$

as n tends to infinity. If $\sigma = 1/2$,

$$\begin{cases} E(|(Z|\pi_{\lambda}X_n)|^{2p}) = O(n^p \log n) \\ E(|(Z|\pi_{\lambda}X_n)|^{2p+1}) = O(n^{p+1} \log n) \end{cases}$$

as n tends to infinity.

In other words, as n tends to infinity,

$$\begin{aligned} \text{if } \Re(\lambda) < 1/2 \text{ then } & \begin{cases} E(|(Z|\pi_{\lambda}X_n)|^{2p}) = O(n^p) \\ E(|(Z|\pi_{\lambda}X_n)|^{2p+1}) = O(n^{p+1}) \end{cases} \\ \text{if } \Re(\lambda) = 1/2 \text{ then } & \begin{cases} E(|(Z|\pi_{\lambda}X_n)|^{2p}) = O(n^p \log n) \\ E(|(Z|\pi_{\lambda}X_n)|^{2p+1}) = O(n^{p+1} \log n) \end{cases} \\ \text{and if } \Re(\lambda) > 1/2 \text{ then } & \begin{cases} E(|(Z|\pi_{\lambda}X_n)|^{2p}) = O(n^{2p\sigma}) \\ E(|(Z|\pi_{\lambda}X_n)|^{2p+1}) = O(n^{2p\sigma+1}). \end{cases} \end{aligned}$$

Remark. Note that we do not know if some value of m leads to $\Re(\lambda) = 1/2$ for some eigenvalue λ . It has no consequence on the final result.

Proof of lemma 2. By induction on the integer p . If $p = 0$, only the assertion on the moment of order $2p + 1$ is nontrivial. If $\|\cdot\|$ denotes the norm associated to the Hermitian form, it follows directly from the definition of the process $(X_n)_n$ that almost surely

$$\|X_{n+1}\| \leq \|X_n\| + \max_{1 \leq i \leq m-1} \|\Delta_i\|$$

for every $n \geq 1$. Therefore, there is some positive constant c depending only on m such that almost surely, for every $n \geq 1$,

$$\|X_n\| \leq cn. \tag{15}$$

The result for $p = 0$ follows from this inequality.

Although it is not needed to make the proof complete, we prove the second moments inequality before presenting the induction; it helps the understanding of the general case, and it is used several times later on. An elementary computation of the conditional expectation, based on the dynamics of the process (X_n) leads to:

$$\begin{aligned}
E^{\mathcal{F}_n} (|(Z|X_{n+1})|^2) &= \sum_{i=1}^{m-1} \frac{1}{n} l_i(X_n) (Z|X_n + \Delta_i) \overline{(Z|X_n + \Delta_i)} \\
&= |(Z|X_n)|^2 + 2\Re \left[\sum_{i=1}^{m-1} \frac{1}{n} l_i(X_n) (Z|X_n) \overline{(Z|\Delta_i)} \right] + \sum_{i=1}^{m-1} \frac{1}{n} l_i(X_n) |(Z|\Delta_i)|^2 \\
&= \Re \left[(Z|X_n) \overline{(Z|(I + \frac{2A}{n})X_n)} \right] + \sum_{i=1}^{m-1} \frac{1}{n} l_i(X_n) |(Z|\Delta_i)|^2.
\end{aligned}$$

Take now the expectation and apply this formula to the vector $\pi_\lambda^* Z$ where u^* denotes the adjoint endomorphism of u relative to the positive definite Hermitian form $(\cdot|\cdot)$. If σ is the real part of λ , one gets

$$E (|(Z|\pi_\lambda X_{n+1})|^2) = \left(1 + \frac{2\sigma}{n}\right) E (|(Z|\pi_\lambda X_n)|^2) + b_n$$

where $b_n = \sum_i l_i(E X_n/n) |(Z|\pi_\lambda \Delta_i)|^2$, the sum being extended to all i between 1 and $m-1$. Since b_n has a limit as n tends to infinity (see (14)), $b_n = O(1)$. We get the explicit form

$$E (|(Z|\pi_\lambda X_n)|^2) = \gamma_n(2\sigma) \left(|(Z|\pi_\lambda X_1)|^2 + \sum_{k=1}^{n-1} \frac{b_k}{\gamma_{k+1}(2\sigma)} \right)$$

and since by (13),

$$\gamma_n(2\sigma) = \frac{n^{2\sigma}}{\Gamma(1+2\sigma)} + O(n^{2\sigma-1}),$$

the above series has not the same behaviour depending on the position of 2σ with respect to 1. This shows the following second moments asymptotics:

$$E (|(Z|\pi_\lambda X_n)|^2) = \begin{cases} O(n) & \text{if } \sigma < 1/2 \\ O(n \log n) & \text{if } \sigma = 1/2 \\ O(n^{2\sigma}) & \text{if } \sigma > 1/2. \end{cases} \quad (16)$$

Suppose now $p \geq 1$. On one hand, if x and y are complex numbers, the binomial formula implies that $|x+y|^{2p} = |x^p + px^{p-1}y + z|^2$ where z is a polynomial in x and y whose degree in x equals $p-2$. Thus

$$|x+y|^{2p} \leq |x|^{2p-2} \Re[x \overline{(x+2py)}] + P(|x|, |y|)$$

where $P(X, Y)$ is a polynomial whose degree in X does not exceed $2p-2$. On the other hand, the inequality (15) provides a positive constant (depending only on m) which bounds from above the number $|l_i(X_n)/n|$ for all i and for all n . The use of the last two facts to bound from above the conditional expectation

$$E^{\mathcal{F}_n} (|(Z|X_{n+1})|^{2p}) = \sum_{i=1}^{m-1} \frac{1}{n} l_i(X_n) |(Z|X_n) + (Z|\Delta_i)|^{2p}$$

leads to the existence of a polynomial Q of degree $\leq 2p - 2$ such that for every $n \geq 1$,

$$E^{\mathcal{F}_n} \left(|(Z|X_{n+1})|^{2p} \right) \leq |(Z|X_n)|^{2p-2} \Re \left[\overline{(Z|X_n)} \left(Z \left| I + \frac{2pA}{n} \right) X_n \right) \right] + Q(|(Z|X_n)|).$$

Now, the same arguments as in the preceding proof for the second moments allow to show the inequality

$$E \left(|(Z|\pi_\lambda X_{n+1})|^{2p} \right) \leq \left(1 + \frac{2p\sigma}{n} \right) E \left(|(Z|\pi_\lambda X_n)|^{2p} \right) + EQ(|(Z|\pi_\lambda X_n)|),$$

which gives the result by induction, assuming the result for all integers $< 2p$. Using (15), the result for the moments of order $2p + 1$ is a straightforward consequence of the inequality

$$E |(Z|X_n)|^{2p+1} \leq \max_{\Omega} |(Z|X_n)| \times E |(Z|X_n)|^{2p},$$

where Ω is the underlying probability space (see beginning of section 2). \square

Lemma 3 is a direct consequence of lemma 2. It makes precise the convergence of the martingale associated to X_n after rescaling with relation (1), establishing that some projections $\pi_\lambda X_n$ have an L^2 and a.s. limit. This lemma is applied for the final result to the first terms of the spectral decomposition of X_n .

Lemma 3 (L^2 -convergence lemma) *Let λ be an eigenvalue of A . If $\Re(\lambda) > 1/2$, then the martingale $\gamma_n^{-1}(\lambda)\pi_\lambda X_n$ converges in L^2 (thus almost surely).*

Proof of lemma 3. The random vector $\gamma_n^{-1}(\lambda)\pi_\lambda X_n$ is a \mathcal{F}_n martingale due to equation (1) and because the restriction of Γ_n to the image of π_λ is the multiplication by $\gamma_n(\lambda)$. Moreover, under the hypothesis on $\Re(\lambda)$, estimation (16) on the second moments implies that $E(\|\pi_\lambda X_n\|^2) = O(n^{2\sigma})$; indeed, it is enough to choose a suitable (orthonormal) basis $(Z_i)_{1 \leq i \leq m-1}$ of $V_{\mathbb{C}}$ such that

$$\|\pi_\lambda X_n\|^2 = \sum_{i=1}^{m-1} |(Z_i|\pi_\lambda X_n)|^2$$

and apply lemma 2 to each vector Z_i . Combining this with (13) gives that

$$E(\|\gamma_n^{-1}(\lambda)\pi_\lambda X_n\|^2) = |\gamma_n^{-2}(\lambda)| E(\|\pi_\lambda X_n\|^2)$$

is a bounded sequence indexed by n so that the martingale $\gamma_n^{-1}(\lambda)\pi_\lambda X_n$ converges in L^2 and thus almost surely by standard theorems on martingales. \square

5 Theorem

Theorem Assume $m \geq 27$. Let $\lambda_2 = \sigma_2 + i\tau_2$ be the eigenvalue of the transition matrix A , having the second greatest real part σ_2 ($\sigma_2 > 1/2$) and a positive imaginary part $\tau_2 > 0$.

For every eigenvalue λ of the transition matrix A , let π_λ be the projection on the eigenspace $\ker(A - \lambda \text{Id})$ associated to λ , relatively to the decomposition of $V_{\mathbb{C}}$ in eigenspaces of A . Let $X := \pi_1 X_1$.

1.

$$X = \lim_{n \rightarrow \infty} \frac{EX_n}{n} = \frac{1}{H_m - 1} \begin{pmatrix} \frac{1}{1 \times 2} \\ \frac{1}{2 \times 3} \\ \vdots \\ \frac{1}{(m-1) \times m} \end{pmatrix} \quad (17)$$

where H_m is the harmonic sum $H_m = \sum_{1 \leq k \leq m} 1/k$.

2. If Λ denotes the limit of the L^2 -convergent martingale $\gamma_n^{-1}(\lambda_2)\pi_{\lambda_2}X_n$, then

$$X_n = nX + 2\Re \left[\frac{n^{\lambda_2} \Lambda}{\Gamma(1 + \lambda_2)} \right] + n^{\sigma_2} \varepsilon_n \quad (18)$$

where the random vector ε_n converges to zero almost surely and in L^2 as n tends to infinity.

Corollary 1 With the same notations as in the theorem, let C and S be the real (and nonrandom) vectors of $V_{\mathbb{C}}$ defined by the relation

$$\pi_{\lambda_2} X_1 = C - iS . \quad (19)$$

Let ρ and φ be respectively the modulus and the argument of the random vector $\frac{2\Lambda}{\Gamma(1+\lambda_2)}$ along the line generated by $\pi_{\lambda_2} X_1$:

$$\rho \exp(i\varphi) \pi_{\lambda_2} X_1 = \frac{2\Lambda}{\Gamma(1 + \lambda_2)}, \quad \rho \geq 0, \varphi \in [0, 2\pi[. \quad (20)$$

Then

$$X_n = nX + n^{\sigma_2} \rho (C \cos(\tau_2 \log n + \varphi) + S \sin(\tau_2 \log n + \varphi)) + n^{\sigma_2} \varepsilon_n,$$

where the random vector ε_n converges to zero almost surely and in L^2 as n tends to infinity.

In other words, the random vector

$$\frac{X_n - nX}{n^{\sigma_2}} - \rho (C \cos(\tau_2 \log n + \varphi) + S \sin(\tau_2 \log n + \varphi))$$

converges to zero almost surely and in L^2 .

The corollary is a straightforward consequence of the theorem. Just write Λ in (20) as the product of a complex random variable and of the nonrandom complex vector $\pi_{\lambda_2} X_1$, and separate the real and imaginary parts of $\pi_{\lambda_2} X_1$ ((19)). Also note that $n^{\lambda_2} = n^{\sigma_2} e^{i\tau_2 \log n}$. Notice that X , C and S are linearly independent vectors of $V_{\mathbb{C}}$ (because $\pi_{\lambda_1} X_1$, $\pi_{\lambda_2} X_1$ and $\pi_{\overline{\lambda_2}} X_1$ are).

Computation of X , C , S . X is the projection of the first vector X_1 on the vector line of the fixed vectors of A . The first equality of (17) has already been shown (see (14)). An easy computation (compute a fixed vector, and add the condition $\sum_i l_i(\pi_1 X_1) = \lim_n \sum_i l_i(E(X_n/n)) = 1$) shows (17).

To express the vectors C and S we sum up how one can compute the projection $\pi_{\lambda} X_1$ of X_1 on the eigenspace $\ker(A - \lambda \text{Id})$ for every eigenvalue λ , and give the result: for each λ , compute first the eigenvector of A associated to λ having 1 as $(m-1)$ -st coordinate. Name it F_{λ} . Decompose $X_1 = \sum_{\lambda \in \text{Sp}(A)} a_{\lambda} F_{\lambda}$ where a_{λ} is the complex number such that $\pi_{\lambda} X_1 = a_{\lambda} F_{\lambda}$. Then, for all $p \geq 0$, one has $A^p X_1 = \sum_{\lambda} a_{\lambda} \lambda^p F_{\lambda}$. With the explicit form of A , one can easily compute the $(m-1)$ -st coordinate of the vectors $A^p X_1$ for $0 \leq p \leq m-2$ (an induction shows that its p -th coordinate is $p!$, its j -th ones are zero for all $j \geq p+1$) and solve the system

$$dx_{m-1} A^p X_1 = \sum_{\lambda \in \text{Sp}(A)} a_{\lambda} \lambda^p, \quad 0 \leq p \leq m-2$$

with Cramer's formula; one writes this way the number a_{λ} as the product of $(m-2)!$ by the quotient of two Vandermonde determinants. After simplification, one gets $a_{\lambda} = (m-2)!/\chi'_A(\lambda)$ where χ_A denotes the characteristic polynomial of A (see (12)). The computation of the logarithmic derivative of $\chi_A + m!$ gives rise to the expression

$$\chi'_A(\lambda) = \prod_{\mu \in \text{Sp}(A) \setminus \{\lambda\}} (\lambda - \mu) = m! \sum_{j=1}^{m-1} \frac{1}{\lambda + j}.$$

The result is now the following: for every eigenvalue $\lambda \in \text{Sp}(A)$,

$$\pi_{\lambda} X_1 = \frac{1}{\chi'_A(\lambda)} \begin{pmatrix} \varpi_1(\lambda) \\ \vdots \\ \varpi_{m-1}(\lambda) \end{pmatrix},$$

where, for every $j \in \{1, \dots, m-1\}$,

$$\varpi_j(\lambda) = (j-1)! \prod_{k=j+1}^{m-1} (k+\lambda) = (j-1)! \frac{\Gamma(m+\lambda)}{\Gamma(j+1+\lambda)} = \frac{m!}{j \gamma_{j+1}(\lambda)}.$$

Proof of the theorem. The proof consists in examining the decomposition

$$X_n = \pi_1 X_n + \pi_{\lambda_2} X_n + \pi_{\overline{\lambda_2}} X_n + \sum_{\Re(\lambda) < \sigma_2} \pi_\lambda X_n, \quad (21)$$

in order to get the expected asymptotic order of magnitude of each term.

The first projection lemma describes the first term, because relation (8) between the number of nodes of each type gives that for every n

$$\sum_{i=1}^{m-1} l_i(X_n) = n,$$

so that

$$\pi_1 X_n = n \pi_1 X_1. \quad (22)$$

For the following two terms in (21), recall that the assumption $m \geq 27$ implies that $\sigma_2 > 1/2$. Let

$$\Lambda = \lim_{n \rightarrow +\infty} \gamma_n^{-1}(\lambda_2) \pi_{\lambda_2} X_n$$

and notice that the random vector Λ is both the L^2 and the almost sure limit of this martingale as guaranteed by the L^2 -convergence lemma (lemma 3) applied to λ_2 . In other words,

$$\gamma_n^{-1}(\lambda_2) \pi_{\lambda_2} X_n = \Lambda + \varepsilon_n$$

where ε_n tends to zero a.s. and in L^2 . In the following, ε_n always denotes a generic random variable which tends to zero a.s. and in L^2 when n tends to infinity, even if it changes from place to place.

Multiply by $\gamma_n(\lambda_2)$ the previous equality and recall the asymptotics of the γ_n given in formula (13) to get:

$$\begin{aligned} \pi_{\lambda_2} X_n &= \gamma_n(\lambda_2) \Lambda + n^{\sigma_2} \varepsilon_n \\ &= \frac{n^{\lambda_2}}{\Gamma(1 + \lambda_2)} \Lambda + n^{\sigma_2} \varepsilon_n. \end{aligned}$$

Summing with $\pi_{\overline{\lambda_2}} X_n$ and noticing that $\pi_{\overline{\lambda_2}} X_n = \overline{\pi_{\lambda_2} X_n}$ gives

$$\pi_{\lambda_2} X_n + \pi_{\overline{\lambda_2}} X_n = 2\Re\left[\frac{n^{\lambda_2} \Lambda}{\Gamma(1 + \lambda_2)}\right] + n^{\sigma_2} \varepsilon_n \quad (23)$$

which provides the second term in (18).

It remains to show that if λ is an eigenvalue of A different from 1, λ_2 and $\overline{\lambda_2}$, then $\pi_\lambda X_n = n^{\sigma} \varepsilon_n$, where ε_n tends to zero a.s. and in L^2 as n tends to infinity. Let σ be the real part of such an eigenvalue, we know that $\sigma < \sigma_2$. The case $\sigma > 1/2$ is easy: lemma 3 of martingale convergence still holds hence L^2 and almost sure convergence are shown together for the martingale $\gamma_n^{-1}(\lambda) \pi_\lambda X_n$. Thus $\pi_\lambda X_n$ is of order n^σ which is negligible to n^{σ_2} .

In case $\sigma \leq 1/2$, let us first prove L^2 -convergence: as in the proof of lemma 3, we have as a corollary of the moments lemma

$$E(\|\pi_\lambda X_n\|^2) = O(n) \text{ or } O(n \log n)$$

hence (recall that $\sigma_2 > 1/2$),

$$\frac{\pi_\lambda X_n}{n^{\sigma_2}} \xrightarrow{L_2} 0 .$$

For the almost sure convergence to zero of $\pi_\lambda X_n/n^{\sigma_2}$, we use Borel-Cantelli lemma: it is sufficient to show that for any $\varepsilon > 0$, the series $\sum_n P(\|\frac{\pi_\lambda X_n}{n^{\sigma_2}}\| > \varepsilon)$ is convergent. By Markov inequality, it is sufficient to show that for some integer p , the moment $E\|\frac{\pi_\lambda X_n}{n^{\sigma_2}}\|^{2p}$ is the general term of a convergent numerical series. It is true, for p large enough, because of the moments lemma: for every positive definite Hermitian form and complex vector Z , for every nonnegative integer p ,

$$E \left(\left| (Z | \frac{\pi_\lambda X_n}{n^{\sigma_2}}) \right|^{2p} \right) = O\left(\frac{1}{n^{p(2\sigma_2-1)}}\right) .$$

Summarizing, for every eigenvalue λ of A different from 1, λ_2 and $\overline{\lambda_2}$,

$$\frac{\pi_\lambda X_n}{n^{\sigma_2}} \xrightarrow[n \rightarrow \infty]{\text{a.s. and in } L^2} 0 . \quad (24)$$

To get the final result, it is now enough to put (21), (22), (23) and (24) together. \square

Corollary 2 *Suppose $m \geq 27$. If χ is any linear form on V , there exist a real number x_χ and real random variables ρ_χ and φ_χ such that*

$$\chi(X_n) = nx_\chi + n^{\sigma_2} \rho_\chi \cos(\tau_2 \log n + \varphi_\chi) + n^{\sigma_2} \varepsilon_n,$$

where ε_n tends to zero almost surely and in L^2 as n tends to infinity .

To prove this, see what happens to $\chi(X_n)$ with corollary 1, and put the sine and cosine terms together to get a new random phase and a new random amplitude.

This corollary describes for example the asymptotic behaviour of the number of nodes of a given type (take $\chi = dx_i$, the i -th coordinate of \mathbb{R}^{m-1}), or of the total (except saturated nodes) number of nodes (take $\chi = \sum dx_i$, where i ranges over all i between 1 and $m-1$).

A natural question arises: what are the laws of the random variables ρ and φ of the theorem ?

6 Simulations

Figure 3 represents simulations for the total number of nodes for $m = 30$. We put the number n of keys inserted in the tree on the x -axis, and $x_n - nx_\chi$ on the y -axis, where x_n is the total number of nodes (except saturated nodes, those with $m - 1$ keys) at time n and x_χ the coefficient $\lim_{+\infty} E(x_n)/n$ of its drift. The graph remains fairly smooth around an “ $n^{\sigma^2} \cos \log n$ ” curve. Note that we only drew one point over one thousand.

Figure 4 illustrates the random amplitude ρ_χ and the random phase φ_χ for the asymptotics of the total number of nodes x_n . On the x -axis, $\log n$; on the y -axis, $(x_n - nx_\chi)/n^{\sigma^2}$ for two simulations. Note the difference between the amplitudes and the phases of both simulations.

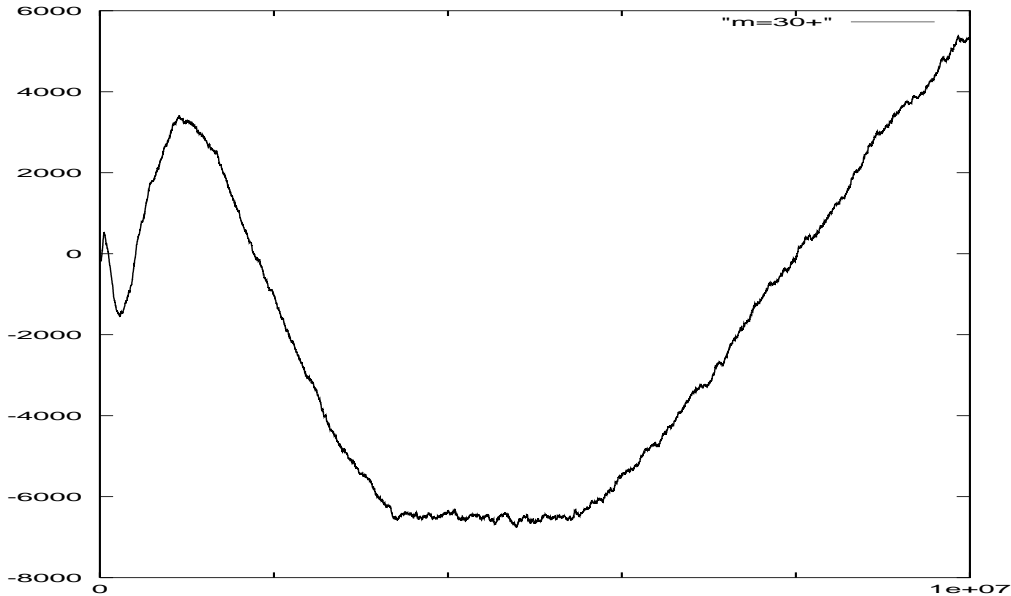


Figure 3: $m = 30$

Acknowledgements

We are indebted to the referees for careful reading and accurate suggestions. We warmly thank Jean-François Marckert for helpful discussions.

References

- [1] K.B. ATHREYA AND P. NEY, Branching processes, Springer, 1972.

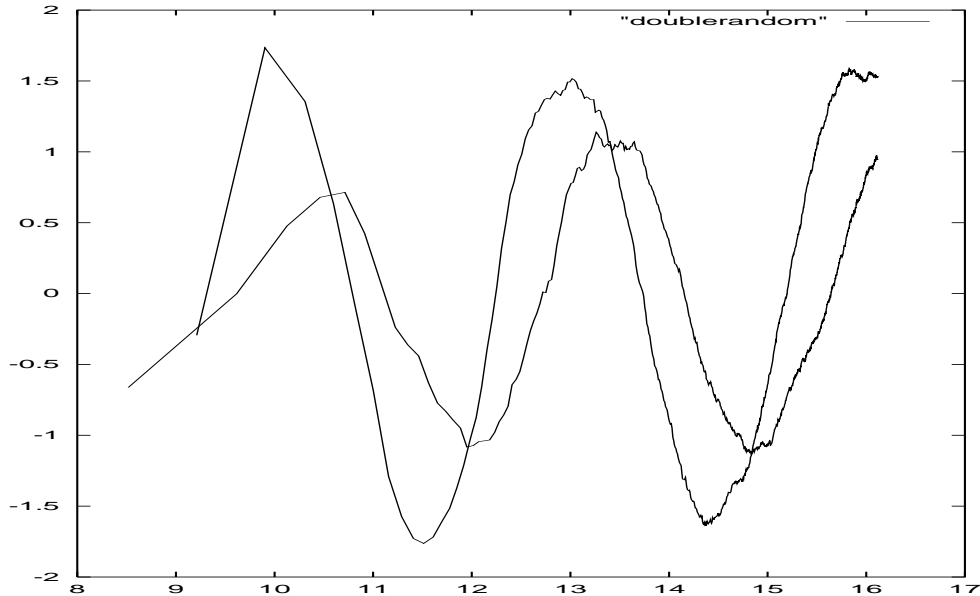


Figure 4: two simulations give two different phases and amplitudes

- [2] J. D. BIGGINS, H. COHN AND O. NERMAN, Multi-type branching in varying environment, *Stochastic Processes and Applications*, **83**, no. 2, 357–400, 1999.
- [3] H.H. CHERN AND H. K. HWANG, Phase changes in random m -ary trees and generalized quicksort, *Random Structures and Algorithms*, **19**, 3-4, 316-358, 2001.
- [4] D. S. DEAN AND S. N. MAJUMDAR, Phase transition in a random fragmentation problem with applications to computer science, *J. Phys. A: Math. Gen.*, **35**, 501–507, 2002.
- [5] H. K. HWANG, Second phase changes in random m -ary trees and generalized quicksort: convergence rates, *Annals of Probability*, **31**, 2, 609–629, 2003.
- [6] J. JABBOUR-HATTAB, Martingales and Large Deviations for Binary Search Trees, *Random Structures and Algorithms*, **19**, no. 2, 112–127, 2001.
- [7] O.D. JONES, On the convergence of multitype branching processes with varying environments, *Ann. Appl. Probab.*, **7**, no. 3, 772–801, 1997.
- [8] W. LEW AND H. MAHMOUD, The joint distribution of elastic buckets in multiway search trees, *Siam J. Comput.* **23**, no. 5, 1050–1074, 1994.

- [9] H. M. MAHMOUD, Evolution of Random Search Trees, John Wiley & Sons, New York, 1992.
- [10] H. M. MAHMOUD, The size of random bucket trees via urn models, *Acta Informatica*, **38**, 813–838, 2002.
- [11] H. M. MAHMOUD AND B. PITTEL, Analysis of the space of search trees under the random insertion algorithm, *Journal of Algorithms*, **10**, 52–75, 1989.
- [12] H. M. MAHMOUD AND R.T. SMYTHE, Probabilistic analysis of bucket recursive trees, *Theoretical Computer Science*, **144**, 180–205, 1995.
- [13] C. J. MODE, Multitype branching processes. Theory and applications. *Modern Analytic and Computational Methods in Science and Mathematics*, **34**, American Elsevier Publishing Co., Inc., New York, 1971.
- [14] R. NEININGER AND L. RUESCHENDORF, A General Limit Theorem for Recursive Algorithms and Combinatorial Structures, *Ann. Appl. Probab.*, to appear, 2003.
<http://www.stochastik.uni-freiburg.de/homepages/rueschendorf/papers/general.ps>
- [15] J. NEVEU, Arbres et processus de Galton-Watson. *Ann. IHP*, **22**, 2, 199–207, 1986.
- [16] R.T. SMYTHE, Central limit theorems for urns models, *Stochastic Processes and Applications*, **65**, 115–137, 1996.