



On discrete-time multiallelic evolutionary dynamics driven by selection

Thierry Huillet

► To cite this version:

Thierry Huillet. On discrete-time multiallelic evolutionary dynamics driven by selection. Journal of Probability and Statistics, 2010, 2010, pp.580762. 10.1155/2010/580762 . hal-00525968

HAL Id: hal-00525968

<https://hal.science/hal-00525968>

Submitted on 13 Oct 2010

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.

ON DISCRETE-TIME MULTIALLELIC EVOLUTIONARY DYNAMICS DRIVEN BY SELECTION

THIERRY E. HUILLET

ABSTRACT. We revisit some problems arising in the context of multiallelic discrete-time evolutionary dynamics driven by fitness. We consider both the deterministic and the stochastic setups and for the latter both the Wright-Fisher and the Moran approaches. In the deterministic formulation, we construct a Markov process whose Master equation identifies with the nonlinear deterministic evolutionary equation. Then, we draw the attention on a class of fitness matrices that plays some role in the important matter of polymorphism: the class of strictly ultrametric fitness matrices. In the random cases, we focus on fixation probabilities, on various conditionings on non-fixation and on (quasi-)stationary distributions.

Keywords: Evolutionary genetics, multiallelic fitness landscape, strictly ultrametric fitness matrix, polymorphism, Moran and Wright-Fisher models, fixation probabilities, multiplicative fitness, conditioning, (quasi-)stationary distribution.

Topics: Evolutionary processes (theory), Population dynamics (Theory).

1. INTRODUCTION

Population genetics aims at elucidating the fate of genotype frequencies undergoing the basic evolutionary processes when various driving ‘forces’ such as fitness, mutation or recombination are at stake in the gene pool. This requires to clarify the updating mechanisms of the gene frequency-distributions over time. Another important additional driving source is the genetic drift whose nature is exclusively random. The corresponding field of interest is the statistical theory arising from this aspect of the gene replacement processes and it requires some use of the Markov chain theory (see [1]).

In this note, we revisit the basics of both the deterministic and stochastic dynamics arising in discrete-time asexual *multiallele* evolutionary genetics driven only by *fitness*. We do not touch at all upon other important mechanisms such as mutation. We start with the haploid case with K alleles before switching to the more interesting diploid case.

Let us summarize and comment the material developed in Section 2. In the deterministic haploid case, a vector of fitness is attached to the alleles. The interest is on the evolution of the allelic frequency distribution over time. The updates of the allele frequency distributions are driven by the relative fitnesses of the alleles,

ending up in a state where only the fittest monomorphic state will survive. This state is also an extremal point of the simplex over which the dynamics takes place. From this dynamics, it appears that the mean fitness increases as time passes by, the rate of increase being the variance in relative fitness (a well-known particular incarnation of the Fisher theorem of natural selection). We recall the background and the evolution equations.

The haploid replicator dynamics, as a nonlinear updating mapping from the simplex to the simplex, may be viewed as the discrete-time nonlinear Master equation of some random Markov process. We supply a construction for this process thereby giving a stochastic interpretation to the original deterministic formulation of the dynamics.

In the diploid case, there is a similar deterministic updating dynamics but now on the full array of the genotype frequencies. It involves the fitness matrix attached to the genotypes. When mating is random so that the Hardy-Weinberg law applies, we may look at the induced marginal allelic frequencies dynamics. The updating dynamics looks quite similar to the one occurring in the haploid case except that the mean fitness is now the mean fitness quadratic form in the current frequencies whereas marginal fitnesses are no longer constant but affine functions in these frequencies. As for the haploid case, it is possible to construct a Markov process whose nonlinear Master equation coincides with the diploid replicator dynamics. We supply this construction which we believe is new.

In the diploid context, the Fisher theorem still holds true but, as a result of the fitness landscape being more complex, there is a possibility for a polymorphic equilibrium state to emerge. Due to its major evolutionary interest, our subsequent concern is to identify examples of diploid dynamics leading to a unique *polymorphic* state on the simplex, either *unstable* or *stable*. We start with the unstable case and draw the attention on a class of fitness matrices leading to a unique unstable polymorphic equilibrium state: the class of strictly *potential* matrices. *Strictly ultrametric matrices* are particular instances of strictly potential matrices ([11]) which therefore will display unstable polymorphism as well. There is a useful canonical representation of strictly ultrametric matrices due to ([14]) which we recall which helps giving specific examples of strictly ultrametric matrices. When dealing with the class of strictly potential fitness matrices, the mean fitness quadratic form is definite-positive; we derive a related class of fitness matrices leading to a definite-negative mean fitness quadratic form. For this class of matrices, there will also be a unique polymorphic equilibrium state for the diploid dynamics and it will be stable. We also draw the attention on a sub-class of the latter: the so-called ‘anti-strictly ultrametric matrices’. For such matrices, among other things, the fitness of each homozygote should not exceed the ones of all the heterozygotes carrying the allele of the homozygote. To the best of our knowledge, the large class of potential fitness matrices as natural candidates for polymorphic states to emerge was not discussed in the literature.

Section 3 is devoted to the stochastic version of these considerations when the transitions in the constitutive allelic population sizes are given by a K -dimensional *Wright-Fisher* model with total constant-size N (see [1] and [12]). It takes into account an additional important driving source of evolution, namely the genetic

drift, whose nature is exclusively random. Under fitness only and in particular in the absence of mutations, the multiallelic Wright-Fisher model is a transient Markov chain on a $\binom{N+K-1}{K-1}$ -dimensional state-space whose absorbing states are the monomorphic states. We give an expression for the fixation probabilities of this process. Then, we develop four conditioning problems: conditioning on fixating in a given monomorphic state, conditioning on avoiding the extremal states before the current instant, conditioning on non-fixation at each transition time and conditioning on avoiding the extremal states in the remote future. For the last three conditioned processes, the equilibrium structure of the fitness mechanism shows up.

Finally we run into similar considerations for the *Moran* model. When dealing with the fixation probabilities in this Moran context, we suggest a new mean-field approximation of these probabilities which is based on a well-known explicit formula for the 2-alleles case ([1]). It concerns the case of *multiplicative fitnesses* only. Finally, we consider the Moran model conditioned on non-fixation at each transition time. We exploit the reversible character of this process to derive a new explicit product formula for its *invariant* probability measure.

2. DETERMINISTIC EVOLUTIONARY DYNAMICS

We start with the haploid case before moving to the more interesting diploid case.

2.1. Single locus: haploid population with K alleles. Consider K alleles A_k , $k = 1, \dots, K$ attached to a single locus. Suppose the current time- t allelic frequency distribution is given by the column vector $\mathbf{x} := x_k$, $k = 1, \dots, K$ ¹. We therefore have $\mathbf{x} \in S_K = \left\{ \mathbf{x} \geq \mathbf{0} : |\mathbf{x}| := \sum_{k=1}^K x_k = 1 \right\}$, the K -simplex as a convex subset of \mathbb{R}^K with dimension $K - 1$. Let $\mathbf{w} := w_k > 0$, $k = 1, \dots, K$ denote the absolute fitnesses of the alleles. Let

$$(1) \quad w(\mathbf{x}) := \sum_l w_l x_l = \mathbf{w}^* \mathbf{x}$$

be the mean fitness of the population at time t . We shall also need

$$(2) \quad \sigma^2(\mathbf{x}) = \sum_{k=1}^K x_k (w_k - w(\mathbf{x}))^2,$$

the variance in absolute fitness and

$$(3) \quad \bar{\sigma}^2(\mathbf{x}) = \sum_{k=1}^K x_k \left(\frac{w_k}{w(\mathbf{x})} - 1 \right)^2 = \sigma^2(\mathbf{x}) / w(\mathbf{x})^2,$$

the variance in relative fitness $w_k / w(\mathbf{x})$.

Dynamics. From the deterministic evolutionary genetics point of view, the discrete-time update of the allele frequency distribution on the simplex S_K is given by²

¹Throughout, a boldface variable, say \mathbf{x} , will represent a column-vector so that its transpose, say \mathbf{x}^* , will be a line-vector.

²The symbol $'$ is a common and useful notation to denote the updated frequency

$$(4) \quad x'_k = p_k(\mathbf{x}) := \frac{x_k w_k}{w(\mathbf{x})}, \quad k = 1, \dots, K.$$

The quantity $\frac{w_k}{w(\mathbf{x})} - 1$ therefore interprets as the frequency-dependent Malthus growth rate parameter of x_k . As required, the vector $\mathbf{p}(\mathbf{x}) := p_k(\mathbf{x})$, $k = 1, \dots, K$, maps S_K into S_K . In vector form, with $D_{\mathbf{x}} := \text{diag}(x_k, k = 1, \dots, K)$, the nonlinear deterministic evolutionary dynamics reads:

$$\mathbf{x}' = \mathbf{p}(\mathbf{x}) = \frac{1}{w(\mathbf{x})} D_{\mathbf{w}} \mathbf{x} = \frac{1}{w(\mathbf{x})} D_{\mathbf{x}} \mathbf{w},$$

or, with $\Delta \mathbf{x} := \mathbf{x}' - \mathbf{x}$, the increment of \mathbf{x}

$$\Delta \mathbf{x} = \left(\frac{1}{w(\mathbf{x})} D_{\mathbf{w}} - I \right) \mathbf{x}.$$

Avoiding the trivial case where fitnesses are all equal, without loss of generality, we can assume that either $w_1 \geq \dots \geq w_K = 1$ or $w_1 \leq \dots \leq w_K = 1$. Thus that allele A_1 or A_K has largest fitness.

Mean fitness increase. According to the dynamical system (4), unless the equilibrium state is attained, the absolute mean fitness $w(\mathbf{x})$ increases:

$$\begin{aligned} \Delta w(\mathbf{x}) &= w(\mathbf{x}') - w(\mathbf{x}) = \sum_k w_k \Delta x_k \\ &= \sum_k w_k x_k \left(\frac{w_k}{w(\mathbf{x})} - 1 \right) = \frac{\sum_k w_k^2 x_k}{w(\mathbf{x})} - w(\mathbf{x}) > 0. \end{aligned}$$

The mean fitness is maximal at equilibrium. The rate of increase of $w(\mathbf{x})$ is:

$$(5) \quad \frac{\Delta w(\mathbf{x})}{w(\mathbf{x})} = \sum_k x_k \left(\frac{w_k}{w(\mathbf{x})} - 1 \right)^2 = \sum_k \frac{(\Delta x_k)^2}{x_k},$$

which is the variance in relative fitness $\bar{\sigma}^2(\mathbf{x})$ defined in (3). These last two facts are sometimes termed the 1930s Fisher fundamental theorem of natural selection. The equilibria of (4) are the extremal states (0-faces) of the boundary of S_K . To make it simple, if there is an allele whose fitness is strictly larger than the ones of the others, the deterministic evolutionary dynamics (4) will attain an equilibrium where only the fittest will survive; starting from any initial state of S_K which is not an extremal (or monomorphic) point, the haploid trajectories will converge to this fittest state.

A stochastic interpretation of the deterministic dynamics (4). A vector \mathbf{x} of S_K can be thought of as a probability vector. The dynamical equation (4), as a nonlinear update mapping from S_K to S_K , may be viewed as the discrete-time nonlinear Master equation of some Markov process whose construction we now give. Suppose we have a population of N haploid individuals each of which can be of one among K types or colors (carrying one among the K possible alleles). We shall need to introduce an extra color-state, say $\partial = \{0\}$, which will be absorbing for the process we shall now construct. Let $\mathbf{K}(t) := K_n(t)$, $n = 1, \dots, N$ be the random color distribution of the population at time t , therefore with enlarged state-space

$\{0, 1, \dots, K\}$. Assume the individuals are indistinguishable leading to the exchangeability property: $K_n(t) \stackrel{d}{=} K_1(t)$, $n = 2, \dots, N$ (equality in distribution). Let $U_{t,m}$, $t = 1, 2, \dots; m = 1, \dots, N$ be N i.i.d. driving sequences of uniformly distributed random variables on $[0, 1]$, independent of $\mathbf{K}(t)$. Let $\bar{w}_k := w_k / \sum_k w_k$, $k = 1, \dots, K$. To decide the allele $K_m(t+1)$ carried by the individual number $m \in \{1, \dots, N\}$ at time $t+1$, with $\mathbf{1}(A)$ the indicator function of the event A , consider the random Markovian dynamical system

$$(6) \quad \mathbf{1}(K_m(t+1) = k, \tau_m > t+1) = \mathbf{1}\left(\frac{\bar{w}_k}{N} \sum_{n=1}^N \mathbf{1}(K_n(t) = k, \tau_n > t) > U_{t+1,m}\right).$$

Here $k \in \{1, \dots, K\}$ and τ_n is the first time that $K_n(t)$ hits the absorbing state ∂ . As a result of $K_n(t) \stackrel{d}{=} K_1(t)$, we naturally assume $\tau_n \stackrel{d}{=} \tau_1$, $n = 2, \dots, N$. For each n , our model therefore attributes a positive probability that $K_n(t) = 0$ for all $t \geq \tau_n$. Although in principle, there is a possibility that the type of the m th particle is the one of the fictitious unobservable allele A_0 , as a result of (6), the sample paths of $\mathbf{K}(t)$ leading to this A_0 are ruled out because focus is on the observable states only.

In words, for the dynamics (6), the observable event $K_m(t+1) = k$ is realized (together then with $\tau_m > t+1$) if the proportion at t of type- k individuals, weighted by the corresponding scaled fitness \bar{w}_k , is large enough (compared to U_{t+1}) and of course if the whole process was not absorbed at $\{0\}$ in the previous step. Taking first the expectation with respect to the driving noise $U_{t+1,m}$ in (6), we get

$$\mathbb{P}(K_m(t+1) = k, \tau_m > t+1 \mid \mathbf{K}(t) > \mathbf{0}) = \frac{\bar{w}_k}{N} \sum_{n=1}^N \mathbf{1}(K_n(t) = k, \tau_n > t).$$

Putting $z_k(t) := \mathbb{P}(K_1(t) = k, \tau_1 > t)$, recalling $(K_n(t); \tau_n) \stackrel{d}{=} (K_1(t); \tau_1)$, $n = 2, \dots, N$, taking the expectation with respect to $\mathbf{K}(t)$, we get an unnormalized version of (4):

$$z_k(t+1) = \bar{w}_k z_k(t), \quad k \in \{1, \dots, K\}.$$

We have $1 > \mathbb{P}(\tau_1 > t) = \sum_{k=1}^K z_k(t) = \sum_{k=1}^K x_k(0) \bar{w}_k^t \rightarrow 0$, geometrically fast. Defining the normalized conditional probabilities

$$x_k(t) = \frac{z_k(t)}{\sum_{k=1}^K z_k(t)} = \mathbb{P}(K_1(t) = k \mid \tau_1 > t),$$

we obtain the normalized haploid dynamics (4)

$$x'_k = \frac{w_k x_k}{\sum_{k=1}^K w_k x_k}, \quad k \in \{1, \dots, K\}.$$

It now may be viewed as the nonlinear Master equation of some stochastic Markov process. Let us make some miscellaneous remarks.

(i) Clearly this construction makes also sense if $N = 1$ (a single particle). (ii) When N is finite, we should stress that the initial condition can be chosen to be deterministic, say with: $x_k(0) = z_k(0) = i_k/N$, for some sequence of integers $i_k \geq 0$ satisfying $i_1 + \dots + i_K = N$ ($i_0 = 0$) and quantifying the initial population sizes. It could also be chosen to be random, with $x_k(0)$ defining the initial probability distribution of the alleles. This occurs in the large N limit if $i_k = \lfloor Nx_k(0) \rfloor$ so that, $i_k/N \rightarrow x_k(0)$. The latter choice may therefore be interpreted as a large N

limit of the former one. (iii) In the stochastic interpretation (6) of the deterministic dynamics (4), $x_k(t)$ can be interpreted either as the probability that the random allele carried by a typical individual is A_k or like the expected proportion of the individuals of type k within the whole population (a frequentist point of view). (iv) The appeal to the coffin state ∂ was a necessary step to understand the normalization $z_k \rightarrow x_k$. (iv) Even though $\mathbf{K}(t)$ is exchangeable, it is not true that, with $n_1 \neq n_2$ any two distinct individuals, their random labels $K_{n_1}(t)$ and $K_{n_2}(t)$ are independent. The random algorithm allowing to update the joint types of $K_{n_1}(t)$ and $K_{n_2}(t)$ could be written down but is much more involved.

2.2. Single locus: diploid population with K alleles. We now run into similar considerations but with diploid populations.

Joint evolutionary dynamics. Let $w_{k,l} \geq 0$, $k, l = 1, \dots, K$ stand for the absolute fitness of the genotypes $A_k A_l$ attached to a single locus. Assume $w_{k,l} = w_{l,k}$ ($w_{k,l}$ being proportional to the probability of an $A_k A_l$ surviving to maturity, it is natural to take $w_{k,l} = w_{l,k}$). Let then W be the symmetric fitness matrix with k, l -entry $w_{k,l}$.

Assume the current frequency distribution at time t of the genotypes $A_k A_l$ is given by $x_{k,l}$. Let X be the frequencies array with k, l -entry $x_{k,l}$. The joint evolutionary dynamics in the diploid case is given by the updating:

$$(7) \quad x'_{k,l} = x_{k,l} \frac{w_{k,l}}{\omega(X)}$$

where the mean fitness ω is now given by: $\omega(X) = \sum_{k,l} x_{k,l} w_{k,l}$. The relative fitness of the genotype $A_k A_l$ is $w_{k,l}/\omega(X)$. The joint dynamics takes the matrix form:

$$X' = \frac{1}{\omega(X)} X \circ W = \frac{1}{\omega(X)} W \circ X$$

where \circ stands for the (commutative) Hadamard product of matrices.

Let J be the flat $K \times K$ matrix whose entries are all 1. Then

$$\Delta X := X' - X = \frac{1}{\omega(X)} (X - J) \circ W = \frac{1}{\omega(X)} W \circ (X - J).$$

We shall also let

$$(8) \quad \sigma^2(X) = \sum_{k,l=1}^K x_{k,l} (w_{k,l} - \omega(X))^2$$

stand for the genotypic variance in absolute fitness and

$$(9) \quad \bar{\sigma}^2(X) = \sum_{k,l=1}^K x_{k,l} \left(\frac{w_{k,l}}{\omega(X)} - 1 \right)^2 = \sigma^2(X) / \omega(X)^2$$

will stand for the diploid variance in relative fitness.

Consider the problem of evaluating the increase of the mean fitness. We have

$$(10) \quad \Delta\omega(X) = \sum_{k,l} \Delta x_{k,l} w_{k,l} = \sum_{k,l} x_{k,l} \left(\frac{w_{k,l}^2}{\omega(X)} - w_{k,l} \right) = \omega(X) \bar{\sigma}^2(X) > 0$$

with a relative rate of increase: $\Delta w(X)/w(X) = \bar{\sigma}^2(X)$. This is the full diploid version of the Fisher theorem.

Marginal allelic dynamics. Assuming a Hardy-Weinberg equilibrium, the frequency distribution at time t , say $x_{k,l}$, of the genotypes $A_k A_l$ is given by: $x_k x_l$ where $x_k = \sum_l x_{k,l}$ is the marginal frequency of allele A_k in the whole genotypic population. The whole frequency information is now enclosed within $\mathbf{x} := x_k$, $k = 1, \dots, K$. For instance, the mean fitness is now given by the quadratic form: $\omega(\mathbf{x}) := \sum_{k,l} x_k x_l w_{k,l} = \mathbf{x}^* W \mathbf{x}$, with \mathbf{x}^* the transposed line vector of the column vector $\mathbf{x} = X \mathbf{1}$ ($\mathbf{1}$ the unit K -vector). We shall also let

$$(11) \quad \sigma^2(\mathbf{x}) = \sum_{k,l=1}^K x_k x_l (w_{k,l} - \omega(\mathbf{x}))^2$$

stand for the genotypic variance in absolute fitness and

$$(12) \quad \bar{\sigma}^2(\mathbf{x}) = \sum_{k,l=1}^K x_k x_l \left(\frac{w_{k,l}}{\omega(\mathbf{x})} - 1 \right)^2 = \sigma^2(\mathbf{x}) / \omega(\mathbf{x})^2$$

will stand for the diploid variance in relative fitness.

Consider now the update of the allelic marginal frequencies \mathbf{x} themselves. If we first define the frequency-dependent marginal fitness of A_k by $w_k(\mathbf{x}) = (W\mathbf{x})_k := \sum_l w_{k,l} x_l$, the marginal dynamics is given as in (4) by:

$$(13) \quad x'_k = x_k \frac{w_k(\mathbf{x})}{\omega(\mathbf{x})} = \frac{1}{\omega(\mathbf{x})} x_k (W\mathbf{x})_k =: p_k(\mathbf{x}), \quad k = 1, \dots, K.$$

This dynamics involves a multiplicative interaction between x_k and $(W\mathbf{x})_k$, the k th entry of the image $W\mathbf{x}$ of \mathbf{x} by W . In (13) there is a normalization by the quadratic form $\omega(\mathbf{x}) = \mathbf{x}^* W \mathbf{x}$. In vector form (13) reads

$$\mathbf{x}' = \frac{1}{\omega(\mathbf{x})} D_{\mathbf{x}} W \mathbf{x} = \frac{1}{\omega(\mathbf{x})} D_{W\mathbf{x}} \mathbf{x} =: \mathbf{p}(\mathbf{x})$$

where \mathbf{p} maps S_K into S_K . Iterating, the time- t frequency distribution is:

$$\mathbf{x}(t) = \mathbf{p}(\mathbf{p}(\dots t \text{ times } \dots (\mathbf{p}(\mathbf{x}_0)))) .$$

Except for the fact that the mean fitness ω in (13) is now a quadratic form in \mathbf{x} and that the marginal fitness of A_k is now frequency-dependent, depending linearly on \mathbf{x} , as far as the marginal frequencies are concerned, the updating formalism (13) in the diploid case looks very similar to the one in (4) describing the haploid case.

In the diploid case, assuming fitnesses to be *multiplicative*, say with $w_{k,l} = w_k w_l$, then $\frac{w_k(\mathbf{x})}{\mathbf{x}^* W \mathbf{x}} = \frac{w_k}{\sum_l w_l x_l}$ and the dynamics (13) boils down to (4). However, the mean fitness in this case is $\omega(\mathbf{x}) = (\sum_l w_l x_l)^2$ and not $w(\mathbf{x}) = \sum_l w_l x_l$ as in the haploid case.

A stochastic interpretation of the deterministic dynamics (13). As for the haploid case, there is a Markov chain governed by the Master equation (13).

Consider a population of diploid individuals. The number of alleles N in this population is therefore twice the number of genes. Each allele can be of one among K types or colors (carrying one among the K possible alleles). As before, we introduce an extra color-state, say $\partial = \{0\}$ which is absorbing for the process to be constructed. For $c = 1, 2$, let $\mathbf{K}^c(t) := K_n^c(t)$, $n = 1, \dots, N$ be two independent copies of the random color distribution of the allelic population at time t . Let $\mathbf{K}(t) = (\mathbf{K}^1(t), \mathbf{K}^2(t))$. Assume the alleles are indistinguishable within each sample, leading to: $K_n^c(t) \stackrel{d}{=} K_1^c(t)$, $n = 2, \dots, N$, $c = 1, 2$. For $c = 1, 2$, let $U_{t,m}^c$, $t = 1, 2, \dots$; $m = 1, \dots, N$ be two mutually independent i.i.d. driving N -sequences of uniformly distributed random variables on $[0, 1]$ and independent of $\mathbf{K}(t)$. To decide the type of the random allele $K_m(t+1)$, $m = 1, \dots, N$, at time $t+1$, consider now the Markovian dynamical system

$$(14) \quad \mathbf{1}(K_m(t+1) = k, \tau_m > t+1) = \mathbf{1}\left(\frac{1}{N} \sum_{n=1}^N \mathbf{1}(K_n^1(t) = k, \tau_n^1 > t) > U_{t+1,m}^1\right) \mathbf{1}\left(\frac{1}{N} \sum_{n=1}^N \bar{W}_{k,K_n^2(t)} \mathbf{1}(\tau_n^2 > t) > U_{t+1,m}^2\right).$$

Here $k \in \{1, \dots, K\}$, τ_n^c are the first hitting times of each $K_n^c(t)$ of the absorbing state ∂ , $c = 1, 2$ and $\bar{W} := W/\|W\|$ for any matrix norm $\|W\|$, say for example: $\|W\| = \sum_{k,l} w_{k,l}$. We assume $\tau_n^c \stackrel{d}{=} \tau_1^c$, $n = 2, \dots, N$, $c = 1, 2$.

In words, for this new dynamics, the observable event $K_m(t+1) = k$ is seen to be realized together with $\tau_m > t+1$ if two natural independent conditions are now satisfied which can be read from the two indicator functions in the right-hand side of (14):

* first, the proportion at t of type- k individuals of the first copy should be large enough (compared to U_{t+1}^1).

* second, for the second sample copy $\mathbf{K}^2(t)$, the expected fitness of the genotypes $A_k A_l$, $l = 1, \dots, K$, containing allele A_k at t should be large enough (compared to U_{t+1}^2).

As for the haploid case, it is necessary that both processes $\mathbf{K}(t)$ should not be absorbed at $\{0\}$ in the previous step. Taking first the expectation with respect to the independent driving noises $U_{t+1,m}^c$ in (14), we get

$$\mathbb{P}(K_m(t+1) = k, \tau_m > t+1 \mid \mathbf{K}(t) > \mathbf{0}) = \frac{1}{N} \sum_{n=1}^N \mathbf{1}(K_n^1(t) = k, \tau_n^1 > t) \cdot \frac{1}{N} \sum_{n=1}^N \bar{W}_{k,K_n^2(t)} \mathbf{1}(\tau_n^2 > t).$$

Putting $z_k(t) = \mathbb{P}(K_1(t) = k, \tau_1 > t)$, taking the expectation with respect to $\mathbf{K}(t)$ and using our independence and exchangeability hypotheses, we get

$$z_k(t+1) = z_k(t) \cdot \mathbb{E} \left[\bar{W}_{k,K_1^2(t)} \mathbf{1}(\tau_1^2 > t) \right] = z_k(t) \cdot \sum_{l=1}^k \bar{W}_{k,l} z_l(t),$$

corresponding to an unnormalized version of (13):

$$z_k(t+1) = z_k(t) (\bar{W} \mathbf{z}(t))_k, \quad k \in \{1, \dots, K\}.$$

Defining the normalized conditional probabilities

$$x_k(t) = \frac{z_k(t)}{\sum_{k=1}^K z_k(t)} = \mathbb{P}(K_1(t) = k \mid \tau_1 > t),$$

we obtain the normalized nonlinear Markovian diploid dynamics (13)

$$x'_k = \frac{x_k(W\mathbf{x})_k}{\sum_{k=1}^K x_k(W\mathbf{x})_k}, \quad k \in \{1, \dots, K\}.$$

Note that this construction makes sense if $N = 2$ (a single individual and its 2 alleles of K possible types). The need for two copies of $\mathbf{K}(t)$ was governed by the quadratic character of the interaction appearing in the numerator of (13).

Increase of mean fitness. Again, the mean fitness $\omega(\mathbf{x})$, as a Lyapunov function, increases as time passes by. We indeed have

$$\Delta\omega(\mathbf{x}) = \omega(\mathbf{x}') - \omega(\mathbf{x}) = \frac{1}{\omega(\mathbf{x})^2} \sum_{k,l} x_k w_k(\mathbf{x}) w_{k,l} x_l w_l(\mathbf{x}) - \sum_{k,l} x_k w_{k,l} x_l > 0,$$

because, defining $0 < X(\mathbf{x}) := \sum_{k,l} x_k \left(1 - \frac{w_k(\mathbf{x})}{\omega(\mathbf{x})}\right) w_{k,l} \left(1 - \frac{w_l(\mathbf{x})}{\omega(\mathbf{x})}\right) x_l$, we have

$$\Delta\omega(\mathbf{x}) = X(\mathbf{x}) + \frac{2}{\omega(\mathbf{x})} \left(\sum_k x_k w_k(\mathbf{x})^2 - \omega(\mathbf{x})^2 \right) > 0.$$

Its partial rate of increase due to frequency shifts only is $\delta\omega(\mathbf{x}) := \sum_k \Delta x_k w_k(\mathbf{x})$. It satisfies

$$(15) \quad \frac{\delta\omega(\mathbf{x})}{\omega(\mathbf{x})} = \sum_k x_k \left(\frac{w_k(\mathbf{x})}{\omega(\mathbf{x})} - 1 \right)^2 = \sum_k \frac{(\Delta x_k)^2}{x_k} = \bar{\sigma}_A^2(\mathbf{x})/2$$

where $\bar{\sigma}_A^2(\mathbf{x})$ is the allelic variance in relative fitness

$$(16) \quad \bar{\sigma}_A^2(\mathbf{x}) := 2 \sum_{k=1}^K x_k \left(\frac{w_k(\mathbf{x})}{\omega(\mathbf{x})} - 1 \right)^2.$$

An alternative representation of the allelic dynamics. There is an alternative vectorial representation of the dynamics (13). Define the symmetric positive-definite matrix $G(\mathbf{x}) = D_{\mathbf{x}}(I - \mathbf{1}\mathbf{x}^*)$ with quadratic entries in the frequencies:

$$G(\mathbf{x})_{k,l} = x_k (\delta_{k,l} - x_l).$$

Introduce the quantity $V_W(\mathbf{x}) = \frac{1}{2} \log \omega(\mathbf{x})$, which is half the logarithm of the mean fitness. Then, (13) may be recast as the gradient-like dynamics:

$$(17) \quad \Delta\mathbf{x} = \frac{1}{\omega(\mathbf{x})} G(\mathbf{x}) W\mathbf{x} = G(\mathbf{x}) \nabla V_W(\mathbf{x}),$$

with $|\Delta\mathbf{x}| = \mathbf{1}^* \Delta\mathbf{x} = 0$ as a result of $\mathbf{1}^* G(\mathbf{x}) = \mathbf{0}^*$. Note

$$\nabla V_W(\mathbf{x})^* \Delta\mathbf{x} = \nabla V_W(\mathbf{x})^* G(\mathbf{x}) \nabla V_W(\mathbf{x}) \geq 0.$$

The dynamics (17) is of pure gradient-type with respect to the Svirezhev-Shashahani distance metric $d_G(\mathbf{x}, \mathbf{x}')$, see [21] and [20]. For this metric, the distance between \mathbf{x} and $\mathbf{x}' = \mathbf{x} + \Delta\mathbf{x}$ of S_K is:

$$d_G(\mathbf{x}, \mathbf{x}') = (\Delta\mathbf{x}^* G^{-1} \Delta\mathbf{x})^{1/2} = \left(\sum_{k=1}^K x_k^{-1} (\Delta x_k)^2 \right)^{1/2}.$$

From (15) and (16), this quantity, which is the length of $\Delta\mathbf{x}$ satisfying $|\Delta\mathbf{x}| = 0$, is also the square-root of half the allelic variance (the standard deviation) in relative fitness.

2.3. Equilibria (diploid case). The mean fitness increase phenomenon occurs till the evolutionary dynamics reaches an equilibrium state. We wish to briefly discuss the questions relative to equilibria in the diploid case.

Preliminaries. In contrast with the haploid case, in the diploid situation, the dynamics (13) can have more complex equilibrium points, satisfying $w_k(\mathbf{x}_{eq}) = w_1(\mathbf{x}_{eq})$, $k = 2, \dots, K$ and $\sum_l x_{eq,l} = 1$. To avoid linear manifolds of equilibria, we first assume that all principal minors of W are non-singular and also that the fitnesses of all homozygotes $w_{k,k}$ are positive. In this case, from the Bézout theorem, the number of equilibria is finite and less or equal than the number $2^K - 1$ of faces of S_K . Note that the K extremal endpoints of S_K (0-faces) are always monomorphic fixed points of (13).

An instructive example fulfilling these preliminary conditions is $W = I$. There are $2^K - 1$ equilibrium points (the barycenters of the $\binom{K}{k+1}$ k -dimensional faces, $k = 0, \dots, K-1$), but only one polymorphic equilibrium which is the barycenter \mathbf{x}_B of S_K . This point is the one with minimal fitness and it is unstable. The 0-faces are stable fixed points whereas the barycenters of the k faces with $k \in \{1, \dots, K-2\}$ are saddle-points. The simplex S_K could be partitioned into pieces each of which being the attraction basins of the stable 0-face states: in contrast with the haploid case, the type of the survivor is not necessarily the one of the fittest; it will depend on the initial condition.

Similar conclusions can be drawn if instead of $W = I$, we start with $W = (I - D\boldsymbol{\lambda})^{-1}$ where $\boldsymbol{\lambda} := (\lambda_k, k = 1, \dots, K)$ satisfies $\mathbf{0} < \boldsymbol{\lambda} < \mathbf{1}$ (meaning $0 < \lambda_k < 1, \forall k$). In this case again, there is only one unstable polymorphic equilibrium which is easily seen to be: $\mathbf{x}_{eq} = (\mathbf{1} - \boldsymbol{\lambda}) / (K - |\boldsymbol{\lambda}|) \in S_K$.

Due to its evolutionary interest, we would like now to discuss the opportunity for a polymorphic state to be stable. Under the above assumptions on W , a unique stable internal (polymorphic) equilibrium state can exist, necessary and sufficient conditions being (i) there is a unique $\mathbf{z} > \mathbf{0}$ for which $W\mathbf{z} = \mathbf{1}$ and (ii) W has exactly one strictly positive dominant eigenvalue and at least one strictly negative eigenvalue or else the sequence of principal minors of W alternates in sign (see Kingman, [7]). If this is the case, the equilibrium polymorphic state is $\mathbf{x}_{eq} = \mathbf{z} / |\mathbf{z}|$. It is stable in the sense that it is a local maximum of the mean fitness $\omega(\mathbf{x}) = \mathbf{x}^* W \mathbf{x}$, with $\omega(\mathbf{x}_{eq}) = 1 / |\mathbf{z}|$. Since $|W| \neq 0$, the linearization of $\mathbf{p}(\mathbf{x})$ at \mathbf{x}_{eq} has no eigenvalue of modulus 1 and so \mathbf{x}_{eq} is hyperbolic and/or isolated (see [15]). A stable isolated polymorphic state is asymptotically Lyapunov stable.

Under these additional assumptions therefore on W , starting from any initial condition in the interior of S_K , all trajectories will be attracted by $\mathbf{x}_{eq} = \mathbf{z}/|\mathbf{z}|$ which is an isolated global maximum of $\omega(\mathbf{x})$.

When there is no such unique globally stable polymorphic equilibrium, all trajectories will still converge but perhaps to a local equilibrium state where some alleles get extinct. Which allele and how many alleles are concerned seems to be an unsolved problem in its full generality.

Special classes of fitness matrices leading to a polymorphic state.

We now draw the attention on a particular class of fitness matrices that lead to a polymorphic state, either unstable or stable. We start with the unstable case, extending the above special diagonal case $W = (I - D_\lambda)^{-1}$ leading to a unique unstable polymorphic state.

(i) *The unstable case.* Let $\Lambda \geq 0$ be a symmetric irreducible strictly substochastic matrix satisfying: $\Lambda \mathbf{1} := \mathbf{q} < \mathbf{1}$: The positive mass-defect vector of Λ is $\mathbf{1} - \mathbf{q}$. Let $\lambda > 0$. Define the symmetric strictly potential matrix: $W = \lambda^{-1}(I - \Lambda)^{-1} \geq 0$, with $W^{-1} = \lambda(I - \Lambda)$ defining a strictly row-diagonally dominant Stieltjes matrix with the properties ([14]): $(W^{-1})_{k,k} > 0$, $(W^{-1})_{k,l} \leq 0$ for $k \neq l$ and $(W^{-1})_{k,k} + \sum_{l \neq k} (W^{-1})_{k,l} > 0 \forall k$. Then

$$(18) \quad \lambda \mathbf{1} > \mathbf{z} = W^{-1} \mathbf{1} = \lambda(I - \Lambda) \mathbf{1} = \lambda(\mathbf{1} - \mathbf{q}) > \mathbf{0}.$$

The vector $\mathbf{z} = \lambda(\mathbf{1} - \mathbf{q})$ is called the equilibrium potential of W . We have: $|\mathbf{z}| = \lambda(K - |\mathbf{q}|)$. For this class of W therefore, $W\mathbf{z} = \mathbf{1}$ admits a positive solution \mathbf{z} .

Conversely, given a non-singular matrix $W \geq 0$ satisfying $W\mathbf{z} = \mathbf{1}$ for some $\mathbf{z} \geq \mathbf{0}$, the matrix $\Lambda = I - \lambda^{-1}W^{-1}$ defines a substochastic matrix if and only W^{-1} satisfies $(W^{-1})_{k,k} > 0$, $(W^{-1})_{k,l} \leq 0$ for $k \neq l$ and $\lambda \geq \max_k (W^{-1})_{k,k}$. Then, W^{-1} is row-diagonally dominant and $W = \lambda^{-1}(I - \Lambda)^{-1}$ is a potential matrix.

Strictly ultrametric (sUm) matrices are special classes of positive-definite and symmetric strictly potential matrices ([11]). A sUm matrix W is symmetric with non-negative entries, satisfying: (i) $w_{k,l} \geq \min\{w_{k,j}, w_{j,l}\}$, $\forall j, k, l$ and (ii) $w_{k,k} > \max_{l \neq k} \{w_{k,l}\}$, $\forall k$ (If in condition (ii), \geq is substituted for $>$, then W is simply an ultrametric matrix and this new condition is implied by (i)). If W is a sUm matrix, the fitness dynamics will admit an unstable polymorphic equilibrium state, as a result of W being positive-definite.

Remark: Suppose Λ is substochastic and primitive. Then $W = \lambda^{-1}(I - \Lambda)^{-1} > 0$ is an ultrametric matrix. If V is the Hadamard reciprocal of W with entries $v_{k,l} = 1/w_{k,l}$, it satisfies: $v_{k,l} \leq \max\{v_{k,j}, v_{j,l}\}$. Therefore V is an ultrametric distance associated to the ultrametric potential W . Tree matrices are ultrametric matrices that are not sUm.

(ii) *The stable case.* To produce a stable equilibrium state from the sUm matrix construction, let $W = \lambda^{-1}(I - \Lambda)^{-1}$ define a symmetric strictly potential matrix as before. Then, there is a $\mathbf{z} = \lambda(\mathbf{1} - \mathbf{q}) > \mathbf{0}$ for which $W\mathbf{z} = \mathbf{1}$. With $\alpha > 1$,

define

$$(19) \quad \widetilde{W} := \frac{\alpha}{|\mathbf{z}|} J - W.$$

With $\widetilde{W} = [\widetilde{w}_{k,l}]$, we have $\min \widetilde{w}_{k,l} = \frac{\alpha}{|\mathbf{z}|} - \max w_{k,l}$ and we can choose $\alpha > 1$ so that $\widetilde{W} \geq 0$. We have $(\alpha - 1)\widetilde{W}\mathbf{z} = \mathbf{1}$ and $\boldsymbol{\delta}^* \widetilde{W} \boldsymbol{\delta} = -\boldsymbol{\delta}^* W \boldsymbol{\delta} < 0$ for all $\boldsymbol{\delta} \neq \mathbf{0}$ satisfying $|\boldsymbol{\delta}| = 0$ showing that $\mathbf{x}_{eq} := \mathbf{z}/|\mathbf{z}| = (\mathbf{1} - \mathbf{q})/(K - |\mathbf{q}|)$ is now a stable polymorphic state for \widetilde{W} . If W is a sUm matrix, then clearly \widetilde{W} satisfies the ‘anti-sUm’ property expressing a fitness domination of the heterozygotes $A_k A_l$ over the homozygotes:

$$(20) \quad \widetilde{w}_{k,l} \leq \max \{\widetilde{w}_{k,j}, \widetilde{w}_{j,l}\}, \forall j, k, l \text{ and } \widetilde{w}_{k,k} < \min_{l \neq k} \{\widetilde{w}_{k,l}\}, \forall k.$$

Non-negative symmetric negative-definite anti-sUm fitness matrices \widetilde{W} will therefore display a stable polymorphic equilibrium state \mathbf{x}_{eq} . Note $\mathbf{x}_{eq}^* \widetilde{W} \mathbf{x}_{eq} = (\alpha - 1)/|\mathbf{z}|$ is now the maximal value of the mean fitness.

Example: When $K = 2$, with $s > -1$, $h > 0$ and $sh > -1$, let

$$W = \begin{bmatrix} 1+s & 1+sh \\ 1+sh & 1 \end{bmatrix} > 0$$

define the fitness matrix with selection parameter s and dominance h . This W is sUm iff $s < 0$ and $h > 1$. The equilibrium state is $\mathbf{x}_{eq}^* := (h/(2h-1); (h-1)/(2h-1))$ and it is unstable. This W is anti-sUm iff $s > 0$ and $h > 1$. The equilibrium state is the same but it is now stable.

Note that a singular multiplicative fitness matrix of the form $W = \mathbf{w}\mathbf{w}^*$ cannot be a strictly potential matrix because its determinant $|W|$ is zero.

A general construction to produce sUm and anti-sUm matrices. Consider the problem consisting in splitting binarily and recursively the set $\{1, \dots, K\}$ till complete reduction to singletons (leaves) which are left behind in the process. For instance, consider the refinement sequence with $K = 6$ of $\{1, \dots, 6\} \equiv (123456)$:

$$(123456) \rightarrow ((23)(1456)) \rightarrow (((2)(3))((16)(45))) \rightarrow (((1)(6))((4)(5))).$$

Starting from the left, there are $2K - 1 = 11$ blocks of symbols (the total number of nodes in the splitting binary tree with K leaves). To each encountered block, numbered from $l = 1$ to $2K - 1$, starting from the left, attach a vector \mathbf{u}_l of size K with i th entry $u_l(i) = 1$ if symbol i is in the block string, 0 otherwise. For instance, from the above sequence: $\mathbf{u}_1^* = (1, 1, 1, 1, 1, 1)$, $\mathbf{u}_2^* = (0, 1, 1, 0, 0, 0)$, $\mathbf{u}_3^* = (1, 0, 0, 1, 1, 1)$, ..., $\mathbf{u}_{10}^* = (0, 0, 0, 1, 0, 0)$, $\mathbf{u}_{11}^* = (0, 0, 0, 0, 1, 0)$. To each such \mathbf{u}_l , attach a number s_l which is > 0 if $|\mathbf{u}_l| = 1$ (the leaves) and ≥ 0 if $|\mathbf{u}_l| > 1$ (the internal nodes, including the root). Then ([14]), for any choice of s_l respecting these constraints

$$W = \sum_{l=1}^{2K-1} s_l \mathbf{u}_l \mathbf{u}_l^* \geq 0$$

is a sUm matrix and any sUm matrix can be represented in this way. Because for the \mathbf{u}_l corresponding to the leaves $s_l > 0$, the diagonal terms of W are necessarily > 0 .

Since for each set $\{1, \dots, K\}$, there are b_K splitting tree sequences where b_K satisfies $b_K = \sum_{k=1}^{K-1} \binom{K}{k} b_k b_{K-k}$, $k \geq 2$, $b_1 = 1$, there are many ways to generate a sUm matrix.

Clearly, for each splitting procedure, with $\lambda^{-1} := \sum_{l=1}^{2K-1} s_l > 0$, W may be written as

$$W = \lambda^{-1} \left(J - \sum_{l=2}^{2K-1} \lambda s_l (J - \mathbf{u}_l \mathbf{u}_l^*) \right) \geq 0$$

where the matrices $J - \mathbf{u}_l \mathbf{u}_l^*$ take values in $\{0, 1\}$.

Now, with $\gamma^{-1} > 0$, any matrix of the form

$$\widetilde{W} = (\lambda^{-1} + \gamma^{-1}) J - W = \gamma^{-1} \left(J + \sum_{l=2}^{2K-1} \gamma s_l (J - \mathbf{u}_l \mathbf{u}_l^*) \right) \geq 0$$

is an anti-sUm matrix. Assuming $\gamma = 1$ and $s_l = s h_l$, \widetilde{W} may be written under the form: $\widetilde{W} = J + sA$ where $A := \sum_{l=2}^{2K-1} h_l (J - \mathbf{u}_l \mathbf{u}_l^*)$. It has at most $2K - 1$ independent parameters, namely the h_l , $l = 2, \dots, 2K - 1$ and s . If the h_l are known, then \widetilde{W} is a one-parameter family of fitness matrices³.

Examples. Assume also that $s_l = s/(2(K-1))$, $\forall l$, with $s > 0$ a selection parameter. Then, defining the $(0, 1]$ -valued matrix

$$(21) \quad A := \frac{1}{2(K-1)} \sum_{l=2}^{2K-1} (J - \mathbf{u}_l \mathbf{u}_l^*) > 0,$$

an anti-sUm matrix of the form $\widetilde{W} = J + sA$ will admit a stable polymorphic equilibrium. Clearly, A itself is anti-sUm. Because of this, there is a $\mathbf{z}_A > \mathbf{0}$ such that $A\mathbf{z}_A = \mathbf{1}$. Thus $(J + sA)\mathbf{z}_A = (|\mathbf{z}_A| + s)\mathbf{1}$ showing that, with $\mathbf{z} = \mathbf{z}_A / (|\mathbf{z}_A| + s)$, $\widetilde{W}\mathbf{z} = (J + sA)\mathbf{z} = \mathbf{1}$. We thus have $|\mathbf{z}| = |\mathbf{z}_A| / (|\mathbf{z}_A| + s)$ and so: $\mathbf{x}_{eq} = \mathbf{z} / |\mathbf{z}| = \mathbf{z}_A / |\mathbf{z}_A|$. Furthermore, the equilibrium mean fitness for such models is $\mathbf{x}_{eq}^* \widetilde{W} \mathbf{x}_{eq} = 1 / |\mathbf{z}| = (|\mathbf{z}_A| + s) / |\mathbf{z}_A| > 1$.

For the following simple sequence example with $K = 4$ (four alleles, say A,C,T,G), $(1234) \rightarrow ((1)(234)) \rightarrow ((24)(3)) \rightarrow ((2)(4))$ we find:

$$A = \frac{1}{6} \begin{bmatrix} 5 & 6 & 6 & 6 \\ 6 & 3 & 5 & 4 \\ 6 & 5 & 4 & 5 \\ 6 & 4 & 5 & 3 \end{bmatrix}$$

which itself clearly is a symmetric anti-sUm matrix, together with $\widetilde{W} = J + sA$. For this example, $\mathbf{x}_{eq}^* = \frac{1}{13} (8, 1, 3, 1)$ and the equilibrium mean fitness is $1 + s/13$. Note that taking $s_l = s/(2(K-1))$ just for the indices l that were initially chosen to satisfy $s_l > 0$ would also lead to anti-sUm matrices $A > 0$ and $\widetilde{W} = J + sA > 0$.

³Fitness matrices of the form $J + sA$ were considered in [3] in the context of the estimation of s problem.

In this case, the sum (21) defining A , should then be restricted to the indices l from 2 to $2K - 1$ for which $s_l > 0$. Proceeding in this extreme way for the above simple example, we get the borderline anti-sUm shapes

$$A = \frac{1}{6} \begin{bmatrix} 3 & 4 & 4 & 4 \\ 4 & 3 & 4 & 4 \\ 4 & 4 & 3 & 4 \\ 4 & 4 & 4 & 3 \end{bmatrix} \text{ and } \widetilde{W} = \begin{bmatrix} \frac{s+2}{2} & \frac{3+2s}{s+2} & \frac{3+2s}{3+2s} & \frac{3+2s}{3+2s} \\ \frac{3+2s}{3+2s} & \frac{s+2}{s+2} & \frac{3+2s}{s+2} & \frac{3+2s}{3+2s} \\ \frac{3+2s}{3+2s} & \frac{3+2s}{3+2s} & \frac{s+2}{s+2} & \frac{3+2s}{3+2s} \\ \frac{3+2s}{3+2s} & \frac{3+2s}{3+2s} & \frac{3+2s}{3+2s} & \frac{s+2}{s+2} \end{bmatrix}$$

Here, \mathbf{x}_{eq}^* is simply the barycenter of the 4-simplex.

Remark: The ultrametric conditions (20) should be compared with the so-called ‘triangle inequality’ conditions (leading to stable polymorphism) pointed out in ([10]), which read:

$$\widetilde{w}_{k,l} < \widetilde{w}_{k,j} + \widetilde{w}_{l,j}, \forall k \neq l \text{ and at least one } j \neq k, j \neq l$$

and

$$\widetilde{w}_{k,l} > \frac{\widetilde{w}_{k,k} + \widetilde{w}_{l,l}}{2}, \forall k \neq l.$$

It is clear that the class of anti-sUm matrices is a particular subclass of the Lewontin one.

3. STOCHASTIC EVOLUTIONARY DYNAMICS

We now switch to the random point of view of multiallelic evolutionary dynamics driven by selection. There are two models of interest: the Wright-Fisher and the Moran models. We start with Wright-Fisher.

3.1. The Wright-Fisher model. The Wright-Fisher model is a discrete space-time model which takes into account another important driving source of evolution, namely the genetic drift whose nature is exclusively random.

The Model and its first properties. Consider an allelic population with constant size N . In the haploid (diploid) case, N is (twice) the number of real individuals. Let $\mathbf{i} := i_k$ and $\mathbf{i}' := i'_k$, $k = 1, \dots, K$ be two vectors of integers quantifying the size of the allelic populations at two consecutive generations t and $t+1$. We shall let $S_{K,N} = \left\{ \mathbf{i} \text{ integers} : |\mathbf{i}| = \sum_{k=1}^K i_k = N \right\}$. Suppose the stochastic evolutionary dynamics is now given by a Markov chain whose one-step transition matrix P from states $\mathbf{I} = \mathbf{i}$ to $\mathbf{I}' = \mathbf{i}'$ is given by the multinomial Wright-Fisher model

$$(22) \quad \mathbb{P}(\mathbf{I}'_{t+1} = \mathbf{i}' \mid \mathbf{I}_t = \mathbf{i}) =: P(\mathbf{i}, \mathbf{i}') = \binom{N}{i'_1 \dots i'_K} \prod_{k=1}^K p_k \left(\frac{\mathbf{i}}{N} \right)^{i'_k}.$$

Therefore, given $\mathbf{I}_t = \mathbf{i}$, to form the next generation, each allele chooses its type independently of the others with probability \mathbf{p} , where $\mathbf{p} = (p_k, k = 1, \dots, K)$ is either given by (4) in the haploid case or by (13) in the diploid case. In the diploid case, the mechanism \mathbf{p} is assumed to present a unique polymorphic state, either stable or unstable.

The state-space dimension of the Markov chain governed by (22) is $n = \binom{N+K-1}{K-1}$ (the number of compositions of integer N into K non-negative parts which is also

the number of ways to assign N indistinguishable balls into K distinguishable boxes). To view $P(\mathbf{i}, \mathbf{i}')$ as a standard transition matrix of some Markov chain, we need first to order the states \mathbf{i} and \mathbf{i}' in (22). Starting from the bottom right corner of P states should be arranged in decreasing order when listing the entries of P moving up and left along the lines and columns respectively; or equivalently, starting from the top left corner of P states should be arranged in increasing order when moving down and right.

For example, we can order the states $\mathbf{i} \in S_{K,N}$ in decreasing order from n to 1, as follows. Let $H(\mathbf{i}) = \frac{1}{N} \sum_{k=1}^K (N+1)^{-(k-1)} i_k$ be a base- $(N+1)$ code of the state \mathbf{i} that will serve as a ranking function. The largest state \mathbf{i}_n for which $H(\mathbf{i})$ is maximal (equal to 1) is $\mathbf{i}_n^* := (N, 0, \dots, 0)$. Given a state \mathbf{i} , define the subsequent state in decreasing order as

$$\sigma_-(\mathbf{i}) = \arg \min_{\mathbf{j}: H(\mathbf{j}) < H(\mathbf{i})} (H(\mathbf{i}) - H(\mathbf{j})).$$

Then $\mathbf{i}_n^* = (N, 0, \dots, 0)$, $\mathbf{i}_{n-1}^* = (N-1, 1, 0, \dots, 0)$, ..., $\mathbf{i}_1^* = (0, \dots, 0, N)$ and to pass from state \mathbf{i} to the next state $\sigma_-(\mathbf{i})$ in this decreasing sequence, there is a unique $\delta_{\mathbf{i}}$ with entries in \mathbb{Z} satisfying $|\delta_{\mathbf{i}}| = 0$ and such that: $\sigma_-(\mathbf{i}) = \mathbf{i} - \delta_{\mathbf{i}}$. This way to order the n states is consistent with the reverse lexicographic order. For instance if $N = 3$, $K = 4$, there are $\binom{6}{3} = 20$ states ordered in decreasing order as follows:

$$\begin{aligned} 3000 &> 2100 > 2010 > 2001 > 1200 > 1110 > 1101 > 1020 > 1011 > 1002 > \\ 0300 &> 0210 > 0201 > 0120 > 0111 > 0102 > 0030 > 0021 > 0012 > 0003. \end{aligned}$$

The way the digits are propagated from left to right is clear and the consecutive $\delta_{\mathbf{i}}$ can easily be obtained. Proceeding in this way to order states, $P(\mathbf{i}, \mathbf{i}')$ is a well-defined conventional object (matrix).

From (22), the marginal transition matrix from \mathbf{i} to $I'_k = i'_k$ is binomial $\text{bin}(N, p_k(\frac{\mathbf{i}}{N}))$ with:

$$P(\mathbf{i}, i'_k) = \binom{N}{i'_k} p_k \left(\frac{\mathbf{i}}{N} \right) \left(1 - p_k \left(\frac{\mathbf{i}}{N} \right) \right)^{N-i'_k}.$$

Given $\mathbf{I} = \mathbf{i}$, the k th component I'_k of the updated state is random with:

$$\mathbb{E}_{\mathbf{i}}(I'_k/N) = p_k \left(\frac{\mathbf{i}}{N} \right) \text{ and } \sigma_{\mathbf{i}}^2(I'_k/N) = p_k \left(\frac{\mathbf{i}}{N} \right) \left(1 - p_k \left(\frac{\mathbf{i}}{N} \right) \right) / N,$$

suggesting that, in the large N population limit, the deterministic evolutionary dynamics should be recovered. Indeed, from (22), the Laplace-Stieltjes transform of the joint law of $\mathbf{I}'_{t+1} \mid \mathbf{I}_t = \mathbf{i}$ reads

$$\mathbb{E}_{\mathbf{i}} \left(e^{-\sum_{k=1}^K \lambda_k I'_k/N} \right) = \left(\sum_{k=1}^K p_k \left(\frac{\mathbf{i}}{N} \right) e^{-\lambda_k/N} \right)^N \underset{N \uparrow \infty}{\sim} e^{-\sum_{k=1}^K \lambda_k p_k(\frac{\mathbf{i}}{N})}.$$

From the strong law of large numbers therefore, if $i_k := \lfloor N x_k \rfloor$, $k = 1, \dots, K$, then, given $\mathbf{x} = x_k$, $k = 1, \dots, K$:

$$\frac{I'_k}{N} \xrightarrow[N \uparrow \infty]{a.s.} x'_k = p_k(\mathbf{x})$$

which is (13): When the population under study is very large, random fluctuations as modelled by (22) can be ignored so that the gene frequencies evolves deterministically.

Fixation probabilities. Let \mathbf{e}_l be the K -null vector except for its l th entry which is 1. The extremal pure states $S_{K,N}^{ex} := \{\mathbf{i}_l^{ex} := N\mathbf{e}_l, l = 1, \dots, K\}$, are all absorbing for this Markov chain because $p_k \left(\frac{\mathbf{i}_l^{ex}}{N} \right) = \delta_{k,l}$ and, from (22), any additional fixed point which \mathbf{p} could have on the boundary-faces of S_K which are not points does not give rise to an absorbing state for P . Under our assumptions on \mathbf{p} , the chain is not recurrent, rather it is transient. Depending on the initial condition, say \mathbf{i}_0 , the chain will necessarily end up in one of the extremal states \mathbf{i}_l^{ex} , with some fixation probability, say $\pi_l(\mathbf{i}_0)$, which can be computed as follows. Let $\boldsymbol{\pi}_l := \pi_l(\mathbf{i})$, $\mathbf{i} \in S_{K,N}$ be the harmonic function of the Wright-Fisher Markov chain which is the smallest solution to the boundary problem:

$$(23) \quad (I - P) \boldsymbol{\pi}_l = \mathbf{0} \text{ if } \mathbf{i} \in S_{K,N} \setminus S_{K,N}^{ex} \text{ and } \boldsymbol{\pi}_l = 1 \text{ (} \mathbf{i} = \mathbf{i}_l^{ex} \text{) if } \mathbf{i} \in S_{K,N}^{ex}.$$

We also have

$$\mathbb{P}(\mathbf{I}_\tau = N\mathbf{e}_l \mid \mathbf{I}_0 = \mathbf{i}_0) = \pi_l(\mathbf{i}_0),$$

where τ ($< \infty$ almost surely) is the random hitting time of $S_{K,N}^{ex}$ for \mathbf{I}_t and the $\pi_l(\mathbf{i}_0)$ s are normalized so as $\sum_l \pi_l(\mathbf{i}_0) = 1$. Thus $\pi_l(\mathbf{i}_0)$ are the searched probabilities to end up in state \mathbf{i}_l^{ex} starting from state \mathbf{i}_0 .

From (23), $\boldsymbol{\pi}_l(\mathbf{i}_0)$ is known if $\mathbf{i}_0 \in S_{K,N}^{ex}$. The remaining unknown restriction, say $\boldsymbol{\pi}_l^Q$, of $\boldsymbol{\pi}_l = \pi_l(\mathbf{i}_0)$ to the non-extremal states is easily seen to be:

$$(24) \quad \boldsymbol{\pi}_l^Q = (I - Q)^{-1} \mathbf{p}_{\mathbf{i}_l^{ex}}, \mathbf{i}_0 \in S_{K,N} \setminus S_{K,N}^{ex}.$$

In (24), Q is obtained from P after erasing the lines and columns corresponding to all the K extremal states and $\mathbf{p}_{\mathbf{i}_l^{ex}}$ is the \mathbf{i}_l^{ex} -column of P where the entries corresponding to the extremal states have been deleted. When dealing with the Wright-Fisher model, Q is a positive matrix and also $\mathbf{p}_{\mathbf{i}_l^{ex}} > \mathbf{0}$, therefore $\pi_l(\mathbf{i}_0) > 0$ for all $\mathbf{i}_0 \in S_{K,N} \setminus S_{K,N}^{ex}$ and this for each l : starting from any state \mathbf{i}_0 which is not an extremal state, there is a positive probability to hit any of the extremal states. Fixation of the state \mathbf{i}_l^{ex} means extinction of the remaining monomorphic states. It would therefore be of interest to understand the structure of the set $\mathcal{A}_l = \{\mathbf{i}_0 : \pi_l(\mathbf{i}_0) > \sum_{k \neq l} \pi_k(\mathbf{i}_0)\} = \{\mathbf{i}_0 : \pi_l(\mathbf{i}_0) > 1/2\}$, for each l , which is the stochastic version of the attraction basin of \mathbf{i}_l^{ex} ; especially when l is the label of the extremal state with largest fitness $w_{k,k}$. If $\mathbf{i}_0 \in \mathcal{A}_l$ indeed, the probability to end up in \mathbf{i}_l^{ex} is larger than the probability to end up in any other extremal state.

Unfortunately, the development of the inverse of $I - Q$ in terms of its adjugate matrix in (24) shows that these fixation probabilities have a very complex determinantal alternating structure and the question of identifying \mathcal{A}_l is very complex.

With $\boldsymbol{\pi}_l$ the solution (24) to the Dirichlet problem with boundary conditions (23), the equilibrium measure of the chain therefore is:

$$\pi_{eq}(\mathbf{i}_0) := \sum_{l=1}^K \pi_l(\mathbf{i}_0) \delta_{\mathbf{i}_l^{ex}},$$

which depends on \mathbf{i}_0 . Necessarily, one allele will fixate and there is no polymorphic equilibrium state even when dealing with diploid populations. Which allele and with what probability will depend on the initial condition. Thanks to fluctuations, the picture therefore looks very different from the one pertaining to the deterministic theory.

Of importance also is the time it takes to get extinct. It relies on similar techniques. For instance, the expected overall fixation time $\alpha(\mathbf{i}_0) := \mathbb{E}_{\mathbf{i}_0}(\tau)$ solves the boundary problem:

$$\begin{aligned} (I - P)\alpha &= \mathbf{1}, \mathbf{i}_0 \in S_{K,N} \setminus S_{K,N}^{ex} \\ \alpha &= 0, \mathbf{i}_0 \in S_{K,N}^{ex} \end{aligned}$$

where $\alpha := \alpha(\mathbf{i}_0)$, $\mathbf{i}_0 \in S_{K,N}$. The restriction α^Q of α to the non-extremal states therefore is

$$\alpha^Q = (I - Q)^{-1} \mathbf{1}, \mathbf{i}_0 \in S_{K,N} \setminus S_{K,N}^{ex}.$$

Conditioning \mathbf{I}_t on non-fixation. There are four places where questions relative to conditioning on fixation are relevant in this context.⁴

(i) Consider the full fixation vector $\pi_l := \pi_l(\mathbf{i}_0)$. Remove from π_l the states \mathbf{i}_0 for which $\pi_l(\mathbf{i}_0) = 0$ only, keeping the one, \mathbf{i}_l^{ex} , for which $\pi_l(\mathbf{i}_l^{ex}) = 1$. The size of this vector, say π_l^R , is $n_R = \binom{N+K-1}{K-1} - K + 1$. Let R be the corresponding $n_R \times n_R$ reduced transition matrix obtained from P after erasing the lines and columns corresponding to all the K monomorphic states except \mathbf{i}_l^{ex} . The Markov chain \mathbf{I}_t conditioned to exit in the extremal state \mathbf{i}_l^{ex} only admits the stochastic transition matrix:

$$R_l := D_{\pi_l^R}^{-1} R D_{\pi_l^R}.$$

It is obtained from R after a diagonal Doob transform based on π_l^R . The chain governed by R_l admits a unique absorbing state which is \mathbf{i}_l^{ex} . The entries of R_l are

$$R_l(\mathbf{i}, \mathbf{i}') = \frac{\pi_l^R(\mathbf{i}')}{\pi_l^R(\mathbf{i})} R(\mathbf{i}, \mathbf{i}')$$

and for this new conditioned Markov chain, transitions to states \mathbf{i}' for which $\pi_l^R(\mathbf{i}') > \pi_l^R(\mathbf{i})$ are favored.

(ii) Let us now consider again the fully reduced transition matrix Q obtained from P after erasing the lines and columns corresponding to all the K monomorphic states. The matrix Q is substochastic and irreducible, with $Q\mathbf{1} < \mathbf{1}$. The law of the process corresponding to \mathbf{I}_t conditioned on avoiding the monomorphic states before t evolves as follows: With $\tau := \wedge_{l=1}^K \tau_l$ the time needed for first hitting one of the extremal states for \mathbf{I}_t , let $\pi_t(\mathbf{i}) = \mathbb{P}(\mathbf{I}_t = \mathbf{i} \mid \tau > t)$. Then, with $\pi_t = \pi_t(\mathbf{i})$,

⁴Similar conditioning problems were considered in [4] in the context of the 2-alleles Wright-Fisher diffusion.

$$\mathbf{i} \in S_{K,N} \setminus S_{K,N}^{ex},$$

$$\pi_{t+1}^* = \frac{\pi_t^* Q}{\pi_t^* Q \mathbf{1}}$$

is the nonlinear Master Equation governing its evolution ([13]). The reduced state-space dimension of this Markov chain is $n_Q = \binom{N+K-1}{K-1} - K$ and $\pi_t \xrightarrow[t \uparrow \infty]{} \pi_\infty$ where π_∞ is the left Perron probability eigenvector of Q associated to the dominant Perron eigenvalue $\rho_Q < 1$, namely: $\rho_Q \pi_\infty^* = \pi_\infty^* Q$.⁵ If the process is started using this limiting quasi-stationary distribution, it remains in the same state over time and the fixation time τ is geometrically distributed with success probability ρ_Q .

(iii) One can define another stochastic process $\bar{\mathbf{I}}_t$ which admits the stochastic transition matrix: $\bar{Q} := D_{Q\mathbf{1}}^{-1} Q$ again defined on the reduced state-space. For each t , we have:

$$\bar{Q}(\mathbf{i}, \mathbf{i}') = \mathbb{P}(\mathbf{I}_{t+1} = \mathbf{i}' \mid \mathbf{I}_t = \mathbf{i}, \tau > 1)$$

and the conditioning on non-fixation occurs at each transition time. This process is an ergodic Markov chain with invariant probability measure solving $\bar{\pi}_{eq}^* = \bar{\pi}_{eq}^* \bar{Q}$. It has the following closed-form determinantal expression (see [17], Section 6 and [16] p. 1559):

$$(25) \quad \bar{\pi}_{eq}(\mathbf{i}) = \frac{|(I - \bar{Q})_{[\mathbf{i}, \mathbf{i}]}|}{\sum_{\mathbf{i}} |(I - \bar{Q})_{[\mathbf{i}, \mathbf{i}]}|}, \quad \mathbf{i} \in S_{K,N} \setminus S_{K,N}^{ex},$$

where $(I - \bar{Q})_{[\mathbf{i}, \mathbf{i}]}$ is the submatrix resulting from the deletion of row \mathbf{i} and column \mathbf{i} of $I - \bar{Q}$. The question as to whether the process governed by \bar{Q} is reversible or not arises. Defining the transition matrix of the time-reversed process \overleftarrow{Q} by:

$$\overleftarrow{Q}^* = D_{\bar{\pi}_{eq}} \bar{Q} D_{\bar{\pi}_{eq}}^{-1},$$

it does not hold that $\overleftarrow{Q} = \bar{Q}$ and so detailed balance does not hold. Indeed, \bar{Q} is similar to the transition matrix of an ergodic Wright-Fisher model and Wright-Fisher chains are not reversible.

(iv) If we condition on non-fixation in the remote future (see [9] for additional details), we get a Markov chain whose stochastic transition matrix is:

$$\tilde{Q} = \rho_Q^{-1} D_{\psi_\infty}^{-1} Q D_{\psi_\infty}.$$

Here ψ_∞ is the positive right Perron eigenvector of Q associated to the Perron eigenvalue $\rho_Q < 1$ satisfying: $\rho_Q \psi_\infty = Q \psi_\infty$. This vector can be chosen so that: $\sum_k \pi_{\infty,k} \psi_{\infty,k} = 1$, where π_∞ is again the left Perron probability eigenvector of Q associated to $\rho_Q < 1$ (See [2]). With $\tilde{\pi}_t(\cdot) = \lim_{s \uparrow \infty} \mathbb{P}(\mathbf{I}_t = \cdot \mid \tau > t + s)$, we have $\tilde{\pi}_{t+1}^* = \tilde{\pi}_t^* \tilde{Q}$. The process governed by \tilde{Q} is an ergodic Markov chain whose invariant probability measure is $\tilde{\pi}_\infty = \pi_\infty \circ \psi_\infty$, the Schur product of π_∞ and ψ_∞ with k th entry $\tilde{\pi}_{\infty,k} = \pi_{\infty,k} \psi_{\infty,k}$.

⁵ π_∞ is called a Yaglom limit (see [22]) or a quasi-stationary distribution..

For the last three conditionings, it is difficult to extract some information on the limiting distribution, either π_∞ or $\bar{\pi}_{eq}$ or $\tilde{\pi}_\infty$, respectively. This is because it would suppose to solve the eigenvalue problems explicitly which is out of reach, at least theoretically. However, assuming a diploid population with a polymorphic equilibrium state \mathbf{x}_{eq} for \mathbf{p} , we expect that these distributions will present a global (local) maximum (minimum) near \mathbf{x}_{eq} if \mathbf{x}_{eq} is stable (unstable). These limiting distributions should be more sharply peaked around the extremum if we consider the conditioning (iv) compared to (ii) because, the latter conditioning being more stringent than the former, it should charge more heavily the sample paths that stay away from the monomorphic states.

Finally, we would like to stress that all these considerations are also relevant in the context of another fundamental stochastic model arising in the context of evolutionary genetics. We shall give some elements of how to proceed with this model presenting very different properties.

3.2. The K -alleles Moran model. We now focus on the Moran model.

The multiallelic Moran model. Let $\alpha, \beta \in \{1, \dots, K\}$. In the Moran version of the stochastic evolution, given $\mathbf{I}_t = \mathbf{I} = \mathbf{i}$, the only accessible values of \mathbf{I}' are the neighboring states: $\mathbf{i}'_{\alpha,\beta} := \mathbf{i} + \mathbf{d}_{\alpha,\beta}$ where $\mathbf{d}_{\alpha,\beta}^* := (0, \dots, 0, -1, 0, \dots, 1, 0, \dots, 0)$. Here -1 is in position α and 1 in position $\beta \neq \alpha$ corresponding to the transfer of an individual from the box α (if non-empty) to the box β . With $n(\mathbf{i}) = \#\{k : i_k > 0\}$ the number of non-empty entries of \mathbf{i} , there are $n(\mathbf{i})(K-1) \leq K(K-1)$ accessible states from \mathbf{i} . The Moran stochastic evolutionary dynamics is now given by a Markov chain whose one-step transition matrix P from states $\mathbf{I} = \mathbf{i}$ to $\mathbf{I}' = \mathbf{i}'$ is:

$$(26) \quad \mathbb{P}(\mathbf{I}_{t+1} = \mathbf{i}' \mid \mathbf{I}_t = \mathbf{i}) = 0 \text{ if } \mathbf{i}' \neq \mathbf{i}'_{\alpha,\beta} \text{ and}$$

$$\mathbb{P}(\mathbf{I}_{t+1} = \mathbf{i}'_{\alpha,\beta} \mid \mathbf{I}_t = \mathbf{i}) =: P(\mathbf{i}, \mathbf{i}'_{\alpha,\beta}) = \frac{i_\alpha}{N} p_\beta \left(\frac{\mathbf{i}}{N} \right),$$

where $\mathbf{p} = (p_\beta \left(\frac{\mathbf{i}}{N} \right), \beta = 1, \dots, K)$ is either given by (4) in the haploid case or by (13) in the diploid case.

Summing $P(\mathbf{i}, \mathbf{i}'_{\alpha,\beta})$ over $\alpha, \beta, \beta \neq \alpha$ in (26), we get the holding probability

$$\mathbb{P}(\mathbf{I}_{t+1} = \mathbf{i} \mid \mathbf{I}_t = \mathbf{i}) = 1 - \sum_{\alpha, \beta: \beta \neq \alpha} \frac{i_\alpha}{N} p_\beta \left(\frac{\mathbf{i}}{N} \right) = \sum_{\alpha} \frac{i_\alpha}{N} p_\alpha \left(\frac{\mathbf{i}}{N} \right),$$

completing the characterization of the K -alleles Moran model. Under our assumptions on \mathbf{p} , the holding probabilities are equal to 1 only for the extremal states $\mathbf{i} \in S_{K,N}^{ex}$ which are therefore the only absorbing states of the Moran chain, just like for the Wright-Fisher model. The drift at state \mathbf{i} is:

$$\mathbb{E}(\mathbf{I}_{t+1} - \mathbf{I}_t \mid \mathbf{I}_t = \mathbf{i}) = \sum_{\alpha} \frac{i_\alpha}{N} \sum_{\beta \neq \alpha} p_\beta \left(\frac{\mathbf{i}}{N} \right) \mathbf{d}_{\alpha,\beta}.$$

Let us compute the Laplace-Stieltjes Transform of \mathbf{I}' in the context of a Moran model. Omitting the argument $\frac{\mathbf{i}}{N}$ in p_β , we get the factorized form:

$$\begin{aligned}
\mathbb{E}_{\mathbf{i}} \left(e^{-\sum_k \lambda_k I'_k} \right) &= \sum_{\alpha, \beta: \alpha \neq \beta} e^{-\sum_k \lambda_k i'_{\alpha, \beta}(k)} P(\mathbf{i}, \mathbf{i}'_{\alpha, \beta}) + e^{-\sum_k \lambda_k i_k} \sum_{\beta} \frac{i_\beta}{N} p_\beta \\
&= e^{-\sum_k \lambda_k i_k} \left(\sum_{\alpha, \beta: \alpha \neq \beta} e^{-\sum_k \lambda_k \mathbf{d}_{\alpha, \beta}(k)} P(\mathbf{i}, \mathbf{i}'_{\alpha, \beta}) + \sum_{\beta} \frac{i_\beta}{N} p_\beta \right) \\
&= e^{-\sum_k \lambda_k i_k} \left(\sum_{\alpha, \beta: \alpha \neq \beta} e^{-(\lambda_\beta - \lambda_\alpha)} \frac{i_\alpha}{N} p_\beta + \sum_{\beta} \frac{i_\beta}{N} p_\beta \right) \\
&= e^{-\sum_k \lambda_k i_k} \left(\sum_{\beta} e^{-\lambda_\beta} p_\beta \sum_{\alpha \neq \beta} \frac{i_\alpha}{N} e^{\lambda_\alpha} + \sum_{\beta} \frac{i_\beta}{N} p_\beta \right) \\
&= e^{-\sum_k \lambda_k i_k} \left(\sum_{\beta} e^{-\lambda_\beta} p_\beta \left(\sum_{\alpha} \frac{i_\alpha}{N} e^{\lambda_\alpha} - \frac{i_\beta}{N} e^{\lambda_\beta} \right) + \sum_{\beta} \frac{i_\beta}{N} p_\beta \right) \\
&= \left(e^{-\sum_k \lambda_k i_k} \right) \left(\sum_{\alpha} \frac{i_\alpha}{N} e^{\lambda_\alpha} \right) \left(\sum_{\beta} e^{-\lambda_\beta} p_\beta \right).
\end{aligned}$$

Putting $\lambda_l = 0$ if $l \neq k$, the k th marginal reads:

$$\mathbb{E}_{\mathbf{i}} \left(e^{-\lambda_k I'_k} \right) = e^{-\lambda_k i_k} \left(1 - \frac{i_k}{N} + e^{\lambda_k} \frac{i_k}{N} \right) (1 - p_k + e^{-\lambda_k} p_k)$$

showing that $I_k(t)$ is of the random walk type. Indeed, we get: $\mathbb{P}_{\mathbf{i}}(I'_k = i'_k) = 0$ if $i'_k \neq i_k \pm 1$ or $i'_k \neq i_k$ and

$$\begin{aligned}
\mathbb{P}_{\mathbf{i}}(I'_k = i_k) &= \left(1 - \frac{i_k}{N} \right) (1 - p_k) + \frac{i_k}{N} p_k \\
\mathbb{P}_{\mathbf{i}}(I'_k = i_k + 1) &= \left(1 - \frac{i_k}{N} \right) p_k = \sum_{l \neq k} P(\mathbf{i}, \mathbf{i}'_{l, k}) \\
\mathbb{P}_{\mathbf{i}}(I'_k = i_k - 1) &= \frac{i_k}{N} (1 - p_k) = \sum_{l \neq k} P(\mathbf{i}, \mathbf{i}'_{k, l}).
\end{aligned}$$

We have

$$\mathbb{E}_{\mathbf{i}} \left(\frac{I'_k}{N} \right) = \frac{i_k}{N} + \frac{1}{N} \left(p_k - \frac{i_k}{N} \right); \quad \sigma_{\mathbf{i}}^2 \left(\frac{I'_k}{N} \right) = \frac{1}{N^2} \left(\frac{i_k}{N} \left(1 - \frac{i_k}{N} \right) + p_k (1 - p_k) \right).$$

There is only a small correction (of order N^{-1}) of the updated mean to its current value and fluctuations around the mean are small too (of order N^{-1}). The evolution process is very slow.

Fixation probabilities. As a random walk model, the Moran model has a much simpler transition matrix P of the Jacobi type. The equilibrium measure of the chain again is:

$$(27) \quad \pi_{eq} := \sum_{l=1}^K \pi_l(\mathbf{i}_0) \delta_{\mathbf{i}_l^{eq}},$$

where π_l again solves the Dirichlet problem (23) but with this new simpler Jacobi P . For the Moran model, the explicit expression (24) of the fixation probability simplifies a little bit because $p_{\mathbf{i}_l^{ex}}(\mathbf{i}) \neq 0$ only for the $K - 1$ neighboring states of \mathbf{i}_l^{ex} that is $\{\mathbf{i} : \mathbf{i} + \mathbf{d}_{\alpha,\beta} = \mathbf{i}_l^{ex}\}$ for some $\mathbf{d}_{\alpha,\beta}$.

When $K = 2$ (2 alleles), the random walk transition probabilities ($p_1 + p_2 = 1$)

$$\mathbb{P}(\mathbf{I}_{t+1} = \mathbf{i}'_{1,2} \mid \mathbf{I}_t = \mathbf{i}) =: P\left(\mathbf{i}, \mathbf{i} + \begin{pmatrix} -1 \\ 1 \end{pmatrix}\right) = \frac{i_1}{N} p_2 \left(\frac{\mathbf{i}}{N}\right),$$

$$\mathbb{P}(\mathbf{I}_{t+1} = \mathbf{i}'_{2,1} \mid \mathbf{I}_t = \mathbf{i}) =: P\left(\mathbf{i}, \mathbf{i} + \begin{pmatrix} 1 \\ -1 \end{pmatrix}\right) = \frac{i_2}{N} p_1 \left(\frac{\mathbf{i}}{N}\right),$$

are the probabilities that the first component $I_{t,1}$ of \mathbf{I}_t moves down and up by one unit respectively. In this case, the Dirichlet problem giving the fixation probabilities solves explicitly. With $\phi(i_0) = 1 + \sum_{i=1}^{i_0-1} \prod_{i_1=1}^i \frac{i_1 p_2(\frac{\mathbf{i}}{N})}{i_2 p_1(\frac{\mathbf{i}}{N})}$ ($\phi(0) = 0$) the harmonic function of the 2-alleles chain, we easily get that

$$\pi_1(i_0, N - i_0) = \frac{\phi(i_0)}{\phi(N)}$$

is the probability that the extremal state $\mathbf{i}_1^{ex} = (N, 0)$ is reached given $\mathbf{i}_0 = (i_0, N - i_0)$. Assuming a model with multiplicative fitnesses: $p_\alpha(\frac{\mathbf{i}}{N}) = \frac{i_\alpha}{N} \omega(\frac{w_\alpha}{N})$, then ϕ takes the simple form ($i_1 + i_2 = N$)

$$\phi(i_0) = 1 + \sum_{i=1}^{i_0-1} \prod_{i_1=1}^i \frac{i_1 p_2(\frac{\mathbf{i}}{N})}{i_2 p_1(\frac{\mathbf{i}}{N})} = 1 + \sum_{i=1}^{i_0-1} \left(\frac{w_2}{w_1}\right)^i$$

showing (See [1], p. 109) that:

$$\pi_1(i_0, N - i_0) = \frac{1 - \left(\frac{w_2}{w_1}\right)^{i_0}}{1 - \left(\frac{w_2}{w_1}\right)^N}.$$

Assuming $w_1 = 1 + s/N$ and $w_2 = 1$, putting $i_0 = [Nx_0]$, for large N , we get [6]

$$\pi_1(Nx_0, N - i_0) \sim \frac{1 - e^{-sx_0}}{1 - e^{-s}}.$$

In the general fitness case:

$$\phi(i_0) = 1 + \sum_{i=1}^{i_0-1} \prod_{i_1=1}^i \frac{i_1 p_2(\frac{\mathbf{i}}{N})}{i_2 p_1(\frac{\mathbf{i}}{N})} = 1 + \sum_{i=1}^{i_0-1} \prod_{i_1=1}^i \frac{(W\mathbf{i}/N)_2}{(W\mathbf{i}/N)_1}$$

where, as usual, $(W\mathbf{i}/N)_k = \sum_{l=1}^2 w_{k,l} i_l / N$, $k, l = 1, 2$, leading to

$$\begin{aligned} (W\mathbf{i}/N)_1 &= w_{1,2} + (w_{1,1} - w_{1,2}) i_1 / N \\ (W\mathbf{i}/N)_2 &= w_{2,2} + (w_{2,1} - w_{2,2}) i_1 / N. \end{aligned}$$

This 2-alleles exact solution can be used in the full diploid K -alleles Moran case with multiplicative fitnesses. Indeed, from this, with $\mathbf{i}_0 = (i_1, \dots, i_K)$, the fixation

probability of A_l can be conjectured to be approximated qualitatively by:

$$(28) \quad \pi_l(\mathbf{i}_l^{ex}) = 1 \text{ and } \pi_l(\mathbf{i}_0) = \frac{1 - \left(\frac{\sum_{k \neq l} i_k w_k / (N - i_l)}{w_l} \right)^{i_l}}{1 - \left(\frac{\sum_{k \neq l} i_k w_k / (N - i_l)}{w_l} \right)^N} \text{ if } \mathbf{i}_0 \neq \mathbf{i}_l^{ex}.$$

This can be justified as follows: mark one particular box with size i_l , corresponding to the allele A_l with fitness w_l . Then clump the $K - 1$ remaining boxes into a single box with size $N - i_l$, corresponding to a fictitious allele with average fitness $\sum_{k \neq l} i_k w_k / (N - i_l)$. We are left with a 2-alleles Moran multiplicative fitness model for which (from the 2-alleles exact solution) the fixation probability of A_l is given by (28). This formula constitutes sort of a mean field approximation to the full Dirichlet problem associated to the Moran model.

Assuming $w_k \sim 1 + s_k/N$, a Kimura-like approximation of (28) would lead for large N to:

$$\pi_l(N\mathbf{x}_0) \sim \frac{1 - e^{\frac{x_l}{1-x_l} \sum_{k \neq l} x_k s_k}}{1 - e^{\frac{1}{1-x_l} \sum_{k \neq l} x_k s_k}},$$

where $\mathbf{x}_0 = (x_1, \dots, x_K)$ is now a point of the continuous K -simplex different from the l -unit vector $\mathbf{e}_l := (0, \dots, 0, 1, 0, \dots, 0)$.

From (28), when the fitness of allele A_l is large (small), compared to the average fitness of the remaining alleles, then the fixation probability of A_l gets close to 1 (respectively to 0). Note also that the larger i_l is, the larger the fixation probability is.

As required also, for all $k \neq l$, $\pi_l(\mathbf{i}_k^{ex}) = \left(1 - \left(\frac{w_k}{w_l}\right)^0\right) / \left(1 - \left(\frac{w_k}{w_l}\right)^N\right) = 0$. As another particular initial configuration case, suppose we start from the 2-alleles type state: $\mathbf{i}_0 =: \mathbf{i}_0(k, l) = (0, \dots, 0, N - 1, 0, \dots, 0, 1, 0, \dots, 0)$ where the 1 is in position l (that is: $i_l = 1$) and the entry $N - 1$ in position $k \neq l$ (that is: $i_k = N - 1$). Although for this choice of the initial state, the fixation of A_k is very likely, there still is a positive probability that allele A_l gets fixed which is seen to be from (28):

$$\pi_l(\mathbf{i}_0(k, l)) = \frac{1 - \frac{w_k}{w_l}}{1 - \left(\frac{w_k}{w_l}\right)^N},$$

depending only on the relative fitnesses of A_l and A_k (See [19] and [18] for a similar expression). As N gets large, this probability gets close to $1 - \frac{w_k}{w_l}$ if $w_k < w_l$ and close to $\left(\frac{w_l}{w_k}\right)^{N-1} \sim 0$ if $w_k > w_l$.

Conversely, assuming the 2-alleles type state \mathbf{i}_0 to be defined by $i_l = N - 1$ and $i_k = 1$,

$$\pi_l(\mathbf{i}_0) = \frac{1 - \left(\frac{w_k}{w_l}\right)^{N-1}}{1 - \left(\frac{w_k}{w_l}\right)^N},$$

which, as required, gets close to 1 as N gets large if $w_k < w_l$ and close to $\frac{w_l}{w_k}$ if $w_k > w_l$.

Conditioning \mathbf{I}_t on non-fixation. The conditioning developments discussed for the Wright-Fisher model are also relevant in the Moran model context substituting the P of Moran for the P of Wright-Fisher. Let us revisit the conditioning (iii).

(iii) With Q now the reduced substochastic matrix of the full Moran transition matrix defined in (26), consider the stochastic process $\bar{\mathbf{I}}_t$ with stochastic transition matrix: $\bar{Q} := D_{Q1}^{-1}Q$ defined on the reduced state-space with dimension n_Q . The process $\bar{\mathbf{I}}_t$ is again an ergodic Markov chain with invariant probability measure solving $\bar{\pi}_{eq}^* = \bar{\pi}_{eq}^* \bar{Q}$. With \bar{Q} the transition matrix of the time-reversed process it now holds that $\bar{Q} = \bar{Q}$ and so detailed balance holds when dealing with the Moran case (see also [5]). This will be proved if we can exhibit an equilibrium probability measure $\bar{\pi}_{eq}$ such that

$$(29) \quad \bar{\pi}_{eq}(\mathbf{i}) \bar{Q}(\mathbf{i}, \mathbf{i}') = \bar{Q}(\mathbf{i}', \mathbf{i}) \bar{\pi}_{eq}(\mathbf{i}'),$$

for all neighboring state $\mathbf{i}' = \mathbf{i} + \mathbf{d}_i$ of \mathbf{i} .

With $\mathbf{j} = (j_1, \dots, j_K)$ any terminal state, suppose we want to use (29) to compute $\bar{\pi}_{eq}(\mathbf{j})$ starting from the smallest available state in the system which is $\mathbf{j}_0 = (0, \dots, 0, 1, N-1)$. This is possible because $\bar{\pi}_{eq}(\mathbf{j})$ may be represented as

$$(30) \quad \bar{\pi}_{eq}(\mathbf{j}) = \bar{\pi}_{eq}(\mathbf{j}_0) \prod_{\mathbf{i}=\mathbf{j}_0}^{\mathbf{j}-\mathbf{d}_{K,1}} \frac{\bar{Q}(\mathbf{i}, \mathbf{i}')}{\bar{Q}(\mathbf{i}', \mathbf{i})},$$

where $\bar{\pi}_{eq}(\mathbf{j}_0)$ can be chosen so that: $\sum_{\mathbf{j}} \bar{\pi}_{eq}(\mathbf{j}) = 1$. Let us give some details.

- Note first, by reversing path, that for two consecutive states $(\mathbf{i}, \mathbf{i}')$, the ratio $\frac{\bar{Q}(\mathbf{i}, \mathbf{i}')}{\bar{Q}(\mathbf{i}', \mathbf{i})}$ can be computed. We have

$$\frac{\bar{Q}(\mathbf{i}, \mathbf{i}')}{\bar{Q}(\mathbf{i}', \mathbf{i})} = \frac{(Q1)_{\mathbf{i}'} Q(\mathbf{i}, \mathbf{i}')}{(Q1)_{\mathbf{i}} Q(\mathbf{i}', \mathbf{i})}$$

where $\mathbf{i}' = \mathbf{i} + \mathbf{d}_i$ for some \mathbf{d}_i of the form $\mathbf{d}_{k,l}$ and therefore

$$(31) \quad Q(\mathbf{i}, \mathbf{i}') = \frac{i_k}{N} p_l \left(\frac{\mathbf{i}}{N} \right) \text{ and } Q(\mathbf{i}', \mathbf{i}) = Q(\mathbf{i}', \mathbf{i}' + \mathbf{d}_{l,k}) = \frac{i_l + 1}{N} p_k \left(\frac{\mathbf{i} + \mathbf{d}_{k,l}}{N} \right).$$

Clearly, from such a structure of the entries of \bar{Q} , for each possible transition $\mathbf{i} \rightarrow \mathbf{i}'$, the ratio $\bar{Q}(\mathbf{i}, \mathbf{i}') / \bar{Q}(\mathbf{i}', \mathbf{i})$ will only depend separately on a ratio involving the terminal and the initial states \mathbf{i}' and \mathbf{i} (the detailed balance condition holds).

- Second, the sequence of \mathbf{i} -s in (30) is governed by the following path starting from \mathbf{j}_0 and ending up in the target state \mathbf{j} : We can use the following sequence of $\mathbf{d}_{k,l}$ s:

$$\left(\mathbf{d}_{K,K-1}^{j_{K-1}-1} \mathbf{d}_{K,K-1} \right) \left(\mathbf{d}_{K,K-2}^{j_{K-2}} \mathbf{d}_{K,K-1} \right) \dots \left(\mathbf{d}_{K,1}^{j_1} \mathbf{d}_{K,1} \right),$$

filling up successively the entries of \mathbf{j} to the left of the last entry of \mathbf{j}_0 by using the $N-1$ individuals of the reservoir state $\mathbf{j}_0 = (0, \dots, 0, 1, N-1)$. By doing so, each intermediate state \mathbf{i} is separated from the next \mathbf{i}' by some clearly identified \mathbf{d}_i , and after evaluating the probability ratio $\bar{Q}(\mathbf{i}, \mathbf{i}') / \bar{Q}(\mathbf{i}', \mathbf{i})$ for each consecutive states of this sequence, we are done. Because there exists a probability distribution $\bar{\pi}_{eq}(\mathbf{j})$ such that the reversibility identity (29) holds, then this Moran process is

reversible with (30) as its stationary distribution. Using (31), this constitutes an exact explicit product-form formula.

4. CONCLUDING REMARKS

This paper studies a classical population genetic model describing a one-locus multiallelic population subject to natural selection, random mating and then random genetic drifts as from the Wright-Fisher and Moran models. Although the two-alleles version of this model is fairly well-studied and understood, this is not so much the case of the multiallelic one, especially in the discrete-time context which we adopt here. Let us summarize our results emphasizing the ones which we believe are new.

Considering first the deterministic updating mechanisms driven by selection, we underline that it has the form of a nonlinear Master equation suggesting that it is possible to construct an underlying Markov process governed by this Master equation. We briefly and intuitively supply such a construction.

In the diploid context, we pay attention on a class of fitness matrices that leads to polymorphism. Would the equilibrium polymorphic state be unstable, we suggest that the class of potential matrices constitute a large such admissible class of fitness matrices. It contains the class of strictly ultrametric matrices which therefore deserves some interest. To the best of our knowledge, there is no discussion of such fitness models in the population genetics context. Would the polymorphic state be stable, we derive a related class of fitness matrices leading to a definite-negative mean fitness quadratic form. Some simple examples are supplied and detailed.

The last Section is devoted to the stochastic version of these considerations taking into account an additional important driving source of evolution, namely the random genetic drift. When driven by selection only and in particular in the absence of mutations, the multiallelic Wright-Fisher model is a transient Markov chain whose absorbing states are the monomorphic states. We give an expression for the fixation probabilities for this process. Then, we develop four conditioning problems: conditioning on fixating in a given monomorphic state, conditioning on avoiding the extremal states before the current instant, conditioning on non-fixation at each transition time and conditioning on avoiding the extremal states in the remote future. Finally we run into similar considerations but for the Moran model. When dealing with the fixation probabilities in this Moran context, we suggest a mean-field approximation of these probabilities which is based on a well-known explicit formula for the 2-alleles case. It concerns the case of multiplicative fitnesses only. Finally, we consider the Moran model conditioned on non-fixation at each transition time. We exploit the reversible character of this process to derive an explicit product formula for its invariant probability measure.

REFERENCES

- [1] Ewens, W. J. *Mathematical population genetics. I. Theoretical introduction*. Second edition. Interdisciplinary Applied Mathematics, 27. Springer-Verlag, New York, 2004.
- [2] Horn, R.A.; Johnson, C.R. *Matrix analysis*. Cambridge University Press, Cambridge, 1985.
- [3] Huillet, T. E. Information and (co)variances in discrete evolutionary genetics involving solely selection. *J. Stat. Mech.* (2009) P090
- [4] Huillet, T. On Wright-Fisher diffusion and its relatives. *J. Stat. Mech.* (2007) P11006

- [5] Khare, K. and Zhou, H. Rates of convergence of some multivariate Markov chains with polynomial eigenfunctions. *Ann. Appl. Probab.* Volume 19, Number 2 (2009), 737-777.
- [6] Kimura, M. On the probability of fixation of mutant genes in a population. *Genetics*, 47, 713-719, (1962).
- [7] Kingman, J. F. C. A mathematical problem in population genetics. *Proc. Cambridge Philos. Soc.* 57, 574-582, 1961.
- [8] Kingman, J. F. C. *Mathematics of genetic diversity*. CBMS-NSF Regional Conference Series in Applied Mathematics, 34. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, Pa., 1980. vii+70 pp. ISBN: 0-89871-166-5
- [9] Lambert, A. Population dynamics and random genealogies. *Stochastic Models* 24, suppl. 1, 45-163, (2008).
- [10] Lewontin, R. C., Ginzburg L. R. and Tuljapurkar, S. D. Heterosis as an explanation for large amounts of genic polymorphism. *Genetics*, 88, 149-170, (1978).
- [11] Martinez, S.; Michon, G. and San Martin J. Inverse of strictly ultrametric matrices are of Stieltjes type. *SIAM J. Matrix Anal. Appl.* 15 (1994), pp. 98-106]
- [12] Maruyama, T. *Stochastic problems in population genetics*. Lecture Notes in Biomathematics, 17. Springer-Verlag, Berlin-New York, 1977.
- [13] McKean, H. P. In: A.K. Aziz, Editor, *Lectures in Differential Equations* vol. 2, Van Nostrand Reinhold Company, New York (1969), p. 177-193.
- [14] Nabben, R.; Varga, R. S. A linear algebra proof that the inverse of a strictly ultrametric matrix is a strictly diagonally dominant Stieltjes matrix. *SIAM J. Matrix Anal. Appl.* 15 (1994), no. 1, 107-113.
- [15] Nagylaki, T.; Lou, Y. Multiallelic selection polymorphism. *Theoretical Population Biology*, Volume 69, Issue 2, Pages 217-229, (2006).
- [16] Rached, Z.; Alajaji, F.; Campbell, L. Rényi's divergence and entropy rates for finite alphabet Markov sources. *IEEE Transactions on Information Theory* 47(4): 1553-1561 (2001).
- [17] Romanovsky, V. I. *Discrete Markov chains*. Translated from the Russian by E. Seneta, Wolters-Noordhoff Publishing, Groningen 1970.
- [18] Sella, G. An exact steady state solution of Fisher's geometric model and other models. *Theoretical Population Biology*, 75(1), 30-34, 2009.
- [19] Sella, G.; Hirsh, A. E. The application of statistical physics to evolutionary biology. *PNAS* 102(27), 9541-9546, 2005.
- [20] Shashahani, S. A new mathematical framework for the study of linkage and selection. *Mem. Amer. Math. Soc.*, 17, 1979.
- [21] Svirezhev, Y. M. Optimum principles in genetics. In: *Studies on Theoretical Genetics*, V. A. Ratner (Ed), pp 86-102, 1972, Novosibirsk, USSR Academy of Science.
- [22] Yaglom, A. M. Certain limit theorems of the theory of branching random processes. *Doklady Akad. Nauk SSSR (N.S.)* 56, 795-798, (1947).