

# A stochastic model for phylogenetic trees

by Thomas M. Liggett\* and Rinaldo B. Schinazi†

University of California at Los Angeles,  
and University of Colorado at Colorado Springs

February 5, 2009

## Abstract

We propose the following simple stochastic model for phylogenetic trees. New types are born and die according to a birth and death chain. At each birth we associate a fitness to the new type sampled from a fixed distribution. At each death the type with the smallest fitness is killed. We show that if the birth (i.e. mutation) rate is subcritical we get a phylogenetic tree consistent with an influenza tree (few types at any given time and one dominating type lasting a long time). When the birth rate is supercritical we get a phylogenetic tree consistent with an HIV tree (many types at any given time, none lasting very long).

## 1 Introduction

The influenza phylogenetic tree is peculiar in that it is very skinny: one type dominates for a long time and any other type that arises quickly dies out. Then the dominating type suddenly dies out and is immediately replaced by a new dominating type. The models proposed so far are very complex and make many assumptions. See for instance Koelle et al. (2006) and van

---

\*Partially supported by NSF grant DMS-0301795

†Partially supported by NSF grant DMS-0701396

*Key words and phrases:* phylogenetic tree, influenza, HIV, stochastic model  
*2000 Mathematics Subject Classification:* 60K35

Nimwegen (2006). We would like to use a simple stochastic model for such a tree. The other motivation for this work comes from the comparison between influenza and HIV phylogenetic trees. An HIV tree is characterized by a radial spread outward from an ancestral node, in sharp contrast with an influenza tree. Moreover, Korber et al. (2001) note that the influenza virus is less diverse worldwide than the HIV virus is in Amsterdam alone. However, both types of trees are supposed to be produced by the same basic mechanism: mutations. Can the same mathematical model produce two trees that are so different? Our simple stochastic model will show a striking difference in behavior depending on the mutation rate.

Our model has a birth and death component and a fitness component. For the death and birth component we do the following. If there are  $n \geq 1$  types at a certain time  $t$  then there is birth of a new type (by mutation) at rate  $n\lambda$ . We think of a birth as the appearance of one new type, not the replacement of one type by two new types. If there are  $n \geq 2$  types then there is death of one type at rate  $n$ . If only one type is left it cannot die. That is,

$$\begin{aligned} n &\longrightarrow n + 1 \text{ at rate } n\lambda \\ n &\longrightarrow n - 1 \text{ at rate } n \text{ if } n \geq 2. \end{aligned}$$

Moreover, each new individual is assigned a fitness value chosen from a fixed distribution, independently each time. Every time there is a death event then the type that is killed is the one with the smallest fitness. Since all that matters is the ranks of the fitnesses, we might as well take their distribution to be uniform on  $[0, 1]$ . For simplicity the process is started with a single type.

We give no specific rule on how to attach a new type after a birth to existing types (in order to construct a tree). Our results do not depend on such a rule. Two natural possibilities are to either attach the new type to the type which has the maximum fitness or to a type taken at random.

**Theorem 1.** *Take  $\alpha \in (0, 1)$ .*

If  $\lambda \leq 1$ , then

$$\lim_{t \rightarrow \infty} P(\text{maximal types at times } \alpha t \text{ and } t \text{ are the same}) = \alpha,$$

while if  $\lambda > 1$ , then this limit is 0.

We see that if  $\lambda < 1$ , 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. On the other hand, if  $\lambda > 1$ , then the dominating type at time  $t$  has likely been present for a time of order shorter than  $t$  and at any given time there will be many types. This is consistent with an HIV tree.

## 2 Proof of Theorem 1

The proof divides into three cases, depending on whether the birth and death chain is positive recurrent, null recurrent, or transient. We present them in order of difficulty.

### 2.1 Case $\lambda < 1$

Let  $\tau_1, \tau_2, \dots$  be the (continuous) times between successive visits of the chain to 1,  $T_n = \tau_1 + \dots + \tau_n$ ,  $\sigma_1, \sigma_2, \dots$  be the number of new types introduced in cycles between successive visits to 1, and  $S_n = 1 + \sigma_1 + \dots + \sigma_n$ . Note that the  $\tau$ 's and  $\sigma$ 's are not independent of each other, but the sequence  $(\tau_1, \sigma_1), (\tau_2, \sigma_2), \dots$  is i.i.d. and independent of the fitness sequence. Define the usual renewal process  $N(t)$  corresponding to the  $\tau$ 's by  $\{N(t) = n\} = \{T_n \leq t < T_{n+1}\}$ .

For  $0 < s < t$ , recalling that  $T_{N(s)} \leq s < T_{N(s)+1}$ , and noting that the maximal type is increasing in time, we see that

$$P(\text{maximal types at times } s \text{ and } t \text{ are the same, } N(s) < N(t)) \tag{1}$$

lies between

$$P(\text{maximal types at times } T_{N(s)} \text{ and } T_{N(t)+1} \text{ are the same, } N(s) < N(t))$$

and

$$P(\text{maximal types at times } T_{N(s)+1} \text{ and } T_{N(t)} \text{ are the same, } N(s) < N(t)).$$

Let  $\mathcal{F}$  be the  $\sigma$ -algebra generated by  $(\tau_1, \sigma_1), (\tau_2, \sigma_2), \dots$ . Then for  $k \leq l$ , since the fitness sequence is i.i.d. and independent of  $\mathcal{F}$ ,

$$P(\text{maximal types at times } T_k \text{ and } T_l \text{ are the same} \mid \mathcal{F}) = \frac{S_k}{S_l}.$$

More precisely, conditional on  $\mathcal{F}$  there are  $S_l$  fitnesses observed by time  $T_l$  and  $S_k$  of them are observed by time  $T_k$ . We claim that in  $n$  i.i.d. observations the probability that the largest occurs among the first  $m$  is  $m/n$ , since any one of  $n$  is equally likely to be the largest. Since  $N(s)$  and  $N(t)$  are  $\mathcal{F}$  measurable, it follows that (1) lies between

$$E \left[ \frac{S_{N(s)}}{S_{N(t)+1}}, N(s) < N(t) \right] \quad \text{and} \quad E \left[ \frac{S_{N(s)+1}}{S_{N(t)}}, N(s) < N(t) \right]. \quad (2)$$

Since  $\lambda < 1$ ,  $E\tau < \infty$ , and the renewal theorem gives

$$N(s)/s \rightarrow 1/E\tau \text{ a.s.},$$

while the strong law of large numbers gives  $S_{N(s)}/N(s) \rightarrow E\sigma$  a.s., so that  $S_{N(s)}/s \rightarrow E\sigma/E\tau$  a.s. It follows by the bounded convergence theorem that

$$\lim_{t \rightarrow \infty} P(\text{maximal types at times } \alpha t \text{ and } t \text{ are the same}) = \alpha. \quad (3)$$

This completes the proof of Theorem 1 in the subcritical case.

## 2.2 Case $\lambda > 1$ .

Define the  $\tau$ 's and  $\sigma$ 's as above, except that now, the cycles used are between the successive times the chain reaches a new high. In other words,  $T_n$  is the hitting time of  $n+1$ ,  $\sigma_n$  is the number of new types born during a first passage cycle from  $n$  to  $n+1$  and  $S_n$  is the number of new types

seen up to time  $T_n$ . Of course, the  $\sigma$ 's and  $\tau$ 's are no longer identically distributed. However,  $(\tau_1, \sigma_1), (\tau_2, \sigma_2), \dots$  are independent. The key to the proof is the following Lemma.

**Lemma 2.** *Assume that  $\lambda > 1$ . Then  $e^{-(\lambda-1)t}N(t)$  is almost surely bounded.*

*Proof of Lemma 2.* Our first step in this proof is to estimate the first two moments of  $\tau_n$ . Following Keilson (1979) (see (5.1.2)) we note that  $\tau_n$  has the same distribution as

$$\frac{X}{(1+\lambda)n} + Y(\tau_{n-1} + \tau'_n) \text{ for } n \geq 2, \quad (4)$$

where  $X$  has a mean 1 exponential distribution,  $\tau'_n$  has the same distribution as  $\tau_n$ ,  $Y$  is a Bernoulli random with  $P(Y = 1) = \frac{1}{\lambda+1}$ , and  $X, Y, \tau_{n-1}$  and  $\tau'_n$  are independent. Letting  $\mu_n = E\tau_n$ , it follows from (4) that

$$\lambda\mu_n = \frac{1}{n} + \mu_{n-1} \text{ for } n \geq 2 \text{ and } \mu_1 = \frac{1}{\lambda}. \quad (5)$$

We will use the following recursion formula, which is easy to prove by induction.

**Lemma 3.** *Let  $a_n$  and  $b_n$  be two sequences of real numbers such that  $a_1 = \lambda^{-1}b_1$  and for  $n \geq 2$ ,  $\lambda a_n = b_n + a_{n-1}$ . Then,*

$$a_n = \sum_{j=1}^n \lambda^{-j} b_{n+1-j}.$$

Applying Lemma 3 to (5) we get

$$\mu_n = \sum_{j=1}^n \frac{\lambda^{-j}}{n+1-j}.$$

Writing  $1/(\lambda-1)$  as a geometric series we have

$$\mu_n - \frac{1}{n(\lambda-1)} = \sum_{j=1}^n \lambda^{-j} \left( \frac{1}{n+1-j} - \frac{1}{n} \right) - \frac{1}{n} \sum_{j=n+1}^{\infty} \lambda^{-j}.$$

Changing the order of summation gives

$$\sum_{n=1}^{\infty} \sum_{j=1}^n \lambda^{-j} \left( \frac{1}{n+1-j} - \frac{1}{n} \right) = \sum_{j=1}^{\infty} \lambda^{-j} \sum_{n=j}^{\infty} \left( \frac{1}{n+1-j} - \frac{1}{n} \right).$$

Note that for  $j \geq 2$

$$\sum_{n=j}^{\infty} \left( \frac{1}{n+1-j} - \frac{1}{n} \right) = \sum_{k=1}^{j-1} \frac{1}{k}$$

and this term is 0 for  $j = 1$ . Hence,

$$\sum_{n=1}^{\infty} \frac{1}{n} \sum_{j=n+1}^{\infty} \lambda^{-j} = \sum_{j=2}^{\infty} \lambda^{-j} \sum_{k=1}^{j-1} \frac{1}{k}.$$

We conclude that

$$\sum_{n=1}^{\infty} \left| \mu_n - \frac{1}{n(\lambda-1)} \right| < \infty$$

and

$$\sum_{n=1}^{\infty} \left( \mu_n - \frac{1}{n(\lambda-1)} \right) = 0.$$

Therefore,

$$E(T_n) - \frac{1}{\lambda-1} \sum_{k=1}^n \frac{1}{k} \text{ converges to } 0. \quad (6)$$

We also need an almost sure result for  $T_n$ , and for this, we will estimate the second moment of  $\tau_n$ . Let  $v_n = \text{Var}(\tau_n)$ . It is easy to check that if  $Y$  is a Bernoulli random variable and is independent of a random variable  $Z$  then

$$\text{Var}(ZY) = E(Y)\text{Var}(Z) + \text{Var}(Y)(EZ)^2.$$

Using this remark and (4) we have for  $n \geq 2$

$$v_n = \frac{1}{(1+\lambda)^2 n^2} + \frac{1}{1+\lambda} (v_n + v_{n-1}) + \frac{\lambda}{(1+\lambda)^2} (\mu_n + \mu_{n-1})^2$$

Therefore, for  $n \geq 2$

$$\lambda v_n = b_n + v_{n-1} \tag{7}$$

where

$$b_n = \frac{1}{(1+\lambda)n^2} + \frac{\lambda}{1+\lambda}(\mu_n + \mu_{n-1})^2.$$

Set  $\mu_0 = 0$ , then  $\lambda v_1 = b_1$ . Hence, Lemma 3 applies to (7), giving

$$v_n = \sum_{j=1}^n \lambda^{-j} b_{n+1-j}. \tag{8}$$

Since  $\mu_n \sim \frac{1}{n(\lambda-1)}$  (that is, the ratio converges to 1),  $b_n \sim \frac{C}{n^2}$  where  $C$  depends on  $\lambda$  only. From (8) we get

$$v_n \sim C' b_n \sim \frac{CC'}{n^2},$$

where  $C'$  depends on  $\lambda$  only. This implies the a.s. convergence of the random series  $\sum(\tau_n - E\tau_n)$  (see for instance Corollary 47.3 in Port (1994)). Therefore the partial sums converge a.s. and

$$T_n - E(T_n) \text{ converges a.s.}$$

Using (6) we get that  $T_n - \frac{1}{\lambda-1} \log n$  converges a.s. and is therefore a.s. bounded. Now use the fact that  $\{N(t) \geq n\} = \{T_n \leq t\}$  to conclude that  $N(t) \exp(-(\lambda-1)t)$  is almost surely bounded. This concludes the proof of Lemma 2.  $\square$

We are now ready to complete the proof of Theorem 1 in the supercritical case. Let  $(Z_i)_{i \geq 1}$  be a discrete time random walk starting at 0 that goes to the right with probability  $\lambda/(\lambda+1)$  and to the left with probability  $1/(\lambda+1)$ . For every  $n \geq 1$ , let  $Z_{i,n}$  be a discrete time random walk starting at 0 with the same rules of evolution as  $Z_i$  except that the random walk  $Z_{i,n}$  has a reflecting barrier at  $-n+1$ . For every  $n \geq 1$ , the two random walks  $Z_i$  and  $Z_{i,n}$  are coupled so that they move together until (if ever) they hit  $-n+1$  and thereafter we still couple them so that  $Z_i \leq Z_{i,n}$  for every  $i \geq 0$ . Let  $U$  and  $U_n$  be the hitting times of 1 for the random walks  $Z_i$  and  $Z_{i,n}$ , respectively.

First note that a new type appears every time there is a birth. Therefore,  $\sigma_n$  is the number of steps to the right of the random walk  $Z_{i,n}$  stopped at 1. That is,  $\sigma_n$  is  $(1 + U_n)/2$ . We now show that  $U_n$  converges a.s to  $U$ . Let  $\delta > 0$  we have

$$P(|U_n - U| > \delta) \leq P(U > U_n) \leq P(Z_i = -n + 1 \text{ for some } i \geq 1).$$

The last probability decays exponentially with  $n$ . Therefore,

$$\sum_{n \geq 1} P(|U_n - U| > \delta) < \infty.$$

An easy application of Borel-Cantelli Lemma implies that  $U_n$  converges a.s. to  $U$ . Since  $U_n \leq U$  the Dominated Convergence Theorem implies that, for every  $k \geq 1$  the  $k$ th moment of  $\sigma_n$  converges to the  $k$ th moment of  $(1 + U)/2$ . In particular,  $Var(\sigma_n)$  is a bounded sequence. This is enough to prove that

$$\frac{1}{n} \sum_{i=1}^n (\sigma_i - E(\sigma_i)) \text{ converges a.s. to } 0;$$

see for instance Proposition 47.10 in Port (1994). Since  $E(\sigma_n)$  is a convergent sequence we get that  $S_n/n$  converges a.s. to the limit of  $E(\sigma_n)$ .

Since  $N(t) \rightarrow \infty$  a.s., this strong law of large numbers gives that  $S_{N(t)}/N(t)$  converges to the limiting expectation of  $\sigma_n$ . This together with Lemma 2 shows that the two terms in (2) converge to 0 when we let  $s = \alpha t$  and  $t$  goes to infinity. The proof of Theorem 1 in the supercritical case is complete.

### 2.3 Case $\lambda = 1$ .

In this subsection we go back to the notation of section 2.1 where  $T_n$  be the time of the  $n$ th visit of the chain to 1.

**Lemma 4.** *Let  $\lambda = 1$ . Then,*

$$\frac{T_n}{n \log n} \rightarrow 1 \text{ in probability.}$$



*Proof of Lemma 4.* When the chain hits 1, it waits a mean 1 exponential time and then jumps to 2. Hence,

$$T_n = \sum_{i=1}^n X_i + \sum_{i=1}^n H_i$$

where the  $X_i$  are independent mean 1 exponential times and  $H_i$  are the hitting times of 1 starting at 2. The  $H_i$  are i.i.d. with distribution function  $F$ . From the backward Kolmogorov equation

$$\int_0^{F(t)} \frac{ds}{1+s^2-2s} = t,$$

we get

$$F(t) = \frac{t}{1+t}.$$

We now use a weak law of large numbers, see Theorem 2 in VII.7 in Feller (1971). It is easier to redo the short proof rather than check the hypotheses of the Theorem. The key is the following consequence of Chebyshev's inequality applied to the truncated random variables:

$$P\left(\left|\frac{1}{nm_n} \sum_{i=1}^n H_i - 1\right| > \epsilon\right) \leq \frac{1}{n\epsilon^2 m_n^2} s_n + n(1 - F(\rho_n)) \quad (9)$$

where

$$m_n = \int_0^{\rho_n} tF'(t)dt \text{ and } s_n = \int_0^{\rho_n} t^2 F'(t)dt,$$

see (7.13) in VII.7 in Feller (1971). We will take  $\rho_n = n\sqrt{\log n}$ . A little Calculus shows that

$$m_n \sim \log \rho_n \sim \log n \text{ and } s_n \sim \rho_n.$$

With our choice of  $\rho_n$ ,  $n(1 - F(\rho_n))$  converges to 0 and

$$\frac{1}{nm_n} \sum_{i=1}^n H_i$$

converges to 1 in probability. This completes the proof of Lemma 4.  $\square$

Since the events  $N(t) \geq n$  and  $T_n \leq t$  are the same, it follows that

$$N(t) \frac{\log t}{t} \rightarrow 1 \tag{10}$$

in probability as  $t \uparrow \infty$ .

Now,  $S_n/n^2$  converges in distribution to a one sided stable law of index  $\frac{1}{2}$  (see Theorem (7.7) in Durrett (2004)). By (10), it follows that  $S_{N(t)}/N(t)^2$  also has this distributional limit. In fact,

$$\frac{S_{N(\alpha t)}}{N^2(t)} \text{ converges to } Y_\alpha$$

in the sense of convergence of finite dimensional distributions, where  $Y_\alpha$  is a stable subordinator (increasing stable process) of index  $1/2$ . (Note that independence between the  $\sigma$ 's and  $\tau$ 's is not required here, which is good since they are highly dependent. All that is needed is that the limit in (10) is constant and that both  $S_n$  and  $N(t)$  are monotone.) So, the limit in (3) is

$$\lim_{t \rightarrow \infty} E\left(\frac{S_{N(\alpha t)}}{S_{N(t)}}\right) = E\left(\frac{Y_\alpha}{Y_1}\right) = \alpha.$$

To check the final equality, it is enough by monotonicity to verify it for rational  $\alpha$ . If  $\alpha = m/n$ , this boils down to the simple fact that if  $V_1, \dots, V_n$  are i.i.d. and positive, then

$$E \frac{V_i}{V_1 + \dots + V_n} = \frac{1}{n}.$$

**Acknowledgements.** We thank an anonymous referee whose remarks helped improve the presentation of the paper and simplify several of our proofs.

## References

R. Durrett (2004) Probability: Theory and Examples (3rd edition). Duxbury press.

W. Feller (1971) *An Introduction to Probability Theory and its Applications, Volume 2* (second edition). Wiley.

J. Keilson (1979) *Markov Chain Models-Rarity and exponentiality*. Springer-Verlag.

B. Korber, B. Gaschen, K. Yusim et al. (2001) Evolutionary and immunological implications of contemporary HIV-1 variation. *British Medical Bulletin* 58, 19-42.

K. Koelle, S. Cobey, B. Grenfell and M. Pascual (2006) Epochal evolution shapes the phylodynamics of interpandemic influenza A (H3N2) in Humans. *Science* vol. 314, 1898-1903.

E. van Nimwegen (2006). Influenza escapes immunity along neutral networks. *Science* vol. 314, 1884-1886.

S. C. Port (1994). *Theoretical Probability for Applications*. Wiley.