## Definition of Evolution

Evolution is the change in the frequency of different types of individuals in a population over time.

Biological evolution is the change in allele frequencies within a gene pool

Evolution generates inheritable traits and is driven by replication, mutation, selection, dispersal, environment, interactions

## Exponential Growth (unrealistic model in the long run)

Discrete case :

$x_t :=$ number of cells in generation $t$

$$X_{t+1} = 2 X_t \quad \Rightarrow \quad X_t = X_0 2^t$$

(doubling)

Continuous case :

$x(t) :=$ continuous number / fractions of cells at time $t$

$T :=$ generation time / time from birth to reproduction

$$T \sim \exp\left(\frac{1}{r}\right), \quad r := \text{growth rate}$$

$$\Rightarrow \mathbb{P}(T \leq t) = 1 - e^{-rt}$$

$$\frac{dX(t)}{dt} = r X(t) \quad \Rightarrow \quad X(t) = X_0 e^{rt}$$

) Cell Death

$d :=$ death rate $\quad \Rightarrow \quad \frac{1}{d} :=$ average life span

$\Rightarrow$ modified differential equation :

$$\frac{dX(t)}{dt} = (r - d) X(t) \quad \Rightarrow \quad X(t) = X_0 e^{(r-d)t}$$

. Basic Reproductive Ratio : expected number of offspring of an individual

$$R_0 = \frac{r}{d}$$

$\hookrightarrow R_0 > 1$ : population expands indefinitely

$\hookrightarrow R_0 < 1$ : population goes extinct

$\hookrightarrow R_0 = 1$ : population size remains constant

but $x^* = x_0$ is not stable

## (4) Logistic Growth

$K :=$ Carrying capacity of the population

- **Logistic Equation** :

$$\frac{dX(t)}{dt} = r X(t) \left( 1 - \frac{X(t)}{K} \right)$$

$$\Rightarrow X(t) = \frac{K X_0 e^{rt}}{K + X_0 (e^{rt} - 1)}$$

- **Equilibrium Points** :

  (1) $x_1^* = 0$ : stable if $r < 0$ and not stable if $r > 0$

  (2) $x_2^* = K$ : stable if $r > 0$ and not stable if $r < 0$

## (5) Stability of Equilibrium Points

- **Discrete case** :

$$x_{t+1} = f(x_t) \quad \Rightarrow \quad x^* = f(x^*)$$

$\Rightarrow x^*$ is attractive if $|f'(x^*)| < 1$

$\quad x^*$ is repelling if $|f'(x^*)| > 1$

- **Continuous case** :

$$\frac{dx(t)}{dt} = f(x(t)) \quad \Rightarrow \quad f(x^*) = 0$$

$\Rightarrow x^*$ is attractive if $f'(x^*) < 0$

$\quad x^*$ is repelling if $f'(x^*) > 0$

## (6) Logistic Difference Equation

$$x_{t+1} = r x_t (1 - x_t) \quad , \quad x_t := \text{fractions of cells}$$

- The number of equilibrium points $x^*$ depends on the value of $r$

  ↳ If $r < 1$, then the point $x^* = 0$ is stable $\Rightarrow$ population goes extinct

  ↳ If $r > 4$, then the system will diverge to $-\infty$ $\Rightarrow$ population goes extinct

  ↳ For $r$ moving from 1 to 4, the number of equilibrium points grows exponentially and eventually becomes chaotic. (around $r = 3.57$)

# (1) Selection

## (1) 2 independent exponentially growing types

Consider A: growth rate $a > 0$, abundance $x(t)$ $\Rightarrow \frac{dx(t)}{dt} = ax(t)$

B: growth rate $b > 0$, abundance $y(t)$ $\Rightarrow \frac{dy(t)}{dt} = by(t)$

### Relative abundance

$$P(t) = \frac{x(t)}{y(t)}$$

$$\Rightarrow \frac{dP(t)}{dt} = \frac{x'(t)y(t) - y'(t)x(t)}{y^2(t)} = (a-b)P(t)$$

$$\Rightarrow P(t) = P_0 \, e^{(a-b)t}$$

↳ If $a > b$, then $P \to \infty$ and selection favors A over B (not extinct)

↳ If $a < b$, then $P \to 0$ and selection favors B over A (not extinct)

↳ If $a = b$, then $P = P_0$.

## (2) 2 Competing types

$x(t) :=$ relative abundance of A, $y(t) :=$ relative abundance of B

### Constraint:

$$x(t) + y(t) = 1 \qquad \forall \, t > 0 \qquad (*)$$

### Dynamical system:

$$x'(t) = x(t)(a - \phi)$$
$$y'(t) = y(t)(b - \phi)$$

where $\phi := ax(t) + by(t)$: average fitness of the population

↳ ensures that constraint $(*)$ is satisfied

By substituting $y(t) = 1 - x(t)$, we get

$$\frac{dx(t)}{dt} = (a-b)x(t)(1 - x(t))$$

### Equilibrium points:

(1) $x^* = 1$: all-A state, stable if $a > b$

(2) $x^* = 0$: all-B state, stable if $a < b$

## 3) n Competing Types

- Probability simplex:

$$S_n = \{(x_1, \ldots, x_n) \mid x_i \geq 0, \ \sum_{i=1}^{n} x_i = 1\}$$

- Type $i$, $i \in \{1, \ldots, n\}$: fitness $f_i$, frequency $x_i(t)$
- The type frequencies are points in the $(n-1)$-dimensional probability simplex $S_n$

$$x_1(t) + \cdots + x_n(t) = 1$$

- Dynamical system:

$$x_i'(t) = x_i(t)(f_i - \phi(t)), \quad i = 1, \ldots, n$$

  where $\phi(t) = x_1(t)f_1 + \cdots + x_n(t)f_n$ : average fitness of the population

- Single Equilibrium Point : starting from any interior point in $S_n$, the fittest type will eventually outcompete all others (survival of fittest)

## (4) Subexponential and superexponential growth

$$x'(t) = ax^c - \phi x$$
$$y'(t) = by^c - \phi y$$

where $x(t) + y(t) = 1$, $A \geq 0$, $\phi(t) = ax(t)^c + by(t)^c$

$$\Leftrightarrow \quad x'(t) = x(t)(1 - x(t)) f(x(t))$$

where $f(x(t)) = ax^{c-1} - bc(1-x)^{c-1}$

- $c = 0$: growth is linear (immigration)
- $c = 1$: exponential growth
- $c > 1$: superexponential growth
- $c < 1$: subexponential growth

Fixed points:

1. $x = 0$
2. $x = 1$
3. For $c \neq 1$, $x^* = \dfrac{1}{1 + \left(\frac{a}{b}\right)^{\frac{1}{c-1}}} \in (0,1)$  (mixed population)

**(arrival of all)** $c < 1$ : $x^*$ is globally stable, whereas all-A and all-B are not (less fit type can invade)

**(arrival of first)** $c > 1$ : $x^*$ is unstable ($x > x^*$: $x \to x^*$, $A$; $x < A$: all-A - B) (less fit can prevent invasion)

# (8) Mutation

- Mutation can occur with /without reproduction.

## (1) Basic Mutation Dynamics :

Assume fitness $a = b = 1$. Let mutation rates $u_1 = P(A \to B)$ and $u_2 = P(B \to A)$

$$x' = x(1 - u_1) + y u_2 - \phi x$$
$$y' = y(1 - u_2) + x u_1 - \phi y$$

where $\phi = ax + by = 1$ and $x + y = 1$

$$\Longleftrightarrow \qquad x' = u_2 - x(u_1 + u_2)$$
$$y' = u_1 - y(u_1 + u_2)$$

$$\Longleftrightarrow \quad x^* = \frac{u_2}{u_1 + u_2} \quad , \quad y^* = \frac{u_1}{u_1 + u_2} \quad \Longrightarrow \quad \frac{x^*}{y^*} = \frac{u_2}{u_1}$$

$\Rightarrow$ Mutation can alone lead to coexistence if $u_1$ and $u_2$ are similar

$\Rightarrow$ if $u_1 >> u_2$, we may assume $u_2 = 0$ and $A$ can go extinct.
Mutation can alone lead to extinction.

## (2) Mutation Dynamics of n types

Let $Q$ be an $n \times n$ matrix with entries

$$q_{ij} = P(\text{type } i \Rightarrow \text{type } j), \quad \sum_{j=1}^{n} q_{ij} = 1$$

$\Rightarrow Q$ is a stochastic matrix

### Dynamical system :

$$x_i' = \sum_{j=1}^{n} x_j q_{ji} - \phi x_i \quad , \quad i = 1, \dots, n$$

$$\Longleftrightarrow \qquad x' = x Q - \phi x$$

- where $\phi = x_1 f_1 + \cdots + x_n f_n$ : average fitness of the population

### Equilibrium point :

$x^* := $ left eigenvector of $Q$ associated with the largest eigenvalue 1.

$\hookrightarrow$ asymmetric mutation can result in selection even without growth differences.

# Chapter 2 Quasispecies

## 1) Central Dogma of Molecular Biology

$$DNA \xrightarrow{\text{transcribes}} RNA \xrightarrow{\text{translates}} Protein$$

## 2) RNA Virus

- use the transcription machinery of the host cell to replicate
- short genome
- very high mutation rate (no proof-reading of reverse transcription)
- high turnover
- exposed to strong selective forces
- extreme evolutionary dynamics

## 3) Sequence Space

- **Sequence of length L** :

$$A^L = \{(a_1, \ldots, a_L) : a_i \in \underset{\downarrow \text{ alphabet}}{A}\}$$

- Sequence space : $A^* = \underset{L \geq 0}{\cup} A^L$

- # of sequences in $A^L := |A|^L$

- **Hamming distance** :

$$\| x - y \| = \sum_{i=1}^{L} \mathbb{1}\{x_i \neq y_i\} \leq L$$

- Evolution is a trajectory of a population in sequence space.

- **Genotype space** :

$$G \subseteq A^*$$

- **Fitness Landscape** :

$$f : G \to \mathbb{R} \quad s.t. \quad f(g) \in \mathbb{R}$$

  ↳ for individual : discrete ; for average population fitness : continuous

  ↳ epistasis : $\varepsilon = (f_{00} + f_{11}) - (f_{01} + f_{10}) \Rightarrow$ measure the strength of additive effects

# (4) The Quasispecies Equation

Let $x(t) = (x_0(t), \ldots, x_n(t))$ be the genotype frequencies for $i = 1, \ldots, n$ genotypes at $t$

Let $Q = (q_{ij}) = (q_{i \to j})$ be a mutation matrix.

Let $f = (f_0, \ldots, f_n)$ be a fitness landscape.

Denote by $\phi = x^T f$ the average fitness of the population.

$$\dot{x}_i = \frac{dx_i}{dt} = \sum_{j=0}^{n} x_j f_j q_{ji} - \phi x_i, \quad i = 1, \ldots, n$$

with $\downarrow$ selection  $\downarrow$ mutation

- **No mutation:** $Q = I$
  $\hookrightarrow$ we recover the selection equation ("survival of the fittest")
- **No selection:** $f = (1, \ldots, 1)$
  $\hookrightarrow$ we recover the mutation equation
- If $Q$ is irreducible, there exists a globally stable equilibrium $x^*$ inside $S_{n+1}$ but $x^*$ does NOT maximize the fitness $\phi$.
- Mutation-Selection Matrix:

$$W = (w_{ij}) = (f_j q_{ji})$$

$$\Rightarrow \quad \dot{x} = Wx - \phi x$$

$\hookrightarrow$ if $W$ is non-negative and irreducible, then the largest eigenvalue of $W$ corresponds to the average fitness $\phi$, and the associated eigenvector after normalization is the global equilibrium point. = mutation-selection balance

# (5) Adaptation (HIV wants to be close to the error threshold but not over it)

- adaptation of a population is localization in sequence space at a local maximum of the fitness landscape
- selection drives the population towards the peak, whereas mutation drives away
- **error threshold:** necessary condition on mutation rate for adaptation to occur
  $\hookrightarrow$ a simplified case: $x_0$ wild type with fitness $f_0 > 1$, $x_1$ other types with fitness 1
  Let $u :=$ mutation rate $\Rightarrow$ wild type is copied error-free $\bar{w}$ $q = (1-u)^L$
  Ignoring back mutation, we get

$$\begin{cases} \dot{x}_0 = x_0 f_0 q - \phi x_0 \\ \dot{x}_1 = x_0 f_0 (1-q) + x_1 - \phi x_1 \end{cases}$$

$$\Rightarrow \quad x_0^* = \begin{cases} \frac{f_0 q - 1}{f_0 - 1} & \text{if } f_0 q > 1 \\ 0 & \text{otherwise} \end{cases}$$

If $u \ll 1$ and $\log f_0 \approx 1$, then
$\log(f_0 q) = 1 + L \log(1-u) \approx 1 - Lu > \log 1 =$
$\Leftrightarrow uL < 1$ (expected mutation per replication)
$\Leftrightarrow u_c = \frac{1}{L}$

If $uL > 1$, have mutational meltdown

# Chapter 3 : Stochastic Models of Finite Populations

## 1) Probability

- **Exponential distribution** : $X \sim \exp(\lambda)$, $f(x) = \lambda e^{-\lambda x}$, $x \geq 0$

$$P(X \leq x) = 1 - e^{-\lambda x}, \quad P(X > x) = e^{-\lambda x}$$

$$E[X] = \frac{1}{\lambda}, \quad Var[X] = \frac{1}{\lambda^2}$$

  ↳ **memoryless property** : $P(X > s + t) = P(X > s) P(X > t)$

$$P(X > s + t \mid X > t) = P(X > s)$$

  ↳ **Competing exponentials** :

$$X \sim \exp(\lambda), \quad Y \sim \exp(\mu), \text{ and } X \perp\!\!\!\perp Y$$

$$\Rightarrow \min(X, Y) \sim \exp(\lambda + \mu) \text{ and } P(X < Y) = \frac{\lambda}{\lambda + \mu}$$

- **Markov chain** :

$$\{ X(t) \mid t \in T = \{0, 1, 2, \ldots\} \}$$

  with

$$P(X(t+1) \mid X(0), \ldots, X(t)) = P(X(t+1) \mid X(t))$$

  ↳ **transition matrix** : $P = (P_{ij})$, $P_{ij} = P(X(t+1) = j \mid X(t) = i)$

  ↳ **time-homogeneous Markov chain** : $P(X(t+1) = j \mid X(t) = i) = P(X(1) = j \mid X(0) = i)$

  ↳ **absorbing state** : $X(t) = x^* \quad \forall t \geq t_0$

  ↳ An **ergodic** Markov chain has a unique **stationary distribution** $\pi$ such that

$$\pi^T P = \pi^T$$

  where
$$\lim_{t \to \infty} P_{ij}(t) = \pi_j, \quad \forall i, j \in S := \text{state space}$$

## 2) The Moran Process

### (1) Model :

- Population size : $N < \infty$.
- 2 types of individuals A and B
- process :

  (1) pick randomly an individual for reproduction

  (2) pick randomly an individual for death

  (3) the offspring of the first individual replaces
     the second individual

- Both types have the **same** probability of reproduction and death.

  $\Rightarrow$ **Neutral drift** : changes in allele frequency are only due to random fluctuations.

(2) Birth-Death Process

The state space is $\{0, 1, \ldots, N\}$. Let $i :=$ # of A individuals.

Let $p = \frac{i}{N} :=$ allele frequency of A.

Then, the transition matrix is given by

$$P_{i,i+1} = p(1-p)$$
$$P_{i,i-1} = (1-p)p$$
$$P_{i,i} = 1 - 2p(1-p) = p^2 + (1-p)^2$$

↳ P is a tri-diagonal matrix with all other entries zero

↳ it is a birth-death process, because the # of A individuals can only change one step at a time

(3) Absorbing states

1. All-A individuals: $P_{N,N} = 1$ and $P_{N,i} = 0 \ \forall i < N$

2. All-B individuals: $P_{0,0} = 1$ and $P_{0,i} = 0 \ \forall i > 0$

(4) Fixation Probabilities:

Let $x_i :=$ probability of ending up in state $N$ when starting from state $i$

Then, in the Moran process, we must have
(neutral)

$$x_i = \frac{i}{N}, \quad i = 0, \ldots, N$$

because each cell has the same chance to produce offspring, then each cell also has the same chance to be the origin of cells eventually in the population, and thus the same chance of fixation.

(5) Stationary Mean:

$$\mathbb{E}[X(t) \mid X(0) = i] = i$$

(6) Time-dependent Variance:

$$\lim_{N \to \infty} \text{Var}(X(t) \mid X(0) = i) = V_1 t$$

where $V_1 = \text{Var}(X(1) \mid X(0) = i) = 2p(1-p) = 2\frac{i}{N}(1 - \frac{i}{N})$

(7) General Birth-Death Process:

$$P_{i,i+1} = \alpha_i, \quad P_{i,i-1} = \beta_i \quad \Rightarrow \quad \gamma_i = \frac{\beta_i}{\alpha_i}, \quad \alpha_0 = \beta_N = 0$$

$$\Rightarrow \quad x_i = \frac{1 + \sum_{j=1}^{i-1} \prod_{k=1}^{j} \gamma_k}{1 + \sum_{j=1}^{N-1} \prod_{k=1}^{j} \gamma_k}, \quad i = 0, \ldots, N$$

:8) <u>Mean Fixation Time</u>

$$-N^2\left[(1-p)\log(1-p) + p\log p\right]$$

:9) <u>Heterozygosity</u>

$H_T :=$ probability that two individuals chosen at random from the population are of *different* types.

$$\mathbb{E}[H_t \mid X(0) = i] = H_0(i)\left(1 - \frac{2}{N^2}\right)^t , \quad H_0(i) = 2\cdot\frac{i}{N}\cdot\frac{N-i}{N-1}$$

↳ decays exponentially at rate $\frac{2}{N^2}$

↳ quantifies the amount of random drift that the population is experiencing

③ <u>Moran Process with Constant Selection</u>

Consider 2 types of individuals A and B with growth rates

$$\lambda_A = r \quad \text{and} \quad \lambda_B = 1$$

Let the waiting time to the next birth of A be

$$T_A \sim \min\underbrace{\{\exp(\lambda_A), \ldots, \exp(\lambda_A)\}}_{i \text{ individuals}} = \exp(ir)$$

Similarly,

$$T_B \sim \exp\left((N-i)\lambda_B\right) = \exp(N-i)$$

Then,

$$P(T_A < T_B) = \frac{ir}{ir + N - i}, \quad P(T_B < T_A) = \frac{N-i}{ir + N - i}$$

The transition probabilities become

$$P_{i,i+1} = \underbrace{\frac{ri}{ri + N - i}}_{P(T_A < T_B)} \cdot \underbrace{\frac{N-i}{N}}_{\text{choose a B to die}}$$

$$P_{i,i-1} = \underbrace{\frac{N-i}{ri + N - i}}_{P(T_A > T_B)} \cdot \underbrace{\frac{i}{N}}_{\text{choose an A to die}}$$

$$P_{i,i} = 1 - P_{i,i+1} - P_{i,i-1}$$

With $\gamma_i = P_{i,i-1}/P_{i,i+1} = \frac{1}{r} \; \forall i = 1, \ldots, N-1, \; \gamma_0 = 0$ and $\gamma_N = 1$, the fixation probabilities become

$$x_i = \frac{1 - 1/r^i}{1 - 1/r^N}, \quad i = 0, \ldots, N$$

In particular,

$$p = x_1, \quad \lim_{r \to 1} p = \frac{1}{N} \quad \text{(neutral Moran process)}$$

# (4) Poisson Process

- A Poisson process is a stochastic counting process, a continuous-time Markov chain with independent Poisson distributions in each interval.

- Model:

$$\{N(t) \mid t \geq 0\} \quad \text{with} \quad N(0) = 0$$

$$N(t+s) - N(s) \sim Poisson(\lambda t) \Rightarrow P(N(t+s) - N(s) = k) = \frac{e^{-\lambda t}(\lambda t)^k}{k!}$$

- Inter-arrival times are exponential:

$$\{T_n \mid n = 1, 2, \ldots\}$$

$$T_n \sim exp(\lambda), \quad n = 1, 2, \ldots$$

sketch proof:

$$P(T_1 > t) = P(N(t) = 0) = e^{-\lambda t}$$

$$P(T_2 > t) = \mathbb{E}_{T_1}[P(T_2 > t) \mid T_1] \quad \text{by law of total expectation}$$

$$= \int_s P(N(s+t) = N(s) \mid T_1 = s) f_{T_1}(s) \, ds$$

$$= \int_s P(N(t) = 0) f_{T_1}(s) \, ds$$

$$= e^{-\lambda t}$$

$$\Rightarrow \text{proof by induction}$$

# ⑤ Rate of Evolution

Consider an all-A population where a B mutant occur at mutation rate $u \ll 1$. Then the number of mutations over time can be modeled by the Poisson process. The waiting time for the first mutant to occur is

$$T_1 \sim exp(\underbrace{Nu}_{\text{total mutation rate}})$$

Suppose B has an selective advantage $r$ and the fixation probability is $P = x_1$.

$\Rightarrow$ Rate of Evolution from all-A to all-B:

$$R = \underbrace{N u}_{\text{prob. of first mutant to occur}} \underbrace{\rho}_{\to \text{ prob. of that mutant to fixate}}$$

↳ If $r = 1$, then $P = \frac{1}{N}$ and $R = u$, which means in the neutral case, the rate of evolution is independent of the population size and depends solely on the mutation rate.

# Chapter 4  Evolutionary Dynamics of Cancer

## 1) Cancer

- Cancer is a breakdown of cellular cooperation for multicellular organisms.
- The somatic evolution of cancer is the uncontrolled, selfish replication of cells, which give rise to tumors.
- Cancer progression: accumulation of mutations in genes, which will increase somatic fitness of cells (apoptosis escape)
- Most cancer cells are aneuploid (strange copy numbers)
- normal $\Rightarrow$ adenoma $\longrightarrow$ cancer $\longrightarrow$ metastasis

## 2) Oncogenes

- Oncogenes increase fitness if one allele is mutated or inappropriately expressed.
- Activation : (1) Point Mutation   (2) Gene Amplification   (3) Chromosomal Fusion
- Fixation : Moran process in a compartment with initially normal cells
  - population size N    - mutation rate $u$    - fitness advantage $r$
  - fixation probability if there is one mutant :

$$\rho = x_1 = \frac{1 - 1/r}{1 - 1/r^N}$$

  - probability that a mutant has been fixed by time $t$ :

$$P(t) = 1 - e^{-Nu\rho t}$$

$\downarrow$ rate of evolution from all-A to all-B

  1) $r > 1$ : advantageous fitness, $P(t)$ increases as N increases
     - large compartments accelerate accumulation of mutations
     - to prevent fixation of a mutant, most tissues with high turnover are organized in many small compartments.
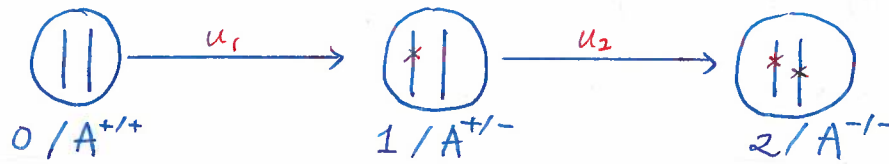  2) $r < 1$ : deleterious fitness, $P(t)$ decreases as N increases
  3) $r = 1$ : fixation of a mutant (new) is independent of N
     - $P(t) = 1 - e^{-ut}$

## 3) Linear Process of Cancer

- Architecture : one or few stem cells that perform asymmetric division into stem cells and cells that are more differentiated. The cells are pushed linearly towards the edge and will eventually undergo apoptosis.
- Only mutations on stem cells can lead to fixation, while all other mutations are "washed ou
- Fixation probability is independent of the fitness advantage $r$ . $P = \frac{1}{N}$ (one stem cell
- Perfect design to protect against mutations in tumour suppressor genes/oncogenes but not instabilit

# ④ Tumour Suppressor Genes (TSG)

- TSGs increase fitness if <u>both alleles</u> are mutated by
  - 1) two point mutations
  - 2) one point mutation + loss of heterozygosity (LOH)
  - ↳ TSGs are inactivated

- <u>Dynamics of TSG Inactivation</u>



$$0 / A^{+/+} \qquad \xrightarrow{u_1} \qquad 1 / A^{+/-} \qquad \xrightarrow{u_2} \qquad 2 / A^{-/-}$$

<u>Goal</u>: In a population of size $N$, what is the probability that at least one cell has been <span style="color:red">hit by 2 mutations</span> by time $t$? $\Rightarrow P(t)$

1) <u>Small Population Size</u>

- <u>Population size</u>: $N$
- <u>Mutation rates</u>: $u_1$, $u_2$

$\Rightarrow$ <u>Fixation probability of the first mutation</u>: (neutral)
$$\rho = x_1 = \frac{1}{N}$$

↳ <u>time until first fixation</u>:
$$T_1 \sim \exp\left(\frac{1}{N}\right)$$
$$\Rightarrow \mathbb{E}[T_1] = N = \text{population size}$$

$\Rightarrow$ <u>Waiting time for the second mutation to occur in any cell</u>:
$$T_2 := \min\{\tilde{T}_1, \ldots, \tilde{T}_N\} \sim \exp(N u_2)$$
$$\Rightarrow \mathbb{E}[T_2] = \frac{1}{N u_2}$$

$\Rightarrow$ <u>Type 1 cells reach fixation before a type 2 cell arises</u>:
$$N \ll \frac{1}{N u_2} \quad \Longleftrightarrow \quad N \ll 1/\sqrt{u_2} \text{ (definition of small population size)}$$

- <u>Dynamical System</u>:

  State 0: all type 0 ; State 1: all type 1 ; State 2: 1 type 2
  $$\dot{X}_0 = -N u_1 \rho X_0 = -N \cdot u_1 \cdot \frac{1}{N} X_0 = -u_1 X_0$$
  $$\dot{X}_1 = u_1 X_0 - N u_2 X_1$$
  $$\dot{X}_2 = N u_2 X_1$$

  where $X_0, X_1, X_2$ are probabilities

  In particular, $P(t) = X_2(t)$

(2) <u>Intermediate Population Size</u>

- Average waiting time for a type 1 cell to occur is

$$\frac{1}{Nu_1}$$

  ↳ we say that it is long if $\frac{1}{Nu_1} > 1 \iff N < \frac{1}{u_1}$

- A type 2 cell is generated before the fixation of type 1 cell if

$$N > 1/\sqrt{u_2}$$

⇒ <u>Intermediate regime for tunneling</u> :

$$\frac{1}{\sqrt{u_2}} << N << \frac{1}{u_1}$$

s.t.

$$0 \longrightarrow 1 \longrightarrow 2$$

(3) <u>Large Population Size</u> :

$$N >> \frac{1}{u_1}$$

  ↳ type 1 cells will be generated immediately (instantaneous waiting time) and they grow linearly according to
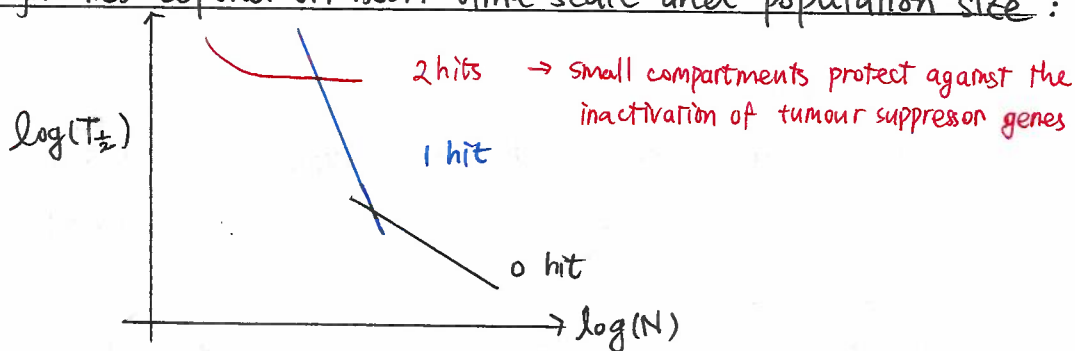
$$x_1(t) = Nu_1 t$$

  ↳ probability of generating a type 2 cell during type 1 growth is

$$P(t) = 1 - \exp\left(-u_2 \underbrace{\int_0^t x_1(t)\, dt}_{\text{abundance}}\right) \quad \text{of type 1 at time } t$$

$$= 1 - \exp\left(-\tfrac{1}{2} Nu_1 u_2 t^2\right)$$

- <u>TSG Inactivation Dynamics depend on both time scale and population size</u> :



$\log(T_{\frac{1}{2}})$    2 hits → small compartments protect against the inactivation of tumour suppressor genes

1 hit

0 hit

→ $\log(N)$

(1) <u>Short time scale</u> : $t << \frac{1}{Nu_2}$ , we have $P(t) \approx Nu_1 u_2 t^2/2$ (2 rate limiting events
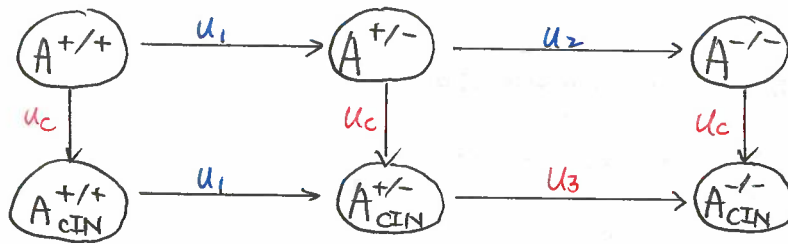
(2) <u>Intermediate time scale</u>: $\frac{1}{Nu_2} << t << \frac{1}{u_1}$ , $P(t) \approx 1 - e^{-u_2 t}$ (first hit is rate limiting

(3) <u>Long time scale</u> : $t >> \frac{1}{u_1}$ , $P(t) = 1 - \exp(-\tfrac{1}{2} Nu_1 u_2 t^2)$ (no rate limiting events

# (5) Chromosomal Instability (CIN)

- CIN leads to increased rate of gaining or losing whole chromosomes (LOH) or large parts of it
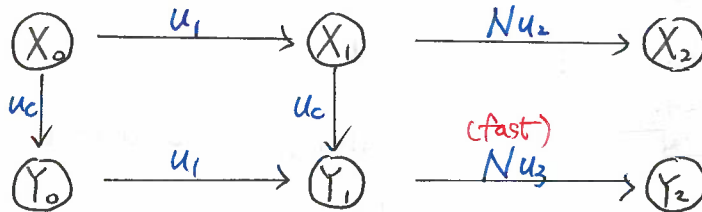  $\Rightarrow$ inactivation of TSGs ($u \approx 10^{-2}$)

- **Model** :



## (1) Neutral CIN :

Consider small compartments: $N \ll \frac{1}{u_1}, \frac{1}{u_2}, \frac{1}{u_c}$

Assume $A^{+/-}$ and CIN cells are neutral and $A^{-/-}$ will be fixed immediately



$\Rightarrow X_2(t) \approx N u_1 u_2 t^2 / 2$ and $Y_2(t) \approx Y_1(t) \approx u_1 u_c t^2$

$\Rightarrow$ waiting time for LOH is negligible since $u_3 \approx 10^{-2}$ !

## (2) Costly CIN in small compartment :

Assume CIN cells have fitness advantage $r < 1$.

$\Rightarrow$ fixation probability of a CIN cell :

$$p = \frac{1 - 1/r}{1 - 1/r^N}$$

$\Rightarrow$ non-CIN cell state to all-CIN cell state rate of evolution :

$$N u_c p$$

$\Rightarrow X_2(t) \approx N u_1 u_2 t^2 / 2$ and $Y_2(t) \approx Y_1(t) \approx N p u_1 u_c t^2 < u_1 u_c t^2$

## (3) Costly CIN in large compartment :

For large $N$, $Np$ becomes vanishingly small, so intermediate CIN types will not fixate. The population tunnels from $X_1$ to $Y_2$ at rate

$$R = \frac{N u_c u_3 r}{1 - r}$$

$\Rightarrow X_2(t) \approx N u_1 u_2 t^2 / 2$ and $Y_2(t) \approx R u_1 t^2 / 2$, $Y_0(t) = Y_1(t) = 0$

$\Rightarrow$ Large compartment protects against the fixation of CIN cells

$\Rightarrow$ Linear design (e.g. colon crypts) is the best design to protect against oncogenes and TSGs but not to genetic instability

# Chapter 5   Cancer Progression

## I) Waiting Times

### (1) Multistage Theory

Assume that tumour progression follows a *linear, multistep* process where each step has mutation rate $u_j \ll 1$.

The waiting time for each step $j$ is exponential, $\exp(u_j)$.

$\Rightarrow$ The waiting time until stage $k$ is reached:

$$T_k \sim \exp(u_1) + \cdots + \exp(u_k)$$

$$\mathbb{E}[T_k] = \mathbb{E}\left[\sum_{j=1}^{k} \exp(u_j)\right] = \sum_{j=1}^{k} \frac{1}{u_j}$$

### (2) Independent mutations:

Assume that each mutation occurs *independently* at time

$$T_j \sim \exp(\lambda_j), \quad j = 1, \ldots, d$$

$\Rightarrow$ The waiting time until *any* $k$ out of $d$ mutations have occurred:

$$T_k \sim \underbrace{\min_{\{j_1,\ldots,j_k\} \subseteq \{1,\ldots,d\}}}_{\substack{\text{pick the earliest} \\ \text{combination of } k}} \underbrace{\max\{T_{j_1}, \ldots, T_{j_k}\}}_{\text{waiting time for all } k \text{ mutations}}$$

$$T_1 \sim \min\{T_1, \ldots, T_d\} = \exp\left(\sum_{j=1}^{d} \lambda_j\right)$$

### (3) Independent mutations + identical mutation rates:

Assume that

$$T_j \sim \exp(\lambda), \quad j = 1, \ldots, d$$

Then, the waiting time until any $k$ out of $d$ mutations:

$$T_1 \sim \exp(d\lambda)$$

$$T_k \sim \underbrace{T_{k-1}}_{\substack{\text{waiting time until} \\ \text{any } (k-1) \text{ out of } d \text{ mutations}}} + \underbrace{\exp\big((d-(k-1))\lambda\big)}_{\text{waiting time until any one out of } d-(k-1) \text{ mutations}}$$

$$\mathbb{E}[T_k] = \mathbb{E}[T_{k-1}] + \frac{1}{(d-(k-1))\lambda}$$

$$= \mathbb{E}[T_{k-2}] + \frac{1}{(d-(k-1))\lambda} + \frac{1}{(d-(k-2))\lambda}$$

$$\vdots$$

$$= \frac{1}{\lambda} \sum_{j=1}^{k} \frac{1}{d-j+1}$$

(4) Mutational Pathways

· A mutational pathway is a genotype lattice can be written as

$$P = j_1 \to j_2 \to \cdots \to j_k$$

where $j_1 < \cdots < j_k$ defines the order of mutations.

· $Exit_i$ := set of all possible mutations in step $i$.

· The probability of $P = j_1 \to \cdots \to j_k$ is

$$Prob(P) = \prod_{i=1}^{k} \frac{\lambda_{j_i}}{\sum_{j \in Exit_i} \lambda_j}$$

· The expected waiting time is

$$E[\tau_P] = \sum_{i=1}^{k} \frac{1}{\sum_{j \in Exit_i} \lambda_j}$$

· The expected waiting time of any $k$ out of $d$ mutations :

$$E[\tau_k] = \sum_{P = j_1 \to \cdots \to j_k} Prob(P) \cdot E[\tau_P]$$

$$\downarrow$$
sum over all possible pathways of length $k$

② The Wright-Fisher Process

· <u>Process</u> : Consider a constant population size $N$ and discrete generations. In each generation $t$, each individual chooses randomly a parent from the previous generation. The individuals can be of different types.

· <u>Model</u> :

Consider 2 types of individuals A and B.

Let $X(t)$ := # of type A individual in generation $t = 0, 1, 2, \ldots$

$$X(t) \in \{0, 1, \ldots, N\}$$

$$\Rightarrow X(t+1) \,|\, X(t) = i \sim Binomial(N, \tfrac{i}{N})$$

                                  # of draws    proportion of type A individuals

$$\Rightarrow P_{ij} = \mathbb{P}(X(t+1) = j \,|\, X(t) = i)$$

$$= \binom{N}{j} \left(\frac{i}{N}\right)^j \left(1 - \frac{i}{N}\right)^{N-j}$$

· <u>Stationary mean</u> :

$$E[X(t) \,|\, X(0) = i] = i$$

\* Both the Moran process and the Wright-Fisher process model the genetic random drift.

## Time-dependent Variance :

$$\lim_{N \to \infty} Var(X(t) | X(0) = i) = i \left(1 - \frac{i}{N}\right) t$$

↳ ratio between WF process and Moran process $= \frac{N}{2}$

↳ Variance increases faster for the WF process because it updates all N individuals simultaneously, whereas the Moran process updates one at a time

· ## Heterozygosity :

$$\mathbb{E}[H_t | X(0) = i] = H_0(i) \left(1 - \frac{1}{N}\right)^t, \quad H_0(i) = 2 \cdot \frac{i}{N} \cdot \frac{N-i}{N-1}$$

↳ the heterozygosity decays exponentially at rate $\frac{1}{N} > \frac{2}{N^2}$, which is faster than in the Moran process.

· ## Fixation Probabilities :

2 absorbing states : $X(t) = 0$ and $X(t) = N$

Let the probability of fixation of type A individuals when starting with $i$ be

$$x_i = \lim_{t \to \infty} \mathbb{P}(X(t) = N | X(0) = i)$$

We have

$$i = \lim_{t \to \infty} \mathbb{E}[X(t) | X(0) = i] = 0 \times (1 - x_i) + N \times x_i$$

(only 2 possibilities →)

$$\iff \qquad x_i = \frac{i}{N}$$

## 3) Wright-Fisher Process with Mutation and Selection

## i) Accumulating Mutations

Consider a binary genome of length $d$.

↳ each locus undergoes independent mutation from 0 to 1 at rate $u$

↳ no back mutations

Let $X_j(t) :=$ # of j-cells (cells with j mutations) in generation $t$

↳ $x_j(t) = \frac{X_j(t)}{N}$

Assume a constant fitness advantage $s$ per mutation. → dominates the waiting time to cance

↳ the fitness of a j-cell $= (1+s)^j$

Assume $X_0(0) = N$, $X_j(0) = 0 \quad \forall j = 1, \dots, d$.

⇒ ## probability of sampling a j-cell

$$\theta_j(t) = \sum_{i=0}^{j} \mathbb{P}(i\text{-cell} \to j\text{-cell})$$

$$= \sum_{i=0}^{j} \mathbb{P}(i\text{-to-}j \text{ mutations}) \, \mathbb{P}(i\text{-cell parent})$$

$$= \sum_{i=0}^{j} \binom{d-i}{j-i} u^{j-i} (1-u)^{d-j} \cdot \frac{(1+s)^i x_i(t)}{\sum_{\ell} (1+s)^\ell x_\ell(t)}$$

The sampling step follows a multinomial distribution with parameters
· $X(t) = (X_0(t), \dots, X_d(t))$
· $\theta(t) = (\theta_0(t), \dots, \theta_d(t))$

# Chapter 6 Diffusion Theory

## ① Structure of the Diffusion Theory

- **Goal**: model the probability distribution of a population over time
- Consider only 2 populations $A_1$ and $A_2$.
- $\psi(p, t) :=$ probability density that $A_1$ has frequency $p$ at time $t$
- $g(p, \varepsilon ; dt) :=$ probability that allele frequency of $A_1$ changes from $p$ to $p+\varepsilon$ in time interval $t$

$$\Rightarrow \quad \psi(p, t+dt) = \int \psi(p-\varepsilon, t) \, g(p-\varepsilon, \varepsilon ; dt) \, d\varepsilon$$

$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\text{marginalizing out the amount of change } \varepsilon}$

<u>Interpretation</u> : if the allele frequency is $p$ at time $t+dt$, then it must have been $p-\varepsilon$ at time $t$ for some amount $\varepsilon > 0$

## ② Two classes of Evolutionary Forces

(1) <u>Directional processes</u>:  $M(p)$
   - ↳ expected change in allele frequency per generation
   - ↳ e.g. mutation, selection, migration, recombination

(2) <u>Nondirectional processes</u>:  $V(p)$
   - ↳ expected variance in the next generation
   - ↳ e.g. genetic drift.

- **example : Moran process**

$$p = p(t) = \frac{X(t)}{N}$$

$$\Rightarrow M(p) = \mathbb{E}[p(t+1) - p(t) \mid p(t)] = p - p = 0$$

$$V(p) = \mathbb{E}[\text{Var}(p(t+1)) \mid p(t)] = \frac{1}{N^2} \mathbb{E}[\text{Var}(X(t+1)) \mid X(t)] = \frac{2p(1-p)}{N^2}$$

   - ↳ the neutral Moran process models only the random drift

③ Kolmogorov Forward Equation / Fokker-Planck Equation / Diffusion Equation

$$\frac{\partial \psi(p,t)}{\partial t} = -\frac{\partial}{\partial p}[\psi(p,t)M(p)] + \frac{1}{2}\frac{\partial^2}{\partial p^2}[\psi(p,t)V(p)] \qquad (*)$$

• Derivation: Let $\psi := \psi(p,t)$, $g := g(p,\varepsilon; \delta t)$
$$\psi(p,t+\delta t) = \int \psi(p-\varepsilon, t)\, g(p-\varepsilon, \varepsilon; \delta t)\, d\varepsilon$$

$$= \int\left[\psi g - \varepsilon\frac{\partial(\psi g)}{\partial p} + \frac{\varepsilon^2}{2}\frac{\partial^2(\psi g)}{\partial p^2} - \frac{\varepsilon^3}{6}\frac{\partial^3(\psi g)}{\partial p^3} + \cdots\right]d\varepsilon \qquad \text{by Taylor's expansion}$$

$$\approx \int\left[\psi g - \varepsilon\frac{\partial(\psi g)}{\partial p} + \frac{\varepsilon^2}{2}\cdot\frac{\partial^2(\psi g)}{\partial p^2}\right]d\varepsilon \qquad \text{by assuming } \varepsilon^2 >> \varepsilon^3$$

$$= \psi\int g\, d\varepsilon - \frac{\partial}{\partial p}\psi\int\varepsilon g\, d\varepsilon + \frac{1}{2}\frac{\partial^2}{\partial p^2}\psi\int\varepsilon^2 g\, d\varepsilon \qquad \text{since } \int p = 1$$

$$\underbrace{\quad}_{=1} \qquad \underbrace{\quad}_{=M(q)\,dt} \qquad \underbrace{\quad}_{=V(q)\,dt}$$

$$= \psi(p,t) - \frac{\partial[\psi(p,t)M(p)]}{\partial p}dt + \frac{1}{2}\frac{\partial^2[\psi(p,t)V(p)]}{\partial p^2}dt$$

→ Subtract both sides by $\psi(p,t)$
→ divide by $dt$ and let $dt \to 0$

• Equilibrium:
Setting (*) to zero gives
$$\frac{1}{2}\frac{\partial}{\partial p}[\psi^*(p,t)V(p)] - \psi^*(p,t)M(p) = 0$$

$$\Rightarrow \qquad \psi^*(p) = \frac{C}{V(p)}\exp\left(\int_0^p \frac{2M(q)}{V(q)}dq\right)$$

↳ In one-dimensional case:
$$\psi^*(p,t) = \frac{1}{\sqrt{2\pi\sigma^2 t}}\exp\left(-\frac{(p-mt)^2}{2\sigma^2 t}\right)$$

which is the pdf of $N(mt, \sigma t)$. As $t \to 0$, $\psi(p,t)$ converges to a point mass $p=0$

④ Combining Mutation, Selection, and Wright-Fisher type Sampling

(1) Selection.
Assume $A_1$ and $A_2$ have frequencies $p$ and $1-p$, with fitness $w_1$ and $w_2$.
Then, the average fitness of the population is
$$\bar{w} = pw_1 + (1-p)w_2 \qquad \text{and} \qquad \frac{d\bar{w}}{dp} = w_1 - w_2$$

In the Wright-Fisher process, we sample a parent at random.
⇒ the allele frequency of $A_1$ in the next generation is
$$p' = \frac{pw_1}{pw_1 + (1-p)w_2} = \frac{pw_1}{\bar{w}}$$

$\Rightarrow$ <u>difference in allele frequency due to selection</u> :

$$\Delta P_{sel} = P' - P$$

$$= \frac{P(w_1 - \bar{w})}{\bar{w}} \quad \text{by definition}$$

$$= \frac{P(1-P)(w_1 - w_2)}{\bar{w}} \quad \text{by definition}$$

$$= P(1-P) \cdot \frac{1}{\bar{w}} \cdot \frac{d\bar{w}}{dp}$$

$$= P(1-P) \frac{d\log(\bar{w})}{dp}$$

$\hookrightarrow$ also called : Wright's equation for an adaptive landscape

## (2) Mutation

Let $u_1 := A_1 - \text{to} - A_2$ mutation rate

$\quad u_2 := A_2 - \text{to} - A_1$ mutation rate

$\Rightarrow \qquad \Delta P_{mut} = -p u_1 + (1-p) u_2$

<u>Combining selection and mutation we get</u> :

$$M(p) = \underbrace{P(1-P) \frac{d\log(\bar{w})}{dp}}_{\Delta P_{sel}} \underbrace{- p u_1 + (1-p) u_2}_{\Delta P_{mut}}$$

## (3) Sampling :

In the Wright-Fisher process :

$$Var[X(t+1) \mid X(t) = i] = Np(1-p) \quad \text{with } p = \frac{i}{N}$$

Thus,

$$V(p) = \mathbb{E}[Var(P(t+1)) \mid p(t)]$$

$$= \mathbb{E}[Var(\tfrac{1}{N} X(t+1)) \mid \tfrac{1}{N} X(t)]$$

$$= \tfrac{1}{N^2} \cdot Np(1-p)$$

$$= \frac{P(1-P)}{N}$$

$$\Rightarrow \int_0^P \frac{M(q)}{V(q)} = \int_0^P N\left( \frac{d\log(\bar{w})}{dq} - \frac{u_1}{1-q} + \frac{u_2}{q} \right) dq$$

$$= N\left( \log \bar{w} + u_1 \log(1-p) + u_2 \log(p) \right)$$

$$\Rightarrow \psi^*(p) = \frac{C}{V(p)} \exp\left( \int_0^P \frac{2M(q)}{V(q)} dq \right)$$

$$= \frac{CN}{P(1-P)} \exp(2N(\log \bar{w} + u_1 \log(1-p) + u_2 \log(p)))$$

$$= C \bar{w}^{2N} (1-p)^{2Nu_1 - 1} p^{2Nu_2 - 1}$$

where $C$ is a normalizing constant s.t. $\int_0^1 \psi^*(p) \, dp = 1$

**(4) Equilibrium under mutation and drift**

For $\bar{w} = 1$ and $u = u_1 = u_2$,

$$\psi^*(p) \propto [p(1-p)]^{2Nu - 1}$$

where $\theta = 2Nu$ : scaled mutation parameter, which captures the impacts of both mutation rate and the population size on the distribution of allele frequency at equilibrium.

↳ for high mutation rate relative to $N$ : coexistence ("unimodal")

↳ for low mutation rate relative to $N$ : fixation to one of the two absorbing states

**(5) Equilibrium under selection and drift**

For $w_1 = 1 + s$ and $w_2 = 1$ for some small $s$ and $u_1 = u_2 = 0$,

$$\Rightarrow \quad \bar{w} = p(1+s) + (1-p)$$
$$= 1 + sp$$
$$\approx e^{sp} \quad \text{for small } s$$

$$\Rightarrow \quad \psi^*(p) \propto \frac{\bar{w}^{2N}}{p(1-p)} \approx \frac{e^{2Nsp}}{p(1-p)}$$

where $\sigma = 2Ns$ : scaled selection parameter, which captures the impacts of both selective advantage and population size on $\psi^*(p)$

↳ for large selective advantage : $\psi^*(p)$ concentrates on large $p$

↳ for small selective advantage : $\psi^*(p)$ concentrates on $\frac{1}{2}$

Putting (4) and (5) together gives the diffusion approximation for the Wright-Fisher process

$$\psi^*(p) \propto [p(1-p)]^{\theta - 1} e^{\sigma p}$$

with $\theta = 2Nu$, $\sigma = 2Ns$

**(6) Compare with the Moran Process with constant selection**

$$M_{WF}(p) \approx N \cdot M_{Moran}(p)$$

$$V_{WF}(p) \approx \frac{N}{2} \cdot V_{Moran}(p)$$

$$P_{WF}\left(\frac{1}{N}\right) = \frac{1 - e^{-2s}}{1 - e^{-2Ns}} \quad , \quad P_{Moran}\left(\frac{1}{N}\right) = \frac{1 - e^{-s}}{1 - e^{-Ns}}$$

$$\approx 2s \qquad\qquad\qquad\qquad \approx s$$

## ⑤ Kolmogorov Backward Equation

$$\frac{\partial \psi(p,t|p_0)}{\partial t} = M(p_0) \frac{\partial \psi(p,t|p_0)}{\partial p_0} + \frac{1}{2} V(p_0) \frac{\partial^2 \psi(p,t|p_0)}{\partial p_0^2}$$

- ### Derivation :

Let $\psi := \psi(p,t|p_0) :=$ probability density that $A_1$ has allele frequency at time $t$, given that the allele frequency was $p_0$ at time $0$.

$$\psi(p, t+dt | p_0) = \int \underbrace{\psi(p,t|p_0+\varepsilon)}_{} \underbrace{g(p_0, \varepsilon; dt)}_{} d\varepsilon$$

       probability that $A_1$      probability to go from
       has frequency $p$ at time     $p_0$ to $p_0+\varepsilon$ in the interval $dt$
       $t+dt$ given that its
       frequency is $p_0+\varepsilon$ at time $dt$

$$= \int \left[ \psi(p,t|p_0) + \varepsilon \frac{\partial \psi}{\partial p_0} + \frac{\varepsilon^2}{2} \frac{\partial^2 \psi}{\partial p_0^2} + \frac{\varepsilon^3}{6} \frac{\partial^3 \psi}{\partial p_0^3} + \cdots \right] g \, d\varepsilon \quad \text{by Taylor's expansion}$$

$$\approx \int \left[ \psi g + \varepsilon g \frac{\partial \psi}{\partial p_0} + \frac{1}{2} \varepsilon^2 g \frac{\partial^2 \psi}{\partial p_0^2} \right] d\varepsilon \quad \text{assuming } \varepsilon^2 >> \varepsilon^3$$

$$= \psi \underbrace{\int g \, d\varepsilon}_{=1} + \frac{\partial \psi}{\partial p_0} \underbrace{\int \varepsilon g \, d\varepsilon}_{M(p_0)\,dt} + \frac{1}{2} \frac{\partial^2 \psi}{\partial p_0^2} \underbrace{\int \varepsilon^2 g \, d\varepsilon}_{V(p_0)\,dt}$$

$$= \psi(p,t|p_0) + \left( M(p_0) \frac{\partial \psi(p,t|p_0)}{\partial p_0} + \frac{1}{2} V(p_0) \frac{\partial^2 \psi(p,t|p_0)}{\partial p_0^2} \right) dt$$

- ### Equilibrium :

$$\frac{\partial \psi^*}{\partial p_0} = C \exp\left( -\int_0^p \frac{2M(q)}{V(q)} dq \right)$$

## ⑥ Fixation probabilities :

$P(p_0) = \psi(1, \infty | p_0) :=$ the fixation probability of $A_1$ given its initial frequency $p_0$

$\Rightarrow$ 2 absorbing states :

$$P(0) = \psi(1, \infty | 0) = 0 \quad , \quad P(1) = \psi(1, \infty | 1) = 1$$

$$\Rightarrow \quad P(p_0) = \frac{\int_0^{p_0} \exp\left( -\int_0^p \frac{2M(q)}{V(q)} dq \right) dp}{\int_0^1 \exp\left( -\int_0^p \frac{2M(q)}{V(q)} dq \right) dp}$$

- ### Wright-Fisher Process with Constant Selection :

With $w_1 = 1+s$ and $w_2 = 1$,

$$\bar{w} = 1+sp \quad , \quad \frac{d\bar{w}}{dp} = s$$

$$\Rightarrow \quad M(p) = p(1-p) \frac{1}{\bar{w}} \cdot \frac{d\bar{w}}{dp} = \frac{sp(1-p)}{1+sp}$$

$$\Rightarrow \quad \frac{2M(p)}{V(p)} = \frac{2sp(1-p)}{1+sp} \Big/ (Np(1-p)) = \frac{2Ns}{1+sp}$$

$$\approx 2Ns \quad \text{for small } p$$

$$\Rightarrow \int_0^p \frac{2M(q)}{V(q)} dq = 2Nsp$$

$$\Rightarrow P(p_0) = \frac{1-e^{-2Nsp_0}}{1-e^{-2Ns}}$$

$$\lim_{s \to 0} P(p_0) = p_0$$

$$P(\tfrac{1}{N}) = \frac{1-e^{-2s}}{1-e^{-2Ns}} \Rightarrow P(\tfrac{1}{N}) \approx 2s$$

for large $N$, small $s$

# 1) Mean Fixation Time

Let $T(p_0)$ := expected waiting time until the fixation of $A_1$, given that it will be fixed

and its initial frequency $p_0$

For the neutral Wright-Fisher process:

$$T(\tfrac{1}{N}) \approx 2N$$

If the mutant has an selective advantage/disadvantage $s$, the larger $s$ is, the smaller $T(\tfrac{1}{N})$ $\Rightarrow$ shorter waiting time for selected mutant

# Chapter 7 Evolutionary Game Theory

## ① Frequency-Dependent Selection (Assuming infinite population)

- Idea: a selective advantage may decrease with increasing population size
- Consider 2 types $A$ and $B$ with frequencies

$$x_A(t) \quad \text{and} \quad x_B(t) \quad , \quad x(t) = (x_A(t), x_B(t))^T$$

and fitness

$$f_A(x(t)) \quad \text{and} \quad f_B(x(t))$$

- The average fitness is

$$\phi(x) = x_A f_A(x) + x_B f_B(x)$$

and the system is

$$\begin{cases} \dot{x}_A = x_A (f_A(x) - \phi(x)) \\ \dot{x}_B = x_B (f_B(x) - \phi(x)) \end{cases}$$

With $x := x_A$ and $1 - x := x_B$, the system is equivalent as

$$\dot{x} = x(1-x)(f_A(x) - f_B(x))$$

- Equilibrium points:

(1) $x^* = 0$ : stable if $f_A(0) < f_B(0)$

(2) $x^* = 1$ : stable if $f_A(1) > f_B(1)$

(3) $x^* \mid x^* \in (0,1), f_A(x^*) = f_B(x^*)$ : stable if

$$\frac{\partial f_A(x^*)}{\partial x} < \frac{\partial f_B(x^*)}{\partial x} \quad \Longleftrightarrow \quad \frac{\partial}{\partial x}(f_A(x^*) - f_B(x^*)) < 0$$

negative slope

## ② Evolutionary Game Theory

- **Definition**: study of frequency-dependent selection

- **Setup**:

(1) A population of players / individuals

(2) Fixed strategy

(3) Random interaction

(4) Goal is to increase the reproductive success / fitness

- **Two-Player Game**

$$\begin{array}{c c} & \begin{array}{cc} A & B \end{array} \\ \begin{array}{c} A \\ B \end{array} & \left( \begin{array}{cc} a & b \\ c & d \end{array} \right) \end{array}$$

← A gets payoff
← B gets payoff

↑ Playing against A    ↑ Playing against B

$$\Rightarrow \dot{x} = x(1-x)[(a-b-c+d)x + (b-d)]$$

$$= \underbrace{ax + b(1-x)}_{f_A(x)} - \underbrace{cx - d(1-x)}_{-f_B(x)}$$

· Strategies :

(1) $a > c$, $b > d$ : A is always the best strategy $\quad$ (A Nash)

(2) $a < c$, $b < d$ : B is always the best strategy $\quad$ (B Nash)

(3) $a > c$, $b < d$ : Playing the same strategy as the opponent is always the best

$$x^* = \frac{d-b}{a-c-b+d} \quad \text{and} \quad \frac{\partial}{\partial x}(f_A(x^*) - f_B(x^*)) = (a-c)+(d-b) > 0$$
$\downarrow$
unstable $\qquad\qquad\qquad\qquad\qquad$ (Both Nash)

(4) $a < c$, $b > d$ : playing the opposite strategy as the opponent is always the best

$$x^* = \frac{d-b}{a-c-b+d} \quad \text{and} \quad \frac{\partial}{\partial x}(f_A(x^*) - f_B(x^*)) = (a-c)+(d-b) < 0$$
$\downarrow$
stable $\qquad\qquad\qquad\qquad\qquad\quad$ (Both not Nash)

(5) $a = c$, $b = d$ : same payoffs for all strategies

③ Nash Equilibrium :

If two players play the same strategy and neither player can increase its payoff by changing strategy, then the strategy is at Nash equilibrium.

④ Evolutionary stable strategy (ESS)

We say that the selection of an all-A population will oppose the invasion of an infinitesimally small amount $\varepsilon$ of B if

$$f_A(1-\varepsilon) > f_B(\varepsilon)$$

$\Leftrightarrow$ $\quad a(1-\varepsilon) + b\varepsilon > c(1-\varepsilon) + d\varepsilon$

A is ESS if either $\qquad a > c \qquad$ or $\qquad a = c$ and $b > d$
$\qquad\qquad\qquad\qquad\qquad \downarrow \qquad\qquad\qquad\qquad \downarrow$
$\qquad$ A strictly Nash $\qquad$ A Nash $\qquad$ B not Nash

⑤ n-Player Evolutionary Game

Suppose there are n strategies.

$$\begin{array}{c} & \begin{array}{cccc} S_1 & S_2 & \cdots & S_n \end{array} \\ \begin{array}{c} S_1 \\ S_2 \\ \vdots \\ S_n \end{array} & \left( \begin{array}{cccc} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right) \end{array}$$

The payoff of playing strategy $S_i$ against strategy $S_j$ is

$$E(S_i, S_j) = a_{ij}$$

(1) <u>Definitions for strategies</u>

$S_k$ is <u>unbeatable</u> if $\forall i \neq k$, $\quad a_{kk} > a_{ik}$ and $\quad a_{ki} > a_{ii}$

$\Downarrow$

$S_k$ is <u>strict Nash</u> if $\forall i \neq k$, $\quad a_{kk} > a_{ik}$

$\Downarrow$

$S_k$ is <u>ESS</u> if $\forall i \neq k$, either $a_{kk} > a_{ik}$ or $a_{kk} = a_{ik}$ and $a_{ki} > a_{ii}$

$\Downarrow$

$S_k$ is <u>weak ESS</u> if $\forall i \neq k$, either $a_{kk} > a_{ik}$ or $a_{kk} = a_{ik}$ and $a_{ki} \geq a_{ii}$

$\Downarrow$

$S_k$ is <u>Nash</u> if $\forall i \neq k$, $\quad a_{kk} \geq a_{ik}$

\* Only the first four are stable against invasion (Evolutionary Stable)

(2) <u>Fitness = Expected Payoff</u>

$$f_i(x) = f_{S_i}(x) = \sum_{j=1}^{n} x_j \, a_{ij}$$

$$\Rightarrow \quad \phi(x) = \sum_{i=1}^{n} x_i f_i(x)$$

(3) <mark>Replicator Equation</mark>

$$\dot{x}_i = x_i \left( f_i(x) - \phi(x) \right), \quad i = 1, \ldots, n$$

where $\quad x = (x_1, \ldots, x_n)^T$, $\quad \underbrace{x_1 + \cdots + x_n = 1}_{\text{simplex } S_n}$

$\hookrightarrow$ the interior of $S_n$ and each face are invariant
$\hookrightarrow$ vertices of $S_n$ are fixed points

(4) <u>3-Player Game / RPC Game</u>

$$A = \begin{pmatrix} 0 & -a_2 & b_3 \\ b_1 & 0 & -a_3 \\ -a_1 & b_2 & 0 \end{pmatrix}$$

- <u>det(A) > 0</u> : unique interior equilibrium, globally stable (damped oscillations)
- <u>det(A) < 0</u> : unique interior equilibrium, unstable (increasing oscillations)
- <u>Special case</u>: zero-sum game

$$A = \begin{pmatrix} 0 & -1 & 1 \\ 1 & 0 & -1 \\ -1 & 1 & 0 \end{pmatrix}$$

$\hookrightarrow$ average fitness is always zero $\Rightarrow \dot{x}_i = x_i f_i(x)$

(5) <u>$n \geq 4$</u> : allows for limit cycles and chaotic attractors

$\hookrightarrow$ at most one isolated equilibrium in the interior : $f_1 = \cdots = f_n$, $x_1 + \cdots + x_n = 1$

## 6) Hawks and Doves

Consider 2 strategies : Hawks (H) escalate fights / Doves (D) retreat.

Benifit of winning $= b$   vs.   Cost of injury $= c$

$$
\begin{array}{c}
 & \begin{array}{cc} H & D \end{array} \\
\begin{array}{c} H \\ D \end{array} &
\left( \begin{array}{cc} \frac{b-c}{2} & b \\ 0 & \frac{b}{2} \end{array} \right)
\end{array}
$$

$\Rightarrow$ __Replicator equation__ :

$$
\dot{x}_H = x_H(1-x_H)\left[ \frac{b-c}{2} x_H + b(1-x_H) - \frac{b}{2}(1-x_H) \right]
$$

$$
= x_H(1-x_H)\left( -\frac{c}{2}x_H + \frac{b}{2} \right)
$$

$\Rightarrow$ __Interior equilibrium__ :

$$
x_H^* = \frac{b}{c}
$$

$\hookrightarrow$ which is stable if $b < c$ $\Rightarrow$ hawks and doves can coexist

## · Mixed Strategies

Consider a strategy that plays H with prob. $p$ and D with prob. $(1-p)$.

The payoff of playing strategy $P_1$ against strategy $P_2$ is

$$
E[P_1, P_2] = P_1 P_2\, E[H,H] + P_1(1-P_2)\, E[H,D] + (1-P_1)P_2\, E[D,H] + (1-P_1)(1-P_2)\, E[D,D]
$$

$$
= P_1 P_2 \frac{b-c}{2} + P_1(1-P_2) b + (1-P_1)(1-P_2) \frac{b}{2}
$$

$$
= \frac{b}{2}P_1 P_2 - \frac{b}{2}\cdot\frac{c}{b}P_1 P_2 + \frac{b}{2}\cdot 2P_1 - \frac{b}{2}\cdot 2 P_1 P_2 + \frac{b}{2} - \frac{b}{2}P_1 - \frac{b}{2}P_2 + \frac{b}{2}P_1 P_2
$$

$$
= \frac{b}{2}\left( 1 + P_1 - P_2 - \frac{c}{b}P_1 P_2 \right)
$$

The strategy $p^* = \frac{b}{c}$ is evolutionary stable since

$$
E[p^*, p^*] = \frac{b}{2}\left( 1 - \frac{b}{c} \right)
$$

$$
E[p^*, p] = \frac{b}{2}\left( 1 + \frac{b}{c} - 2p \right)
$$

$$
E[p, p^*] = \frac{b}{2}\left( 1 - \frac{b}{c} \right)
$$

$$
E[p, p] = \frac{b}{2}\left( 1 - \frac{c}{b}p^2 \right)
$$

$\Rightarrow$
$$
E[p^*, p^*] = E[p, p^*]
$$
$$
\Rightarrow E[p^*, p] > E[p, p] \quad \forall p \neq p^*.
$$

$\Rightarrow p^*$ is Nash, weak ESS, ESS but not strictly Nash or unbeatable

· If we consider only pure strategies, then there may not be a Nash equilibrium.

But if we consider all mixed strategies, there is always a Nash equilibrium.

# 1) Prisoner's Dilemma

- **2 strategies**: (1) Cooperation (C), (2) Defection (D)

$$\begin{array}{c} \\ C \\ D \end{array} \begin{pmatrix} \phantom{C} C & \phantom{D} D \\ R & S \\ T & P \end{pmatrix}$$

$$T > R > P > S \quad \text{and} \quad R > \frac{T+P}{2}$$

- CC: Reward for mutual cooperation
- DC: Temptation to defect
- DD: Punishment for mutual defection
- CD: Sucker's payoff

- **Direct Reciprocity**: the game is repeated $m$ times

(1) **GRIM vs ALLD**:

- GRIM: cooperate until the opponent defects
- ALLD: always defect

$$\begin{array}{c} \\ GRIM \\ ALLD \end{array} \begin{pmatrix} \phantom{GRIM} GRIM & \phantom{ALLD} ALLD \\ mR & S+(m-1)P \\ T+(m-1)P & mP \end{pmatrix}$$

For $m > \frac{T-P}{R-P}$, GRIM and ALLD are both strictly Nash

$\Rightarrow$ Direct Reciprocity can stabilize cooperation but not initiate.

(2) **GRIM vs GRIM\***:

- GRIM\*: given $m$, defect on the $m^{th}$ round

$$\begin{array}{c} \\ GRIM \\ GRIM^* \end{array} \begin{pmatrix} \phantom{GRIM} GRIM & \phantom{GRIM^*} GRIM^* \\ mR & (m-1)R+S \\ (m-1)R+T & (m-1)R+P \end{pmatrix}$$

$\Rightarrow$ GRIM $\to$ GRIM\* $\to$ GRIM\*\* $\to$ $\cdots$ $\to$ ALLD : only Nash equilibrium

- **Variable number of rounds**

(1) $\omega$ := probability that another round will be played

$\Rightarrow$ probability of playing $k$ rounds : $\omega^{k-1}(1-\omega)$

$\Rightarrow$ **expected # of rounds**:

$$\bar{m} = \sum_{k=1}^{\infty} \omega^{k-1}(1-\omega) = \frac{1}{1-\omega}$$

$$\begin{array}{c} \\ GRIM \\ ALLD \end{array} \begin{pmatrix} \phantom{GRIM} GRIM & \phantom{ALLD} ALLD \\ \bar{m}R & S+(\bar{m}-1)P \\ T+(\bar{m}-1)P & \bar{m}P \end{pmatrix}$$

$\Rightarrow$ For $\bar{m} > \frac{T-P}{R-P}$, GRIM is evolutionary stable

$\Rightarrow$ defecting in the last round is no longer possible

(2) TFT vs. ALLD

- TFT: start with cooperation and do what the opponent have done in the previous round

$$
\begin{array}{c}
\phantom{xxxxx} \text{TFT} \phantom{xxxxxxx} \text{ALLD} \\
\begin{array}{c} \text{TFT} \\ \text{ALLD} \end{array}
\left(
\begin{array}{cc}
\bar{m}R & S + (\bar{m}-1)P \\
T + (\bar{m}-1)P & \bar{m}P
\end{array}
\right)
\end{array}
$$

↳ same payoff matrix as GRIM vs. ALLD

↳ For $\bar{m} > \dfrac{T-P}{R-P}$ , TFT is evolutionary stable against ALLD

↳ it can resume cooperation while GRIM cannot

↳ not robust against error, if one player starts to defect, then both defeet → ∞

↳ TFT can be easily invaded by ALLC

(3) TFT vs. ALLC

- ALLC : always cooperate

$$
\begin{array}{c}
\phantom{xxxx} \text{TFT} \phantom{xxxxx} \text{ALLC} \\
\begin{array}{c} \text{TFT} \\ \text{ALLC} \end{array}
\left(
\begin{array}{cc}
\bar{m}R & \bar{m}R \\
\bar{m}R & \bar{m}R
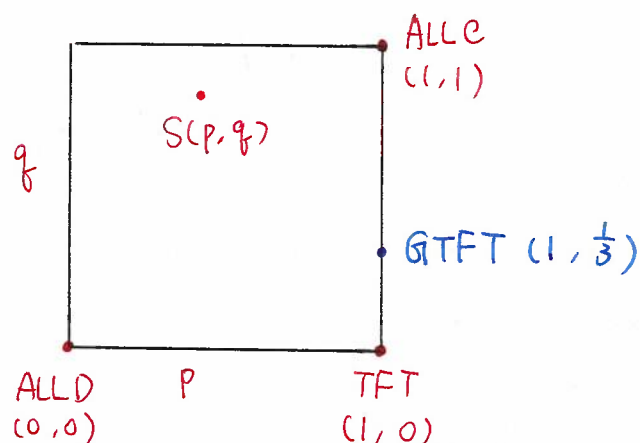\end{array}
\right)
\end{array}
$$

↳ TFT is not evolutionary stable against random drift.

# 3) Reactive Strategies

Strategy $S(p, q)$ cooperates with

(1) prob. $p$ if the opponent cooperated in the previous move

(2) prob. $q$ if the opponent defected in the previous move

- Markov chain for repeated Prisoner's Dilemma

  Consider 2 reactive strategies $S_1(p_1, q_1)$ and $S_2(p_2, q_2)$

  <u>State space</u> : $\{CC, CD, DC, DD\}$

  <u>Transition matrix</u> : $M$

  $$
  \begin{array}{c|cccc}
   & CC & CD & DC & DD \\
  \hline
  CC & p_1 p_2 & p_1(1-p_2) & (1-p_1)p_2 & (1-p_1)(1-p_2) \\
  CD & q_1 p_2 & q_1(1-p_2) & (1-q_1)p_2 & (1-q_1)(1-p_2) \\
  DC & p_1 q_2 & p_1(1-q_2) & (1-p_1)q_2 & (1-p_1)(1-q_2) \\
  DD & q_1 q_2 & q_1(1-q_2) & (1-q_1)q_2 & (1-q_1)(1-q_2)
  \end{array}
  $$

  <u>Probability distribution of the game after $t$ rounds</u> :

  $$x(t) = (x_{CC}(t), x_{CD}(t), x_{DC}(t), x_{DD}(t))$$

  <u>Dynamical system</u> :

  $$x(t+1) = x(t) \cdot M$$

  <u>Payoff at stationary distribution</u> .

  $$E(S_1, S_2) = R\, s_1 s_2 + S\, s_1(1-s_2) + T(1-s_1)s_2 + P(1-s_1)(1-s_2)$$

  where $S_1(p_1, p_2, q_1, q_2)$ and $S_2(p_1, p_2, q_1, q_2)$ are stationary distributions of cooperation

- <u>Generous TFT (GTFT)</u>

  $$GTFT = S(1, q), \quad q = \min\left\{1 - \frac{T-R}{R-S}, \frac{R-P}{T-P}\right\}$$

  $$ALLD \to TFT \dashrightarrow GTFT$$

  because GTFT can correct mistakes stochastically

- <u>Win-Stay Lose-Shift (WSLS)</u>

  Cooperate when CC or DD , defect when CD or DC

  ↳ deterministic corrector for errors

  ↳ WSLS dominates ALLC, whereas GTFT does not .

## 1) Two Players in a finite population

Consider 2 strategies : A , B

Population size : $N$

Payoff matrix :

$$
\begin{array}{c c}
 & \begin{array}{cc} A & \qquad B \end{array} \\
\begin{array}{c} A \\ B \end{array} &
\left( \begin{array}{cc} a & b \\ c & d \end{array} \right)
\end{array}
$$

Let $i :=$ # of A individuals

The expected payoff for A and B :

A : $\qquad F_i = \dfrac{(i-1)\,a + (N-i)\,b}{N-1}$

freq of type B

freq of other type A

B : $\qquad G_i = \dfrac{i\,c + (N-i-1)\,d}{N-1}$

We say that __selection opposes A invading B__ if

$$F_1 < G_1$$

$\Longleftrightarrow \qquad (N-1)\,b < c + (N-2)\,d$ , independent of $a$

For $N=2$, we have simply $b < c$.

$\hookrightarrow$ playing A against B gets smaller payoff than playing B against A

## 2) Intensity of Selection

Let $w :=$ intensity of selection.

Define the __frequency-dependent fitness__ as

$$f_i = 1 - w + w F_i$$
$$g_i = 1 - w + w G_i$$

$\hookrightarrow$ if $w=0$, no force of selection / the game does not contribute

$\hookrightarrow$ if $w=1$, fitness is determined entirely by the payoff

$\hookrightarrow$ __weak selection__ : send $w \to 0$

$\hookrightarrow$ $w$ gets cancelled out in the deterministic replicator equation

but is important in the stochastic process in finite population

## 3) Moran Process with Constant Selection

Replace the original fitness values by frequency-dependent fitness:

$$\underbrace{P_{i,i+1} = \frac{i f_i}{i f_i + (N-i) g_i}}_{P(T_A < T_B)} \cdot \underbrace{\frac{N-i}{N}}_{\text{choose a } B \text{ to die}}$$

$$\underbrace{P_{i,i-1} = \frac{(N-i) g_i}{i f_i + (N-i) g_i}}_{P(T_A > T_B)} \cdot \underbrace{\frac{i}{N}}_{\text{choose an } A \text{ to die}}$$

$$P_{i,i} = 1 - P_{i,i+1} - P_{i,i-1}$$

and $P_{0,0} = P_{N,N} = 1$.

· **Fixation Probability**

$$\gamma_i = \frac{P_{i,i-1}}{P_{i,i+1}} \implies P_A = \frac{1}{1 + \sum\limits_{j=1}^{N-1} \prod\limits_{k=1}^{j} \frac{g_i}{f_i}} = \frac{1}{1 + \sum\limits_{j=1}^{N-1} \prod\limits_{k=1}^{j} \frac{1 - w + w G_i}{1 - w + w F_i}}$$

## ④ Weak Selection Limit

As $w \to 0$, $P_A \approx \frac{1}{N} \frac{1}{1 - (\alpha N - \beta) w / 6}$

where $\alpha = a + 2b - c - 2d$, $\beta = 2a + b + c - 4d$

↳ selection <span style="color:red">favors</span> the fixation of $A$ if

$$P_A > \frac{1}{N} \iff \alpha N > \beta \iff a(N-2) + b(2N-1) > c(N+1) + d(2N-4)$$

↳ For $N = 2$, we have

$$P_A > \frac{1}{N} \iff b > c \iff F_1 > G_1$$

· **Weak Selection and Large Population**

As $N \to \infty$, we have

$$P_A > \frac{1}{N} \iff a + 2b > c + 2d \iff a - c > 2(d - b)$$

In the case where $a > c$, $b < d$:

· Both $A$ and $B$ are Nash. <span style="color:red">Whichever has higher frequency gets higher fitness.</span>

· The unstable interior equilibrium:

$$x^* = \frac{d - b}{a - b - c + d} = \frac{d - b}{(a - c) + (d - b)}$$

seems contradictory

$\Rightarrow$ The $\frac{1}{3}$ Law:

$$P_A > \frac{1}{N} \iff x^* < \frac{1}{3} \leftarrow \text{-----}$$

* it can happen that selection opposes the invasion of $A$ but favours the fixation of $A$

## 5) Evolutionary Stability in Finite Population

B is ESS$_N$ if

$$\begin{array}{c} & \begin{array}{cc} A & B \end{array} \\ \begin{array}{c} A \\ B \end{array} & \begin{pmatrix} a & b \\ c & d \end{pmatrix} \end{array}$$

(1) Selection protects against invasion:

$$F_1 < G_1, \qquad \Longleftrightarrow \qquad (N-1)\,b < C + (N-2)\,d$$

(2) Selection protects against replacement:

$$P_A < \frac{1}{N} \qquad \Longleftrightarrow \qquad a(N-2) + b(2N-1) < C(N+1) + d(2N-4)$$
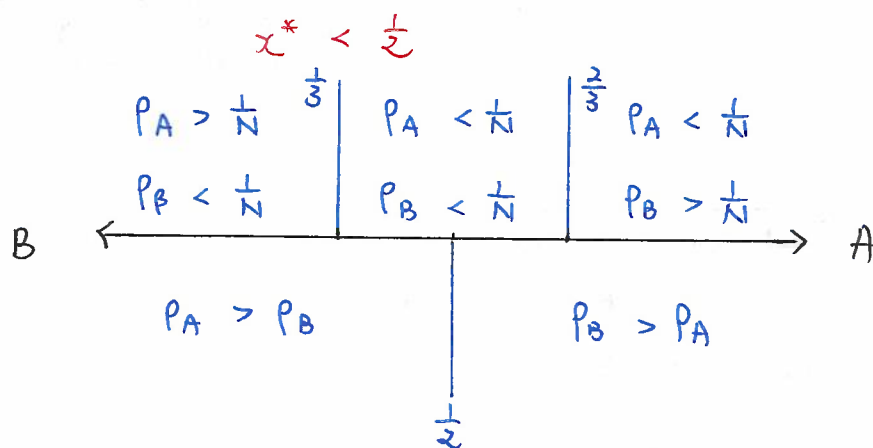
↳ For $N=2$, B is ESS$_N$ if $b < c$, ESS is neither necessary nor sufficient

↳ For large $N$, B is ESS$_N$ if $b < d$ and $x^* > \frac{1}{3}$, ESS is necessary but not sufficient.

## 6) Risk Dominance

A is risk dominant over B if $P_A > P_B$.

For large $N$, we have



## 7) Prisoner's Dilemma in Finite Population

- In the case of an infinite population, both TFT and ALLD are stable against invasion by each other if $m > \dfrac{T-P}{R-P}$

- In the finite population, selection favors TFT to replace ALLD if $P_{TFT} > \frac{1}{N}$

# Chapter 8 Evolutionary Graph Theory
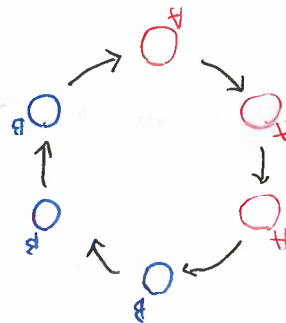
## ① Evolutionary graph

- $G = (V, E)$ where $V = \{1, 2, \ldots, N\}$ are individuals
  and $i \to j \in E$ means offspring of $i$ can replace $j$

- In each generation, one individual $i \in \{1, 2, \ldots, N\}$ is chosen randomly
  and its offspring replaces $j$ with probability $w_{ij}$.
  The matrix $W = (w_{ij})$ is a stochastic matrix

- **Moran process**

  $$w_{ij} = \frac{1}{N} \quad \forall \, i,j, \qquad \rho = \frac{1 - 1/r}{1 - 1/r^N}$$

  → a complete graph with identical weights

- **Markov chain on a directed cycle:**

$$
\begin{pmatrix}
0 & 0 & 0 & 0 & 0 & 1 \\
1 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 \\
0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0
\end{pmatrix}
$$

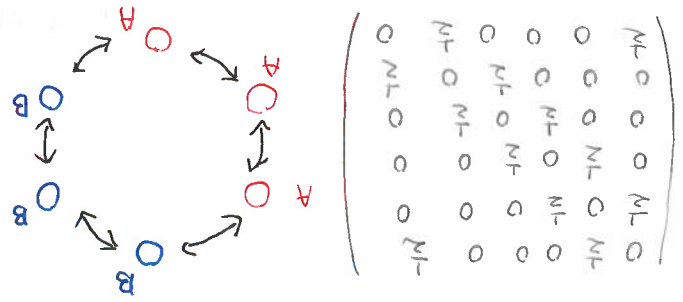→ Starting from one B mutant, only one connected cluster of B can emerge

→ Let $m :=$ # of B individuals with fitness $r$. A has fitness 1.

$$P_{m, m+1} = \frac{r}{N - m + rm}$$

$$P_{m, m-1} = \frac{1}{N - m + rm}$$

$$\Rightarrow \quad \gamma_m = \frac{1}{r}$$

$$\Rightarrow \quad \rho = \frac{1 - 1/r}{1 - 1/r^N} \qquad (\text{same as the Moran process})$$

- **Bidirected Cycle**

$$
\begin{pmatrix}
0 & \tfrac{1}{2} & 0 & 0 & 0 & \tfrac{1}{2} \\
\tfrac{1}{2} & 0 & \tfrac{1}{2} & 0 & 0 & 0 \\
0 & \tfrac{1}{2} & 0 & \tfrac{1}{2} & 0 & 0 \\
0 & 0 & \tfrac{1}{2} & 0 & \tfrac{1}{2} & 0 \\
0 & 0 & 0 & \tfrac{1}{2} & 0 & \tfrac{1}{2} \\
\tfrac{1}{2} & 0 & 0 & 0 & \tfrac{1}{2} & 0
\end{pmatrix}
$$

$$\Rightarrow \quad \rho = \frac{1 - 1/r}{1 - 1/r^N}$$

→ Now, each cell has 2 directions to place offspring
but only one direction will increase its frequency.

· **Linear Process**

$$O \leftarrow \underset{\text{stem cell}}{O} \leftarrow O \leftarrow O \leftarrow O \leftarrow O$$

→ expected fraction of a randomly placed mutant:

$$\rho = \frac{1}{N}$$

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

· **Burst**



→ again, $\rho = \frac{1}{N}$

$$\begin{pmatrix} 0 & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

## ② Amplifiers and Suppressors

· A graph $G$ is an **amplifier** of selection if for $r > 1$

$$\rho_G > \rho_{moran}$$

· A graph $G$ is a **suppressor** of selection if for $r < 1$

$$\rho_G < \rho_{moran}$$

· The strongest suppressor of selection is independent of $r$ :

$$\rho_G = \frac{1}{N}$$

## ③ Isothermal Theorem

· The temperature of a vertex $j$ :

$$T_j = w_{ij} + \cdots + w_{nj}$$

· Hot : many incoming edges → **Isothermal** := all vertices have the same temperature

· $\rho_G = \rho_{moran}$ ⟺ $G$ is isothermal ⟺ $G$ is doubly stochastic

· Isothermal graphs

  i) Directed cycle
  ii) Cycle
  a) All symmetric graphs : $w_{ij} = w_{ji}$ $\forall i,j$

· All one-rooted graphs : $\rho = \frac{1}{N}$
· **Multiple-rooted** graphs prevent **fixation**
· Star (Bidirected burst) is an amplifier of selection
· Superstar is a strong amplifier of selection

## ④ Games on evolutionary graph (slides)

# Chapter 9 Spatial Models of the Evolution of Solid Tumours

① Four processes of population genetics

　(1) selection　(2) mutation　(3) drift　(4) gene flow
　　　　　　　　　　　　　　　　　　　　　↑
　　　　　　　　　　　　　　　　　spatial structure

② Cellular Automata

- regular grid of sites, which are associated with sets of states
- each site is in a neighbourhood
- rules: depend on current state of site and its neighbourhood
  ↳ can be probabilistic

③ Eden Growth Model

- 2 states: unoccupied ($S_0$), occupied ($S_1$)
- von Neumann neighbourhood: adjacent sites



- <u>Available site-focused rule</u>: (smooth boundary
  randomly choose an $S_0$ site that adjoins at least one $S_1$ site and switch to $S_1$
- <u>Bond-focused rule</u>:
  randomly choose an $S_1$ site with probability proportional to the number of $S_0$ sites
  then randomly choose an $S_0$ neighbour and switch to $S_1$　(rough boundary)
- <u>Cell-focused rule</u>:
  randomly choose an $S_1$ site that adjoins at least one $S_0$ site, then randomly
  choose an $S_0$ neighbour and switch to $S_1$　(smooth boundary)

  ↳ all resemble a disc in 2D and a ball in 3D in the long run

④ Eden Growth Model with Mutation

- Multiple occupied states $\{S_1, S_2, \ldots\}$
- Mutation rates: $S_i \to S_j$
- Mutation at division
- Neutral model: a disc
- With selection: e.g. $S_i : (1+s)^i$ ⇒ new mutations grow faster

## 5) Eden Growth Model w̄ Cell Death and Migration

- Cell death facilitates selection
- Migration facilitates growth by increasing the surface-to-volume ratio

## 6) Deme-Based Models

- each deme contains multiple cells
- assume cells within demes are well-mixed
- cells can migrate between demes
- When each deme contains only one cell, back to cellular automaton
- As migration rate → 0, reduce to a set of independent non-spatial processes

## 7) Spatial Moran Model

- <u>Assumptions</u>:

  1) Offspring of one cell replaces another cell
  2) replacement is chosen with prob. ∝ cell fitness (death-birth model)
  3) prob. of being replaced by local parent = $1 - m$

     " " " " " neighbouring " = $m$ = dispersion probability

- Consider a mutant invading an infinite row of demes

  Let $n = \{\cdots, n_{i-1}, n_i, n_{i+1}, \cdots\}$ be the vector of mutant population sizes

  Let $\mu :=$ death rate, $s :=$ fitness advantage, $N :=$ deme population size

  - <u>Probability that the number of mutants $n_i$ in deme $i$ increases by one:</u> (decreases)

$$W_i^+(n) = \underbrace{\mu(N - n_i)}_{\substack{\text{prob. that a} \\ \text{WT cell dies}}} \times \underbrace{(1+s)}_{\substack{\text{fitness} \\ \text{of mutant}}} \times \underbrace{\left((1-m) \times \frac{n_i}{N} + \frac{m}{2} \times \frac{n_{i-1}}{N} + \frac{m}{2} \times \frac{n_{i+1}}{N}\right)}_{\text{replacement + dispersion}}$$

$$W_i^-(n) = \mu n_i \times \left((1-m)\frac{N-n_i}{N} + \frac{m}{2} \times \frac{N-n_{i-1}}{N} + \frac{m}{2} \times \frac{N-n_{i+1}}{N}\right)$$

  - <u>Diffusion Approximation / Fisher's Equation</u>

$$\frac{\partial u}{\partial t} = \underbrace{D[1 + s(1-u)]\frac{\partial^2 u}{\partial x^2}}_{\substack{\text{speed of mutant spreading} \\ \text{through space}}} + \underbrace{\mu s u(1-u)}_{\substack{\text{growth of a mutant} \\ \text{within a deme}}} \qquad \propto D\frac{\partial^2 u}{\partial x^2} + ru(1-u)$$

  ↳ $u = \frac{\langle n_i \rangle}{N}$, $x \approx li$ (distance along the row of demes)

  ↳ approximation only works when $N$ is large and $s$ is large

  ↳ should also consider clonal interference, environmental heterogeneity

# Chapter 10  Branching Processes in Biology

## ① Galton-Watson Process

- **Process**: A *single* ancestor lives for *one* unit of time, after which it produces a *random* number of offspring $Z \sim P$ fixed. Each offspring i.i.d. $\sim P$

  Let $Z_n :=$ # of individuals in generation $n$. $Z_0 = 1$ and $Z_1 = Z$.

  The Galton-Watson process is
  $$\{ Z_n \mid n = 0, 1, 2, \ldots \}$$

- defined on the nonnegative integers

  $$\Rightarrow \quad Z_{n+1} = \left( \sum_{j=1}^{Z_1} Z_n^{(j)} \right) \longrightarrow$$ # of offspring that ancestor $j$ in generation 1 produces in generation $n$

  $\downarrow$
  random variable

- **Transition Probabilities**:

  Let $$P_k = \text{Prob}(Z = k), \quad k \in \mathbb{N}$$

  $$\Rightarrow \quad P(i,j) = \text{Prob}(Z_{n+1} = j \mid Z_n = i)$$

  $$\Rightarrow \quad P(1,k) = P_k, \quad k \in \mathbb{N}$$

  **In general**:
  $$P(0,j) = \delta_{0j} = \begin{cases} 1 & \text{if } j = 0 \\ 0 & \text{else} \end{cases}$$

  For $i \geq 1$:
  $$P(i,j) = P_j^{*i} = \sum_{k_1 + k_2 + \cdots + k_i = j} P_{k_1} \cdots P_{k_i}$$

  $\underbrace{\phantom{k_1 + k_2 + \cdots + k_i = j}}$
  all combinations of producing $j$ offsprings

  $\hookrightarrow \{ P_k^{*i} \}_{k \geq 0} := i$-fold convolution of $\{ P_k \}_{k \geq 0}$

## ② Probability Generating Functions

For $Z \sim \{ P_k \}_{k \geq 0}$, the pgf of $Z$ is
$$f(s) = \mathbb{E}[s^Z] = \sum_{k=0}^{\infty} P_k s^k, \quad s \in [0,1] \quad \Rightarrow f(1) = \sum_{k=0}^{\infty} P_k = 1$$

The pgf generates the distribution $P$ through derivatives:
$$\frac{d^k f}{ds^k}(0) = k! \, P_k, \quad k \geq 0$$

(1) **Moments of $Z$**:
$$f'(s) = \sum_{k=0}^{\infty} k P_k s^{k-1}, \quad f''(s) = \sum_{k=0}^{\infty} k(k-1) P_k s^{k-2} = \sum_{k=0}^{\infty} k^2 P_k s^{k-2} - \sum_{k=0}^{\infty} k P_k s^{k-2}$$

$$\Rightarrow \quad \mathbb{E}[Z] = f'(1), \quad \text{Var}(Z) = \mathbb{E}[Z^2] - (\mathbb{E}[Z])^2 = f'(1) + f''(1) - f'(1)^2$$

(2) <u>Powers of $f$</u>:

$$f(s) = \sum_{j=0}^{\infty} P_j \, s^j = \sum_{j=0}^{\infty} P(1,j) \, s^j$$

$$\Rightarrow [f(s)]^k = \sum_{j=0}^{\infty} P(k,j) \, s^j \quad , \quad k \geq 1$$

(3) <u>Iterations</u>:

$$f^{(0)}(s) = s$$
$$f^{(1)}(s) = f(s)$$
$$f^{(n+1)}(s) = f(f^{(n)}(s)) \quad , \quad n \geq 1$$

(4) <u>n-Step Transitions</u>

$P_n(i,j) :=$ transition $\underset{\wedge}{\text{from}}$ $i$ individuals to $j$ individuals in $n$ steps
$\phantom{P_n(i,j) :=}$ probability

(5) <u>Chapman-Kolmogorov Equation</u>

$$P_{n+m}(i,j) = \sum_{k=0}^{\infty} P_n(i,k) \, P_m(k,j)$$

(6) <u>Proposition : $f_n = f^{(n)}$</u>

Let $f_n$ be the pgf of $Z_n$. Then $f_n$ is equivalent as applying the pgf of $Z$ for $n$ times.

<u>proof</u>:
$$f_{n+1}(s) = \sum_j P_{n+1}(1,j) \, s^j \qquad \text{by definition of } f_{n+1}$$

$$= \sum_j \sum_k P_n(1,k) P(k,j) \, s^j \qquad \text{by Chapman-Kolmogorov equation}$$

$$= \sum_k P_n(1,k) \sum_j P(k,j) \, s^j \qquad \text{by rearranging}$$

$$= \sum_k P_n(1,k) \, [f(s)]^k \qquad \text{by definition of power of } f$$

$$= f_n(f(s)) \qquad \text{by definition of } f_n$$

$$\vdots$$

$$= f^{(n+1)}(s)$$

7) <u>Moments of $Z_n$</u>

Assume that $P_0 + P_1 < 1$ and $P_j \neq 1$ for all $j$
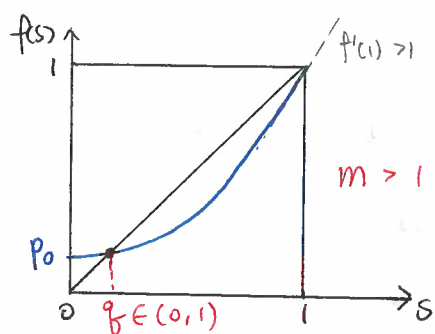Set $m = \mathbb{E}[Z]$ and $\sigma^2 = Var[Z]$. Then

$$\mathbb{E}[Z_n] = m^n \quad , \quad Var[Z_n] = \begin{cases} \dfrac{\sigma^2 m^{n-1}(m^n - 1)}{m-1} & \text{if } m \neq 1 \\[2mm] n\sigma^2 & \text{if } m = 1 \end{cases}$$

$\downarrow$

exponential in expected # of offsprings

## ③ Extinction

Given that $Z_n = 0$ is an absorbing state, the probability of extinction is

$$P = \text{Prob}(Z_i = 0 \text{ for some } i \geq 0)$$
$$= \lim_{n \to \infty} \text{Prob}(Z_i = 0 \text{ for some } 1 \leq i \leq n)$$
$$= \lim_{n \to \infty} \text{Prob}(Z_n = 0)$$
$$= \lim_{n \to \infty} f_n(0) = \lim_{n \to \infty} f^{(n)}(0)$$



- ## Theorem

The extinction probability of the Galton-Watson process $\{Z_n\}$ is the smallest non-negative root $q$ of the equation $f(s) = s$.
If $E[Z] = m \leq 1$, then $q = 1$. If $m > 1$, then $q < 1$.

- ## Criticality

(1) Supercritical

$$m > 1, \quad E[Z_n] \to \infty, \quad q < 1$$

(2) Critical

$$m = 1, \quad E[Z_n] = 1, \quad q = 1$$

(3) Subcritical

$$m < 1, \quad E[Z_n] \to 0, \quad q = 1$$

- ## Instability :

$$\lim_{n \to \infty} \text{Prob}(Z_n = k) = 0, \quad k \geq 1$$

$$\text{Prob}(\lim_{n \to \infty} Z_n = 0) = q$$

$$\text{Prob}(\lim_{n \to \infty} Z_n = \infty) = 1 - q$$

Interpretation: In the Galton-Watson process, the population cannot stay in any state indefinitely. It will either go extinct with probability $q$, or keep on growing forever with probability $(1-q)$ $\Rightarrow$ useful for small population

# 4) Multi-Type Galton Watson Process

- Consider 2 types : type 0 (wild type) , type 1 (mutant)
  with counts $Z_0(t)$ and $Z_1(t)$, $t \in \{0, 1, 2, \dots\}$

- Each cell gives 2 offsprings.

- Each offspring of a type 0 cell can mutate to type 1 with rate $\alpha$, irreversible

- **Probability generating functions**

$$F = (F_0, F_1)$$

$$F_0(s_0, s_1 ; t) = E\left[ s_0^{Z_0(t)} s_1^{Z_1(t)} \mid Z_0(0) = 1 , Z_1(0) = 0 \right] =: F_0(s; t)$$

$$F_1(s_0, s_1 ; t) = E\left[ s_0^{Z_0(t)} s_1^{Z_1(t)} \mid Z_0(0) = 0 , Z_1(0) = 1 \right] := F_1(s; t)$$

$$\Rightarrow F_0(s; t) = \left[ (1-\alpha) F_0(s; t-1) + \alpha F_1(s; t-1) \right]^2$$

$$F_1(s; t) = \left[ F_1(s; t-1) \right]^2$$

$$\Rightarrow E[Z_0(t) \mid Z_i(0) = \delta_{1i}] = 2\, E[Z_0(t-1) \mid Z_i(0) = \delta_{1i}] = 0$$

$$E[Z_0(t) \mid Z_i(0) = \delta_{0i}] = 2(1-\alpha) E[Z_0(t-1) \mid Z_i(0) = \delta_{0i}]$$

$$\Rightarrow E[Z_0(t) \mid Z_i(0) = \delta_{0i}] = [2(1-\alpha)]^t = \text{expected number of WT at time } t$$

## Expected total # of cells

$$N(t) = E[Z_0(t) + Z_1(t) \mid Z_i(0) = \delta_{0i}] = 2^t$$

## Expected # of mutant cells

$$r(t) = E[Z_1(t) \mid Z_i(0) = \delta_{0i}] = 2^t - (2(1-\alpha))^t = 2^t(1 - (1-\alpha)^t)$$

## Probability of a mutant-free population

$$P_0(t) = F_0(1, 0 ; t) = E\left[ 1^{Z_0(t)} 0^{Z_1(t)} \mid Z_i(0) = \delta_{0i} \right] = E\left[ \mathbb{1}_{\{Z_1(t) = 0\}} \mid Z_i(0) = \delta_{0i} \right]$$

$$P_1(t) = F_1(1, 0 ; t) = \text{prob. of being WT-free}$$

$$\Rightarrow \text{For } t = 0, 1, 2, \dots$$

$$P_0(t) = (1-\alpha)^{2^{t+1} - 2} , \qquad P_1(t) = 0$$

For each fixed $N$, we have

$$P_0(r) = \left(1 - \frac{r}{N}\right)^{\frac{2(N-1)}{\log_2 N}}$$

$$\downarrow$$

can be measured and evaluated

# Chapter 11 : Evolutionary Escape

## ① Posets and Distributive Lattices

### (1) Binary Sequence Space

↳ "0" : unmutated ; "1" : mutated (irreversible)

↳ Genotype := binary string of length $n$

↳ Wild Type :
$$0 = 00 \cdots 0$$

↳ Escape Type :
$$\underline{1} = 11 \cdots 1$$

### (2) Partially Ordered Sets (Posets) (there are pairs of elements that cannot be compared)

A poset is a set $\mathcal{E}$ together with a binary relation "$\leq$", which is

- reflexive : $\forall e \in \mathcal{E}, \ e \leq e$

- antisymmetric : $\forall e_1, e_2 \in \mathcal{E}$, if $e_1 \leq e_2$ and $e_2 \leq e_1$, then $e_1 = e_2$

- transitive : $\forall e_1, e_2, e_3 \in \mathcal{E}$, if $e_1 \leq e_2$ and $e_2 \leq e_3$, then $e_1 \leq e_3$

↳ Write $e_1 < e_2$ if $e_1 \leq e_2$ and $e_1 \neq e_2$, $e_1 < e_2$ is called a cover relation if there is no $e' \in \mathcal{E}$ s.t. $e_1 < e' < e_2$.

↳ Hasse Diagram : $G = (\mathcal{E}, E)$ and $e_1 \to e_2 \in E \iff e_1 < e_2$ is a cover relation

### (3) Order Ideal

An order ideal, $g$, in a poset $\mathcal{E}$ is a subset of $\mathcal{E}$ that is closed downward
i.e. if $e_2 \in g$ and $e_1 \leq e_2$, then $e_1 \in g$.

### (4) Distributive Lattice

The set of all order ideals of $\mathcal{E}$ forms a distributive lattice $J(\mathcal{E})$ under inclusion.

↳ $(J(\mathcal{E}), \subseteq)$ is a poset

↳ every pair of order ideals $(g_1, g_2)$ has a unique supremum, $g_1 \cup g_2$, and a unique infimum, $g_1 \cap g_2$.

### (5) Genotype Lattice

- Let $\mathcal{E}$ be a set of $n = |\mathcal{E}|$ irreversible genetic events
- Since evolution may not proceed in any random order, the posets $(\mathcal{E}, \leq)$ encode constraints on the order in which the mutations can accumulate.
- The order ideals $g$ of $J(\mathcal{E})$ are the genotype that can evolve subject to order constraints

→ $G = J(\mathcal{E}) = $ Genotype Lattice

## 6) The Empty Poset

The empty poset is defined as the set $\mathcal{E} = \{1, 2, \ldots, n\}$ with no relations

Then, the genotype lattice $\mathcal{G}$ is simply the hypercube. $\Rightarrow$ no constraints
$$\{0, 1\}^n$$

## 7) Chains

A chain in $\mathcal{G} = (J(\mathcal{E}), \subseteq)$ of length $k$ is a collection of $k$ totally ordered subsets

$$g_1 \subset g_2 \subset \cdots \subset g_k$$

one must be a subset of the next

The chains in $\mathcal{G}$ are mutational pathways consistent with the poset $\mathcal{E}$.

## 8) Fitness Landscapes

A fitness landscape is a mapping $f : \mathcal{G} \to \mathbb{R}$

$\hookrightarrow$ for every genotype $g \in \mathcal{G}$, there is a corresponding fitness

## ) Evolutionary Escape

We consider a setup where only the escape type has fitness $f > 1$, whereas other genotypes have fitness $f < 1$. In this case, only the escape type has the chance to replicate in the long run, while others will die out eventually with probability 1.

Our question is : under selective pressure, what is the probability that the wild type reaches the escape state before extinction?

## 11) Mutational Neighbourhood

The mutational neighbourhood of a genotype $g \in \mathcal{G}$ is the set of genotypes $h \in \mathcal{G}$ that can be reached by mutation:

$$N(g) = \{h \in \mathcal{G} : g \subset h\}$$

Let $m = |\mathcal{G}| \leq 2^n$ be the number of mutations / genotypes.

Let $\mu_e$ be the mutation rate of event $e \in \mathcal{E}$.

Assume that mutations are independent and fix a total order of $\mathcal{G}$, we can write the mutation matrix $\underset{m \times m}{U} = (u_{gh})_{g, h \in \mathcal{G}}$ by

$$u_{gh} = \begin{cases} \prod_{e \in h \backslash g} \mu_e & , \text{ if } h \in N(g) \\ 0 & , \text{ otherwise} \end{cases}$$

(2) **k-Step Offspring**

Let $f: G \to \mathbb{R}$ be a fitness landscape and set $\underset{m \times m}{F} = \text{diag}(f)$

$\hookrightarrow$ $(UF)_{g,h} :=$ probability of genotype $g$ producing offspring of type $h$ in one step

$\hookrightarrow$ $(UF)_{g,h}^{k} :=$ probability of genotype $g$ producing offspring of type $h$ along any
mutational pathway of length $k$ in $G$

Define a matrix $B$ such that

$$B = UF + (UF)^2 + \cdots + (UF)^n = (I - UF)^{-1} - I$$

where $B = (bgh)_{g,h \in G}$ :

$$bgh = \begin{cases} u_{gh} \, f(h) \, P_{gh}(f) & , \text{if } g \subset h \\ \\ 0 & , \text{else} \end{cases}$$

where $P_{gh}(f)$ is a polynomial of degree $|h \backslash g| - 1$ on $\mathbb{R}^G$

$\hookrightarrow$ generating function for all chains from $g$ to $h$ in $G$
pathways

(3) **Risk Polynomial**

Consider $g = 0 :=$ wild type and $h = 1 :=$ escape type.
Write $f(g) = f_g \; \forall \; g \in G$ .

Then, the risk Polynomial is defined as

$$R(G, f) := P_{01}(f) = \underbrace{\sum_{0 = g_0 \subset g_1 \subset \cdots \subset g_{k-1} \subset g_k = 1}}_{\text{sum over all chains of length } k} f_{g_1} f_{g_2} \cdots f_{g_{k-1}}$$

(4) **Invasion**

Let $R_g :=$ basic reproductive ratio of an invading pathogen $g$
        $= \#$ of offspring a single individual will produce

We are interested in the case: $R_1 > 1$ and $R_g < 1$ for all $g \neq 1$

Define the fitness landscape by :

$$f_g = \frac{R_g}{1 - R_g} = R_g + R_g^2 + R_g^3 + \cdots$$

$$\iff \quad R_g = \frac{f_g}{1 + f_g}$$

For $g \neq 1$, $f_g \approx R_g$

# ) Multitype Branching Process

Consider a branching process on the type space $G$ with a Poisson offspring distribution
The probability that a single individual of type $g$ produces $k$ offsprings of type $h$ is

$$p_{gh}^k = \frac{(u_{gh} R_g)^k \, e^{-u_{gh} R_g}}{k!} \qquad p_{gh} \sim \text{Poisson}(u_{gh} R_g)$$

Let $\xi_g :=$ probability of escape starting with one individual of type $g$ (risk of escape)

$\Rightarrow \quad 1 - \xi_g :=$ probability of extinction

Since for extinction of type $g$, all lineages of type $g$ must go extinct, We get a recursive formula

$$1 - \xi_g = \prod_{h \supseteq g} \sum_{k=0}^{\infty} p_{gh}^k \, (1 - \xi_h)^k$$

$\underset{\text{multiply over the}}{\downarrow}$ mutational neighbourhood of $g$     $\underset{\text{chance of }g}{\downarrow}$ producing $k$ offsprings of $h$     $\searrow$ all $k$ offsprings go extinct

Substituting the Poisson distribution gives :

$$1 - \xi_g = \prod_{h \supseteq g} \sum_{k=0}^{\infty} \frac{(u_{gh} R_g)^k}{k!} e^{-u_{gh} R_g} (1 - \xi_h)^k$$

$$= \prod_{h \supseteq g} e^{-u_{gh} R_g} \sum_{k=0}^{\infty} \frac{[u_{gh} R_g (1 - \xi_h)]^k}{k!}$$

$$= \prod_{h \supseteq g} e^{-u_{gh} R_g} \, e^{u_{gh} R_g (1 - \xi_h)} \qquad \text{by Taylor series } e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

$$\Rightarrow \quad \log(1 - \xi_g) = -\sum_{h \supseteq g} u_{gh} R_g \xi_h$$

For $g \neq 1$, we have $\xi_g \ll 1$ and $(R_g)^2 \approx 0$. Thus

$$-\xi_g \cong -R_g \sum_{h \supseteq g} u_{gh} \xi_h \qquad \text{by } \log(1 - \xi_g) \approx -\xi_g$$

$$\Rightarrow \quad \xi_g \approx R_g \left( \xi_g + \sum_{h > g} u_{gh} \xi_h \right)$$

$$\Rightarrow \quad (1 - R_g) \xi_g \approx R_g \sum_{h > g} u_{gh} \xi_h$$

$$\Rightarrow \quad \xi_g \approx f_g \sum_{h > g} u_{gh} \xi_h \qquad \text{by } f_g := \frac{R_g}{1 - R_g}$$

In particular,

$$\xi_0 \approx f_0 \sum_{h \in G} u_{0h} \xi_h$$

Solving it yields

$$\xi_0 \approx \underset{\substack{\text{probability of escape} \\ \text{of one wild type}}}{\nwarrow} \xi_1 \cdot \underset{\substack{\text{probability of} \\ \text{escape of one} \\ \text{escape type}}}{\nwarrow} f_0 \cdot \underset{\substack{\text{fitness} \\ \text{of WT}}}{\downarrow} \prod_{e \in \mathcal{E}} u_e \cdot \underset{\substack{\text{accumulate all} \\ \text{necessary mutations}}}{\underbrace{R(G; f)}}$$

$\longrightarrow$ generates all mutational pathways from the wild type to the escape type

By i.i.d. assumption, the probability of escape of $N$ wild-type pathogens is

$$1 - (1 - \xi_0)^N \approx 1 - e^{-N \xi_0}$$

We define the critical population size as

$$N^* = \frac{1}{\zeta_0} \quad \Leftrightarrow \quad 1 - e^{-N^*\zeta_0} \approx 1 - e \approx \frac{1}{e}$$

↳ If $N \gg N^*$, then the escape is almost certain

↳ If $N = N^*$, then the probability of successful intervention is $\frac{1}{e} \approx \frac{1}{e}$

↳ If $N \ll N^*$, then escape is almost impossible

⇒ Successful treatment depends on the tumour size!

\* The more constraints are imposed on $g$, the larger the critical population size

# Chapter 12  Coalescent Theory

## ① The Coalescent

- In the Wright-Fisher Process, each individual in the new generation chooses a parent cell from the previous generation uniformly at random.
  - ⤷ there is a certain probability that two or more individuals have one common ancestor

- The coalescent process = thinking the Wright-Fisher Process <u>backward in time</u>
  - ⤷ there is a certain probability that two or more individuals coalesce

- Coalescent events = branching points in the process, where two or more lineages meet/coalesce

- Coalescent times = the waiting time between $j$ and $(j-1)$ lineages , $T(j)$
  /the waiting time until the next coalescent event
  /branch lengths in the tree

- The probability that $j$ individuals have no common ancestor in the previous generation

$$j = 2 : \quad 1 - \frac{1}{N}$$

$$j = 3 : \quad \left(1 - \frac{1}{N}\right)\left(1 - \frac{2}{N}\right)$$

$$\vdots$$

$$\prod_{i=1}^{j-1} \left(1 - \frac{i}{N}\right) = 1 - \binom{j}{2} N^{-1} + \underline{\mathcal{O}(N^{-2})}$$

$$\qquad\qquad\qquad\qquad\qquad \to 0 \text{ as } N \to \infty$$

- We measure time in units of $N$ generations.
  Let $T(j)$ be the coalescent time between $j$ and $j-1$ lineages
  = time in which $j$ individuals have no common ancestor.

$$\mathbb{P}(T(j) > t) = \left( \prod_{i=1}^{j-1} \left(1 - \frac{i}{N}\right) \right)^{Nt}$$

$$= \left( 1 - \binom{j}{2} N^{-1} + \mathcal{O}(N^{-2}) \right)^{Nt}$$

$$\to \exp\left(-\binom{j}{2} t\right) \qquad \text{as } N \to \infty$$

⟹ only pairwise coalescent events occur in the limit

⟹ $T(j) \sim$ exponential $\left( \binom{j}{2} \right)$

⟹ The stochastic process that models the coalescent time is called the coalescent

## 2) Time to the Most Recent Common Ancestor (MRCA)

- **Definition** : the most recent common ancestor refers to the root of the smallest tree where the leaves are the set of individuals in consideration

- **Time to MRCA** : for a sample of size $n$ ($\neq$ population size $N$)

$$T_{MRCA}(n) = \sum_{j=2}^{n} T(j)$$

- **Expectation** :

$$E[T_{MRCA}(n)] = \sum_{j=2}^{n} E[T(j)] \quad \text{by linearity of expectation}$$

$$= \sum_{j=2}^{n} 1/\binom{j}{2} \quad \text{since } T(j) \sim \exp(\binom{j}{2})$$

$$= \sum_{j=2}^{n} \frac{2}{j(j-1)}$$

$$= \sum_{j=2}^{n} 2\left(\frac{1}{j-1} - \frac{1}{j}\right)$$

$$= 2\left(1 - \frac{1}{n}\right) \quad \text{by telescoping sum}$$

* $E[T(2)] = 1$ : expected time for the last coalescent event is 1 in the unit of $N$ generations

* $\lim_{n \to \infty} E[T_{MRCA}(n)] = 2 = 2 \cdot E[T(2)]$ : twice as long as if there are only 2 individuals ! (50%)

- **Variance** :

$$Var[T_{MRCA}(n)] = \sum_{j=2}^{n} Var(T(j)) \quad \text{since } T(j)\text{'s are non-overlapping}$$

$$= \sum_{j=2}^{n} 1/\binom{j}{2}^{2} \quad \text{since } T(j) \sim \exp(\binom{j}{2})$$

$$= \sum_{j=2}^{n} \left(\frac{2}{j(j-1)}\right)^{2}$$

$$= 8\sum_{j=1}^{n} \frac{1}{j^{2}} + \frac{4}{n^{2}} - 8\left(1 - \frac{1}{n}\right) - 4$$

* $Var[T(2)] = 1$ : variance of the waiting time for the last coalescent event is also 1

* $\lim_{n \to \infty} Var[T_{MRCA}(n)] = \frac{8\pi^{2}}{6} - 12 \approx 1.16$ : the large proportion of this variance is explained by 2-individual coalescent time variance (86%)

## (6) Detecting Selection

· Simulation of mutation process

(1) Simulate a tree according to the coalescent process

(2) Superimpose a Poisson process that puts down mutations independently on all branches at rate $\frac{θ}{2}$ where $θ = 2Nμ$ (scaled mutation rate)

· Infinite sites model

· Assume an infinite number of sites and each mutation to affect a different nucleotide site (unique mutation)

↳ only appropriate for long DNA sequences with uniform mutation rate across sites

· Number of segregating sites : # of sites where not all alleles are identical

Under the infinite sites model,

S = total number of mutations

The total branch length is

$$T_{tot}(n) = \sum_{j=2}^{n} j\, T(j)$$

⇒ $E[S] = \frac{θ}{2} E[T_{tot}(n)]$  by the Poisson process

$= \frac{θ}{2} \sum_{j=2}^{n} j \cdot (1/\binom{j}{2})$  since $T(j) \sim \exp(\binom{j}{2})$

$= θ \underbrace{\sum_{j=2}^{n} \frac{1}{j-1}}_{C_n : \text{constant that depends only on } n}$

⇔ $θ = C_n^{-1} E[S]$

· Average Pairwise Nucleotide Distance

$$K = \sum_{\substack{i,j \in [n] \\ i<j}} \| g_i - g_j \|_0$$

⇒ average pairwise Hamming distance between 2 sequences in the sample

⇒ $E[K] = \frac{θ}{2} \cdot 2\, \underbrace{E[T(2)]}_{=1} = \frac{θ}{2} \cdot 2 = θ$

> $E[K] = θ = C_n^{-1} E[S]$
> under neutral infinite sites model

· Effects of selections on S and K

(1) S is not sensitive to **allele frequency** but sensitive to **low-frequency alleles**

(2) K is sensitive to **allele frequency** but not sensitive to **low-frequency alleles**

- <u>Tajima's D</u>

$$D = \frac{\hat{K} - c_n^{-1}\hat{S}}{\sqrt{\hat{V}}}$$

- <u>Null hypothesis</u>: there is no selection / the process is neutral

- the distribution of D under the null can be obtained through simulation of the Coalescent process many times

- can calculate the p-value of D using the null distribution and the data

## 4) <u>Inference under the coalescent</u>

- model parameters $\theta$ : mutation rate, population size, selective advantage, ...

- likelihood:

$$L(\theta) = P(D \mid \theta) = \int \underbrace{P(D \mid T, \theta)}_{\substack{\text{statistical phylogenetic} \\ \text{tree model}}} \underbrace{P(T \mid \theta)}_{\text{coalescent}} dT$$

- use MCMC to approximate integral

- model parameters are estimated by MLE or Bayesian inference.

## ① Tumor Evolution

- <u>Intra-tumour heterogeneity</u> : there could be multiple clones within a tumour
- <u>Clone:</u> a group of tumour cells that share a highly similar genotype and mutational profile
- <u>Subclone</u> : a group of tumour cells that diverge from an ancestral clone by acquiring additional mutations
- <u>Clone expansion</u> : process in which one genotype with higher fitness expands in frequency in the tumour mass
- <u>Selective sweep</u> : process in which a genotype with a very high fitness emerges and outcompetes all other clones in the tumour (e.g. cancer)
- <u>Driver mutations</u> : mutations that confer a fitness advantage
- <u>Passenger mutations</u> : mutations that have no effect on fitness
- <u>Truncal mutations</u> : ancestral mutations in the trunk of the phylogenetic tree that are shared by all clones
- <u>Subclonal mutations</u> : mutations in a lineage that has diverged from the trunk
- <u>Models of tumour evolution</u> :
  (1) Linear    (2) Neutral    (3) Branching    (4) Punctuated

## ② Variant-Allele-Frequency-Spectrum

- Binary leaf-labeled tree T : cells are the leaves and mutations occur on branches
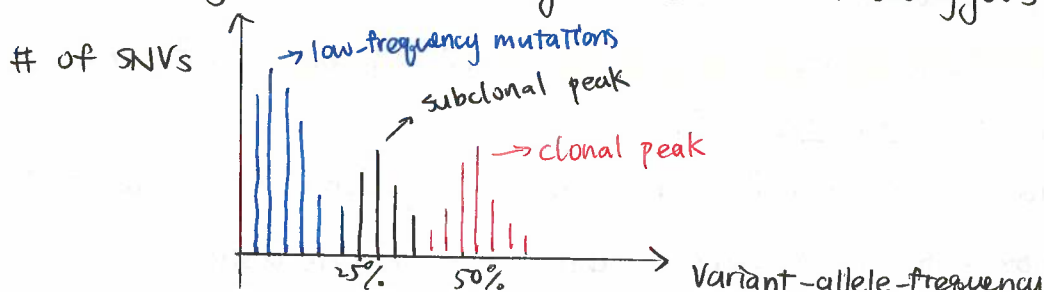- Relation matrix M :

$$M_{ij} = \begin{cases} 1 & \text{if cell } j \text{ is located below mutation } i \text{ in } T \\ 0 & \text{otherwise} \end{cases}$$

- Mutation frequencies:

$$\text{freq}(i) = \sum_j M_{ij}$$

- Data is usually from a mixture of cells.
- Mutations usually occur on a single strand ⇒ heterozygous ⇒ 50% frequencies

5) Inference under the neutral model

- Starting with a single cell, the number of tumour cells at time t under exponential growth is

$$N(t) = e^{\lambda \beta t}$$

where $\lambda$: cell division rate

$\beta$: successful division rate

- Assumptions: Infinite Sites Model

(1) Founding cell has acquired all mutations that give fitness advantage

(2) Subclonal mutations are neutral

- Expected number of new mutations per time interval:

$$\frac{dM(t)}{dt} = \mu \pi \lambda N(t)$$

where $\mu$: mutation rate at cell division

$\pi$: ploidy (# of chromosome sets : 2 for human)

- Total number of mutations in interval $[t_0, t]$:

$$M(t) = \mu \pi \lambda \int_{t_0}^{t} N(t) \, dt = \frac{\mu \pi}{\beta} (e^{\lambda \beta t} - e^{\lambda \beta t_0})$$

- Mutation age t vs. allele frequency f:

$$f = \frac{1}{\pi N(t)} = \frac{1}{\pi e^{\lambda \beta t}}$$

$$f_{max} = \frac{1}{\pi e^{\lambda \beta t_0}} = \frac{1}{2} \quad \text{for } \pi = 2 \text{ and } t_0 = 0$$

$$\Leftrightarrow \quad e^{\lambda \beta t} = \frac{1}{\pi f}$$

↳ older mutations = higher frequency

→ right slope of neutral cluster

↓ proportional to

$$\Rightarrow M(f) = \frac{\mu}{\beta} \left( \frac{1}{f} - \frac{1}{f_{max}} \right) \quad \text{or} \quad \frac{dM}{df} = -\mu \pi \lambda \frac{1}{f^2}$$

↓

mutation rate per effective cell division

for frequency f, $M(f)$ is the area under the VAF spectrum for all frequencies > f.

$\Rightarrow$ there should be a *linear relationship* b/w $M(f)$ and $\frac{1}{f}$

- Inference for Neutral Evolution:

↳ If the relationship is nonlinear, then the process is not neutral.

↳ If the relationship is linear, then it is inconclusive.

## (4) Inference in the presence of selection

- Assume 2 cell populations: host tumour and subclone with different growth rates

$$\lambda_{sub} \gtrsim \lambda_{host}$$

- **Selective advantage**:

$$s = \frac{\lambda_{sub} - \lambda_{host}}{\lambda_{host}}$$

- **Estimating subclone properties from the VAF spectrum**

↳ mutation rate $\mu$: estimated from the right slope of the neutral cluster

↳ subclone frequency $f_{sub}$: estimated from the mean of the subclone peak

↳ number of mutations in the subclone $M_{sub}$ at the time $t_1$ (when subclone appears): estimated from the area under the VAF curve of the subclone cluster

= the number of mutations acquired between $t_0$ and $t_1$

$$\Rightarrow \quad M_{sub} = \mu T' = \mu (2 \log(2) t_1)$$

where $T' :=$ mean # of successful cell divisions b/w $t_0$ and $t_1$

$$\Rightarrow \quad t_1 \simeq \frac{M_{sub}}{2 \log(2) \hat{\mu}} = \text{age of subclone in terms of doubling}$$

↳ Can also estimate the selective advantage:

$$s = \frac{\log\left(\frac{f_{sub}}{1 - f_{sub}}\right) + \lambda t_1}{\lambda (t_{end} - t_1)}$$

where $t_{end}$ can be estimated from tumour size

$$2^{t_{end}} = (1 - f_{sub}) 10^{10}$$

↳ size of tumour

and $\lambda \triangleq \log(2)$

## (5) Confounders of the VAF spectrum

(1) Contamination with normal cells can shift VAF spectrum to the left

(2) Copy number changes can shift VAF spectrum to the right

(3) Mutation losses destroy connection b/w allele frequency and age