

On models of evolution of species

Rahul Roy

Indian Statistical Institute, New Delhi

Joint work with Hideki Tanemura (Keio University, Yokohama)

arXiv: 1909.09759

Structure of the talk

The first part of the talk is an introduction to some mathematical models on the evolution of species.

In the second part we present our contribution.

“Those who are horrified at Mr. Darwin’s theory, may comfort themselves with the assurance that, if we are descended from the ape, we have not descended so far as to preclude all hope of return.”

Ambrose Bierce 1874, *The Fiend’s delight*.

The Bak-Sneppen model

Consider a population consisting of N species located on a circle. Also let $\{f_i(0) : 1 \leq i \leq N\}$ be i.i.d. $\text{Uniform}[0, 1]$ random variables. $f_i(0)$ represents the **fitness** of the i th species initially.

A simple evolutionary mechanism

Let $\{f_i(n) : 1 \leq i \leq N\}$ be the fitness values at time n , the fitness values at time $n + 1$ is given by

$$f_i(n+1) := \begin{cases} f_i(n) & \text{if } f_i(n) \neq \mathbf{\min}\{f_l(n) : 1 \leq l \leq N\} \\ g_i(n) & \text{if } f_i(n) = \mathbf{\min}\{f_l(n) : 1 \leq l \leq N\} \end{cases}$$

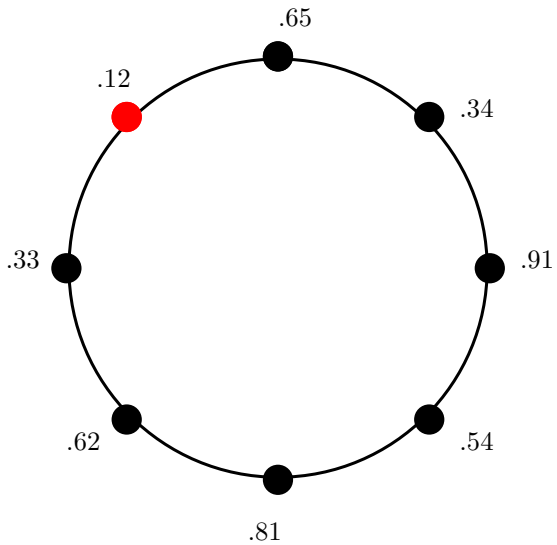
where $\{g_i(n) : 1 \leq i \leq N, n \geq 0\}$ is another collection of i.i.d. $\text{Uniform}[0, 1]$ random variables, independent of the previous collection.

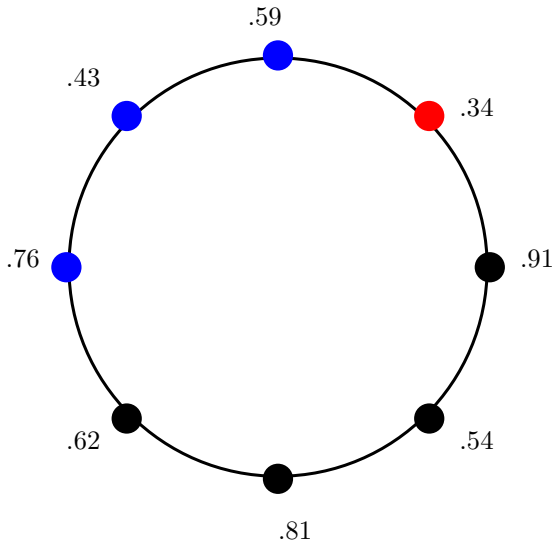
This is really a trivial model and it can be easily shown that $f_i^N(n) \rightarrow 1$ in distribution as $n, N \rightarrow \infty$.

So we need to add a twist to it.

We do an i.i.d. Uniform[0, 1] sampling not only at the site of the minimum fitness, but also its two neighbours, i.e.

$$f_i(n+1) := \begin{cases} g_i(n) & \text{if } f_i^{(n)} = \mathbf{\min}\{f_l(n) : 1 \leq l \leq N\} \\ g_j(n) & \text{if } j = i \pm 1 \text{ and } f_i(n) = \mathbf{\min}\{f_l(n) : 1 \leq l \leq N\} \\ f_i(n) & \text{for all other sites.} \end{cases}$$





Bak and Sneppen (1993)¹ observed that for large N that the 1-dimensional marginals are uniform on $(f_c, 1)$ for some $f_c \sim 2/3$.

Power law behaviour

Assuming the existence of f_c , physicists study **avalanches**.

Fix $0 < q < 1$ and let

$$\tau_1 = \inf\{n : f_i^N(n) > q \text{ for all } 1 \leq i \leq N\}$$

and

$$\tau_2 = \inf\{n > \tau_1 : f_i^N(n) > q \text{ for all } 1 \leq i \leq N\}.$$

The time $\tau_2 - \tau_1$ is the **length of an avalanche of fitnesses below q** .

For $q = f_c$, its distribution is supposed to have a power law behaviour.

¹Bak, P. and Sneppen, K. (1993). Punctuated equilibrium and criticality in a simple model of evolution. Phys. Rev. Lett., 74, 4083–4086.

q=0,56

0.6390	0.8958	0.6686	0.9367	0.8230	0.8034	0.5628	0.5991	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	n
0.6390	0.8958	0.6686	0.9367	0.8230	0.1752	0.4742	0.4726	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	n+1
0.6390	0.8958	0.6686	0.9367	0.1697	0.7738	0.7412	0.4726	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	
0.6390	0.8958	0.6686	0.2409	0.1259	0.3001	0.7412	0.4726	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	
0.6390	0.8958	0.6686	0.5477	0.6114	0.0782	0.7412	0.4726	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	
0.6390	0.8958	0.6686	0.5477	0.3931	0.3915	0.1207	0.4726	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	
0.6390	0.8958	0.6686	0.5477	0.3931	0.8757	0.5937	0.7068	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	n+6
0.6390	0.8958	0.6686	0.3636	0.8600	0.5634	0.5937	0.7068	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	
0.6390	0.8958	0.6777	0.8808	0.5634	0.5937	0.7068	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059		
0.6390	0.7437	0.5188	0.0971	0.8808	0.5634	0.5937	0.7068	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	
0.6390	0.7437	0.6038	0.6972	0.5377	0.5634	0.5937	0.7068	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	n+10
0.6390	0.7437	0.6038	0.5772	0.7767	0.6691	0.5937	0.7068	0.9981	0.6765	0.7533	0.7663	0.7680	0.6059	

Range=7

Duration=11

Taken from

<http://membres-timc.imag.fr/Herve.Guiol/activites/ParisSlidessimple.pdf>

Rigorous results

Let F_N be the distribution function of the 1-dimensional marginal in the stationary regime with N sites

Theorem

(Meester and Znamenski, 2003)

There exists $0 < q < 1$ and $c_q > 0$ such that, uniformly in N , we have

$$F_N(q) > c_q.$$

This establishes the non-triviality of the Bak-Sneppen model.

Avalanches

First a triviality – let $N = 3$.

Let $X_k := \min\{f_i(k) : i = 1, 2, 3\}$. Since the system is flushed at every time point, X_0, X_1, \dots is a sequence of i.i.d. random variables. So, by symmetry,

$$P(X_0 > \max\{X_1, \dots, X_n\}) = 1/n.$$

So, if L is the length of an avalanche of fitness X_0 , then

$$P(L > n) = 1/n, \text{ and so, } E(L) = \infty.$$

In case $N > 3$ the above argument fails because we do not flush the system at every time point.

However, for any N if $X_0 := \min\{f_i(0) : i = 1, \dots, N\}$, then

Theorem

(Gillett, Meester and van der Wal, 2006)

Let L be the length of an avalanche of fitness X_0 , then $E(L) = \infty$.

This lends validity to the ‘power law’ belief of the physicists.

Before we end this discussion on the Bak-Sneppen model we mention a curious result from this same paper.

Let $n > N$ and $Y_0 := \min\{f_1(0), \dots, f_N(0), U_1, \dots, U_{n-N}\}$, where $\{U_1, U_2, \dots\}$ is another collection of i.i.d. Uniform[0, 1] random variables. Then, for the model with N sites, we have

$E(L) < \infty$, where L is the length of an avalanche of fitness Y_0 ;

which again lends credence to the hypothesis of physicists that the power law phenomenon holds only for avalanches of fitness $f > f_c$.

Meester and Sarkar (2012)

Meester and Sarkar (2012) considered a modified Bak-Sneppen model, which retained the imprimatur of the original Bak-Sneppen model..

Instead of replacing the least fit and its two neighbours, they considered a model where, besides the least fit, a random one is chosen uniformly from the remaining $N - 1$.

They show that, if the initial configuration has an individual of fitness 0 and the rest have a fitness $> f$, then the length L_f^N of an avalanche at level f satisfies

Theorem

1. For $f < 1/2$ and all N ,

$$P(L_f^N > n) \leq \mathbf{exp}(-c_1(f)n) \text{ for some positive constant } c_1(f),$$

2. For $f > 1/2$,

$$\lim_{N \rightarrow \infty} P(L_f^N > n) \geq c_2(f) \text{ for some positive constant } c_2(f),$$

3. For $f = 1/2$,

$$\lim_{n \rightarrow \infty} \lim_{N \rightarrow \infty} \sqrt{n}P(L_f^N > n) = 2/\sqrt{\pi}.$$

Liggett and Schinazi (2009)

Liggett and Schinazi (2009) wanted to study the following question (which is very topical)

“Korber, et al.² noted that the influenza virus is less diverse worldwide than the HIV virus is in Amsterdam alone. However, both types of (phylogenetic) tree are supposed to be produced by the same mechanism: mutations. Can the same mathematical model produce two trees that are so different?”

²Evolutionary and immunological implications of contemporary HIV-1 variation. *British Med. Bull.* 58, 19–42.

Their model is a continuous time birth and death chain, with every birth having a fitness component which is an i.i.d. random variable and each death is of the individual which has the smallest fitness.

The birth rate is

$$n \rightarrow n + 1 \text{ at rate } \lambda n$$

and the death rate

$$n \rightarrow n - 1 \text{ at rate } n \text{ provided } n \geq 1.$$

Since only the ordering of the fitness is needed, w.l.o.g assume that the fitness is Uniform[0, 1] distributed.

Let $M(t) :=$ maximal fitness at time t .

They show

Theorem

(Liggett and Schinazi (2009))

For $\alpha \in (0, 1)$, we have

$$\lim_{t \rightarrow \infty} P(M(\alpha t) = M(t)) = \begin{cases} \alpha & \text{if } \lambda \leq 1 \\ 0 & \text{if } \lambda < 1. \end{cases}$$

So for $\lambda < 1$, the dominating type (i.e. the fittest type) at time t has likely been present for a time of order t and at any given time there will not be many types. This is consistent with the observed structure of an influenza tree.

While, for $\lambda > 1$ the dominating type at time t has likely been present for a time of order smaller than t and at any given time there will be many types. This is consistent with an HIV tree.

Guiol, Machado and Schinazi (2010)

The Bak-Sneppen model suffers from two major shortcomings:–

1. The number of species is fixed, although we study the asymptotics as $N \rightarrow \infty$.
2. It could be the case that along with the least fit species, we also remove the most fit, because it happened to be the neighbour of the least fit species.

Guiol, Machado and Schinazi (2010) proposed a model where

1. The number of species present at any instant of time is random.
2. Only the least fit species die.

The GMS model is a ‘birth and death’ model defined as follows:

- (i) At time 0 there is a particle at 0.
- (ii) With probability p there is a birth and the individual is born with a **fitness** $f \sim \text{Unif}[0, 1]$, independent of other births.
- (iii) With probability $1 - p$ a death takes place and the individual with the smallest fitness is removed.

So if

$$L_n^f := \#\{\text{individuals with fitness } \leq f \text{ at time } n\}$$

$$R_n^f := \#\{\text{individuals with fitness } > f \text{ at time } n\}$$

$$N_n := L_n^f + R_n^f = \text{population size at time } n$$

then, for $k, \ell \geq 1$,

$$\begin{aligned} & P\left((L_{n+1}^f, R_{n+1}^f) = (k, \ell)\right) \\ &= \begin{cases} pf & \text{if } (L_n^f, R_n^f) = (k-1, \ell) \\ 1-p & \text{if } (L_n^f, R_n^f) = (k+1, \ell) \\ p(1-f) & \text{if } (L_n^f, R_n^f) = (k, \ell-1) \end{cases} \end{aligned}$$

Similarly we may write the transition probabilities when either k or ℓ equals 0.

So for $f = f_c := \frac{1-p}{p}$ (with $p > 1/2$), we have

$$L_{n+1}^f - L_n^f = \begin{cases} +1 & \text{w.p. } 1 - p \\ -1 & \text{w.p. } 1 - p \\ 0 & \text{w.p. } 2p - 1. \end{cases}$$

So L_n^f is a symmetric random walk with reflecting boundary at 0 and hence

$$P(L_n^f = 0 \text{ infinitely often}) = 1 \text{ for } f \leq f_c.$$

Also, a simple SLLN argument shows,

$$\text{for } f_c < a < b \leq 1, \text{ we have } \frac{R_n^a - R_n^b}{n} \rightarrow p(b - a) \text{ almost surely.}$$

This is from Guiol, Machado and Schinazi (2010), building on the work of Liggett and Schinazi (2006).

GMS– continuous version

Guiol, Machado and Schinazi (2013) has a continuous version of the previous model.

Here a new species is born at a rate λ and an existing species dies at rate μ .

At birth, a new species has a random (positive) fitness determined by a distribution F .

And the death is of the species with the least fitness.

Assume that the distribution function F admits a density.

Suppose at time 0 there are k species with fitnesses

$$0 < f_1 < \dots < f_k = f.$$

Let $\tau_f^k :=$ time of survival of the species with fitness f .

Theorem

τ_f^k has a Bessel distribution given by

$$P(\tau_f^k \leq t) = \left(\frac{\mu}{\lambda_f}\right)^{k/2} \int_0^t e^{-s(\mu+\lambda_f)} \frac{k}{s} I_k(2\sqrt{s\mu\lambda_f}) ds,$$

where $\lambda_f = \lambda F(f)$ and I_k is the modified Bessel function of the first kind with index k defined by

$$I_k(x) = \sum_{l=0}^{\infty} \frac{1}{(l+k)! l!} (x/2)^{2l+k}.$$

From the theorem we observe the following asymptotic behaviour

1. If $\lambda_f < \mu$ then

$$P(\tau_f^k > t) \sim C_k \frac{e^{-\gamma t}}{t^{3/2}}$$

where $\gamma = (\sqrt{\mu} - \sqrt{\lambda_f})^2$ and C_k depends on k , λ_f and μ .

2. If $\lambda_f > \mu$ then

$$P(\tau_f^k = \infty) = 1 - \left(\frac{\mu}{\lambda_f}\right)^k,$$

$$P(t < \tau_f^k < \infty) \sim C_k \frac{e^{-\gamma t}}{t^{3/2}}.$$

3. If $\lambda_f = \mu$

$$P(\tau_f^k > t) \sim k \frac{1}{\sqrt{t\pi\mu}}.$$

Thus there is a phase transition:

If $\lambda > \mu$ then taking $f_c = F^{-1}(\mu/\lambda)$, then species with fitness larger than f_c have a positive probability of eternal survival.

Also as in the discrete model:

The number of species at time t with fitnesses less than f_c is a null recurrent birth and death process,

while the number of species at time t with fitnesses in (a, b) with $f_c < a < b$ behaves like $t^{\frac{\lambda(F(b)-F(a))}{\lambda+\mu}}$ asymptotically almost surely.

Ben-Ari and Schinazi (2016):

Returning to the discrete GMS model, suppose at birth (an event which occurs with probability p), either

- a) with probability r a mutant is born with a fitness $f \sim \text{Unif}[0, 1]$, independent of other births.
- b) with probability $1 - r$ the individual born has a fitness chosen uniformly at random among the fitnesses of the existing individuals at that time.

A 'death' removes all the individuals with the smallest fitness.

Here again we have a similar phase transition at $f_c = \frac{1-p}{pr}$ (we assume $pr > 1 - p$).

Also, for a given n, k, f , let

$U_n^k(f) := \#\{s \in [f, 1] : \text{there are exactly } k \text{ individuals with fitness } s\}$.

i.e., number of sites in $[f, 1]$ with a population of exactly k at time n .

For $A \subset \mathbb{N}^* \times [0, 1]$ Borel, consider the empirical distribution

$$H_n(A) := \begin{cases} \sum_{(k,f) \in A} [U_n^k(f) - U_n^k(f+)] & \text{if } N_n(0) > 0 \\ \delta_{(0,0)}(A) & \text{if } N_n(0) = 0 \end{cases}$$

Ben-Ari and Schinazi (2016) showed that

H_n converges weakly to a product measure of $\text{Geom}\left(\frac{pr-(1-p)}{p-(1-p)}\right)$ and $\text{Unif}[f_c, 1]$

Michael and Volkov (2012)

Our work is on the GMS model, however before we talk about our work, we discuss a variant of the GMS model introduced by Michael and Volkov (2012).

Let X_1, X_2, \dots and Z_1, Z_2, \dots be two independent collections of i.i.d. positive integer valued random variables.

Let $0 < p < 1$ and $T_0 = 0$. At time n , the state of the system T_n is a finite subset of $[0, 1]$. The Markov process is as follows:

At time $n + 1$, with probability p , we have $\#T_{n+1} = \#T_n + Z_n$, with each of the new Z_n individuals being assigned $\text{Unif}[0, 1]$ fitnesses, independent of other individuals and independent of other random variables;

with probability $1 - p$, we have $\#T_{n+1} = \mathbf{max}\{\#T_n - X_n, 0\}$, and we remove all species with the smallest X_n fitness.

Let $\mu_Z = EZ$ and $\mu_X = EX$, then we have

1. Suppose $\mu_Z = \infty$ and $\mu_X < \infty$ then, as $n \rightarrow \infty$, we have that T_n approaches a random sample from $\text{Unif}[0, 1]$, in the sense that if B_n is the set of all species born until time n and D_n the set of all species removed until time n , then $T_n = B_n \setminus D_n$ and

$$\limsup_n \frac{\#(T_n \Delta B_n)}{B_n} = 0 \text{ almost surely .}$$

2. Suppose $\mu_Z < \infty$ and $\mu_X < \infty$ and let $p > p_c := \frac{\mu_X}{\mu_X + \mu_Z}$. Then T_n approaches a random sample from $\text{Unif}[f, 1]$, where $f = \frac{(1-p)\mu_X}{p\mu_Z} \in (0, 1)$.

3. Suppose $\mu_Z < \infty$ and (i) $\mu_X < \infty$ with $p < p_c$ or (ii) $\mu_X = \infty$, then $T_n = \emptyset$ for infinitely many n .

The preferential attachment model

At time 0 there is one individual of fitness 0. At time n , there is either a birth or a death of an individual from the existing population with probability p or $1 - p$ respectively, and independent of any other random mechanism considered earlier.

- (P1) In case of a birth, there are two possibilities.
- (i) with probability r , a mutant is born and has a fitness parameter f uniformly at random in $[0, 1]$, or
 - (ii) with probability $1 - r$ the individual born has a fitness f with a probability proportional to the number of individuals with fitness f among the entire population present at that time. Here we have a caveat that, if there is no individual present at the time of birth, then the fitness of the individual is sampled uniformly in $[0, 1]$.
- (P2) In case of a death, an individual from the population at the site closest to 0 is eliminated.

The formal structure

Let $X_n = \{(k_i, x_i) : k_i \geq 1, x_i \in [0, 1], i = 1, \dots, \ell\}$, where the total population at time n is divided in exactly ℓ sites x_1, \dots, x_ℓ , with the size of the population at site x_i being exactly k_i . In case there is no individual present at time n we take $X_n = \emptyset$.

The process X_n is Markovian on the state space

$$\mathcal{S} := \{\emptyset\} \cup \{ \{(k, x)\}_{x \in \Lambda} : (k, x) \in \mathbb{N} \times [0, 1], \#\Lambda < \infty, \}.$$

For a given $f \in (0, 1)$, let L_n^f denote the size of the population at time n at sites in $[0, f]$,

$$L_n^f := \sum k_s : \text{sum over } s \in [0, f] \text{ and } (k_s, s) \in X_n,$$

R_n^f denote the size of the population at time n at sites in $(f, 1]$,

$$R_n^f := \sum k_s : \text{sum over } s \in (f, 1] \text{ and } (k_s, s) \in X_n,$$

and N_n denote the size of the population at time n ,

$$N_n := L_n^f + R_n^f.$$

The pair (L_n^f, R_n^f) is Markovian:

(1-1) If $(L_n^f, R_n^f) = (0, 0)$

$$(L_{n+1}^f, R_{n+1}^f) = \begin{cases} (1, 0) & \text{w. p. } fp \\ (0, 1) & \text{w. p. } (1-f)p \\ (0, 0) & \text{w. p. } 1-p \end{cases} \quad (1)$$

(1-2) If $(L_n^f, R_n^f) \in \{0\} \times \mathbb{N}$

$$(L_{n+1}^f, R_{n+1}^f) = \begin{cases} (1, R_n^f) & \text{w. p. } fpr \\ (0, R_n^f + 1) & \text{w. p. } (1-f)pr + p(1-r) \\ (0, R_n^f - 1) & \text{w. p. } 1-p \end{cases} \quad (2)$$

(1-3) If $(L_n^f, R_n^f) \in \mathbb{N} \times \{0\}$

$$(L_{n+1}^f, R_{n+1}^f) = \begin{cases} (L_n^f + 1, 0) & \text{w. p. } fpr + p(1 - r) \\ (L_n^f, 1) & \text{w. p. } (1 - f)pr \\ (L_n^f - 1, 0) & \text{w. p. } 1 - p \end{cases} \quad (3)$$

(1-4) If $(L_n^f, R_n^f) \in \mathbb{N} \times \mathbb{N}$

$$(L_{n+1}^f, R_{n+1}^f) = \begin{cases} (L_n^f + 1, R_n^f) & \text{w. p. } fpr + p(1 - r) \frac{L_n^f}{N_n} \\ (L_n^f, R_n^f + 1) & \text{w. p. } (1 - f)pr + p(1 - r) \frac{R_n^f}{N_n} \\ (L_n^f - 1, R_n^f) & \text{w. p. } 1 - p. \end{cases} \quad (4)$$

Unlike in other cases, the transition probabilities here are not spatially homogeneous.

The model exhibits a phase transition at a critical position f_c defined as

$$f_c := \frac{1-p}{pr}$$

as given in the following theorem:

Theorem

- (1) In case $p \leq 1 - p$, the population dies out infinitely often a.s., in the sense that

$$P(N_n = 0 \text{ for infinitely many } n) = 1$$

- (2) In case $1 - p < rp$, the size of the population goes to infinity as $n \rightarrow \infty$, and most of the population is distributed at sites in the interval $[f_c, 1]$, in the sense that

$$P\left(\lim_{n \rightarrow \infty} \frac{R_n^{f_c}}{N_n} = 1\right) = 1 \text{ and } P\left(\liminf_{n \rightarrow \infty} \frac{R_n^{f_c} - R_n^f}{N_n} > 0\right) = 1 \text{ for any } f > f_c.$$

- (3) In case $rp \leq 1 - p < p$, the size of the population goes to infinity as $n \rightarrow \infty$, and most of the population is concentrated at sites near 1, in the sense that

$$P\left(\lim_{n \rightarrow \infty} N_n = \infty\right) = 1 \text{ and, for any } \epsilon > 0, P\left(\lim_{n \rightarrow \infty} \frac{R_n^{1-\epsilon}}{N_n} = 1\right) = 1.$$

Let $F_n(f)$ denote the empirical distribution of sites at time n , i.e.

$$F_n(f) := \frac{\#\{s \in [0, f] : (k, s) \in X_n \text{ for some } k \geq 1\}}{\#\{s \in [0, 1] : (k, s) \in X_n \text{ for some } k \geq 1\}},$$

we have a Glivenko-Cantelli type result

Corollary:

If $1 - p < rp$ (i.e., $f_c < 1$), then

$$F_n(f) \rightarrow \frac{\max\{f - f_c, 0\}}{1 - f_c} \quad \text{uniformly a.s.}$$

For a given n, k, f let $U_n^k(f) := \#\{s \in [f, 1] : (k, s) \in X_n\}$ denote the number of sites in $[f, 1]$ at time n which has a population of size exactly k .

Clearly the total number of sites is $S_n = \sum_k U_n^k(0)$.

Taking $U_n^k(f+) = \lim_{s \downarrow f} U_n^k(s)$, for $A \subseteq \mathbb{X}$, define the empirical distribution of size and fitness on $\mathbb{N} \times [0, 1]$ by

$$H_n(A) := \begin{cases} \frac{\sum_{(k,f) \in A} U_n^k(f) - U_n^k(f+)}{S_n} & \text{if } S_n > 0, \\ \delta_{(0,0)}(A) & \text{if } S_n = 0. \end{cases}$$

Theorem

For $pr > 1 - p$, as $n \rightarrow \infty$, H_n converges weakly to a product measure on $\mathbb{N} \times [0, 1]$ whose density is given by

$$p_k \frac{1_{[f_c, 1]}(x)}{1 - f_c} dx, \quad \text{for } (k, x) \in \mathbb{N} \times [0, 1]$$
$$\text{with } p_k = \frac{2p - 1}{p(1 - r)} B\left(1 + \frac{2p - 1}{p(1 - r)}, k\right) \text{ for } k \in \mathbb{N},$$

where $B(a, b) := \int_0^1 t^{a-1} (1 - t)^{b-1} dt$ is the Beta function with parameters $a, b > 0$.

For k large and s fixed, $B(s, k) \sim \Gamma(s)k^{-s}$ as $k \rightarrow \infty$.

So the probability density

$$p_k \sim k^{-\left(1 + \frac{2p-1}{p(1+r)}\right)}, \text{ i.e. } \sum_{k \geq j} p_k \sim j^{-\frac{2p-1}{p(1+r)}}.$$

Since $p(1+r) > p + (1-p) = 2p - 1$, we have a power law behaviour as is to be expected from preferential attachment models.

Number of individuals of a fixed fitness

Fix $f \in [0, 1]$ and let N_n^f denote the number of individuals with fitness f at time n . When $rp > 1 - p$, i.e. $f_c < 1$, from the first Theorem we know that, $P(L_n^f = 0 \text{ infinitely often}) = 0$ for $f \in (f_c, 1)$. Thus, if a mutant with fitness $f \in (f_c, 1)$ is born at some large time ℓ , then the chances of the mutant dying is small, and so a natural question is ‘for some $n > \ell$, how many individuals did this mutant attract by time n ’, i.e., what is the value of N_n^f ?

Theorem

Fix $f \in (f_c, 1)$, we have, for $\ell < n$, as $\ell, n \rightarrow \infty$

$$\begin{aligned} & E[N_n^f | \text{a mutant with fitness } f \text{ is born at time } \ell] \\ & \sim \frac{\Gamma((2p - 1)\ell + 1)\Gamma((2p - 1)n + 1 + p(1 - r))}{\Gamma((2p - 1)\ell + 1 + p(1 - r))\Gamma((2p - 1)n + 1)} \\ & \sim \left(\frac{n}{\ell}\right)^{p(1-r)}. \end{aligned}$$

A crucial result used in the proofs is

Lemma

(1) Let $f_c = \frac{1-p}{rp} < 1$.

(i) For $f < f_c$ and for any $\eta \in (0, 1)$ we have

$$P\left(\exists T > 0 \text{ such that } \rho_n^f \equiv \frac{L_n^f}{N_n} \leq \eta \forall n \geq T\right) = 1, \quad (5)$$

$$\text{and } P(L_n^f = 0 \text{ infinitely often}) = 1. \quad (6)$$

(ii) Let $f > f_c$. Then

$$P(L_n^f = 0 \text{ infinitely often}) = 0. \quad (7)$$

(2) Let $1 \leq f_c = \frac{1-p}{rp} < \frac{1}{r}$.

(i) For $f < 1$ and for any $\eta \in (0, 1)$ we have (9) and (10).

(ii) Let $f = 1$. Then we have (11).

First we prove the following theorem:

Theorem

For $pr > 1 - p$, as $n \rightarrow \infty$, H_n converges weakly to a product measure on $\mathbb{N} \times [0, 1]$ whose density is given by

$$p_k \frac{1_{[f_c, 1]}(x)}{1 - f_c} dx, \quad \text{for } (k, x) \in \mathbb{N} \times [0, 1]$$

$$\text{with } p_k = \frac{2p - 1}{p(1 - r)} B\left(1 + \frac{2p - 1}{p(1 - r)}, k\right) \text{ for } k \in \mathbb{N},$$

where $B(a, b) := \int_0^1 t^{a-1} (1 - t)^{b-1} dt$ is the Beta function with parameters $a, b > 0$.

To prove the above theorem, for $k, t_1, n \in \mathbb{N}$ let

$A_k(t_1, n)$, be the event that a mutant born at time t_1 gets $k - 1$ attachments until time n , and let

$$q_k(t_1, n) := P(A_k(t_1, n)).$$

First we show that

Lemma

For the preferential attachment model with $p = 1$, i.e., no deaths, we have

$$E \left[\left\{ \frac{1}{n} \sum_{t_1=1}^n (1_{A_k(t_1, n)} - q_k(t_1, n)) \right\}^2 \right] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

The left hand side above is

$$\begin{aligned} & \frac{1}{n^2} \sum_{t_1=1}^n \sum_{s_1=1}^n [\mathbb{P}(A_k(s_1, n) \cap A_k(t_1, n)) - \mathbb{P}(A_k(s_1, n))\mathbb{P}(A_k(t_1, n))] \\ &= \frac{1}{n^2} \sum_{t_1=1}^n \sum_{s_1=1}^n \mathbb{P}(A_k(s_1, n)) [\mathbb{P}(A_k(t_1, n)|A_k(s_1, n)) - \mathbb{P}(A_k(t_1, n))] . \end{aligned}$$

We show that for any $x_1, y_1 \in (0, 1)$ with $x_1 < y_1$

$$\mathbb{P}(A_k(y_1 n, n)|A_k(x_1 n, n)) - \mathbb{P}(A_k(y_1 n, n)) \rightarrow 0, \quad n \rightarrow \infty,$$

which suffices to prove the lemma.

This is done by breaking up $A_k(t_1, n)$ according to the arrival times of the $k - 1$ mutants. □

Lemma Let $p = 1$. For each $k \in \mathbb{N}$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t_1=1}^n q_k(t_1, n) = \frac{1}{1-r} B\left(\frac{2-r}{1-r}, k\right) = p_k.$$

To prove the above lemma observe that,
for $k = 1$, we have

$$q_1(t_1, n) = r \prod_{j=t_1+1}^n \left(1 - \frac{1-r}{j}\right),$$

since the number of individuals till time $j - 1$ is j and the probability that the mutant who arrived at time t_1 gets an attachment at time j is $\frac{1-r}{j}$.

For $k = 2$

$q_2(t_1, n)$

$$= r \sum_{t_2=t_1+1}^n \left\{ \prod_{j=t_1+1}^{t_2-1} \left(1 - \frac{1-r}{j} \right) \right\} \frac{1-r}{t_2} \left\{ \prod_{j=t_2+1}^n \left(1 - \frac{2(1-r)}{j} \right) \right\},$$

where t_2 is the time of the second attachment.

Similarly for each $k \in \mathbb{N}$

$$q_k(t_1, n) = r \sum_{t_1 < t_2 < \dots < t_k \leq n} \prod_{\ell=1}^k \prod_{j=t_\ell+1}^{t_{\ell+1}} \left(1 - \frac{\ell(1-r)}{j}\right) \prod_{\ell=1}^{k-1} \frac{\ell(1-r)}{t_{\ell+1} - \ell(1-r)}.$$

By using Stirling's formula we see that

$$\prod_{j=t_\ell+1}^{t_{\ell+1}} \left(1 - \frac{\ell(1-r)}{j}\right) \sim \left(\frac{t_\ell}{t_{\ell+1}}\right)^{\ell(1-r)}, \quad t_\ell, t_{\ell+1} \rightarrow \infty.$$

Now letting $n \rightarrow \infty$ and taking $t_\ell = nx_\ell$ we have

$$\begin{aligned}
 & \frac{1}{n} \sum_{t_1=1}^n q_k(t_1, n) \\
 & \sim r \int_{0 < x_1 < \dots < x_k < 1} dx_1 \cdots dx_k \prod_{\ell=1}^k \left(\frac{x_\ell}{x_{\ell+1}} \right)^{\ell(1-r)} \prod_{\ell=1}^{k-1} \frac{\ell(1-r)}{x_{\ell+1}} \\
 & = r(1-r)^{k-1} \int_0^1 dx_1 x_1^{1-r} \prod_{\ell=2}^k \int_{x_1}^1 dx_\ell x_\ell^{-r} = r \int_0^1 dx_1 x_1^{1-r} (1-x_1^{1-r})^{k-1} \\
 & = \frac{r}{1-r} \int_0^1 dy y^{\frac{1}{1-r}} (1-y)^{k-1} = \frac{r}{1-r} B\left(\frac{2-r}{1-r}, k\right).
 \end{aligned}$$

When $p = 1$ we have

$$\frac{1}{n} \sum_{t_1=1}^n 1_{A_k(t_1, n)} \rightarrow \frac{r}{1-r} B\left(\frac{2-r}{1-r}, k\right) \text{ as } n \rightarrow \infty, \text{ in probability.}$$

Noting that

$$\lim_{n \rightarrow \infty} \frac{S_n}{n} = r, \quad \text{a.s.} \quad (\text{Recall } S_n \text{ is the number of sites till time } n)$$

we have

$$\lim_{n \rightarrow \infty} \frac{\sum_{f \in (0,1)} U_n^k(f) - U_n^k(f+)}{S_n} = \frac{1}{1-r} B\left(\frac{2-r}{1-r}, k\right) = p_k \text{ in probability.} \quad (8)$$

Noting that the sites are uniformly distributed on $[0, 1]$ independently, and preferential attachment does not depend on the position of sites, we obtain the Theorem for $p = 1$.

Next we consider the case where $p \in (0, 1)$. We introduce another Markov process \hat{X}_n , $n \in \mathbb{N} \cup \{0\}$, which is a pure birth process, as follows:

1. At time 0 there exists one individual at a site uniformly distributed on $(f_c, 1)$.
2. with probability $\hat{r} := \text{pr}(1 - f_c)$ a mutant is born with a fitness uniformly distributed in $[f_c, 1]$,
3. with probability $p(1 - r)(1 - f_c)$ the individual born has a fitness f with a probability proportional to the number of individuals of fitness f and we increase the corresponding population of fitness f individuals by 1.
3. With probability $1 - p(1 - f_c)$ nothing happens, i.e. neither a birth nor a death occurs.

For the Markov process \hat{X}_n , $n \in \mathbb{N} \cup \{0\}$, we define \hat{q}_k , \hat{S}_n and \hat{U}_n in the same manner as q_k , S_n and U_n for X_n , $n \in \mathbb{N} \cup \{0\}$. Then by the same argument as above we see that

$$\frac{1}{n} \sum_{t_1=1}^n \tilde{q}_k(t_1, n) \sim p(1 - rf_c) \frac{\hat{r}}{1 - \hat{r}} B\left(\frac{2 - \hat{r}}{1 - \hat{r}}, k\right)$$

and

$$\lim_{n \rightarrow \infty} \frac{\hat{S}_n}{n} = p(1 - f_c).$$

Hence

$$\lim_{n \rightarrow \infty} \frac{\sum_{f \in (0,1)} \hat{U}_n^k(f) - \hat{U}_n^k(f+)}{\hat{S}_n} = \frac{1}{1 - \hat{r}} B\left(\frac{2 - \hat{r}}{1 - \hat{r}}, k\right) = p_k,$$

From the main Lemma, we know that deletions of individuals in $(f_c, 1)$ occur finitely often and $\frac{R_n^f}{L_n^f + R_n^f} \rightarrow 1$ almost surely as $n \rightarrow \infty$. Thus we have

$$\begin{aligned}
 & \lim_{n \rightarrow \infty} \frac{\sum_{f \in (0,1)} U_n^k(f) - U_n^k(f+)}{S_n} \\
 &= \lim_{n \rightarrow \infty} \frac{\sum_{f \in (0,1)} \hat{U}_n^k(f) - \hat{U}_n^k(f+)}{\hat{S}_n} \quad \text{a.s.} \\
 &= \frac{1}{1 - \hat{r}} B\left(\frac{2 - \hat{r}}{1 - \hat{r}}, k\right) \\
 &= \frac{2p - 1}{p(1 - r)} B\left(1 + \frac{2p - 1}{p(1 - r)}, k\right)
 \end{aligned}$$

and so (8) for $p \in (0, 1]$. Noting that the sites are uniformly distributed on $[0, 1]$ independently, and preferential attachment does not depend on the position of sites, we obtain the theorem. \square

Now we prove

Theorem

Fix $f \in (f_c, 1)$, we have, for $\ell < n$, as $\ell, n \rightarrow \infty$

$$\begin{aligned} & E[N_n^f | \text{a mutant with fitness } f \text{ is born at time } \ell] \\ & \sim \frac{\Gamma((2p-1)\ell+1)\Gamma((2p-1)n+1+p(1-r))}{\Gamma((2p-1)\ell+1+p(1-r))\Gamma((2p-1)n+1)} \\ & \sim \left(\frac{n}{\ell}\right)^{p(1-r)}, \end{aligned}$$

where N_n^f denotes the number of individuals with fitness f at time n .

Proof. Since we are interested in the region $f > f_c$ and also, for the calculation of the expectation, we just need to factor out the death rate $(1 - p)$, so we modify the Markov process \hat{X}_n introduced earlier, by removing the times when ‘nothing happens’, i.e. the process does not move. This is done as follows: let \hat{N}_n be the number of individuals of the process \hat{X}_n at time n , we define a new Markov process \check{X}_n , for $n \geq 0$, by

$$\hat{X}_n = \check{X}_{\hat{N}_n - 1}.$$

Since $\hat{N}_0 = 1$, we see that $\check{N}_\ell = \ell + 1$, where \check{N}_ℓ is the number of individuals of the process \check{X} at time ℓ .

Let \check{N}_m^f be the number of individuals of \check{X} of fitness f at time m .

We have

$$\begin{aligned} E[\check{N}_m^f | \check{N}_{m-1}^f] &= \{1 - p(1 - r)\} \check{N}_{m-1}^f \\ &\quad + p(1 - r) \left\{ (\check{N}_{m-1}^f + 1) \frac{\check{N}_{m-1}^f}{m} + \check{N}_{m-1}^f \left(1 - \frac{\check{N}_{m-1}^f}{m} \right) \right\} \\ &= \left(1 + \frac{p(1 - r)}{m} \right) \check{N}_{m-1}^f. \end{aligned}$$

If $\check{N}_0^f = \check{N}_0 = 1$ then we have

$$E[\check{N}_m^f | \check{N}_0^f = 1] = \prod_{k=1}^m \left(\frac{k + p(1-r)}{k} \right) = \frac{\Gamma(m+1+p(1-r))}{\Gamma(1+p(1-r))\Gamma(m+1)},$$

while, if $\check{N}_\ell^f = 1$ then we have

$$E[\check{N}_m^f | \check{N}_\ell^f = 1] = \prod_{k=\ell+1}^m \left(\frac{k + p(1-r)}{k} \right) = \frac{\Gamma(\ell+1)\Gamma(m+1+p(1-r))}{\Gamma(\ell+1+p(1-r))\Gamma(m+1)}.$$

Since $\frac{\hat{N}_n}{n} \rightarrow \text{pr}(1 - f_c) + p(1 - r) = 2p - 1$, if $\hat{N}_0^f = 1$ then we have

$$\begin{aligned} & \mathbb{E}[\hat{N}_n^f | \hat{N}_0^f = 1] \\ & \sim \prod_{k=1}^{(2p-1)n} \left(\frac{k + p(1-r)}{k} \right) = \frac{\Gamma((2p-1)n + 1 + p(1-r))}{\Gamma(1 + p(1-r))\Gamma((2p-1)n + 1)}. \end{aligned}$$

$$\begin{aligned} \mathbb{E}[\hat{N}_n^f | \hat{N}_\ell^f = 1] &= \prod_{k=(2p-1)\ell+1}^{(2p-1)n} \left(\frac{k + p(1-r)}{k} \right) \\ &= \frac{\Gamma((2p-1)\ell + 1)\Gamma((2p-1)n + 1 + p(1-r))}{\Gamma((2p-1)\ell + 1 + p(1-r))\Gamma((2p-1)n + 1)}. \end{aligned}$$

From the main Lemma we have $\mathbb{E}[N_n^f | N_\ell^f = 1] \sim \mathbb{E}[\hat{N}_n^f | \hat{N}_\ell^f = 1]$, we have the theorem. □

Now we prove the main lemma

Lemma

(1) Let $f_c = \frac{1-p}{rp} < 1$.

(i) For $f < f_c$ and for any $\eta \in (0, 1)$ we have

$$P\left(\exists T > 0 \text{ such that } \rho_n^f \equiv \frac{L_n^f}{N_n} \leq \eta \forall n \geq T\right) = 1, \quad (9)$$

$$\text{and } P(L_n^f = 0 \text{ infinitely often}) = 1. \quad (10)$$

(ii) Let $f > f_c$. Then

$$P(L_n^f = 0 \text{ infinitely often}) = 0. \quad (11)$$

(2) Let $1 \leq f_c = \frac{1-p}{rp} < \frac{1}{r}$.

(i) For $f < 1$ and for any $\eta \in (0, 1)$ we have (9) and (10).

(ii) Let $f = 1$. Then we have (11).

The idea of the proof is that, since for $f < f_c \wedge 1$, R_n^f will be much larger than L_n^f , we stochastically bound the non-spatially homogeneous Markov chain by a spatially homogeneous Markov chain, and study the modified Markov chain. As such, for $\varepsilon \in [0, 1]$, we introduce a Markov chain $(L_n^f(\varepsilon), R_n^f(\varepsilon))$ as follows:

If $(L_n^f(\varepsilon), R_n^f(\varepsilon)) \in \mathbb{N} \times \mathbb{N}$

$$\begin{aligned}
 & (L_{n+1}^f(\varepsilon), R_{n+1}^f(\varepsilon)) \\
 = & \begin{cases} (L_n^f(\varepsilon) + 1, R_n^f(\varepsilon)) & \text{w. p. } fpr + p(1-r)\varepsilon \\ (L_n^f(\varepsilon), R_n^f(\varepsilon) + 1) & \text{w. p. } (1-f)pr + p(1-r)(1-\varepsilon) \\ (L_n^f(\varepsilon) - 1, R_n^f(\varepsilon)) & \text{w. p. } 1-p. \end{cases}
 \end{aligned}$$

Similarly for other cases.

Recall the original transition probabilities were

$$(L_{n+1}^f, R_{n+1}^f) = \begin{cases} (L_n^f + 1, R_n^f) & \text{w. p. } fpr + p(1-r) \frac{L_n^f}{N_n} \\ (L_n^f, R_n^f + 1) & \text{w. p. } (1-f)pr + p(1-r) \frac{R_n^f}{N_n} \\ (L_n^f - 1, R_n^f) & \text{w. p. } 1-p. \end{cases}$$

For $\varepsilon \in [0, 1]$, we couple the processes $\{(L_n^f(\varepsilon), R_n^f(\varepsilon)) : n \geq 1\}$ such that

$$L_n^f(\varepsilon) \leq L_n^f(\varepsilon'), \quad R_n^f(\varepsilon) \geq R_n^f(\varepsilon') \quad \text{for } \varepsilon \leq \varepsilon' \text{ and all } n \geq 1.$$

For $\rho_n^f := \frac{L_n^f}{N_n}$,

$$\begin{aligned} L_{n+1}^f &= L_{n+1}^f(\rho_n^f), & R_{n+1}^f &= R_{n+1}^f(\rho_n^f), \\ L_n^f(0) &\leq L_n^f \leq L_n^f(1), & R_n^f(1) &\leq R_n^f \leq R_n^f(0). \end{aligned}$$

Also, death rates don't change, so

$$N_n(\varepsilon) := L_n^f(\varepsilon) + R_n^f(\varepsilon) = N_n.$$

By the law of large numbers we have

$$\lim_{n \rightarrow \infty} \frac{L_n^f(\varepsilon)}{n} = [\text{fpr} + p(1-r)\varepsilon - 1 + p]_+, \text{ almost surely.}$$

$$\lim_{n \rightarrow \infty} \frac{N_n}{n} = 2p - 1, \text{ almost surely,}$$

and so, for $\rho_n^f(\varepsilon) := \frac{L_n^f(\varepsilon)}{N_n}$, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \rho_n^f(\varepsilon) &= \left[\frac{\text{fpr} + p(1-r)\varepsilon - 1 + p}{2p - 1} \right]_+ \\ &= \left[\frac{\text{fpr} - 1 + p}{2p - 1} + \frac{p(1-r)\varepsilon}{2p - 1} \right]_+. \end{aligned}$$

To study $\frac{L_n^f(\epsilon)}{N_n}$ we introduce the linear function defined by

$$h(x) = \frac{fpr - 1 + p}{2p - 1} + \frac{p(1 - r)}{2p - 1}x.$$

Note that $\frac{p(1-r)}{2p-1} > 0$. By a simple calculation we see that if $f \leq 1$

$$h(0) \leq \frac{pr - 1 + p}{2p - 1} < 0 \quad \text{if } pr < 1 - p \quad \text{and} \quad h(1) \leq 1 - \frac{pr(1 - f)}{2p - 1} < 1.$$

Then we see that

$$\lim_{n \rightarrow \infty} h^n(1) < 0,$$

where $h^2(x) = h(h(x))$ and $h^{n+1}(x) = h(h^n(x))$. Hence, for any $\eta > 0$

$$P \left(\exists T > 0 \text{ s.t. } \alpha_n^f \equiv \frac{L_n^f}{N_n} \leq \eta \quad \forall n \geq T \right) = 1,$$

which will give the 2nd part of the theorem.

Mean field heuristics

We now present some mean field heuristics about the location of the leftmost site x_t at time t in the case when $pr < 1 - p < p$, i.e. $f_c > 1$. Let $y_t = 1 - x_t$. The number of individuals to enter the interval $(x_t, 1]$ is approximately

$$pr y_t dt + p(1 - r)dt,$$

where the first term counts the births which are mutants and the second term counts the births which are not mutants. While the number of individuals deleted in the interval $(x_t, 1]$ is approximately

$$-\frac{dy_t}{y_t} \{p - (1 - p)\}t,$$

this being the absolute value of the deletions since $\frac{dy_t}{dt} < 0$.

Thus we consider the following differential equation:

$$p r y_t dt + p(1-r)dt + \frac{dy_t}{y_t} \{2p-1\}t = (2p-1)dt,$$

Solving we have, for $\gamma := \frac{pr(f_c-1)}{2p-1} = \frac{1-p-pr}{2p-1}$

$$y_t = \frac{f_c - 1}{Ct^\gamma - r} \\ \sim C't^{-\gamma}, \quad t \rightarrow \infty.$$

and the number of sites is

$$rpty_t \sim C'rpt^{1-\gamma}.$$

For $f_c > 1$ we have $\gamma = \gamma(p, r) > 0$, and $\gamma(p, r)$ is a decreasing function of p . Also

- (i) when $p = 1 - p$, i.e., $p = \frac{1}{2}$, then $\gamma = \infty$; this corresponds to the case when the process dies out repeatedly,
- (ii) when $pr = 1 - p$, i.e., $f_c = 1$, then $\gamma = 0$; this corresponds to the case when there are only a bounded number of sites surviving,
- (iii) when $p = \frac{2}{3+r} \in \left(\frac{1}{2}, \frac{1}{1+r}\right)$, then $\gamma = 1$; this too corresponds to the case when there are only a bounded number of sites surviving.

From the above, we see that there are three critical values

$$p_c^{(0)} := \frac{1}{2} < p_c^{(1)} := \frac{2}{3+r} < p_c^{(2)} := \frac{1}{r+1} < 1$$

and four phases:

1. For $p \in (p_c^{(2)}, 1)$, $\gamma = -\infty$ and individuals exist in the interval $(f_c, 1]$.
2. For $p \in (p_c^{(1)}, p_c^{(2)}]$, $\gamma \in (0, 1)$ and the number of sites are increasing with the order $t^{1-\gamma}$ and the average number of individuals per site is of order t^γ .
3. For $p \in (p_c^{(0)}, p_c^{(1)}]$, $\gamma \in (1, \infty)$, that is, $1 - \gamma$ is negative, and the number of sites is finite, with the average number of individuals being of order t .
4. For $p \in (0, p_c^{(0)}]$ the process dies out infinitely often.

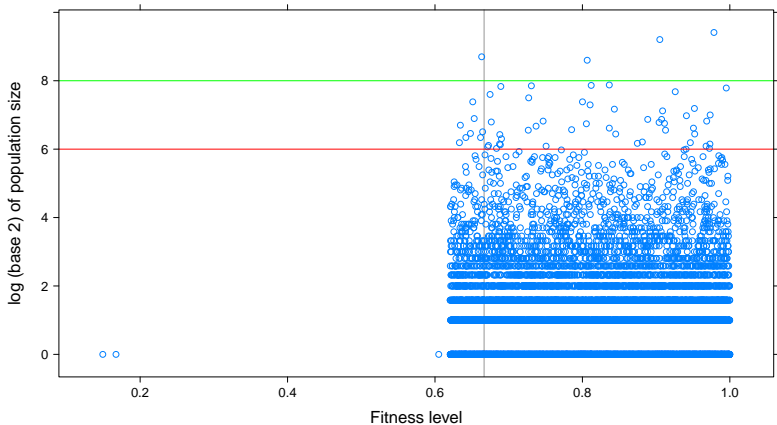


Figure: Population (in \log_2 scale) at various fitness levels.

$p = 3/4$, $r = 1/2$, so that $f_c = 2/3$. The simulation has been conducted with $n = 100,000$.

Source: Deepayan Sarkar

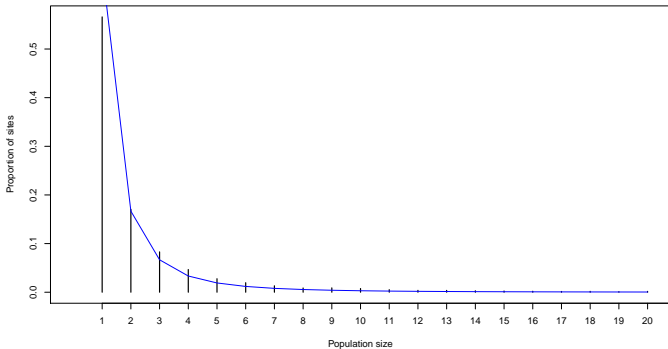


Figure: Population (in \log_2 scale) at various fitness levels.

Source: Deepayan Sarkar

“TALK, v.t. To commit an indiscretion without temptation, from an impulse without purpose.”

Ambrose Bierce (1906) *The Devil's Dictionary*.