

# On the use of Markov chains and Perron Frobenius Theorem in Population Genetics

Ola Hössjer  
Dept. of Mathematics  
Stockholm University

June 2015

# Motivation and Outline



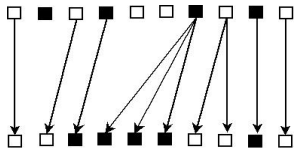
- ▶ Conservation biology: Protect genetic diversity in order to
  - ▶ **Prevent inbreeding**
  - ▶ Keep species viable
  - ▶ Improve environmental adjustment of species
- ▶ Outline:
  - ▶ Populations with substructure
  - ▶ Formulas for long term rate of loss of genetic variants

# Effective population size $N_e$

- ▶ Population of size  $N$
- ▶  $N_e$ : Size of ideal population with same rate of loss of genetic variants as studied population (smaller  $N_e \rightarrow$  faster rate)
- ▶ Short term protection rule:  $N_e \geq 50$
- ▶  $N_e/N$  varies between species

## Genetic drift for size $N$ population

- ▶ Discrete time  $t = 0, 1, 2, \dots$
- ▶ Genes  $g = 1, \dots, 2N$ .
- ▶ Two variants white and black
- ▶  $\nu_{tg}$  nr. of offspring of  $g$ , time  $t$ .
- ▶ Freq. of black, time  $t + 1$ , is



$$\begin{aligned} X_{t+1} &= \frac{1}{2N} |\{g; g \text{ is black at } t+1\}| \\ &= \frac{1}{2N} \sum_{g; g \text{ black at time } t} \nu_{tg}. \end{aligned}$$

$$\begin{aligned} 2N &= 10, \\ X_t &= 0.4, \\ X_{t+1} &= 0.5, \\ \nu_{t2} &= \nu_{t5} = \nu_{t6} = 0, \\ \nu_{t1} &= \nu_{t3} = \nu_{t4} = \nu_{t9} = \nu_{t,10} = 1, \\ \nu_{t8} &= 2, \\ \nu_{t7} &= 3. \end{aligned}$$

- ▶  $\{X_t\}$  Markov chain, state space

$$\begin{aligned} \mathcal{X} &= \left\{0, \frac{1}{2N}, \dots, \frac{2N-1}{2N}, 1\right\} \\ &= \{0\} \cup \{1\} \cup \left\{\frac{1}{2N}, \dots, \frac{2N-1}{2N}\right\} \\ &= \mathcal{X}_0 \cup \mathcal{X}_1 \cup \mathcal{X}_2. \end{aligned}$$

Absorbing states:  $\mathcal{X}_0, \mathcal{X}_1$

Transient states:  $\mathcal{X}_2$

# Ideal population: Wright-Fisher (WF) model<sup>1</sup>

Children pick parental genes independently;

$$\{\nu_{tg}\}_{g=1}^{2N} \sim \text{Mult} \left( 2N; \frac{1}{2N}, \dots, \frac{1}{2N} \right),$$

so that the transition kernel of  $\{X_t\}$  is

$$\begin{aligned} \mathbf{P} &= (P(x, y))_{x, y \in \mathcal{X}}, \\ P(x, y) &= P(X_{t+1} = y | X_t = x) = \binom{2N}{2Ny} x^{2Ny} (1-x)^{2N(1-y)}. \end{aligned}$$

Feller (1951) showed that the eigenvalues of  $\mathbf{P}$  are

$$\lambda_1 = \lambda_2 = 1, \lambda_j = \frac{(2N-1)(2N-2)\cdots(2N-j+2)}{(2N)^{j-2}}, \quad j = 3, \dots, 2N+1,$$

so that the largest non-unit eigenvalue is

$$\lambda = \lambda_3 = 1 - \frac{1}{2N}. \quad (1)$$

It gives the asymptotic rate of fixation of one variant;

$$\lim_{t \rightarrow \infty} \frac{P(X_t \in \mathcal{X}_2)}{\lambda^t} = C, \quad (0 < C < \infty).$$

---

<sup>1</sup>Fisher (1921), Wright (1931).

## Cannings model

Cannings (1974) showed more generally that

$$\lambda_1 = \lambda_2 = 1, \lambda_j = E \left( \prod_{g=1}^{j-1} \nu_{tg} \right), \quad j = 3, \dots, 2N + 1,$$

if  $\{\nu_{tg}\}_{g=1}^{2N}$  are exchangeable, so that in particular,

$$\lambda = \lambda_3 = E(\nu_{t1}\nu_{t2}) = 1 + \text{Cov}(\nu_{t1}, \nu_{t2}) = 1 - p,$$

with coalescence probability

$$\begin{aligned} p &= -\text{Cov}(\nu_{t1}, \nu_{t2}) \\ &= \frac{E[\nu_{t1}(\nu_{t1}-1)]}{2N-1} \\ &= 2N \cdot E \left[ \binom{\nu_{t1}}{2} \right] / \binom{2N}{2} \\ &= P(\text{two offspring have the same parent}). \end{aligned}$$

Eigenvalue effective size<sup>2</sup>  $N_{eE} =$  size of a WF population with fixation rate  $\lambda$ :

$$\lambda = 1 - \frac{1}{2N_{eE}} \implies N_{eE} = \frac{1}{2(1-\lambda)} = \frac{N - \frac{1}{2}}{E[\nu_{t1}(\nu_{t1}-1)]}.$$

---

<sup>2</sup>Crow (1954), Ewens (1982).

## Gene diversities

Introduce (predicted) gene diversities at time  $t = 0, 1, 2, \dots$

$$\begin{aligned}H_t &= 2X_t(1 - X_t), \\ &= P(\text{two genes picked with repl. have diff. variants} | X_t) \\ h_t &= E(H_t) \\ &= P(\text{two genes picked with repl. have diff. variants}).\end{aligned}$$

It can be shown that

$$h_{t+1} = \lambda h_t.$$

Hence gene diversities

$$h_t = \lambda^t h_0 = (1 - p)^t h_0, \quad (2)$$

tend to zero at the same multiplicative rate as non-fixation probabilities  $P(X_t \in \mathcal{X}_2)$ .

But gene diversities and coalescence probabilities ( $p$ ) are easier to analyze theoretically!!

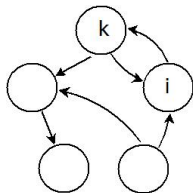
# Structured population

Divide into  $s$  subpopulations  $i = 1, \dots, s$ . Let

$2N_i$  = number of genes in subpop.  $i$ , ( $\sum_i N_i = N$ )

$X_{ti}$  = fraction of black variants in subpop.  $i$  at time  $t$ ,

$\nu_{tkig}$  = nr. of offspring of gene  $g$  of subpop.  $k$  at time  $t$  that end up in subpop.  $i$  at time  $t + 1$ ,  
exchangeable for  $g = 1, \dots, 2N_k$ .



This gives dynamics

$$X_{t+1,i} = \frac{1}{2N_i} \sum_{k=1}^s \sum_{\substack{g: g \in \text{subpop. } k \text{ at time } t, \\ g \text{ is black}}} \nu_{tkig}.$$

Find, under suitable conditions,

$$N_{eE} = \frac{1}{2(1 - \lambda)},$$

where  $\lambda$  is rate of fixation, so that for some  $0 < C < \infty$ ,

$$\lim_{t \rightarrow \infty} \frac{P(\text{non-fixation at time } t)}{\lambda^t} = C.$$



# Structured population, contd.

Let

$$\mathbf{X}_t = (X_{t1}, \dots, X_{ts}).$$

If reproduction is time invariant,  $\{\mathbf{X}_t\}$  is Markov chain with state space

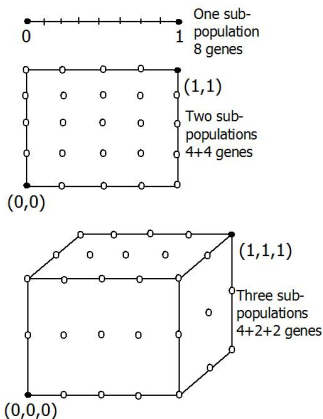
$$\begin{aligned}\mathcal{X} &= \left\{0, \frac{1}{2N_1}, \dots, \frac{2N_1-1}{2N_1}, 1\right\} \times \dots \times \left\{0, \frac{1}{2N_s}, \dots, \frac{2N_s-1}{2N_s}, 1\right\} \\ &= \mathcal{X}_0 \cup \mathcal{X}_1 \cup \mathcal{X}_2.\end{aligned}$$

where

$$\begin{aligned}\mathcal{X}_0 &= \{(0, \dots, 0)\}, \\ \mathcal{X}_1 &= \{(1, \dots, 1)\}, \\ \mathcal{X}_2 &= \mathcal{X} \setminus (\mathcal{X}_0 \cup \mathcal{X}_1),\end{aligned}$$

and

$$\begin{aligned}\mathbf{P} &= (P(\mathbf{x}, \mathbf{y}))_{\mathbf{x}, \mathbf{y} \in \mathcal{X}}, \\ P(\mathbf{x}, \mathbf{y}) &= P(\mathbf{X}_{t+1} = \mathbf{y} | \mathbf{X}_t = \mathbf{x}), \\ \lambda &= \text{3rd largest eigenvalue of } \mathbf{P}.\end{aligned}$$



## Gene diversities for structured population

Introduce (predicted) gene diversities<sup>3</sup> at time  $t = 0, 1, 2, \dots$

$$\begin{aligned}H_{tij} &= X_{ti}(1 - X_{tj}) + (1 - X_{ti})X_{tj}, \\ &= P(\text{two genes picked with repl. from } i \text{ and } j \text{ have diff. variants} | \mathbf{X}_t) \\ h_{tij} &= E(H_{tij}) \\ &= P(\text{two genes picked with repl. from } i \text{ and } j \text{ have diff. variants})\end{aligned}$$

between all pairs of subpopulations  $i$  and  $j$ . The column vector

$$\mathbf{h}_t = \text{vec} \left( (h_{tij})_{i,j=1}^s \right)$$

with  $s^2$  predicted gene diversities satisfies

$$\mathbf{h}_{t+1} = \mathbf{A}\mathbf{h}_t \implies \mathbf{h}_t = \mathbf{A}^t \mathbf{h}_0, \quad (3)$$

where

$$\mathbf{A} = (A_{ij,kl})_{ij,kl \in \{1, \dots, s\} \times \{1, \dots, s\}}$$

is a square matrix of order  $s^2$ , and<sup>4</sup>

$$\lambda = \lambda_{\max}(\mathbf{A}). \quad (4)$$

---

<sup>3</sup>Maruyama (1970), Felsenstein (1972), Nei (1973).

<sup>4</sup>By Perron-Frobenius Theorem applied to  $\mathbf{P}$ .

# Backward migration and coalescence theory to find **A**

If children pick parental subpopulations independently, with

$B_{ik}$  = probability by which genes in  $i$  pick  
parental subpopulations from  $k \in \{1, \dots, s\}$ ,

$p_{ijk}$  = coalescence probability within  $k$   
=  $P(\text{two genes from } i, j \text{ with parents in } k, \text{ have same parent})$   
=  $\frac{N_k}{2N_i N_j B_{ik} B_{jk}} \left( \frac{E(\nu_{tki1}(\nu_{tki1}-1))}{1 - \frac{1}{2N_i}} \right)^{\{i=j\}} E(\nu_{tki1} \nu_{tkj1})^{\{i \neq j\}}$ .

Then (3) holds, with

$$A_{ij,kl} = \left(1 - \frac{1}{2N_i}\right)^{\{i=j\}} \left(\frac{1 - p_{ijk}}{1 - \frac{1}{2N_k}}\right)^{\{k=l\}} B_{ik} B_{jl}. \quad (5)$$

## Motivation of (3) and (5)

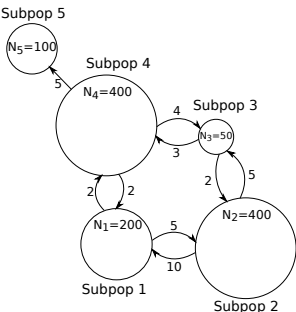
We have that

$$\begin{aligned}h_{t+1,ij} &= P(\text{genes from } i \text{ and } j \text{ at time } t+1 \text{ have different variants}) \\&= P(\text{different genes picked from } i \text{ and } j \text{ at time } t+1) \\&\quad \sum_{k,l} \cdot P(\text{gene from } i \text{ has parent from } k \text{ at time } t) \\&\quad \cdot P(\text{gene from } j \text{ has parent from } l \text{ at time } t) \\&\quad \cdot P(\text{different parents from } k \text{ and } l) \\&\quad \cdot P(\text{different variants of the different parents at time } t) \\&= \left(1 - \frac{1}{2N_i}\right)^{\{i=j\}} \sum_{k,l} B_{ik} B_{jl} (1 - p_{ijk})^{\{k=l\}} \cdot h_{t,kl} / \left(1 - \frac{1}{2N_k}\right)^{\{k=l\}} \\&= \sum_{k,l} A_{ij,kl} h_{t,kl}.\end{aligned}$$

with  $A_{ij,kl}$  as in (5). In vector form this writes

$$\mathbf{h}_{t+1} = \mathbf{A}_t \mathbf{h}_t.$$

# Computation of $N_{eE}$



$$(N_i)_{i=1}^s = (200, 400, 50, 400, 100),$$

$$N = \sum_{i=1}^s N_i = 1150,$$

$$(B_{ik})_{i,k=1}^s = \begin{pmatrix} 0.94 & 0.05 & 0 & 0.01 & 0 \\ 0.0125 & 0.9825 & 0.005 & 0 & 0 \\ 0 & 0.1 & 0.82 & 0.08 & 0 \\ 0.005 & 0 & 0.0075 & 0.9875 & 0 \\ 0 & 0 & 0 & 0.05 & 0.95 \end{pmatrix}$$

$$\begin{aligned} \text{Repr.} &= \{(\nu_{tkig})_{g=1}^{2N_k}\}_{k=1}^s \\ &\sim \text{Mult}\left(2N_i; \frac{B_{i1}}{2N_1}, \dots, \frac{B_{i1}}{2N_1}, \dots, \frac{B_{is}}{2N_s}, \dots, \frac{B_{is}}{2N_s}\right) \\ &\text{independently for } i = 1, \dots, s, \end{aligned}$$

$$p_{ijk} = 1/(2N_k),$$

$$N_{eE} = 970,$$

## Gene diversity effective size $N_{eG}$

Let

$$W_{ij} = P(\text{choose gene pair from } i \text{ and } j),$$

and collect them into row vector of length  $s^2$ :

$$\mathbf{W} = \text{vec}((W_{ij})_{1 \leq i, j \leq s})'.$$

Predicted gene diversity for two randomly sampled genes at time  $t$ , is

$$\begin{aligned} h_t &= P(\text{the two genes have different variants}) \\ &= \sum_{1 \leq i, j \leq s} W_{ij} h_{tij} \\ &= \mathbf{W} \mathbf{h}_t \\ &= \mathbf{W} \mathbf{A}^t \mathbf{h}_0. \end{aligned}$$

It follows from (1) and (2), that for Wright-Fisher model

$$h_t = \left(1 - \frac{1}{2N}\right)^t \cdot h_0. \quad (6)$$

Gene diversity effective size over time interval  $[0, t]$  solves (6), i.e.

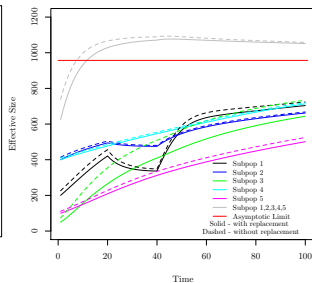
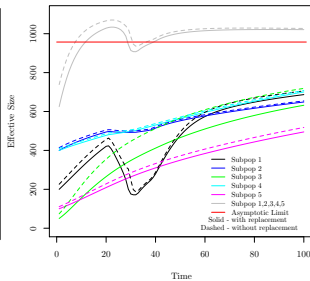
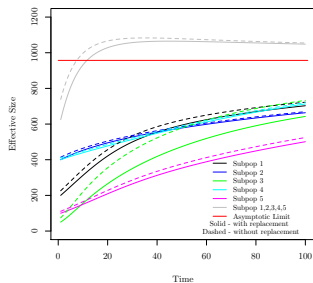
$$N_{eG}([0, t]) = \frac{1}{2 \left[1 - \left(\frac{h_t}{h_0}\right)^{1/t}\right]} \xrightarrow{t \rightarrow \infty} \frac{1}{2(1 - \lambda)} = N_{eE}.$$

# Local/global $N_{eG}$ and $N_{eE}$

Constant subpop. sizes

Local bottleneck in 1

Blocked migration 1-2



Horizontal:  $N_{eE}$

Solid:  $t \rightarrow N_{eG}([0, t])$  for whole population and subpopulations

Global weights:  $W_{ij} = 1/s^2$

Local weights, subpopulation  $k$ :  $W_{ij} = 1_{\{(i,j)=(k,k)\}}$

## Proof of (4)

We have that

$$h_t = \mathbf{W}h_t = \mathbf{W}\mathbf{A}^t\mathbf{h}_0 = C\lambda_{\max}(\mathbf{A})^t + o(\lambda_{\max}(\mathbf{A})^t) \quad (7)$$

as  $t \rightarrow \infty$ . But also

$$\begin{aligned} h_t &= \sum_{ij} W_{ij} h_{tij} \\ &= \sum_{ij} W_{ij} E[X_{ti}(1 - X_{tj}) + X_{tj}(1 - X_{ti})] \\ &= E[\phi(\mathbf{X}_t)] \\ &= E[E(\phi(\mathbf{X}_t)|\mathbf{X}_0)] \\ &= \sum_{\mathbf{x}, \mathbf{y}} \pi(\mathbf{x}) P^{(t)}(\mathbf{x}, \mathbf{y}) \phi(\mathbf{y}), \end{aligned} \quad (8)$$

where

$$\begin{aligned} \phi(\mathbf{x}) &= \sum_{i,j} W_{ij} [x_i(1 - x_j) + x_j(1 - x_i)], \\ \pi(\mathbf{x}) &= P(\mathbf{X}_0 = \mathbf{x}), \\ \mathbf{P}^t &= (P^{(t)}(\mathbf{x}, \mathbf{y}); \mathbf{x}, \mathbf{y} \in \mathcal{X}). \end{aligned}$$



## Perron-Frobenius

Block decompose transition matrix as

$$\mathbf{P} = \begin{pmatrix} 1 & 0 & \mathbf{0} \\ 0 & 1 & \mathbf{0} \\ \mathbf{P}_{20} & \mathbf{P}_{21} & \mathbf{P}_{22} \end{pmatrix} \implies \mathbf{P}^t = \begin{pmatrix} 1 & 0 & \mathbf{0} \\ 0 & 1 & \mathbf{0} \\ \mathbf{P}_{20}^{(t)} & \mathbf{P}_{21}^{(t)} & \mathbf{P}_{22}^t \end{pmatrix},$$

Since  $\mathbf{P}_{22}$  is non-negative, irreducible and aperiodic, it has unique largest eigenvalue  $\lambda = \lambda_3$ , and therefore

$$P^{(t)}(\mathbf{x}, \mathbf{y}) = \lambda^t r(\mathbf{x}) l(\mathbf{y}) + o(\lambda^t), \quad \mathbf{x}, \mathbf{y} \in \mathcal{X}_2, \quad (9)$$

where

$$\begin{aligned} \mathbf{l} &= (l(\mathbf{x}); \mathbf{x} \in \mathcal{X}_2) \\ \mathbf{r} &= (r(\mathbf{x}); \mathbf{x} \in \mathcal{X}_2)' \end{aligned}$$

are left and right eigenvectors of  $\mathbf{P}_2$  with eigenvalue  $\lambda$  and components

$$\begin{aligned} l(\mathbf{x}) &> 0, \\ r(\mathbf{x}) &> 0, \end{aligned} \quad (10)$$

for all  $\mathbf{x} \in \mathcal{X}_2$ .

## Proof of (4), contd.

Use (8), (9), (10) and the fact that

$$\begin{aligned}\phi(\mathbf{x}) &= 0, & \mathbf{x} \in \mathcal{X}_0 \cup \mathcal{X}_1, \\ \phi(\mathbf{x}) &> 0, & \mathbf{x} \in \mathcal{X}_2 \text{ (if all } W_{ij} > 0), \\ \pi(\mathbf{x}) &> 0, & \text{for some } \mathbf{x} \in \mathcal{X}_2,\end{aligned}$$

to conclude

$$\begin{aligned}h_t &= \lambda^t \sum_{\mathbf{x}, \mathbf{y} \in \mathcal{X}_2} \pi(\mathbf{x}) r(\mathbf{x}) l(\mathbf{y}) \phi(\mathbf{y}) + o(\lambda^t) \\ &= \lambda^t \sum_{\mathbf{x} \in \mathcal{X}_2} \pi(\mathbf{x}) r(\mathbf{x}) \cdot \sum_{\mathbf{y} \in \mathcal{X}_2} l(\mathbf{y}) \phi(\mathbf{y}) + o(\lambda^t) \\ &= C \lambda^t + o(\lambda^t),\end{aligned}$$

with  $C > 0$ . Combining this with (7), we finally deduce

$$\lambda = \lambda_{\max}(\mathbf{A}).$$

## Other topics

- ▶ Various types of structure (geographic, age, sex, combinations, ...).
- ▶ Computer program GESP (Olsson et al, 2015).
- ▶ Large population asymptotics:

$$N_{eE} = \frac{N}{C_1} + o(N) = N_{eC} + o(N) \text{ as } N \rightarrow \infty,$$

where  $N_{eC}$  is coalescence effective size<sup>5</sup> and  $C_1$  a coalescence rate.

- ▶ Small migration asymptotics:

$$N_{eE} = \frac{N}{C_2 B} + o(B^{-1}) \text{ as } B \rightarrow 0,$$

where  $B$  = long term rate of subpopulation change in ancestral line.

- ▶ Leading right eigenvector of  $\mathbf{A}$  to assess subpopulation differentiation<sup>6</sup>

---

<sup>5</sup>Nordborg and Krone (2002), Sjödin et al. (2005).

<sup>6</sup> $F_{ST}$  of Wright (1943),  $G_{ST}$  of Nei (1973).

# References

Hössjer, O. (2011). Coalescence theory for a general class of structured populations with fast migration. *Advances in Probability Theory* **43**(4), 1027-1047.

Hössjer, O., Jorde, P.E. and Ryman, N. (2013). Quasi equilibrium approximations of the fixation index under neutrality: The island model. *Theoretical Population Biology* **84**, 9-24.

Ryman, N., Allendorf, F.W., Jorde P.E., Laikre, L. and Hössjer, O. (2013). Samples from subdivided populations yield biased estimates of effective size that overestimate the rate of loss of genetic variation. *Molecular Ecology Resources* **14**, 87-99.

Olsson, F. Hössjer, O., Laikre, L. and Ryman, N. (2013). Characteristics of the variance effective population size over time using an age structured model with variable size. *Theoretical Population Biology* **90**, 91-103.

Hössjer, O. (2014). Spatial autocorrelation for subdivided populations with invariant migration schemes. *Methodology and Computing in Applied Probability* **16**(4), 777-810. DOI 10.1007/s11009-013-9321-3.

Hössjer, O. and Ryman, N. (2014). Quasi equilibrium, variance effective population size and fixation index for models with spatial structure. *Journal of Mathematical Biology* **69**(5), 1057-1128. DOI 10.1007/s00285-013-0728-9.

Hössjer, O., Olsson, F., Laikre, L. and Ryman, N. (2014). A new general analytical approach for modeling patterns of genetic differentiation and effective size of subdivided populations over time. *Mathematical Biosciences* **258**, 113-133. DOI: 10.1016/j.mbs.2014.10.001.

Hössjer, O., Olsson, F., Laikre, L. and Ryman, N. (2015). Metapopulation Inbreeding Dynamics, Effective Size and Subpopulation Differentiation - a General Analytical Approach for Diploid Organisms. *Theoretical Population Biology* **102**, 40-59. DOI:10.1016/j.tpb.2015.03.006.

**Hössjer, O. (2015). On the eigenvalue effective size in structured populations. To appear in *Journal of Mathematical Biology*. Available online, DOI 10.1007/s00285-014-0832-5.**

THANKS!