

Supporting Text for Prevolutionary Dynamics

Martin A. Nowak & Hisashi Ohtsuki

Program for Evolutionary Dynamics, Department of Organismic and Evolutionary Biology, Department of Mathematics, Harvard University, Cambridge MA 02138, USA

1 Prolife

Prolife dynamics are given by

$$\dot{x}_i = a_i x_{i'} - (d + a_{i0} + a_{i1}) x_i. \quad (1)$$

The index i represents all binary strings (sequences). Longer strings are produced from shorter ones by adding 0 or 1 on the right side. Each string, i , has one precursor, i' , and two followers, $i0$ and $i1$. For example, the precursor of string 0101 is 010; the two followers are 01010 and 01011. For the precursors of strings 0 and 1 we set $x_{0'} = x_{1'} = 1$. The constants a_i denote the rate at which string i arises from i' by addition of an activated monomer (which is either 0* or 1*). Eq.(1) assumes that the concentration of activated monomers is constant. All strings are removed (die) at rate d .

Prolife dynamics define a tree with the activated monomers at the root. The tree of prolife has infinitely many lineages. A lineage is a sequence of strings that follow each other. For example, one such lineage is 0, 00, 000,

At equilibrium, the right hand side of Eq.(1) is zero, so we obtain

$$x_i = b_i x_{i'}, \quad (2)$$

where b_i is given by

$$b_i = \frac{a_i}{d + a_{i0} + a_{i1}}. \quad (3)$$

Using Eq.(2) recursively gives us

$$x_i = b_i b_{i'} b_{i''} \cdots b_\sigma, \quad (4)$$

where σ is the ancestral monomer (0 or 1) of sequence i .

Let us consider super-symmetric prelife with $a_0 = a_1 = \alpha/2$ and $a_i = a$ for all other sequences, i . From Eq.(4), we obtain the following results.

The abundance of a sequence of length n is

$$x_n = \frac{\alpha}{2a} \left(\frac{a}{2a+d} \right)^n. \quad (5a)$$

The total abundance of all sequences of length n is

$$X_n = 2^n x_n = \frac{\alpha}{2a} \left(\frac{2a}{2a+d} \right)^n. \quad (5b)$$

The total abundance of all sequences is

$$X = \sum_{n=1}^{\infty} X_n = \frac{\alpha}{d}. \quad (5c)$$

The total abundance of all sequences in one lineage is

$$\tilde{X} = \sum_{n=1}^{\infty} x_n = \frac{\alpha}{2(a+d)}. \quad (5d)$$

The average sequence length is

$$\bar{n} = \frac{\sum_{n=1}^{\infty} n X_n}{X} = 1 + \frac{2a}{d}. \quad (5e)$$

Although there are infinitely many lineages, the abundance of any one lineage is a considerable fraction of the entire population. The reason is that short sequences belong to many lineages and they are much more abundant than long sequences.

2 Prelife landscape

Let us consider a random prelife landscape where reaction rates of sequences of length more than two are randomly given by

$$a_i = \begin{cases} a + s & \text{(with prob. } p) \\ a & \text{(with prob. } 1 - p). \end{cases} \quad (6)$$

The other parameters are the same as before: $a_0 = a_1 = \alpha/2$.

From Eq.(4), at equilibrium we obtain the following results. The average abundance of a sequence of length n is

$$\bar{x}_n = \frac{\alpha}{2} AB^n, \quad (7)$$

where

$$A = \frac{(2a + d)^2 + (2a + d)(3 - 2p)s + 2(1 - p)^2 s^2}{a(2a + d)^2 + (2a + d)(3a + pd)s + \{2a + p(2 - p)d\}s^2} \quad (8)$$

and

$$B = \frac{a(2a + d)^2 + (2a + d)(3a + pd)s + \{2a + p(2 - p)d\}s^2}{(2a + d)(2a + d + s)(2a + d + 2s)}. \quad (9)$$

A sequence is selected if its equilibrium abundance is not vanishing as $s \rightarrow \infty$. For sequence i of length n , rewriting Eq.(4) yields

$$x_i = \frac{1}{d + a_{i0} + a_{i1}} \cdot \underbrace{\frac{a_i}{d + a_{i'0} + a_{i'1}} \cdot \frac{a_{i'}}{d + a_{i''0} + a_{i''1}} \cdots \frac{a_{\sigma\rho}}{d + a_{\sigma0} + a_{\sigma1}}}_{n-1 \text{ terms}} \cdot \frac{\alpha}{2}, \quad (10)$$

where $\sigma\rho$ represents the first two digits of sequence i . The first term in the right hand side of Eq.(10) is

$$\begin{cases} \frac{1}{(a+s)+(a+s)+d} \xrightarrow{s \rightarrow \infty} 0 & \text{(with prob. } p^2) \\ \frac{1}{(a+s)+a+d} \xrightarrow{s \rightarrow \infty} 0 & \text{(with prob. } 2p(1 - p)) \\ \frac{1}{a+a+d} \xrightarrow{s \rightarrow \infty} \frac{1}{a+a+d} & \text{(with prob. } (1 - p)^2). \end{cases} \quad (11)$$

The first term does not vanish with probability $(1 - p)^2$.

For each of the next $n - 1$ terms on the right hand side of Eq.(10) we have

$$\left\{ \begin{array}{ll} \frac{a+s}{(a+s)+(a+s)+d} \xrightarrow{s \rightarrow \infty} \frac{1}{2} & \text{(with prob. } p^2) \\ \frac{a+s}{(a+s)+a+d} \xrightarrow{s \rightarrow \infty} 1 & \text{(with prob. } p(1-p)) \\ \frac{a}{(a+s)+a+d} \xrightarrow{s \rightarrow \infty} 0 & \text{(with prob. } p(1-p)) \\ \frac{a}{a+a+d} \xrightarrow{s \rightarrow \infty} \frac{1}{a+a+d} & \text{(with prob. } (1-p)^2). \end{array} \right. \quad (12)$$

Each term does not vanish with probability $1 - p(1 - p)$. Therefore, the probability that a sequence of length n is selected (does not vanish) is given by

$$(1-p)^2[1-p(1-p)]^{n-1}. \quad (13)$$

The expected number of sequences of length n that are selected is

$$2^n(1-p)^2[1-p(1-p)]^{n-1}. \quad (14)$$

For example, if $a = 1, d = 1, \alpha = 1$ and $p = 1/2$ as in Figure 2, we obtain from Eq.(7) for the average abundance of sequences of length n

$$\bar{x}_n = \frac{18 + 12s + s^2}{36 + 42s + 11s^2} \left(\frac{36 + 42s + 11s^2}{12(3+s)(3+2s)} \right)^n. \quad (15)$$

Note that $\bar{x}_n(s)$ a monotonically decreasing function (of s) for $n \leq 3$, a one-humped function for $3 < n < 12$, and a monotonically increasing function for $n \geq 12$. From Eq.(14), the expected number of sequences of length n that survive for large s is given by $(1/3)(3/2)^n$.

3 Master sequence

In this section, we study the case where all reactions leading to one particular sequence (the master sequence) occur at the increased rate b , while all other reactions occur at rate a .

Suppose $0^n = \underbrace{00 \cdots 0}_n$ is the master sequence. The reaction rates are given by

$$\begin{aligned} a_0 &= a_1 = \alpha/2 \\ a_i &= b \quad \text{for } i = 00, \cdots, 0^n \\ a_i &= a \quad \text{for other } i. \end{aligned} \quad (16)$$

From the general formula, Eq.(4), the abundances of sequences $i = \underbrace{0 \cdots 0}_\ell \underbrace{1 * \cdots *}_m$ at equilibrium are given by

$$x_i = \begin{cases} \frac{\alpha}{2a} \left(\frac{a}{2a+d} \right)^m & \text{if } \ell = 0 \\ \frac{\alpha}{2b} \left(\frac{b}{a+b+d} \right)^\ell \left(\frac{a}{2a+d} \right)^m & \text{if } 1 \leq \ell \leq n-1 \\ \frac{\alpha}{2a} \left(\frac{b}{a+b+d} \right)^{n-1} \left(\frac{a}{2a+d} \right)^{\ell+m+1-n} & \text{if } \ell \geq n. \end{cases} \quad (17)$$

In particular, we are interested in the abundances of all sequences that have the same length as the master sequence. Let x_i denote the abundance of a sequence of the form $\underbrace{0 \cdots 0}_i \underbrace{1 * \cdots *}_{n-i}$. In this notation, x_n represents the abundance of the master sequence. From eq.(17), we obtain

$$x_i = \begin{cases} \frac{\alpha}{2a} \left(\frac{a}{2a+d} \right)^n & \text{if } i = 0 \\ \frac{\alpha}{2b} \left(\frac{b}{a+b+d} \right)^i \left(\frac{a}{2a+d} \right)^{n-i} & \text{if } 1 \leq i \leq n-1 \\ \frac{\alpha}{2a} \left(\frac{b}{a+b+d} \right)^{n-1} \left(\frac{a}{2a+d} \right) & \text{if } i = n. \end{cases} \quad (18)$$

Since $b > a$, we find

$$x_0 > x_1 < x_2 < \cdots < x_{n-1} < x_n \quad \text{and} \quad x_0 < x_n. \quad (19)$$

The master sequence is most abundant among all sequences of length n .

If $b \rightarrow \infty$, then the abundance of the master sequence converges to

$$x_{n,max} = \lim_{b \rightarrow \infty} x_n = \frac{\alpha}{2(2a + d)}. \quad (20)$$

Let us now calculate the condition for the abundance of the master sequence, x_n , to exceed a fraction, $1/k$, of the maximum value, $x_{n,max}$. From Eqs.(18) and (20), we have

$$\frac{\alpha}{2a} \left(\frac{b}{a + b + d} \right)^{n-1} \left(\frac{a}{2a + d} \right) > \frac{1}{k} \cdot \frac{\alpha}{2(2a + d)}. \quad (21)$$

This condition is rewritten as

$$b > \frac{a + d}{k^{\frac{1}{n-1}} - 1} \approx \frac{a + d}{\ln k} n \quad (n \gg 1). \quad (22)$$

Hence, for a master sequence of length n to make up a significant fraction of the population, the rate constant b must grow as a linear function of n .

4 Master sequence with mutation

As before, we assume that all reactions leading to the master sequence occur at an increased rate, b , but there is a probability u of incorporating the wrong monomer. The rate of those reactions that stay within the lineage leading to the master sequence is given by $b(1 - u)$, while the reactions that come off the lineage occur at rate $a + bu$. We have

$$\begin{aligned} a_0 &= a_1 = \alpha/2 \\ a_i &= b(1 - u) \quad \text{for } i = 00, \dots, 0^n \\ a_i &= a + bu \quad \text{for } i = 01, \dots, 0^{n-1}1 \\ a_i &= a \quad \text{for all other } i. \end{aligned} \quad (23)$$

Consider sequences of the form $i = \underbrace{0 \dots 0}_\ell \underbrace{1 * \dots *}_m$. As always the asterisks represent either 0 or 1. From the general formula, Eq.(4), the equilibrium abundance of

sequence i is given by

$$x_i = \begin{cases} \frac{\alpha}{2a} \left(\frac{a}{2a+d} \right)^m & \text{if } \ell = 0 \\ \frac{\alpha}{2b(1-u)} \left(\frac{b(1-u)}{a+b+d} \right)^\ell & \text{if } 1 \leq \ell \leq n-1, m = 0 \\ \frac{\alpha}{2b(1-u)} \cdot \frac{a+bu}{a} \left(\frac{b(1-u)}{a+b+d} \right)^\ell \left(\frac{a}{2a+d} \right)^m & \text{if } 1 \leq \ell \leq n-1, m \geq 1 \\ \frac{\alpha}{2a} \left(\frac{b(1-u)}{a+b+d} \right)^{n-1} \left(\frac{a}{2a+d} \right)^{\ell+m+1-n} & \text{if } \ell \geq n. \end{cases} \quad (24)$$

Let us now compare the abundances of all sequences of length n . Let x_i denote the abundances of sequences of the form $\underbrace{0 \cdots 0}_i \underbrace{1 * \cdots *}_{n-i}$. In this notation, the abundance of the master sequence is given by x_n . From eq.(24), we obtain

$$x_i = \begin{cases} \frac{\alpha}{2a} \left(\frac{a}{2a+d} \right)^n & \text{if } i = 0 \\ \frac{\alpha}{2a} \cdot \frac{a+bu}{b(1-u)} \left(\frac{b(1-u)}{a+b+d} \right)^i \left(\frac{a}{2a+d} \right)^{n-i} & \text{if } 1 \leq i \leq n-1 \\ \frac{\alpha}{2a} \left(\frac{b(1-u)}{a+b+d} \right)^{n-1} \left(\frac{a}{2a+d} \right) & \text{if } i = n. \end{cases} \quad (25)$$

In order to understand the relative ranking of the equilibrium abundances of all sequences of length n , we must distinguish three cases.

Case (i) $u < \frac{a}{2a+d}$:

(i-a) If $b < \frac{a(a+d)}{(a+d)-u(2a+d)}$ then $x_0 < x_1 > x_2 > \cdots > x_{n-1} > x_n$.

(i-b) If $\frac{a(a+d)}{(a+d)-u(2a+d)} < b < \frac{a}{1-2u}$ then $x_0 < x_1 < x_2 < \cdots < x_{n-1} > x_n$.

(i-c) If $\frac{a}{1-2u} < b < \frac{a^2}{a-u(2a+d)}$ then $x_0 < x_1 < x_2 < \cdots < x_{n-1} < x_n$.

(i-d) If $b > \frac{a^2}{a-u(2a+d)}$ then $x_0 > x_1 < x_2 < \cdots < x_{n-1} < x_n$ and $x_0 < x_n$.

Case (ii) $\frac{a}{2a+d} \leq u < \frac{a+d}{2a+d}$:

(ii-a) If $b < \frac{a(a+d)}{(a+d)-u(2a+d)}$ then $x_0 < x_1 > x_2 > \cdots > x_{n-1} > x_n$.

(ii-b) If $\frac{a(a+d)}{(a+d)-u(2a+d)} < b < \frac{a}{1-2u}$ then $x_0 < x_1 < x_2 < \cdots < x_{n-1} > x_n$.

(ii-c) If $b > \frac{a}{1-2u}$ then $x_0 < x_1 < x_2 < \cdots < x_{n-1} < x_n$.

Case (iii) $u \geq \frac{a+d}{2a+d}$:

(iii-a) If $b < \frac{a}{1-2u}$, then $x_0 < x_1 > x_2 > \cdots > x_{n-1} > x_n$.

(iii-b) If $b > \frac{a}{1-2u}$, then $x_0 < x_1 > x_2 > \cdots > x_{n-1} < x_n$ and $x_1 > x_n$.

In summary, the equilibrium abundance of the master sequence is

$$x_n = \frac{\alpha}{2(2a+d)} \left(\frac{b(1-u)}{a+b+d} \right)^{n-1}. \quad (26)$$

The master sequence is most abundant among all sequences of length n if

$$u < \frac{a+d}{2a+d} \quad \text{and} \quad b > \frac{a}{1-2u}. \quad (27)$$

If $b \rightarrow \infty$, then the abundance of the master sequence converges to

$$x_{n,max} = \lim_{b \rightarrow \infty} x_n = \frac{\alpha}{2(2a+d)} (1-u)^{n-1}. \quad (28)$$

For x_n to exceed a fraction, $1/k$, of this maximum value, $x_{n,max}$, we need

$$\frac{\alpha}{2(2a+d)} \left(\frac{b(1-u)}{a+b+d} \right)^{n-1} > \frac{1}{k} \cdot \frac{\alpha}{2(2a+d)} (1-u)^{n-1}, \quad (29)$$

which is simplified to

$$b > \frac{a+d}{k^{\frac{1}{n-1}} - 1} \approx \frac{a+d}{\ln k} n. \quad (n \gg 1). \quad (30)$$

If $b \rightarrow \infty$ and $u \rightarrow 0$, then the abundance of the master sequence converges to

$$\hat{x}_{n,max} = \lim_{\substack{b \rightarrow \infty \\ u \rightarrow 0}} x_n = \frac{\alpha}{2(2a+d)}. \quad (31)$$

For x_n to exceed a fraction, $1/k$, of this maximum value, $\hat{x}_{n,max}$, we need

$$\frac{\alpha}{2(2a+d)} \left(\frac{b(1-u)}{a+b+d} \right)^{n-1} > \frac{1}{k} \cdot \frac{\alpha}{2(2a+d)}, \quad (32)$$

which is rewritten as

$$\left(\frac{a+b+d}{b(1-u)} \right)^{n-1} < k. \quad (33)$$

When $b \gg a+d$, $u \ll 1$ and $n \gg 1$, the left hand side of Eq.(33) is approximated by

$$\left[\left(1 + \frac{a+d}{b} \right) (1+u) \right]^n \approx \left(1 + \frac{a+d}{b} + u \right)^n \approx \exp \left[n \left(\frac{a+d}{b} + u \right) \right]. \quad (34)$$

Therefore condition (33) is simplified to

$$\frac{a+d}{b} + u < \frac{\ln k}{n}. \quad (35)$$

For $u = 0$ we obtain the previous condition on b . For $b \rightarrow \infty$ we obtain the error-threshold

$$u < \frac{\ln k}{n}. \quad (36)$$

The mutation rate of prelife must be less than the inverse of the sequence length, for the master sequence to reach a significant abundance in the population.

5 Replication

Let us assume that some sequences have the ability to replicate. Incorporating replication into prelife dynamics leads to the following differential equation:

$$\dot{x}_i = a_i x_{i'} - (d + a_{i0} + a_{i1})x_i + r x_i (f_i - \phi). \quad (37)$$

The first part of this equation describes prelife as before. The second part represents the standard selection equation. The coefficient, r , measures the relative contribution of selection dynamics in Eq.(37). The fitness of sequence i is given by f_i . The quantity, ϕ , is an additional death rate, which cancels out the additional production of sequences by replication. From

$$\sum_i r x_i (f_i - \phi) = 0, \quad (38)$$

we have

$$\phi = \frac{\sum_i f_i x_i}{\sum_i x_i}. \quad (39)$$

In other words, ϕ represents the average fitness of the population.

For $r = 0$, replication is absent and we recover prelife dynamics, Eq.(1). For $r \rightarrow \infty$, replication dominates and we obtain the standard selection dynamics.

We define the net reproductive rate of sequence i as

$$g_i \equiv r(f_i - \phi) - (d + a_{i0} + a_{i1}). \quad (40)$$

As in the main text, the sign of the net reproductive rate predicts a phase transition between prelife and life.

6 Replication with mutation

Imagine that sequence i of length n is the unique replicator, but its replication is susceptible to errors. In each elongation step a wrong monomer is attached with

probability u .

Let f_i be the fitness of the replicator in the absence of errors. As replication is error-free with probability $(1 - u)^n$, the realized fitness of the replicator becomes

$$(1 - u)^n f_i. \quad (41)$$

For the replicator to be selected, the net reproductive rate must be positive:

$$g_i = r\{(1 - u)^n f_i - \phi\} - (d + a_{i0} + a_{i1}) > 0. \quad (42)$$

By using

$$(1 - u)^n \approx \exp(-un) \quad (u \ll 1 \text{ and } n \gg 1), \quad (43)$$

and by neglecting ϕ (which is very small at the error threshold), condition (42) can be rewritten as

$$u < \frac{1}{n} \log \left[\frac{r f_i}{d + a_{i0} + a_{i1}} \right]. \quad (44)$$

Therefore, the replicator is selected if the mutation rate is less than the inverse of the sequence length.