

3. 遺伝子系図の合祖過程

集団遺伝の古典的理論として遺伝子頻度の変化を記述するマルコフ連鎖モデルおよび、拡散過程モデルについて紹介してきた。近年集団遺伝学においては遺伝子の系図を表現する合祖過程(Coalescent Process)と呼ばれる確率モデルが導入され広範に研究されている。この章ではこの遺伝子の系図に関するマルコフ過程モデルについて紹介する。

3. 1 合祖モデル

近交係数(inbreeding coefficient)、同祖的 (identical by descent)など遺伝子系図に関連する概念は古くから集団遺伝学では用いられている。Felsenstein(1971)は遺伝子系図の視点からN個の半数体生物から成る有限集団の遺伝的変異の減衰速度を論じた。突然変異および自然選択は無いものと仮定しよう。集団からランダムに選んだ r 個の個体(遺伝子)がちょうど s 個の異なる親(親遺伝子)を持つ確率を $G_{r,s}$ とする。当然 $r < s$ のとき $G_{r,s} = 0$ である。第 t 世代にランダムに取り出した k 個の遺伝子が異なる j 種のアレルを含む確率を

$P_t(k, j)$ とすると、1世代遡ることにより漸化式 $P_t(k, j) = \sum_{s=1}^N G_{k,s} P_{t-1}(s, j)$ を得る。

$N \times N$ 行列を $G = (G_{r,s})$, $P(t) = (P_t(k, j))$ とすると、 $P(t) = GP(t-1)$ と表される。これより、 $P(t) = G^t P(0)$ を得る。ここで行列の積 G^t の (k, j) 成分は k 個のサンプル遺伝子が t 世代遡った祖先集団で j 個の異なる祖先に由来する確率を表す。すなわちサンプル遺伝子の系図を表すマルコフ連鎖の推移行列であることに注意しよう。行列 G の性質を調べることにより、Felsensteinは十分大きな t に対して漸近的に $P_t(N, k) \sim A_k G_{kk}^t$ (A_k は定数)が成り立つことを示し、全集団から k 種のアレルの遺伝的多様性が失われる率が漸的に毎世代 G_{kk} で表されることを証明した。集団の遺伝的多様性の消失速度については第3. 3. 2節で詳しく論じる。さらに、Gladstien(1978)は可換モデルについて詳しい研究をしているが、ここでは簡単に次の定理を紹介しておこう。

定理 3. 1

可換モデルにおいては上記の確率 $G_{r,s}$ は次式で与えられる。

$$G_{rs} = \left[\frac{\binom{N}{s}}{\binom{N}{r}} \right] \sum_{R \in A(r,s)} E \left[\prod_{i=1}^s \binom{\theta_i}{r_i} \right]$$

ただし、 $A(r, s) = \{R = (r_1, r_2, \dots, r_s); r_1, \dots, r_s > 0, r_1 + \dots + r_s = r\}$ 、

$\sum_{R \in A(r,s)}$ は $A(r,s)$ に属す全ての $R = (r_1, \dots, r_s)$ についての和を表す。

(証明) s 個の親から生まれた r 個のサンプル中で i 番目の親から生まれた子供の数を r_i ($1 \leq i \leq s$) で表す。 $G_{r,s}$ はランダムに選んだ r 個の遺伝子が s 個の親遺伝子を持つ確率なので、まず全集団からランダムに r 個の遺伝子を選ぶ方法が $\binom{N}{r}$ 通りあり、前の世代から s 個の親個体を選ぶ方法は $\binom{N}{s}$ 通り。さらに、これらの間に親子関係の結び方を考えると、 i 番目の親が生む子供の数は θ_i 、その中から r_i 個の子供を選ぶ方法が各々独立に $\binom{\theta_i}{r_i}$ 通りあるので、可換モデルに従って $(\theta_1, \dots, \theta_s)$ に関する確率分布で期待値を取ると定理の結果を得る。

幾つかのモデルで具体例を計算してみよう。

例1 ライト・フィッシャーモデル

$$G_{kj} = \frac{N(N-1)\dots(N-j+1)}{N^k} S_k^{(j)} \quad (3.1)$$

$S_k^{(j)}$ は第2種スターリング数で異なる k 個のものを j 個のブロックに分ける場合の数を表す (詳しくは付録 (C) 参照)。

例2 モランモデル

$$G_{kk} = 1 - \frac{k(k-1)}{N^2} = 1 - G_{k,k-1} \quad (3.2)$$

Kingman(1982a,b,c)は可換モデルにおいてある条件の下で、集団のサイズ N を無限大にする拡散近似と同じ時間スケールの極限操作により、上記の遺伝子系図のプロセスが合祖過程(Coalescent Process)と呼ばれている連続時間のマルコフ連鎖(死滅過程)に収束することを示した。同じ合祖過程モデルはほとんど同時期に田嶋(Tajima(1983))によっても発表されている。

まず、可換モデルの一つライト・フィッシャーモデルを例にとってみよう。

N 個の半数体生物から成る集団を考えよう。この集団から 2 個の遺伝子をサンプルしたとき、異なる親を持つ確率 $P(2,2)$ を考えると、ある一つの遺伝子の親は N 通りあるので

$$P(2,2) = \frac{N(N-1)}{N^2} = 1 - \frac{1}{N}。さらに t 世代遡っても異なる祖先に由来する確率を $P_t(2,2)$$$

とすると、 $P_t(2,2) = \left(1 - \frac{1}{N}\right)^t$ となる。N 世代を単位時間 $t = 1$ とする時間スケールを取り

$[Mt]$ 世代遡っても異なる祖先に由来する確率を $P_t^N(2,2)$ とすると $P_t^N(2,2) = \left(1 - \frac{1}{N}\right)^{[Mt]}$

さらに $N \rightarrow \infty$ とすると、 $P_t(2,2) = \lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^{[Mt]} = \exp(-t)$ 、平均 1 の指数分布に従う。

より一般に n 個のサンプル遺伝子が k 個の親に由来する確率は(3.1)より

$$P(n,k) = G_{n,k} = \frac{N(N-1)\cdots(N-k+1)}{N^n} S_n^{(k)}$$

$$= \begin{cases} \frac{n(n-1)}{2N} + O\left(\frac{1}{N^2}\right) & \text{if } k = n-1 \\ 1 - \frac{n(n-1)}{2N} + O\left(\frac{1}{N^2}\right) & \text{if } k = n \\ O\left(\frac{1}{N^2}\right) & \text{if } k \leq n-2 \end{cases} \quad (3.3)$$

N 個のサンプル遺伝子が $[Mt]$ 世代遡っても異なる n 個の祖先をもつ確率 $P_t^N(n,n)$ は

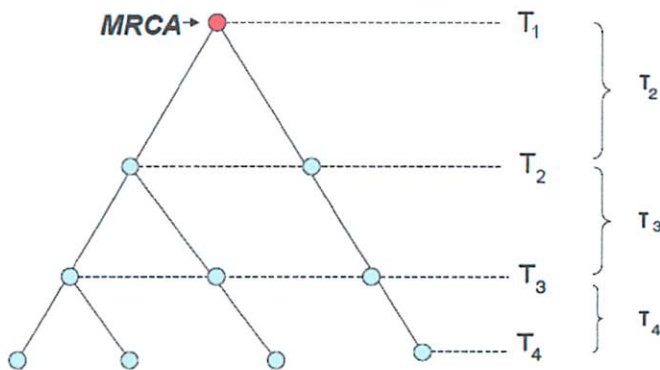
$$P_t^N(n,n) = \left(1 - \frac{n(n-1)}{2N} + O\left(\frac{1}{N^2}\right)\right)^{[Mt]}$$

よって $N \rightarrow \infty$ の極限をとると、

$$P_t(n,n) = \lim_{N \rightarrow \infty} \left(1 - \frac{n(n-1)}{2N} + O\left(\frac{1}{N^2}\right)\right)^{[Mt]} = \exp\left(-\frac{n(n-1)}{2}t\right)$$

$P_t(n,n)$ は平均が $2/n(n-1)$

の指数分布に従う。状態 n への滞在時間は平均 $2/n(n-1)$ の指数分布に従い、その後状態 $n-1$ へ推移する。



サンプル数 4 の場合の図
状態 4,3,2 の滞在時間を
 τ_4, τ_3, τ_2 とする。
 $T_4 = 0, T_1 = \tau_4 + \tau_3 + \tau_2$

n 個のサンプル遺伝子は最終的に確率 1 で一つの共通祖先 (Most Recent Common Ancestor) に到達する。その待ち時間を T_1 とすると平均 $E[T_1] = \sum_{k=2}^n \frac{2}{k(k-1)} = 2 \left(1 - \frac{1}{n}\right)$ となる。この確率過程を合祖過程 (Coalescent Process) と呼ぶ。

Kingman(1982)はより一般の可換モデルにおいて同様に合祖過程が導かれることを証明した。N個の半数体生物から成る集団から n 個の個体 I_1, I_2, \dots, I_n をサンプルしたとする。その祖先の共有関係によってサンプル遺伝子の集合に次のような同値関係 $R_r^{(N)}$ を定義する。

すなわち、二つの個体 I_k, I_j が r 世代前に祖先を共有するとき $(I_k, I_j) \in R_r^{(N)}$ と書き、この

時刻で同じ同値類に属すと言う。 $R_r^{(N)}$ は r 世代前に祖先を共有するか否かによって

I_1, I_2, \dots, I_n 上に定義された同値類である。初期状態は全てが異なるという状態でこれを

$R_0^{(N)} = \Delta$ で表す、また r 世代前にこの n 個の個体が単一の共通祖先を持つならばこれを

$R_r^{(N)} = \Theta$ で表すことにする。 $R_r^{(N)}$ は I_1, I_2, \dots, I_n 上に定義された同値関係を状態空間とする

マルコフ連鎖である。また $R_r^{(N)}$ に含まれるクラス数を $A_r^{(N)} = |R_r^{(N)}|$ で定義する。例えば

$n=5$ のとき、

$$R_0^{(N)} = \{(I_1)(I_2)(I_3)(I_4)(I_5)\}, R_1^{(N)} = \{(I_1 I_2)(I_3 I_4)(I_5)\},$$

$$R_2^{(N)} = \{(I_1, I_2)(I_1, I_2, I_5)\}, R_3^{(N)} = \{(I_1, I_2, I_3, I_4, I_5)\} = \Theta \text{ とすると}$$

$$A_0^{(N)} = 5, A_1^{(N)} = 3, A_2^{(N)} = 2, A_3^{(N)} = 1 \text{ となる。}$$

集合 $\{I_1, I_2, \dots, I_n\}$ 上に定義される全ての同値関係の集合を E_n とする。 $\alpha, \beta \in E_n$ に対して、

α が β の細分であるとき $\alpha \subseteq \beta$ で表す。例えば $\alpha = \{(I_1, I_2)(I_3, I_4)(I_5)\}$,

$\beta = \{(I_1, I_2)(I_3, I_4 I_5)\}$ のとき $\alpha \subseteq \beta$ である。また $\alpha = \{(I_1, I_2)(I_3, I_4)(I_5)\}$,

$\beta = \{(I_1, I_3)(I_2, I_4 I_5)\}$ のとき、 α は β の細分ではない。 α に含まれるクラス数を $|\alpha|$ とす

ると、 $\alpha \subseteq \beta$ のとき $|\alpha| \geq |\beta|$ である。 n 個のサンプル遺伝子が r 世代前の祖先によって α

$\in E_n$ のとき、 $r+1$ 世代前に状態 $\beta \in E_n$ である推移確率 $P_{\alpha, \beta}$ を求める。同値関係による

$\{I_1, I_2, \dots, I_n\}$ の分割を

$\beta : \{C_1, C_2, \dots, C_b\}$, $|\beta| = b$ とする。また、 α は β の細分なので、

$\alpha : \{C_{11}, \dots, C_{1a_1}, C_{21}, \dots, C_{2a_2}, \dots, C_{b1}, \dots, C_{ba_b}\}$ ただし $\sum_{i=1}^b a_i = |\alpha| = a$ 、

$$\bigcup_{i=1}^{a_j} C_{ji} = C_j \quad (j=1, 2, \dots, b) \text{ である。}$$

N 個の個体から成る集団からランダムに a 個の異なる個体を取り出したとき、その親が異なる b 個体である確率なので、 (j_1, j_2, \dots, j_b) を $\{1, 2, \dots, N\}$ の N 個の自然数から選んだ異なる b 個の自然数とすると、各親個体 j_i が生む k_i 個の子供から a_i 個の子供を選べば良いので、可換モデルであることに注意すると次式が成り立つ。

$$\begin{aligned} P_{\alpha, \beta} &= \frac{1}{N_{[a]}} E \left[\sum_{(j_1, j_2, \dots, j_b)} (v_{j_1})_{[a_1]} (v_{j_2})_{[a_2]} \dots (v_{j_b})_{[a_b]} \right] \\ &= \frac{1}{N_{[a]}} \sum_{(j_1, j_2, \dots, j_b)} \sum_{(k_1)_{[a_1]} (k_2)_{[a_2]} \dots (k_b)_{[a_b]}} P(v_{j_1} = k_1, \dots, v_{j_b} = k_b) \\ &= \frac{N_{[b]}}{N_{[a]}} \sum_{(k_1)_{[a_1]} \dots (k_b)_{[a_b]}} P(v_1 = k_1, \dots, v_b = k_b) = \frac{N_{[b]}}{N_{[a]}} E[(v_1)_{[a_1]} \dots (v_b)_{[a_b]}] \quad (3.4) \end{aligned}$$

ただし $x_{[k]} = x(x-1)(x-2)\dots(x-k+1)$

$\{P_{\alpha, \beta}; \alpha, \beta \in E_n\}$ はマルコフ連鎖 $R_r^{(N)}$ の推移確率である。

定義: $\alpha, \beta \in E_n$ に対して、 $\alpha \subseteq \beta$, $|\alpha| = |\beta| + 1$ のとき、 $\alpha \prec \beta$ で表す。すなわち

β は α の中の二つのクラスが合体(合祖)したものである。このとき次の定理を得る。

定理 3. 2 (Kingman(1982c))

極限值 $\lim_{N \rightarrow \infty} \text{Var}(v_1) = \sigma^2 > 0$ が存在し、すべての $p \geq 1$ に対して $\text{Sup}_N E[v_1^p] < \infty$ のとき

$$P_{\alpha, \beta} = \begin{cases} \frac{\sigma^2}{N} + o\left(\frac{1}{N}\right) & \text{if } \alpha \prec \beta \\ 1 - \frac{a(a-1)}{2} \times \frac{\sigma^2}{N} + o\left(\frac{1}{N}\right) & \text{if } \alpha = \beta \\ o\left(\frac{1}{N}\right) & \text{otherwise} \end{cases} \quad \text{ただし } a = |\alpha|$$

(証明)

(i) $\alpha \prec \beta$ のとき、 $|\alpha| = |\beta| + 1$ より $a_1 = 2, a_2 = 1, \dots, a_b = 1$ の場合を示せば十分である。

$$P_{\alpha,\beta} = \frac{1}{N_{[a]}} E \left[\sum_{(j_1, \dots, j_b)} \nu_{j_1} (\nu_{j_1} - 1) \nu_{j_2} \dots \nu_{j_b} \right] = \frac{1}{N_{[a]}} E \left[\sum_{j=1}^N \{ \nu_j (\nu_j - 1) \left(\sum_{\substack{(j_2, \dots, j_b) \\ j_i \neq j (2 \leq i \leq b)}} \nu_{j_2} \dots \nu_{j_b} \right) \} \right]$$

$$\leq \frac{1}{N_{[a]}} E \left[\sum_{j=1}^N \nu_j (\nu_j - 1) \{ \nu_1 + \nu_2 + \dots + \nu_N \}^{b-1} \right]$$

$$\sum_{i=1}^N \nu_i = N, b = a-1 \text{ より}$$

$$P_{\alpha,\beta} \leq \frac{N^{a-2}}{N_{[a]}} E \left[\sum_{j=1}^N \nu_j (\nu_j - 1) \right] = \frac{N^{a-1}}{N_{[a]}} E [\nu_1 (\nu_1 - 1)] = \frac{\sigma^2}{N} + o\left(\frac{1}{N}\right) \quad (3.5)$$

$$S = \sum_{j=1}^N \left\{ \nu_j (\nu_j - 1) \left(\sum_{\substack{(j_2, \dots, j_b) \\ j_i \neq j (2 \leq i \leq b)}} \nu_{j_2} \dots \nu_{j_b} \right) \right\} \text{ とすると、 } \sum_{\substack{(j_2, \dots, j_b) \\ j_i \neq j}} \nu_{j_2} \nu_{j_3} \dots \nu_{j_b} \text{ は } (\nu_1 + \dots + \nu_N - \nu_j)^{b-1}$$

の展開式から、いずれかの i について ν_i の 2 次以上を含む項を全て引いたものに等しい

$$\sum_{\substack{(j_2, \dots, j_b) \\ j_i \neq j}} \nu_{j_2} \nu_{j_3} \dots \nu_{j_b} \geq (N - \nu_j)^{b-1} - \sum_{i \neq j} \{ (\nu_i + \sum_{k \neq i, j} \nu_k)^{b-1} \text{ の中で } \nu_i \text{ の 2 次以上のを含む項の和} \}$$

{ } 内は $b-1 = a-2$ 個の $\left(\nu_i + \sum_{k \neq i, j} \nu_k \right)$ の積なので、その中から 2 個は ν_i 、他は ν_j を除

く任意の変数を選んで掛けた項全体の和より小さく、 $\sum_{k \neq j} \nu_k \leq N$ なので結局

$$\sum_{\substack{(j_2, \dots, j_b) \\ j_i \neq j}} \nu_{j_2} \nu_{j_3} \dots \nu_{j_b} \geq (N - \nu_j)^{b-1} - \binom{a-2}{2} \sum_{i \neq j} \nu_i^2 N^{a-4}$$

$p \geq 1, 0 \leq x \leq 1$ のとき、 $(1-x)^p \geq 1-px$ より

$$(N - \nu_j)^{b-1} = (N - \nu_j)^{a-2} = N^{a-2} \left(1 - \frac{\nu_j}{N} \right)^{a-2} \geq N^{a-2} \left\{ 1 - \frac{(a-2)\nu_j}{N} \right\} \quad \text{これより}$$

$$S = \sum_{j=1}^N \left\{ \nu_j (\nu_j - 1) \left(\sum_{\substack{(j_2, \dots, j_b) \\ j_i \neq j (2 \leq i \leq b)}} \nu_{j_2} \dots \nu_{j_b} \right) \right\}$$

$$\geq \sum_{j=1}^N \nu_j (\nu_j - 1) \left[N^{a-2} \left\{ 1 - \frac{(a-2)\nu_j}{N} \right\} - \binom{a-2}{2} \sum_{i \neq j} \nu_i^2 N^{a-4} \right]$$

$$\begin{aligned}
&= N^{a-2} \sum_{j=1}^N v_j(v_j-1) - (a-2)N^{a-3} \sum_{j=1}^N v_j^2(v_j-1) - \binom{a-2}{2} N^{a-4} \sum_j \sum_{i \neq j} v_j(v_j-1)v_i^2 \\
P_{\alpha,\beta} &= \frac{E[S]}{N_{[a]}} \geq \frac{N^{a-1}}{N_{[a]}} \left\{ E[v_1(v_1-1)] - \frac{(a-2)}{N} E[v_1^3] - \frac{\binom{a-2}{2}}{N} E[v_1^2 v_2^2] \right\} \\
&= \frac{\sigma^2}{N} + o\left(\frac{1}{N}\right) \dots\dots\dots \tag{3.6}
\end{aligned}$$

(3.5)(3.6)より $P_{\alpha,\beta} = \frac{\sigma^2}{N} + o\left(\frac{1}{N}\right)$ 。

(ii) $\alpha \subseteq \beta, |\alpha| \geq |\beta| + 2$ すなわち $a \geq b + 2$ のとき

$$\begin{aligned}
P_{\alpha,\beta} &= \frac{N_{[b]}}{N_{[a]}} E[(v_1)_{[a_1]}(v_2)_{[a_2]} \dots (v_b)_{[a_b]}] \leq \frac{N_{[b]}}{N_{[a]}} E[v_1^{a_1} v_2^{a_2} \dots v_b^{a_b}] \\
&\leq \frac{N_{[b]}}{N_{[a]}} (E[v_1^a])^{\frac{a_1}{a}} (E[v_2^a])^{\frac{a_2}{a}} \dots (E[v_b^a])^{\frac{a_b}{a}} \quad (\text{ヘルダーの不等式(注)}) \\
&= \frac{N_{[b]}}{N_{[a]}} E[v_1^a] = o\left(\frac{1}{N}\right)
\end{aligned}$$

(注) Appendix(E) ヘルダーの不等式で $f_i = v_i^{a_i}, \frac{1}{p_i} = \frac{a_i}{a}$ ($i = 1, 2, \dots, b$) と置く。

(iii) $\alpha = \beta$ のとき

$\alpha \not\subseteq \beta$ のとき、明らかに $P_{\alpha,\beta} = 0$ 。 $|\alpha| = a$ のとき $a = b + 1$ となる組み合わせは $\binom{a}{2}$ 通り

よって (i) (ii) より $P_{\alpha,\alpha} = 1 - \frac{a(a-1)}{2} \times \frac{\sigma^2}{N} + o\left(\frac{1}{N}\right)$ を得る。

(証明終わり)

$$Q = (Q_{\alpha,\beta}), \quad \text{ただし } Q_{\alpha,\beta} = \begin{cases} 1 & \text{if } \alpha < \beta \\ -\frac{a(a-1)}{2} & \text{if } \alpha = \beta, |\alpha| = a \\ 0 & \text{その他} \end{cases} \tag{3.7}$$

とすると推移確率は $P = (P_{\alpha,\beta}) = I + c_N Q + o(c_N), \quad c_N = \frac{\sigma^2}{N}$ と表される。 (3.8)

集団から二つの遺伝子をランダムに取り出したとき、同じ親を持つ確率を c_N とすると

$$c_N = \frac{1}{N(N-1)} \sum_{i=1}^N E[v_i(v_i-1)] = \frac{E[v_1(v_1-1)]}{N-1} = \frac{Var(v_1)}{N-1}$$

(1) Wright-Fisher model

$$v_1 \text{ は二項分布 } B(N, \frac{1}{N}) \text{ に従うので、} v_1 \text{ の母関数は } g(x) = E[e^{xv_1}] = (\frac{e^x}{N} + 1 - \frac{1}{N})^N$$

$$Var(v_1) = N \times \frac{1}{N} \times (1 - \frac{1}{N}) = 1 - \frac{1}{N}, \quad \lim_{N \rightarrow \infty} Var(v_1) = 1 \quad c_N = \frac{Var(v_1)}{N-1} = \frac{1}{N}$$

$$\sup_N E[v_1^p] = \sup_N \lim_{x \rightarrow 0} \frac{\partial^p}{\partial x^p} g(x) < \infty$$

(2) Moran model

$$P(v_1 = 0) = \frac{N-1}{N^2}, \quad P(v_1 = 1) = \frac{1}{N^2} + (\frac{N-1}{N})^2 = \frac{N^2 - 2N + 2}{N^2}, \quad P(v_1 = 2) = \frac{N-1}{N^2}$$

$$E[v_1] = 1, \quad Var(v_1) = \frac{2(N-1)}{N^2}, \quad c_N = \frac{2}{N^2}$$

定理 3. 2 の結果より N 世代を単位時間とする時間のスケーリングを行い、 N を無限に大きくすると遺伝子の系図を表す合祖過程(Coalescent Process)と呼ばれる連続時間マルコフ連鎖が得られる。

定理 3. 3

一般に行列 $A = (a_{\xi\eta})$ は $a_{\xi\eta} \geq 0, \sum_{\eta} a_{\xi\eta} = 1$ を満たすとき確率行列と呼ばれる。このとき

ノルム $\|A\| = \max_{\xi} \sum_{\eta} |a_{\xi\eta}|$ に関して縮小作用素 ($\|A\| \leq 1$) であるので行列 $A_i, B_i (i = 1, 2, \dots, r)$ が

縮小作用素のとき、 r に関する数学的帰納法により $\|A_1 A_2 \dots A_r - B_1 B_2 \dots B_r\| \leq \sum_{i=1}^r \|A_i - B_i\|$

が成り立つ。これより

$$\left\| P_N^{[t/c_N]} - (I + c_N Q)^{[t/c_N]} \right\| \leq \frac{t}{c_N} \left\| P_N - (I + c_N Q) \right\| = t \left\| \frac{P_N - I}{c_N} - Q \right\|。$$

故に $\lim_{N \rightarrow \infty} P_N^{[t/c_N]} = \lim_{N \rightarrow \infty} (I + c_N Q)^{[t/c_N]} = e^{tQ}$ 。これより時間スケールされたプロセス

$(R_{[t/c_N]}^{(N)})_{t \geq 0}$ の有限次元分布収束が導かれる。

サンプル α の定理 3. 3 で得られた連続時間マルコフ連鎖を $\alpha_t = \lim_{N \rightarrow \infty} R_{t/c_N}^{(N)}$ とし n -合祖過程と呼ぶ。 $|\alpha_t|$ は時刻 t での祖先遺伝子数を表し、初期条件は $|\alpha_0| = n$ で $|\alpha_t|$ は時間と共に減少する死滅過程である。連続時間マルコフ連鎖 α_t が状態 $|\alpha_t| = k$ に滞在する時間を τ_k とすると、

$$P(\tau_k > t) = \exp\left(-\frac{k(k-1)}{2}t\right), \quad k = n, n-1, \dots, 2 \quad (3.9)$$

$$\text{分布密度 } \frac{d}{dt}P(\tau_k \leq t) = \frac{k(k-1)}{2} \exp\left(-\frac{k(k-1)}{2}t\right), \quad \tau_n, \tau_{n-1}, \dots, \tau_2 \text{ は独立。}$$

祖先遺伝子数 $|\alpha_t|$ の推移確率を $P_t(k, j) = P(|\alpha_t| = j | |\alpha_0| = k)$ とすると次の方程式を満たす。

$$\frac{d}{dt}P_t(k, j) = \frac{k(k-1)}{2} \{P_t(k-1, j) - P_t(k, j)\} \quad P_0(k, j) = \delta_{k,j} \quad (3.10)$$

この解については次の節で突然変異を含むより一般的な形で求める。

n -合祖過程 α_t のジャンプ過程を $\{\mathfrak{R}_k; k = n, n-1, \dots, 1\}$ とすると、 $\beta \in E_n, |\beta| = k$ のとき

$$P(\alpha_t = \beta | |\alpha_0| = n) = P(|\alpha_t| = k | |\alpha_0| = n)P(\mathfrak{R}_k = \beta | \mathfrak{R}_n = \Delta), \text{ すなわち二つの過程 } |\alpha_t| \text{ と } \mathfrak{R}_k \text{ は独立である。}$$

3. 2 遺伝的多様性と遺伝子系図

生物集団内の遺伝的多様性は DNA 塩基配列の多様性であり、それは DNA に生じる突然変異に起因する。この節ではまず突然変異を含む遺伝子系図について考えよう。

突然変異はそれまでに存在しない常に新しいタイプであるという無限対立遺伝子モデルを仮定する。このとき、系図上に現れる突然変異は常に新しいアレルタイプの起源と見なされる。サンプル遺伝子の系図を過去に遡って行くと、そのアレルタイプの起源となる祖先遺伝子に到達する。同一の突然変異を起源にもつ遺伝子を同祖的 (identical by descent) と呼ぶことにすると、サンプル遺伝子の祖先遺伝子に含まれる家系数 (Lines of Descent) に関するマルコフ過程を得る。突然変異は毎世代 $\frac{\theta}{2N}$ の率で起こり、生じた突然変異は常に新しいタイプとする (infinite-allele model)。

k 個のサンプル遺伝子について、時間を t だけ遡った集団において突然変異を経ることなしに j 個の祖先に由来する確率を $P_t(k, j)$ とすると、

$$1 \text{ 世代当たりの変化は } 1 \text{ 世代を } \Delta t = \frac{1}{N} \text{ の時間スケールにとると}$$

$$P_{i+\Delta t}(k, k) = \left(1 - \frac{\vartheta}{2N}\right)^k \left(1 - \frac{k(k-1)}{2N}\right) P_i(k, k) + O\left(\frac{1}{N^2}\right)$$

$$P_{i+\Delta t}(k, j) = \left(1 - \frac{\vartheta}{2N}\right)^k \left(1 - \frac{k(k-1)}{2N}\right) P_i(k, j) + \left(\frac{k\vartheta}{2N} + \frac{k(k-1)}{2N}\right) P_i(k-1, j) + O\left(\frac{1}{N^2}\right)$$

これより次の微分方程式を得る。

定理 3. 4

$$\frac{d}{dt} P_i(k, j) = \frac{k(k+\vartheta-1)}{2} \{P_i(k-1, j) - P_i(k, j)\} \quad (k > j)$$

$$\frac{d}{dt} P_i(j, j) = -\frac{j(j+\vartheta-1)}{2} P_i(j, j), \quad P_0(k, j) = \delta_{k,j}$$

$\vartheta = 0$ とすると定理 3. 3 の合祖過程が得られる。

この方程式の解は Tavaré(1984, p129) で解説されているが、ここではラプラス変換を用いた方法を紹介しよう。

$Q_\lambda(k, j) = \int_0^\infty e^{-\lambda t} P_i(k, j) dt$ としよう。上式の両辺をラプラス変換すると

$$\lambda Q_\lambda(k, j) - \delta_{k,j} = \alpha_i \{Q_\lambda(k-1, j) - Q_\lambda(k, j)\} \quad \text{ただし } \alpha_i = \frac{i(i+\vartheta-1)}{2}.$$

$$\text{これより } k > j \text{ のとき } Q_\lambda(k, j) = \frac{\alpha_k}{\lambda + \alpha_k} Q_\lambda(k-1, j) + \frac{\delta_{k,j}}{\lambda + \alpha_k} = \frac{\alpha_k}{\lambda + \alpha_k} Q_\lambda(k-1, j)$$

$$\text{これを繰り返すと } Q_\lambda(k, j) = \left\{ \prod_{i=j+1}^k \frac{\alpha_i}{\lambda + \alpha_i} \right\} Q_\lambda(j, j) = \frac{1}{\alpha_j} \prod_{i=j}^k \frac{\alpha_i}{\lambda + \alpha_i},$$

$$\text{ただし } P_i(j, j) = \exp\left\{-\frac{j(j+\vartheta-1)}{2} t\right\} \text{ より } Q_\lambda(j, j) = \frac{1}{\lambda + \alpha_j} \text{ である。}$$

さらにこの式を部分分数に分解すると

$$Q_\lambda(k, j) = \frac{1}{\alpha_j} \sum_{i=j}^k \frac{A_i}{\lambda + \alpha_i}, \quad \text{ただし } A_i = \left(\prod_{r=j}^k \alpha_r \right) / \prod_{\substack{r=j \\ r \neq i}}^k (\alpha_r - \alpha_i), \quad i \geq j$$

ラプラス逆変換により

$$P_i(k, j) = \sum_{i=j}^k \frac{A_i}{\alpha_j} \exp\left[-\frac{i(i+\vartheta-1)}{2} t\right] \quad (3.11)$$

と表現される。ただし

$$\begin{aligned} \frac{A_t}{\alpha_j} &= \left(\prod_{r=j+1}^k \alpha_r \right) / \prod_{\substack{r=j \\ r \neq i}}^k (\alpha_r - \alpha_i) = \prod_{r=j+1}^k \{r(r+\vartheta-1)\} / \prod_{\substack{r=j \\ r \neq i}}^k \{(r-i)(r+i+\vartheta-1)\} \\ &= \frac{k!(2i+\vartheta-1) \prod_{r=j+1}^k (r+\vartheta-1)}{j!(-1)^{t-j}(k-i)! \prod_{r=j}^k (r+i+\vartheta-1)} = \frac{(-1)^{t-j}(2i+\vartheta-1)k_{[j]}(j+\vartheta)_{(t-1)}}{j!(i-j)!(k+\vartheta)_{(t)}} \end{aligned}$$

ここで $x_{(j)} = x(x+1)\dots(x+j-1)$, $x_{[j]} = x(x-1)\dots(x-j+1)$

さらにサンプル数 k を無限に大きくすると、

$$\lim_{k \rightarrow \infty} \frac{A_t}{\alpha_j} = \lim_{k \rightarrow \infty} \frac{(-1)^{t-j}(2i+\vartheta-1)k_{[j]}(j+\vartheta)_{(t-1)}}{j!(i-j)!(k+\vartheta)_{(t)}} = \frac{(-1)^{t-j}(2i+\vartheta-1)(j+\vartheta)_{(t-1)}}{j!(i-j)!}$$

$k \rightarrow \infty$ とすると、 $P_t(j) = \lim_{k \rightarrow \infty} P_t(k, j)$ は全集団の遺伝子が t 遡った祖先集団で突然変異を経ることなく j 個の遺伝子と同祖的(lines of descent)につながっている確率となる。

$$P_t(j) = \lim_{k \rightarrow \infty} P_t(k, j) = \sum_{i=j}^{\infty} \frac{(-1)^{t-j}(2i+\vartheta-1)(j+\vartheta)_{(t-1)}}{j!(i-j)!} \exp\left[-\frac{i(i+\vartheta-1)}{2}t\right] \quad (3.12)$$

特に $\vartheta = 0$ とすると、全集団の遺伝子が t 遡った j 個の祖先遺伝子の子孫である確率は

$$P_t(j) = \sum_{i=j}^{\infty} \frac{(-1)^{t-j}(2i-1)(j)_{(t-1)}}{j!(i-j)!} \exp\left[-\frac{i(i-1)}{2}t\right] \quad (3.13)$$

$j = 1$ とすると、 t 遡った祖先集団で一つの共通祖先に由来する確率が

$$P_t(1) = \sum_{i=1}^{\infty} \frac{(-1)^{t-1}(2i-1)(1)_{(t-1)}}{(i-1)!} \exp\left[-\frac{i(i-1)}{2}t\right] = 1 - \sum_{i=2}^{\infty} (-1)^t(2i-1) \exp\left[-\frac{i(i-1)}{2}t\right] \quad (3.14)$$

遺伝子はアデニン(A)、チミン(T)、シトシン(C)、グアニン(G)の4種の塩基の一次配列として表現される。集団から3本の遺伝子をサンプルし、次のようなDNA塩基配列が得られたと仮定しよう。

ATTGCCTAGGTCAACTGGACCTGA (サンプル1)
 ATTGTCTAGCTCAACTGGACCTGA (サンプル2)
 ATTGCCTAGCTCAACTGGACCTGA (サンプル3)

全部で24個の塩基から成るDNA領域で、ほとんど同じ配列であるが、例えば左から5番目の塩基サイトはサンプル1と3はシトシン(C)であるがサンプル2はチミン(T)であり、同様に10番目もサンプルによって塩基の種類(G, C)が異なる。この様に、サ

ンプルによって異なる塩基を持つようなサイトを分離サイト(segregating site)という。一般に二つのサンプル遺伝子について、その共通祖先までの合祖時間が長いと、それだけ多くの突然変異を蓄積し分離サイトの数も増加する。従って、逆に分離サイトの数が多いと、共通祖先までの合祖時間も長いと予想される。そこで、サンプル中の分離サイトの情報から共通祖先までの総合祖時間を推定する問題を考えてみよう。サンプル遺伝子の系図過程で j 個の祖先遺伝子である期間の長さを τ_j ($2 \leq j \leq n$) , とする。 $T = \sum_{j=2}^n \tau_j$ は合祖時

間の合計、 $L = \sum_{j=2}^n j\tau_j$ は系図の全ての枝(branch)の長さの合計である。

サンプル遺伝子数が 2 個のとき、田島(Tajima(1983))は分離サイトの数が $S = k$ という条件の下で、同一祖先までの合祖時間の分布が次のガンマ分布で与えられることを示した。合祖時間の分布は指数分布 $P(T = t) = \exp(-t)$ に従い、 $T = t$ の条件下で分離サイト

の数は平均 $\frac{\theta}{2} \times 2t = \theta t$ のポアソン分布 $P(S = k|T = t) = \frac{(\theta t)^k}{k!} \exp(-\theta t)$ に従うので

$$\begin{aligned} P_2(T = t|S = k) &= \frac{P_2(T = t, S = k)}{P_2(S = k)} = \frac{P_2(T = t)P_2(S = k|T = t)}{P_2(S = k)} \\ &= \frac{\exp(-t) \times \frac{(\theta t)^k}{k!} \exp(-\theta t)}{\frac{1}{(1+\theta)} \left(\frac{\theta}{1+\theta}\right)^k} = \frac{(1+\theta)^{k+1}}{k!} t^k \exp[-(1+\theta)t] = \text{Gamma}(k+1, \frac{1}{1+\theta}) \end{aligned} \quad (3.15)$$

ここで $\theta = 2Nu$, u は 1 世代 1 遺伝子当たりの突然変異率である。

ここでは一般に n 個の遺伝子をサンプルしたとき、その分離サイトから共通な祖先までの総合祖時間を推定する問題を考える。そのために、 (T, S) の同時母関数を定義する。

$$G(\lambda, z) = E[\exp(-\lambda T) \times z^S] = \sum_{k=0}^{\infty} z^k \left\{ \int_0^{\infty} P(T = t, S = k) e^{-\lambda t} dt \right\} \text{ ただし } \lambda \geq 0, 0 \leq z \leq 1.$$

このとき、次のようにして容易に母関数を求めることができる(Tavare et al(1997))。

補助定理 3. 5

$$G(\lambda, z) = \prod_{j=2}^n \left\{ \frac{j(j-1)}{j(j-1) + \theta j(1-z) + 2\lambda} \right\}$$

(証明) 突然変異は系図の全長 L のとき、平均 $\frac{\theta L}{2}$ のポアソン分布で発生し、滞在時間 τ_j は分布密度が指数分布 $\frac{d}{dt} P(\tau_j \leq t) = \frac{j(j-1)}{2} \exp\left(-\frac{j(j-1)}{2} t\right)$ に従うことより、

$$\begin{aligned}
G(\lambda, z) &= E[E[\exp(-\lambda T)z^S | \tau_2, \dots, \tau_n]] = E[\exp(-\lambda T)E[z^S | \tau_2, \dots, \tau_n]] \\
&= E[\exp(-\lambda T) \left\{ \sum_{k=0}^{\infty} z^k P(S = k | L) \right\}] = E \left[\exp(-\lambda T) \left\{ \sum_{k=0}^{\infty} z^k \times \frac{\left(\frac{\vartheta L}{2}\right)^k}{k!} \exp\left(-\frac{\vartheta L}{2}\right) \right\} \right] \\
&= E[\exp(-\lambda T) \exp\left(\frac{\vartheta L z}{2}\right) \exp\left(-\frac{\vartheta L}{2}\right)] = E[\exp(-\lambda T - \frac{\vartheta}{2} L(1-z))] \\
&= E[\exp\{-\lambda \sum_{j=2}^n \tau_j - \frac{\vartheta}{2} (1-z) \sum_{j=2}^n j \tau_j\}] = \prod_{j=2}^n E[\exp\{-\lambda \tau_j - \frac{\vartheta}{2} (1-z) j \tau_j\}] \\
&= \prod_{j=2}^n \left[\int_0^{\infty} \exp\{-\lambda t - \frac{\vartheta}{2} (1-z) j t\} \times \frac{j(j-1)}{2} \exp\left(-\frac{j(j-1)}{2} t\right) dt \right] \\
&= \prod_{j=2}^n \left(\frac{j(j-1)}{j(j-1) + \vartheta j(1-z) + 2\lambda} \right)
\end{aligned}$$

ここで $E[\bullet | \tau_2, \dots, \tau_n]$ は τ_2, \dots, τ_n を与えた下での条件付き期待値を表す (付録(H)参照)。

母関数 $G(\lambda, z)$ を z のベキ級数に展開すると

$$G(\lambda, z) = \sum_{k=0}^{\infty} z^k \left\{ \sum_{|r|=k} \left\{ \prod_{j=2}^n \left(\frac{\vartheta}{\vartheta + j - 1} \right)^{r_j} \left(\frac{j-1}{\vartheta + j - 1} \right) \right\} \left\{ \prod_{j=2}^n \left(\frac{j(\vartheta + j - 1)}{2\lambda + j(\vartheta + j - 1)} \right)^{r_j + 1} \right\} \right\}$$

ここで $\sum_{|r|=k}$ は $|r| = \sum_{i=2}^n r_i = k$ を満たす全ての $n-1$ 次元ベクトル $r = (r_2, \dots, r_n)$ についての和を表す。これより

$$\int_0^{\infty} P(T = t, S = k) e^{-\lambda t} dt = \sum_{|r|=k} \left\{ \prod_{j=2}^n \left(\frac{\vartheta}{\vartheta + j - 1} \right)^{r_j} \left(\frac{j-1}{\vartheta + j - 1} \right) \right\} \left\{ \prod_{j=2}^n \left(\frac{j(\vartheta + j - 1)}{2\lambda + j(\vartheta + j - 1)} \right)^{r_j + 1} \right\}.$$

最後にラプラス逆変換により次の定理を得る。

定理 3. 6

$$P(T = t, S = k) = \sum_{|r|=k} \left\{ \prod_{j=2}^n \left(\frac{\vartheta}{\vartheta + j - 1} \right)^{r_j} \left(\frac{j-1}{\vartheta + j - 1} \right) \right\} \left\{ \prod_{j=2}^n \text{Gamma} \left(r_j + 1, \frac{2}{j(\vartheta + j - 1)} \right) \right\}$$

$$\text{ここで } \text{Gamma} \left(r_j + 1, \frac{2}{j(\vartheta + j - 1)} \right) = \frac{\left(\frac{j(\vartheta + j - 1)}{2} \right)^{r_j + 1}}{r_j!} t^{r_j} \exp \left[-\frac{j(\vartheta + j - 1)}{2} t \right]$$

また $\ast \prod_{j=2}^n \text{Gamma}(\cdot)$ はガンマ分布の畳み込みを意味する。さらに

$$P(S=k) = \sum_{|r|=k} \prod_{j=2}^n \left\{ \left(\frac{\vartheta}{\vartheta+j-1} \right)^{\vartheta} \left(\frac{j-1}{\vartheta+j-1} \right) \right\}$$

$$P_n(T=t|S=k) = \frac{P(T=t, S=k)}{P(S=k)}$$

定理 3. 6 の内容は第 1 章、補助定理 1. 17 に基づき次のように説明できる。

一つの遺伝子の系図上で突然変異率は $\frac{\vartheta}{2}$ である。故に k 個遺伝子の系図のいずれかで mutation が起こるまでの待ち時間を σ_k とすると指数分布

$$P(\sigma_k > t) = \left\{ \exp\left[-\frac{\vartheta}{2}t\right] \right\}^k = \exp\left[-\frac{\vartheta k}{2}t\right]$$

$$P(\tau_k > t) = \exp\left[-\frac{k(k-1)}{2}t\right]$$

k 個のサンプル遺伝子について最初に突然変異または合祖が起こるまでの待ち時間 $\text{Min}(\sigma_k, \tau_k)$ の分布は

$$P(\text{Min}(\sigma_k, \tau_k) > t) = \exp\left[-\frac{k\vartheta}{2}t\right] \exp\left[-\frac{k(k-1)}{2}t\right] = \exp\left[-\left(\frac{k\vartheta+k(k-1)}{2}\right)t\right]$$

変化が起きたとき、それが突然変異、合祖である確率はそれぞれ

$$P(\text{突然変異}) = P(\text{Min}(\sigma_k, \tau_k) = \sigma_k) = \frac{k\vartheta/2}{\{k\vartheta+k(k-1)\}/2} = \frac{\vartheta}{\vartheta+k-1}$$

$$P(\text{合祖}) = P(\text{Min}(\sigma_k, \tau_k) = \tau_k) = \frac{k(k-1)/2}{\{k\vartheta+k(k-1)\}/2} = \frac{k-1}{\vartheta+k-1}$$

これより k 個のサンプル遺伝子の系図について r 回突然変異が生じ最後に coalesce が起こる確率は $\left(\frac{\vartheta}{\vartheta+k-1}\right)^r \left(\frac{k-1}{\vartheta+k-1}\right)$ 、 r 回の突然変異と 1 回の合祖が起こるまでの待ち時間は指数分布 $\text{Exp}\left(\frac{2}{k(\vartheta+k-1)}\right)$ の $r+1$ 重畳み込みで与えられるので、ガンマ分布

$$\text{Gamma}\left(r+1, \frac{2}{k(\vartheta+k-1)}\right)$$

となる (ガンマ分布については付録 A 参照)。
 $|r| = \sum_{j=2}^n r_j = k$, を満たすベクトル $r = (r_2, \dots, r_n)$ について $\left(\frac{\vartheta}{\vartheta+j-1}\right)^{\vartheta} \left(\frac{j-1}{\vartheta+j-1}\right)$ は j 個

の系統が $j-1$ 個の系統に合祖する前に r_j 回突然変異が起こる確率であり、

$\prod_{j=2}^n \left(\frac{\vartheta}{\vartheta+j-1} \right)^{r_j} \left(\frac{j-1}{\vartheta+j-1} \right)$ は共通な一つの祖先に到達するまでに、 $r = (r_2, \dots, r_n)$ に対応

した突然変異と合祖が起こる確率である。状態 j から $j-1$ へ合祖するまでに r_j 回の突

然変異が起こるとき状態 j の滞在時間はガンマ分布 $\text{Gamma}(r_j + 1, \frac{2}{j(\vartheta+j-1)})$ に従う。

$\eta^r_j (j=2, \dots, n)$ をガンマ分布 $\text{Gamma}\left(r_j + 1, \frac{2}{j(\vartheta+j-1)}\right)$ に従う独立な確率変数とす

ると、総計 $\eta^r = \sum_{j=2}^n \eta_j$ の分布はこれらのガンマ分布の畳み込みで表される。

$r = (r_2, \dots, r_n)$ に対して、 $A(r) = \prod_{j=2}^n \left\{ \left(\frac{\vartheta}{\vartheta+j-1} \right)^{r_j} \left(\frac{j-1}{\vartheta+j-1} \right) \right\} / P(S=k)$ とすると、明

らかに $\sum_{|r|=k} A(r) = 1$ 。T の条件付分布は $P_n(T=t|S=k) = \sum_{|r|=k} A(r) P(\eta^r = t)$ となる。

ガンマ分布 $\text{Gamma}(\alpha, \beta)$ の平均、分散はそれぞれ $\alpha\beta$ および $\alpha\beta^2$ である。よって

$$E[\eta^r] = \sum_{j=2}^n \frac{2(r_j+1)}{j(\vartheta+j-1)} \quad \text{and} \quad \text{Var}[\eta^r] = \sum_{j=2}^n \frac{4(r_j+1)}{j^2(\vartheta+j-1)^2}.$$

条件 $S=k$ の下での合祖時間 T の条件付平均および分散は次式のようにになる。

$$\begin{cases} E\{T|S=k\} = \sum_{|r|=k} A(r) E[\eta^r] \\ \text{Var}[T|S=k] = \sum_{|r|=k} A(r) \text{Var}[\eta^r] + \sum_{|r|=k} A(r) \{E[\eta^r] - E\}^2 \quad \text{where } E = E\{T|S=k\} \end{cases} \quad (3.16)$$

分離したサイト数 S について別の表現式を紹介しよう。系図の全枝の長さ $L = \sum_{j=2}^n j\tau_j$ が

与えられたとき、突然変異の数はパラメーター $\frac{\vartheta}{2}L$ のポアソン分布に従うので

$$P(S=k|L) = \frac{(\vartheta L/2)^k}{k!} \exp\left(-\frac{\vartheta L}{2}\right). \quad L \text{ の分布密度を } f_L(t) \text{ とすると}$$

$P(S=k) = \int_0^\infty P(S=k|t) f_L(t) dt$ で与えられる。 $L_j = j\tau_j (j=2, 3, \dots, n)$ の分布は

$$P(L_j = j\tau_j > t) = P\left(\tau_j > \frac{t}{j}\right) = \exp\left(-\frac{j(j-1)}{2} \times \frac{t}{j}\right) = \exp\left(-\frac{(j-1)t}{2}\right) \text{ より}$$

分布密度は $f_{L_j}(t) = \frac{j-1}{2} \exp\left(-\frac{(j-1)t}{2}\right)$ となる。 L_2, L_3, \dots, L_n は互いに独立で

$L = \sum_{j=2}^n L_j$ なので、 L の分布密度は $\{f_{L_j}(t); j=2, 3, \dots, n\}$ の畳み込みで得られるので

$$\text{最終的に、 } f_L(t) = \sum_{j=2}^n (-1)^j \binom{n-1}{j-1} \frac{j-1}{2} \exp\left(-\frac{(j-1)t}{2}\right) \text{ となる。} \quad (3.17)$$

また、これは $f_L(t) = \frac{n-1}{2} e^{-t/2} (1 - e^{-t/2})^{n-2}$ という表示も可能である。

これより、分離サイト数 S の確率分布は

$$\begin{aligned} P(S = k) &= \int_0^\infty \frac{(g t / 2)^k}{k!} \exp\left(-\frac{g t}{2}\right) \times \left(\sum_{j=2}^n (-1)^j \binom{n-1}{j-1} \frac{j-1}{2} \exp\left(-\frac{(j-1)t}{2}\right)\right) dt \\ &= \left(\frac{g}{2}\right)^k \sum_{j=2}^n (-1)^j \binom{n-1}{j-1} \frac{j-1}{2} \int_0^\infty \frac{t^k}{k!} \exp\left(-\frac{(g+j-1)t}{2}\right) dt \\ &= \sum_{j=2}^n (-1)^j \binom{n-1}{j-1} \left(\frac{j-1}{g+j-1}\right) \left(\frac{g}{g+j-1}\right)^k. \end{aligned} \quad (3.18)$$

3. 3 合祖過程の種々の性質

合祖過程を用いて集団内の遺伝子の系図や多様性に関する様々な性質を導くことができる。ここでは Kingman(1982), Tavaré(1984) から幾つかの結果を紹介する。

この節では突然変異は仮定しない Kingman の合祖過程から導かれる性質を紹介する。

3. 3. 1 合祖過程のジャンプ過程とファミリーサイズ過程

合祖過程 α_t のジャンプ過程 $\{\mathfrak{R}_k; k = n, n-1, \dots, 1\}$ の推移確率を求める。 $\{\mathfrak{R}_k\}$ の

$$1 \text{ 回の推移確率は } \xi, \eta \in E_n \text{ に対し } P(\mathfrak{R}_{k-1} = \eta | \mathfrak{R}_k = \xi) = \begin{cases} \frac{2}{k(k-1)} & \xi > \eta \text{ のとき} \\ 0 & \text{その他} \end{cases}.$$

これより次の定理を得る。

定理 3. 7

n -合祖過程 α_t において、祖先サイズ過程 $|\alpha_t|$ とジャンプ過程 $\{\mathfrak{R}_k\}$ は独立であり

$\alpha_t = \mathfrak{R}_{|\alpha_t|}$ と表される。また \mathfrak{R}_k の分布は次式で与えられる。

$$P(\mathfrak{R}_k = \xi | \mathfrak{R}_n = \Delta) = \frac{(n-k)! k! (k-1)!}{n! (n-1)!} \lambda_1! \lambda_2! \dots \lambda_k! \quad (3.19)$$

ここで $\lambda_1, \lambda_2, \dots, \lambda_k$ は ξ に含まれる k 個の同値類の各クラスサイズで、 $\lambda_1 + \dots + \lambda_k = n$ を満たす。

(証明) $\alpha_i = \mathfrak{R}_{|\alpha_i|}$ は 3. 1 節で述べたので、式 (3.19) を証明する。 n から始め、逆向きの数学的帰納法で示す。

(i) $k = n$ のとき、 $\lambda_1 = \dots = \lambda_n = 1$ なので、(3.19) の右辺 = 1 が成り立つ。

(ii) k のとき成り立つと仮定して、 $\xi \in E_n, |\xi| = k$ のとき $P_k(\xi) = P(\mathfrak{R}_k = \xi | \mathfrak{R}_n = \Delta)$

と書くことにする。

(iii) $k-1$ のとき、 $\eta \in E_n, |\eta| = k-1$ に対して

$$P_{k-1}(\eta) = \sum_{\xi} P(\mathfrak{R}_{k-1} = \eta | \mathfrak{R}_k = \xi) P_k(\xi) = \frac{2}{k(k-1)} \sum_{\xi \succ \eta} P_k(\xi)$$

η のクラスサイズを $\lambda_1, \lambda_2, \dots, \lambda_{k-1}$ (ただし $\lambda_1 + \dots + \lambda_{k-1} = n$) とする。 ξ に含まれる二つのクラスが 1 クラスに合併して η の $k-1$ 個のクラスが出来るので ξ のクラスサイズは $1 \leq i \leq k-1, 1 \leq \nu \leq \lambda_i - 1$ に対して、 $\lambda_1, \dots, \lambda_{i-1}, \nu, \lambda_i - \nu, \lambda_{i+1}, \dots, \lambda_{k-1}$ と η の i 番目のクラスを二つに分割する形で与えられる。故に全ての可能な i, ν について加えればよいので、 λ_i 個の遺伝子を ν と $\lambda_i - \nu$ の二つのクラスに分割する方法は $\frac{1}{2} \binom{\lambda_i}{\nu}$ 通りに注意して

$$\begin{aligned} P_{k-1}(\eta) &= \sum_{i=1}^{k-1} \sum_{\nu=1}^{\lambda_i-1} \frac{2}{k(k-1)} \times \frac{(n-k)! k! (k-1)!}{n!(n-1)!} \lambda_1! \dots \lambda_{i-1}! \nu! (\lambda_i - \nu)! \dots \lambda_{k-1}! \times \frac{1}{2} \binom{\lambda_i}{\nu} \\ &= \frac{(n-k)! (k-1)! (k-2)!}{n!(n-1)!} \sum_{i=1}^{k-1} \lambda_1! \dots \lambda_{i-1}! \lambda_{i+1}! \dots \lambda_{k-1}! \left\{ \sum_{\nu=1}^{\lambda_i-1} \nu! (\lambda_i - \nu)! \binom{\lambda_i}{\nu} \right\} \\ &= \frac{(n-k+1)! (k-1)! (k-2)!}{n!(n-1)!} \lambda_1! \lambda_2! \dots \lambda_{k-1}! \end{aligned}$$

よって $k-1$ のときも成り立ち証明された。

n -合祖過程から派生するマルコフ過程としてファミリーサイズ過程 $\{F_t; t \geq 0\}$ がある。
 $\alpha \in E_n$ がサイズ i の同値クラスを m_i 個 ($i = 1, 2, \dots, n$) 含むとき $f(\alpha) = (m_1, m_2, \dots, m_n)$

と書くことにする。 $\sum_{i=1}^n m_i$ は α に含まれる同値クラスの数であり、 $\sum_{i=1}^n i m_i = n$ である。

このとき、 n -合祖過程 α_t に対し、 $F_t = f(\alpha_t)$ によってファミリーサイズ過程 F_t を定義する。 α_t は n 個の異なるクラスに始まり、最終的に 1 つのクラスに合祖するのでファミリーサイズ過程は $F_0 = (n, 0, \dots, 0)$ に始まり $F_{\infty} = (0, \dots, 0, 1)$ に終わる。ファミリーサイ

ズ過程の推移確率 $P(F_t = M | F_0)$ (ただし $M = (m_1, m_2, \dots, m_n)$ 、 $F_0 = (n, 0, \dots, 0)$) はクラス

数を $\sum_{i=1}^n m_i = k$ とすると

$$\begin{aligned} P(F_t = M | F_0) &= P(|\alpha_t| = k | |\alpha_0| = n) \left\{ \sum_{\xi, f(\xi)=M} P(\mathfrak{R}_k = \xi) \right\} \\ &= P_t(n, k) \frac{(n-k)! k! (k-1)!}{n! (n-1)!} \lambda_1! \lambda_2! \dots \lambda_k! \left\{ \sum_{\xi, f(\xi)=M} 1 \right\} \\ f(\xi) = M \text{ を満たす } \xi \text{ の数は } & \frac{n!}{(m_1! \dots m_n!) (\lambda_1! \dots \lambda_k!)} \text{ なので} \end{aligned}$$

$$P(F_t = M | F_0) = P_t(n, k) \times \frac{(n-k)! (k-1)!}{(n-1)!} \times \frac{k!}{m_1! \dots m_n!} = P_t(n, k) \times \left\{ \frac{\binom{k}{M}}{\binom{n-1}{k-1}} \right\} \quad (3.20)$$

ここで $\frac{\binom{k}{M}}{\binom{n-1}{k-1}} = \frac{k!}{m_1! \dots m_n!}$ である。 (Tavare(1984))

3. 3. 2 集団からアレルが消失する速度

突然変異が存在しないとき、集団内の遺伝的多様性は時間とともに消失してゆく。この消失の速度を Felsenstein(1971)に従い遺伝子系図を利用して導いた Tavare(1984)の結果を紹介する。

時刻 t にサンプルした n 個の遺伝子が j 種のアレルタイプを含む確率を $Q_t(n, j)$ とする。

このサンプル遺伝子の系図を遡ると

$$Q_t(n, j) = \sum_{k=j}^n P_t(n, k) Q_0(k, j) \quad (3.21)$$

ここで $P_t(n, k)$ は n の遺伝子が k 個の祖先に由来する n -合祖過程の推移確率である。

$n \rightarrow \infty$ として全集団が j 種のアレルタイプを含む確率を $Q_t(j) = \lim_{n \rightarrow \infty} Q_t(n, j)$ とすると

(3.10)を使って

$$Q_t(j) = \lim_{n \rightarrow \infty} Q_t(n, j) = \sum_{k=j}^{\infty} P_t(k) Q_0(k, j) \quad (3.22)$$

ここで $Q_0(k, j)$ は時刻 $t = 0$ に取り出した k 個のサンプルがちょうど j 種のアレルタイプを持つ確率である。このとき次の補助定理が成り立つ。

補助定理 3. 8 (Tavare(1984), (7.4)式)

$t = 0$ の初期集団に K 種のアレルタイプが存在したと仮定し、頻度を p_1, p_2, \dots, p_K (ただ

し $\sum_{i=1}^K p_i = 1$) とする。このとき

$$Q_0(k, j) = \sum_{s=1}^j (-1)^{j-s} \binom{K-s}{j-s} \sum_s (p_{i_1} + \dots + p_{i_s})^k, \quad j = 1, 2, \dots, \min(k, K) \quad (3.23)$$

ここで \sum_s は与えられた各 s に対して $1 \leq i_1 < i_2 < \dots < i_s \leq K$ を満たすすべての (i_1, \dots, i_s)

についての和を表す。

(証明) アレルタイプを $\{1, 2, \dots, K\}$ とし、この中から j 種のアレル $1 \leq i_1 < i_2 < \dots < i_j \leq K$ を固定して考える。この j 種の中から r 種の部分集合を取り出し $I = \{i_1, i_2, \dots, i_r\}$, $r \leq j$ とする。集団からランダムに k 個の遺伝子をサンプルするとき、 $A(I) = A(i_1, \dots, i_r)$ をこの k 個のサンプルに含まれるアレルタイプは i_1, \dots, i_r の部分集合であるという事象、また

$B(I) = B(i_1, \dots, i_r)$ を含まれるアレルタイプがちょうど i_1, \dots, i_r であるという事象とする。

二つのアレルタイプの部分集合 $J_1 = (k_1, \dots, k_s), J_2 = (j_1, \dots, j_u)$ について $J_1 \neq J_2, J_1 \subseteq I, J_2 \subseteq I$ ならば、 $B(J_1) \cap B(J_2) = \emptyset$ 、すなわち $B(J_1)$ と $B(J_2)$ は排反事象である。よって

$A(I) = \sum_{J \subseteq I} B(J)$ なので $P(A(I)) = \sum_{J \subseteq I} P(B(J))$ 。このとき包含と排除の原理 (付録 D 参

照) により $P(B(I)) = \sum_{J \subseteq I} (-1)^{|I|-|J|} P(A(J))$ と表せる。ここで $i = |I|, j = |J|$ (各集合の要素数)。

これより $Q_0(k, j) = \sum_j P(B(i_1, i_2, \dots, i_j))$ 、ここで \sum_j は K 種の中の全ての j 種のアレルタイ

$1 \leq i_1 < i_2 < \dots < i_j \leq K$ についての和を表す。包含と排除の原理より

$$Q_0(k, j) = \sum_j P(B(i_1, i_2, \dots, i_j)) = \sum_j \sum_{L \subseteq J} (-1)^{j-s} P(A(L)) \quad (\text{ただし } s = |L| \leq |J| = j)$$

$$= \sum_{s=1}^j \sum_{L: |L|=s} P(A(L)) (-1)^{j-s} \left(\sum_{J(\supseteq L)} 1 \right)$$

ここで $\sum_{J(\supseteq L)} 1$ は集合 $L = (k_1, \dots, k_s)$ を部分集合として含む集合 J (ただし $|J| = j$) の数である。

これは K 種のアレルタイプから L に含まれる s 種のアレルを除いた $K-s$ 種のアレルから $j-s$ 種のアレルを選ぶ場合の数に等しいので $\sum_{J(\supseteq L)} 1 = \binom{K-s}{j-s}$ 。また $L = (k_1, \dots, k_s)$ のとき

$P(A(L)) = (p_{k_1} + \dots + p_{k_s})^k$ なので(3.23)を得る。

$k < j$ のとき明らかに $Q_0(k, j) = 0$ であるが、このとき(3.23)の右辺 = 0 も成り立つ。
これに注意し、(3.23)を((3.22)に代入すると

$$\begin{aligned} Q_i(j) &= \sum_{k=1}^{\infty} P_i(k) Q_0(k, j) = \sum_{k=1}^{\infty} P_i(k) \left\{ \sum_{s=1}^j (-1)^{j-s} \binom{K-s}{j-s} \sum_S (p_{i_1} + \dots + p_{i_s})^k \right\} \\ &= \sum_{s=1}^j (-1)^{j-s} \binom{K-s}{j-s} \left\{ \sum_S \sum_{k=1}^{\infty} P_i(k) (p_{i_1} + \dots + p_{i_s})^k \right\} \end{aligned} \quad (3.24)$$

二つのアレル A, a を仮定しアレル A の初期頻度が p のとき、 $f(p; t) = \sum_{k=1}^{\infty} P_i(k) p^k$ とすると $f(p; t)$ は時刻 t における集団がアレル A に固定している確率を表す。以上をまとめると

定理 3. 9 (Tavare(1984), (7.6)式)

時刻 t に集団に存在するアレルの種類数がちょうど j である確率は次式で与えられる。

$$Q_i(j) = \sum_{s=1}^j (-1)^{j-s} \binom{K-s}{j-s} \sum_S f(p_{i_1} + \dots + p_{i_s}; t) \quad (3.25)$$

これらの結果は Littler(1975), Griffiths(1979)でも得られている。

さらに集団内の平均のアレル数を求めると

$$\sum_{j=1}^{\infty} j Q_i(j) = \sum_{j=1}^{\infty} j \left\{ \sum_{k=1}^{\infty} P_i(k) Q_0(k, j) \right\} = \sum_{k=1}^{\infty} \left\{ \sum_{j=1}^{\infty} j Q_0(k, j) \right\} P_i(k)$$

アレルタイプを A_1, A_2, \dots, A_K 、初期頻度を p_1, p_2, \dots, p_K とすると

$$\sum_{j=1}^{\infty} j Q_0(k, j) = \sum_{j=1}^K \{1 - (1 - p_j)^k\} = K - \sum_{j=1}^K (1 - p_j)^k$$

$$\sum_{j=1}^{\infty} j Q_i(j) = \sum_{k=1}^{\infty} \left\{ K - \sum_{j=1}^K (1 - p_j)^k \right\} P_i(k) = K - \sum_{j=1}^K f(1 - p_j; t) \quad (3.26)$$

と表される。

3. 4 モランモデルとサブサンプルの遺伝子系図

Kingman の合祖過程は定理 3. 2、定理 3. 3 で示した様に Cannings の可換モデルからある条件の下で導かれるが、この範疇に入らないモデルとしてモランモデルがある。まず Watterson(1984)によるモランモデルの遺伝子系図について紹介し、さらにサンプル遺伝子の中からさらに遺伝子をサンプルしたサブサンプルを考えたとき、サンプルとそのサブサンプルの入れ子構造の遺伝子系図を考察した Saunderson et al.(1984)の結果を紹介する。

3. 4. 1 モランモデルの遺伝子系図

第2. 1節で紹介したN個の半数体生物集団のモランモデルを仮定する。すなわち毎世代1個の親が選ばれ $1-u$ の確率で親と同じタイプの子供を、確率 u で親と異なる全く新しいタイプの子供を1個生む。突然変異は無限対立遺伝子モデルを仮定する。その後、N個の親の中からランダムに選ばれた1個体が集団から除かれ元の集団サイズに戻る。以上を1世代とするモランモデルで、時刻 $t=0$ に集団からランダムに n 個の遺伝子をサンプルし、この遺伝子サンプルを $\{I_1, I_2, \dots, I_n\}$ とする。これらの遺伝子の中に次のような二種類の同値関係 $\{D_t, F_t; t=1, 2, 3, \dots\}$ を定義する。二つの遺伝子 I_k, I_j が突然変異を経由する事無く t 世代前に共通な祖先を持つとき、 $(I_k, I_j) \in D_t$ と書き、二つの遺伝子は D_t 同値という。

D_t は突然変異を経ずに t 世代前に同じ祖先を共有するもので分類されたクラスの集合である。 D_t に含まれるクラス数を $|D_t|$ で、 D_t に含まれるクラス ξ_i ($i=1, 2, \dots, |D_t|$)のサイズを λ_i とする。二つの遺伝子 I_k, I_j が1世代前から t 世代前までの何れかの時点で、ある突然変異

遺伝子を共通祖先として持つとき $(I_k, I_j) \in F_t$ と書き、二つの遺伝子は F_t 同値という。同様にクラス数を $|F_t|$ で、 F_t に含まれるクラス η_j ($j=1, 2, \dots, |F_t|$)のサイズを μ_j とする。 ξ_i は n 個のサンプル遺伝子の t 世代前まで突然変異を経由しない1つの祖先を共有するクラスを表わし、 η_j は突然変異を起こした1つの祖先を共有するクラスを表している。

二つの同値類による n 個の遺伝子 $\{I_1, I_2, \dots, I_n\}$ の分割を

$$R_t = \{\xi_i, i=1, 2, \dots, |D_t|; \eta_j, j=1, 2, \dots, |F_t|\} = \{\xi; \eta\} \text{ とする。 } \sum_{i=1}^p \lambda_i + \sum_{j=1}^q \mu_j = n \text{ (ただし}$$

$p=|D_t|, q=|F_t|$) である。モランモデルの遺伝子系図過程 R_t の推移確率を考えよう。

初期状態は $|D_0|=n, |F_0|=0, R_0 = \{\xi_i = (I_i), i=1, 2, \dots, n; \emptyset\}$ である。

(i) 第 t 世代の状態を $R_t = \{\xi; \eta\}$ 、 $R_{t+1} = \{\xi^*; \eta\}$ とする。Kingman の合祖過程と同様に ξ^* が ξ の中の二つのクラスを一つのクラスに融合することで得られるとき $\xi \prec \xi^*$ と書くことにする。 ξ の中のクラス ξ_i と ξ_j が1世代前に共通な親を持ち ξ^* でひとつのクラス $\xi_i \cup \xi_j$ となったと仮定すると、これは親 $\xi_i \cup \xi_j$ が突然変異を起こさずに子供 ξ_i を残した

場合と子供 ξ_j を残した場合の二通りあるので、その確率は $\xi < \xi^*$ のときに限り

$$P(R_{t+1} = \{\xi^*; \eta\} | R_t = \{\xi; \eta\}) = \frac{2(1-u)}{N^2} \text{となる。} \quad (2.2 \text{節 例2参照}) \quad (3.27)$$

(ii) t 世代前の祖先 ξ_i が、その1世代前のある親個体の突然変異を起こした子供であるとき、第 $t+1$ 世代では ξ_i はF同値類に属すクラスになる。これを $R_t = \{\xi; \eta\}$ および $R_{t+1} = \{\xi - \xi_i; \eta + \xi_i\}$ と表わすことにするとその推移確率は

$$P(R_{t+1} = \{\xi - \xi_i; \eta + \xi_i\} | R_t = \{\xi; \eta\}) = \frac{u}{N} \text{。} \quad (3.28)$$

(iii) 変化を起こさない確率は

$$P(R_{t+1} = R_t | R_t = \{\xi; \eta\}) = 1 - \frac{p(p-1)(1-u)}{N^2} - \frac{pu}{N}, \text{ ただし } p = |D_t| \quad (3.29)$$

これは(i)の場合の数が $\frac{p(p-1)}{2}$ 通り、(ii)が p 通りあることによる。

1世代当たりの推移確率が(3.27)(3.28)(3.29)で与えられる遺伝子系図のマルコフ連鎖 $\{R_t; t=0,1,2,\dots\}$ の推移はD同値類の中の二つのクラスの合祖かまたは突然変異による遷移なので、1世代の $\{R_t\}$ の推移確率はD同値類の数 $|D_t|$ にしか依存しない。Kingmanの合祖過程と同様に $\{R_t\}$ のジャンプ過程を $\{\mathfrak{R}_k; k=n, n-1, \dots, 0\}$ とすると $R_t = \mathfrak{R}_{|D_t|}$ が成り立

つ。 $R_t = \{\xi; \eta\}$ 、 $p = |D_t|$ 、 $q = |F_t|$ とすると、

$$P(R_t = \{\xi; \eta\}) = P(|D_t| = p)P(R_t = \{\xi; \eta\} | |D_t| = p) = P(|D_t| = p)P(\mathfrak{R}_p = \{\xi; \eta\})$$

が成り立つ。このとき次の定理が成り立つ。

定理3. 10

初期条件 $|D_0| = n, |F_0| = 0$ 、 $R_0 = \{\xi_i = (I_i), i=1,2,\dots,n; \phi\}$ のとき

$$P(|F_t| = q, R_t = \{\xi; \eta\} | |D_t| = k) = \frac{(n-k)!k!g^q}{n!(k+g)_{(n-k)}} \prod_{i=1}^k \lambda_i! \prod_{j=1}^q (\mu_j - 1)! \quad (3.30)$$

ただし $k=0,1,\dots,n$; $q=0,1,\dots,n-k$; $g = \frac{Nu}{1-u}$, $\sum_{i=1}^k \lambda_i + \sum_{j=1}^q \mu_j = n$ 。

(証明)

$P(\mathfrak{R}_p = \{\xi; \eta\}) = P(|F_t| = q, R_t = \{\xi; \eta\} \mid |D_t| = p)$ なので、ジャンプ過程の確率分布

$P(\mathfrak{R}_p = \{\xi; \eta\})$ が(3.30)の右辺で与えられることを p に関する逆向きの帰納法で示す。

(i) $p = n$ のとき、 $\mathfrak{R}_n = \{\xi_1, \dots, \xi_n; \emptyset\}$, $|\xi_i| = \lambda_i = 1$ ($1 \leq i \leq n$), $q = 0$ なので(3.30)の右辺

$$= \frac{0!n!\mathcal{G}^0}{n!(n+\mathcal{G})_{(0)}} = 1 \text{ が成り立つ。}$$

(ii) $p = k+1$ のとき成り立つと仮定して、ジャンプ過程の状態を R , W で表わすと

$$P(\mathfrak{R}_k = R) = \sum_W P(\mathfrak{R}_{k+1} = W)P(\mathfrak{R}_k = R \mid \mathfrak{R}_{k+1} = W) \quad (3.31)$$

ジャンプ過程 $\{\mathfrak{R}_k\}$ の W から R への状態変化は合祖又は突然変異によるので、場合分けをして考える。

(a) 合祖による遷移: W の D 同値類に属す二つのクラス ξ_{i_1} と ξ_{i_2} が合祖して R の一つのク

ラス $\xi_i = \xi_{i_1} \cup \xi_{i_2}$ になったと仮定しよう。 $|\xi_i| = \lambda_i$, $|\xi_{i_1}| = \lambda_{i_1}$ とすると $|\xi_{i_2}| = \lambda_i - \lambda_{i_1}$ 。

ξ_i を二つのクラス ξ_{i_1}, ξ_{i_2} に分割する方法は $\frac{1}{2} \binom{\lambda_i}{\lambda_{i_1}}$ 通り。 $1 \leq \lambda_{i_1} \leq \lambda_i - 1$ なので $W \rightarrow R$ が

合祖による遷移である確率は

$$P(\mathfrak{R}_k = R \mid \mathfrak{R}_{k+1} = W) = \frac{2(1-u)/N^2}{((k+1)k(1-u)/N^2) + (k+1)u/N} = \frac{2}{(k+1)(k+\mathcal{G})}$$

(3.26)で合祖による遷移の部分だけの和は帰納法の仮定より

$$\begin{aligned} \sum_W P(\mathfrak{R}_{k+1} = W) \frac{2}{(k+1)(k+\mathcal{G})} &= \frac{1}{(k+1)(k+\mathcal{G})} \sum_{i=1}^k \left[\sum_{s=1}^{\lambda_i-1} \binom{\lambda_i}{s} \left\{ \frac{(n-k-1)!(k+1)!\mathcal{G}^q}{n!(k+1+\mathcal{G})_{(n-k-1)}} \right. \right. \\ &\quad \left. \left. \times \lambda_1! \dots s! (\lambda_i - s)! \dots \lambda_k! \left(\prod_{j=1}^q (\mu_j - 1)! \right) \right\} \right] \\ &= \frac{(n-k-1)!\mathcal{G}^q k!}{n!(k+1+\mathcal{G})_{(n-k-1)}(k+\mathcal{G})} \left(\prod_{j=1}^q (\mu_j - 1)! \right) \left\{ \sum_{i=1}^k \sum_{s=1}^{\lambda_i-1} \binom{\lambda_i}{s} \lambda_1! \dots s! (\lambda_i - s)! \dots \lambda_k! \right\} \\ &= \frac{(n-k-1)!\mathcal{G}^q k!}{n!(k+\mathcal{G})_{(n-k)}} \left(\prod_{j=1}^q (\mu_j - 1)! \right) \left(\prod_{i=1}^k \lambda_i! \right) \left(\sum_{i=1}^k (\lambda_i - 1) \right) \quad (3.32) \end{aligned}$$

(b) 突然変異による遷移

$W = \{\xi_1, \dots, \xi_k, \eta_j; \eta_1, \dots, \eta_{j-1}, \eta_{j+1}, \dots, \eta_q\}$, $R = \{\xi_1, \dots, \xi_k; \eta_1, \dots, \eta_q\}$ とする。すなわち R の中の

F 同値類の中の $\eta_j (1 \leq j \leq q)$ が W の中の D 同値類に属していた一つのクラスの突然変異に

よって遷移した確率を考える。(a)と同様にして $P(\mathfrak{R}_k = R | \mathfrak{R}_{k+1} = W) = \frac{\mathcal{G}}{(k+1)(k+\mathcal{G})}$ 、

突然変異による変化の部分の和を取ると

$$\begin{aligned} \sum_W P(\mathfrak{R}_{k+1} = W) \frac{\mathcal{G}}{(k+1)(k+\mathcal{G})} &= \frac{\mathcal{G}}{(k+1)(k+\mathcal{G})} \sum_{j=1}^q \frac{(n-k-1)!(k+1)!\mathcal{G}^q}{n!(k+1+\mathcal{G})!} \left\{ \left(\prod_{i=1}^k \lambda_i! \right) \mu_j! \right\} \\ &\quad \times \prod_{\substack{s=1 \\ s \neq j}}^q (\mu_s - 1)! \\ &= \frac{(n-k-1)!k!\mathcal{G}^q}{n!(k+\mathcal{G})_{(n-k)}} \left(\prod_{i=1}^k \lambda_i! \right) \left(\prod_{j=1}^q (\mu_j - 1)! \right) \left(\sum_{j=1}^q \mu_j \right) \end{aligned} \quad (3.33)$$

(3.27)(3.28)を加え合わせると $\sum_{i=1}^k (\lambda_i - 1) + \sum_{j=1}^q \mu_j = n - k$ なので(3.30)を得る。

よって、全ての $k = n, n-1, \dots, 0$ について成り立つ。

さらに $|D_t|$ の分布について次の定理が成り立つ。

定理 3. 1 1

$$P(|D_t| = k) = \sum_{i=k}^n \left\{ 1 - \frac{2d_i}{N(N+\mathcal{G})} \right\}^t (-1)^{t-k} (2i+\mathcal{G}-1) \frac{(k+\mathcal{G})_{(t-1)} n_{[t]}}{k!(i-k)!(n+\mathcal{G})_{(t)}} \quad (3.34)$$

ただし、 $k = 0, 1, 2, \dots, n$; $d_i = \frac{i(i+\mathcal{G}-1)}{2}$, $x_{(n)} = x(x+1)\dots(x+n-1)$,

$x_{[n]} = x(x-1)\dots(x-n+1)$ 。 $k = 0$ のときは、 $i = 0$ の項は 1 とする。

(証明) 数学的帰納法を用いる。

(i) $k = n$ のとき (3.34) は $P(|D_t| = n) = \left(1 - \frac{2d_n}{N(N+\mathcal{G})} \right)^t$ 、これは第 t 世代まで全く変

化がない確率であり明らかに正しい。

(ii) $t = 0$ のとき(3.34)は(3.11)の $t = 0$ の時の右辺と一致するので $P(|D_0| = k) = \delta_{n,k}$ 。

$(k, t) = (n, 0)$ を出発点として k は $n, n-1, n-2, \dots, 0$; t は $0, 1, 2, \dots$ と二変数の数学的帰納法で証明する。全ての $\{(k, 0); k = n, n-1, \dots, 0\}$ および $\{(n, t); t = 0, 1, 2, \dots\}$ については成り立っているので、 $t-1$ のとき、 $k = n, n-1, \dots, 0$ について成り立ち、 t のとき $k = n, (n-1), \dots, p$ まで成り立つと仮定する。

(iii) 時刻 t および $k = p-1$ のとき

$$\begin{aligned} P(|D_t| = p-1) &= P(|D_{t-1}| = p)P(|D_t| = p-1 \mid |D_{t-1}| = p) \\ &\quad + P(|D_{t-1}| = p-1)P(|D_t| = p-1 \mid |D_{t-1}| = p-1) \\ &= P(|D_{t-1}| = p) \frac{2d_p}{N(N+g)} + P(|D_{t-1}| = p-1) \left\{ 1 - \frac{2d_{p-1}}{N(N+g)} \right\} \end{aligned}$$

$P(|D_{t-1}| = p), P(|D_{t-1}| = p-1)$ は帰納法の仮定より (3.34) で与えられるので代入し整理すると (3.34) を得る。

(3.34) で世代を $\left[\frac{N^2}{2}t\right]$ と置き、 $N^2/2$ 世代を単位時間とするタイムスケールを取り $N \rightarrow \infty$

の極限を取ると $\lim_{N \rightarrow \infty} \left\{ 1 - \frac{2d_i}{N(N+g)} \right\}^{[N^2t/2]} = \exp(-d_i t)$, $d_i = \frac{i(i+g-1)}{2}$ なので

$$\lim_{N \rightarrow \infty} P(|D_{[N^2t/2]}| = k) = \sum_{i=k}^n (-1)^{i-k} \frac{(2i+g-1)(k+g)_{(i-1)} n_{[t]}}{k!(i-k)!(n+g)_{(i)}} \exp\left[-\frac{i(i+g-1)}{2}t\right] \quad (3.35)$$

となり (3.11) と一致する。また、突然変異率 $u = 0$ のとき、遺伝子系図過程 R_t は合祖過程となり、常に $|F_t| = 0$ 。 (3.30) において $q = 0, g = 0$ とすると

$$P(R_t = \{\xi; \phi\} \mid |D_t| = k) = \frac{(n-k)!k!}{n!k_{(n-k)}} \lambda_1! \dots \lambda_k! = \frac{(n-k)!k!(k-1)!}{n!(n-1)!} \lambda_1! \dots \lambda_k! \quad (3.36)$$

これは (3.19) と一致する。突然変異率 $u > 0$ のとき、 $t \rightarrow \infty$ において $|D_t| \rightarrow 0$ となり、全てのサンプルはある突然変異を祖先とする遺伝子となり遺伝子系図過程 R_t は次の定常分布に収束する。

$$P(|F_\infty| = q, R_\infty = \{\phi; \eta\} \mid D_\infty = \phi) = \frac{g^q}{g_{(n)}} (\mu_1 - 1)! (\mu_2 - 1)! \dots (\mu_q - 1)! \quad (3.37)$$

これは第 4 章で論じる Ewens のサンプリング公式と呼ばれる式 (4.3) に対応し、定常状態の集団からランダムに取り出した遺伝子が q 種類のアレルタイプを含み、それぞれ μ_i 個

($i = 1, 2, \dots, q$) のサンプルを含んでいる確率を表わす。

3. 4. 2 サブサンプルの遺伝子系図

N 個の半数体生物から成る離散時間 Moran モデルを考える。集団から n 個の遺伝子をサンプルし、さらにこの n 個のサンプルから m 個 ($m \leq n$) のサブサンプルを取り出す。このサンプルとサブサンプルの遺伝子系図に関する Saunderson, Tavare and Watterson (1984) の結果を紹介する。突然変異はないものとする。 $A_1^N(t)$ を n 個のサンプルの t 世代前の祖先遺伝子の数、 $A_2^N(t)$ を m 個のサブサンプルの祖先遺伝子の数とする。前節の記法では $A_1^N(t) = |D_t|$ である。 $(A_1^N(t), A_2^N(t))$ の同時推移確率を

$$P_i((n, m), (k, j)) = P(A_1^N(t) = k, A_2^N(t) = j \mid A_1^N(0) = n, A_2^N(0) = m) \quad (3.38)$$

$A_1^N(t)$ の推移確率を $g_i(n, k) = P(A_1^N(t) = k \mid A_1^N(0) = n)$ とすると、定理 3. 1 1 より

$$g_i(n, k) = \sum_{i=k}^n \left\{ 1 - \frac{i(i-1)}{N^2} \right\}^t \frac{(-1)^{i-k} (2i-1) k_{(t-1)} n_{[t]}}{k!(i-k)! n_{(t)}} \quad (3.39)$$

二つの遺伝子系図過程 $A_1^N(t), A_2^N(t)$ についてそれぞれ祖先の数が r 個になる最初の時刻を $T_1(r) = \min\{t; A_1^N(t) = r\}$, $T_2(r) = \min\{t; A_2^N(t) = r\}$ とし、 $A_2^{N*}(r) = A_2^N(T_1(r))$ とする。 $A_2^{N*}(r)$ はサンプルの系図過程が r 個の祖先に成った時点でのサブサンプルの祖先遺伝子の数を表わす。このとき、次の補助定理が成り立つ。

補助定理 3. 1 2

$A_1^N(0) = n, A_2^N(0) = m$ とする。このとき、 $\{A_2^{N*}(n-i); i = 0, 1, 2, \dots, n-1\}$ は推移確率が次式で与えられる時間的に非一様なマルコフ連鎖である。

$$P(A_2^{N*}(r-1) = k-1 \mid A_2^{N*}(r) = k) = 1 - P(A_2^{N*}(r-1) = k \mid A_2^{N*}(r) = k) = \frac{k(k-1)}{r(r-1)} \quad (3.40)$$

(証明)

これは遺伝子系図過程でサブサンプルの祖先の数が減少するときは必ずサンプルの祖先の数が減少するので、(3.40) はサンプルの祖先数が r になった時点でサブサンプルの祖先が k のとき、サンプルの系図過程が合祖を起し $r-1$ になった時点でそれがサブサンプルの祖先の合祖でもある確率なので、明らかに場合の数の比より ${}_k C_2 / {}_r C_2 = \frac{k(k-1)}{r(r-1)}$ を得る。ま

た

$A_2^{N*}(r) = k$ のとき、明らかに $A_2^{N*}(r-1) = k$ または $k-1$ なので (3.40) を得る。

この結果は Moran モデルに限らず全ての可換モデルで成り立つ。

これより $A_2^N(t)$ の条件付分布が得られる。

定理 3. 1 3

$$P(A_2^N(t) = j | A_1^N(t) = k, A_1^N(0) = n, A_2^N(0) = m) \\ = \frac{(n-m)!(n-k)!m!(m-1)!k!(k-1)!}{(m-j)!(k-j)!n!(n-1)!j!(j-1)!} \times \frac{(n+j-1)!}{(k+m-1)!(n+j-k-m)!} \quad (3.41)$$

(証明)

$A_1^N(t) = k, A_2^N(t) = j, A_1^N(t+1) = k$ ならば $A_2^N(t+1) = j$ なので、左辺は

$P(A_2^{N*}(k) = j | A_1^N(0) = n, A_2^N(0) = m)$ に等しい。すなわち、時刻 $T_1(k)$ において(3.41)が成

り立つことを示せば十分である。 $P(A_2^{N*}(k) = j | A_1^N(0) = n, A_2^N(0) = m) = \phi(k, j)$ と書くと

補助定理 3. 1 2 より $k \leq n-1, j \leq m-1$ のとき

$$\phi(k, j) = \phi(k+1, j+1) \frac{j(j+1)}{k(k+1)} + \phi(k+1, j) \left\{ 1 - \frac{j(j-1)}{k(k+1)} \right\} \quad (3.42)$$

境界条件は $\phi(n, m) = 1, j > m$ のとき $\phi(n, j) = 0$ である。まず $j = m$ のとき(3.41)が成り立つことを証明する。 $k+1$ まで成り立つと仮定すると

$$\phi(k, m) = \phi(k+1, m) \left\{ 1 - \frac{m(m-1)}{k(k+1)} \right\} \\ = \frac{(n-m)!(n-k-1)!m!(m-1)!(k+1)!k!}{0!(k+1-m)!n!(n-1)!m!(m-1)!} \times \frac{(n+m-1)!}{(k+m)!(n-k-1)!} \times \frac{(k+m)(k-m+1)}{k(k+1)} \\ = \frac{(n-m)!(n-k-1)!m!(m-1)!k!(k-1)!}{0!(k-m)!n!(n-1)!m!(m-1)!} \times \frac{(n+m-1)!}{(k+m-1)!(n-k-1)!}$$

これより $j = m$ のとき全ての k について(3.41)が成り立つことが示された。最後に (k, j) ま

$$\phi(k, j-1) = \phi(k+1, j) \frac{j(j-1)}{k(k+1)} + \phi(k+1, j-1) \left\{ 1 - \frac{(j-1)(j-2)}{k(k+1)} \right\}$$

数学的帰納法により上式を整理すると証明される。

定理 3. 1 3 と(3.39)より $(A_1^N(t), A_2^N(t))$ の結合分布も得られる。

$$P_t((n, m), (k, j)) = P(A_1^N(t) = k, A_2^N(t) = j | A_1^N(0) = n, A_2^N(0) = m) \\ = P(A_2^N(t) = j | A_1^N(t) = k, A_1^N(0) = n, A_2^N(0) = m) P(A_1^N(t) = k | A_1^N(0) = n) \\ = g_t(n, k) \phi(k, j) \quad (3.43)$$

$g_1(n, k) = \begin{cases} n(n-1)/N^2 & (k = n-1 \text{ のとき}) \\ 1 - n(n-1)/N^2 & (k = n \text{ のとき}) \end{cases}$ より 1 世代当たりの $(A_1^N(t), A_2^N(t))$ の推

移確率 $P((n, m), (k, j))$ は

$$P((n, m), (k, j)) = \begin{cases} \frac{n(n-1)}{N^2} \times \frac{m(m-1)}{n(n-1)} = \frac{m(m-1)}{N^2} & (k = n-1, j = m-1) \\ \frac{n(n-1)}{N^2} \times \frac{(n-m)(n+m-1)}{n(n-1)} = \frac{(n-m)(n+m-1)}{N^2} & (k = n-1, j = m) \\ \left\{1 - \frac{n(n-1)}{N^2}\right\} \times 1 = 1 - \frac{n(n-1)}{N^2} & (k = n, j = m) \end{cases} \quad (3.44)$$

$A_1^N(T_2(j))$ はサブサンプルの系図過程が j 個の祖先に到達した最初の時刻におけるサンプル遺伝子の祖先の数である。故に $A_1(T_2(j)) = \max\{k; A_2^*(k) = j\}$ である。定理 3. 1 3 より次の定理を得る。

定理 3. 1 4

$$\begin{aligned} P(A_1^N(T_2(j) = k | A_1^N(0) = n, A_2^N(0) = m) &= P(A_2^{N*}(k+1) = j+1, A_2^{N*}(k) = j) \\ &= \frac{(n-m)!(n-k-1)!}{(m-j-1)!(k-j)!} \times \frac{m!(m-1)!k!(k-1)!}{n!(n-1)!j!(j-1)!} \times \frac{(n+j)!}{(k+m)!(n+j-k-m)!} \end{aligned}$$

(証明)

$A_1^N(T_2(j)) = k$ となるのは $A_1^N = k+1$ のとき $A_2^N = j+1$ 、かつ $A_1^N = k$ のとき $A_2^N = j$ 、すなわち合祖がサブサンプルの祖先内で起こることなので補助定理 3. 1 2 より

$$\begin{aligned} P(A_1^N(T_2(j) = k | A_1^N(0) = n, A_2^N(0) = m) &= P(A_1^{N*}(k+1) = j+1, A_2^{N*}(k) = j) \\ &= \phi(k+1, j+1) \times \frac{j(j+1)}{k(k+1)} \end{aligned}$$

これを整理して上の結果を得る。

特に $j=1$ とするとサブサンプルの遺伝子系図 $A_2(t)$ が共通な一つの祖先に到達したときのサンプル遺伝子の系図 $A_1(t)$ の条件付分布が得られる。

系 3. 1 5

$$P(A_1^N(T_2(1)) = k | A_1^N(0) = n, A_2^N(0) = m) = (m-1)(n+1) \frac{k!m!(n-k-1)!(n-m)!}{(m+k)!(n-1)!(n-k-m+1)!}$$

特に $k=1$ とすると

$$P(A_1^N(T_2(1))=1 | A_1^N(0)=n, A_2^N(0)=m) = P(T_1(1)=T_2(1)) = \frac{(m-1)(n+1)}{(m+1)(n-1)}$$

高い確率でサンプルとサブサンプルが同時に共通な祖先に到達することが分かる。 $n \rightarrow \infty$ とすると、全集団の遺伝子と m 個のサンプルが同時に共通な祖先に到達する確率は

$$P(A_1^N(T_2(1))=1 | A_1^N(0)=\infty, A_2^N(0)=m) = \frac{m-1}{m+1}$$

の数 j が 1 個ずつ減少する死滅過程であり、サンプルとサブサンプルの共通祖先の集合は高い確率である世代で一致することを意味する。これをサンプルとサブサンプル遺伝子系図過程のカップリング(coupling)と呼ぶ。ある世代でカップリングが生じるとそれ以後は二つの遺伝子系図は同一の系図過程となる。カップリングが最初に起こった世代におけるサンプルおよびサブサンプルの祖先の数を L とすると、初期条件 $A_1(0)=n, A_2(0)=m$ のとき、

その確率 $P(L=j)$ は $\phi(j, j) = P(A_2^{N*}(j)=j | A_1^N(0)=n, A_2^N(0)=m)$ を用いて

$$P(L=j) = \phi(j, j) - \phi(j+1, j+1) = 2j(n-m) \frac{m!(m-1)!(n-j-1)!(n+j+1)!}{n!(n-1)!(m+j)!(m-j)!} \quad (3.45)$$

突然変異を考慮に入れたモランモデルを考える。サンプル、サブサンプルの遺伝子系図で 3. 4. 1 節で定義した突然変異を経由しない祖先遺伝子の同値クラス類 D_i に着目したブ

ロセスを考えよう。すなわち $A_1^N(t) = |D_i|$ であり、 $A_1^N(t), A_2^N(t)$ はサンプル、サブサンプ

ルの突然変異を経由していない祖先遺伝子の数を表わす。 $A_1^N(t)$ の 1 世代あたり推移確率

を $h_{n,k} = P(A_1^N(1)=k | A_1^N(0)=n)$ とすると、(3.27)(3.28)(3.29)より

$$h_{n,k} = \begin{cases} 1 - \frac{n(n+g-1)}{N(N+g)} & (k=n) \\ \frac{n(n+g-1)}{N(N+g)} & (k=n-1) \\ 0 & \text{その他} \end{cases} \quad \text{ただし } g = \frac{Nu}{1-u}. \quad (3.46)$$

$A_1^N(t)$ の推移確率を $h_t(n, k) = P(A_1^N(t)=k | A_1^N(0)=n)$ とすると $h_t(n, k) = P(|D_t|=k)$ よ

り(3.34)式で与えられる。 $(A_1^N(t), A_2^N(t))$ の同時分布を

$P_t((n, m), (k, j)) = P(A_1^N(t)=k, A_2^N(t)=j | A_1^N(0)=n, A_2^N(0)=m)$ とすると、1 世代当たり

の推移確率は

$$P_1((n, m), (k, j)) = \begin{cases} 1 - \frac{n(n+g-1)}{N(N+g)} & (k = n, j = m) \\ \frac{(n-m)(n+m+g-1)}{N(N+g)} & (k = n-1, j = m) \\ \frac{m(m+g-1)}{N(N+g)} & (k = n-1, j = m-1) \\ 0 & \text{その他} \end{cases} \quad (3.47)$$

例えば $(k, j) = (n-1, m-1)$ となるのはサンプル遺伝子の祖先 D_i が合祖または突然変異によって n から $n-1$ に減少しそれがサブサンプルの系図に含まれる場合なので、その確率は $\frac{n(n+g-1)}{N(N+g)} \times \frac{m(m+g-1)}{n(n+g-1)} = \frac{m(m+g-1)}{N(N+g)}$ となる。

$(A_1^N(t), A_2^N(t))$ の状態空間は $\mathfrak{S} = \{(k, j) \mid k = 0, 1, \dots, n; j = 0, 1, \dots, \min(m, k)\}$ である。

定理 3. 1 3 の拡張として次の定理を得る。

定理 3. 1 6

$$\begin{aligned} P(A_2^N(t) = j \mid A_1^N(t) = k, A_1^N(0) = n, A_2^N(0) = m) \\ = \frac{(n-m)! m! (n-k)! \Gamma(m+g) \Gamma(n+j+g) \Gamma(k+g)}{n! (n-k-m+j)! j! (k-j)! (m-j)! \Gamma(n+g) \Gamma(j+g) \Gamma(m+k+g)} \end{aligned} \quad (3.48)$$

(証明)

補助定理 3. 1 2 と同様にして $\{A_2^{N^*}(n-i); i = 0, 1, 2, \dots, n\}$ は推移確率が次式で与えられる非斉次のマルコフ連鎖である。

$$P(A_2^{N^*}(k-1) = j-1 \mid A_2^{N^*}(k) = j) = 1 - P(A_2^{N^*}(k-1) = j \mid A_2^{N^*}(k) = j) = \frac{j(j+g-1)}{k(k+g-1)}$$

定理 3. 1 3 と同様にして(3.48)の右辺は $\phi(k, j) = P(A_2^{N^*}(k) = j \mid A_1(0) = n, A_2(0) = m)$

に等しい。(3.42)と同様に次式を得る。

$$\phi(k, j) = \phi(k+1, j+1) \frac{(j+g)(j+1)}{(k+g)(k+1)} + \phi(k+1, j) \left\{ \frac{(k-j+1)(k+j+g)}{(k+g)(k+1)} \right\} \quad (3.49)$$

境界条件 $\phi(n, j) = \begin{cases} 1 & (j = m) \\ 0 & (j > m) \end{cases}$ の下で、これより(3.43)を得る。

$(A_1^N(t), A_2^N(t))$ の同時分布 $P_i((n, m), (k, j))$ は(3.48)と $A_1(t)$ の推移確率 $h_i(n, k)$ を用いて

$$P_i((n, m), (k, j)) = h_i(n, k) \phi(k, j) \quad (3.50)$$

またサブサンプルの系図過程が $A_2^N(t) = j$ となる最初の時刻でのサンプルの系図過程

$A_1^N(T_2(j))$ について、定理3.14と同様に次の定理が成り立つ

定理3.17

$$P(A_1^N(T_2(j) = k | A_1^N(0) = n, A_2^N(0) = m) = P(A_2^{N^*}(k+1) = j+1, A_2^{N^*}(k) = j) \\ = \frac{(n-m)!(n-k-1)!m!k!\Gamma(m+\mathcal{G})\Gamma(n+j+1+\mathcal{G})\Gamma(k+\mathcal{G})}{(m-j-1)!(k-j)!j!(n-k+j-m)!n!\Gamma(n+\mathcal{G})\Gamma(j+\mathcal{G})\Gamma(m+k+1+\mathcal{G})}$$

ただし $0 \leq j \leq k \leq n-m+j$ 。

3.4.3 サブサンプルの連続時間遺伝子系図過程

離散時間 Moran モデルを用いて考察してきたが、 $N^2/2$ 世代を単位時間とし集団のサイズ N を無限に大きくすると連続時間の系図過程が得られる。まず突然変異がないとき、(3.34)より

$$g_t(n, k) = \lim_{N \rightarrow \infty} P(A_1^N(\frac{N^2 t}{2}) = k | A_1^N(0) = n) = \sum_{i=k}^n \exp[-\frac{i(i-1)}{2}t] \frac{(-1)^{i-k} (2i-1)k_{(i-1)}n_{(i)}}{k!(i-k)!n_{(i)}}$$

定理3.13はサンプルの系図過程 $A_1(\cdot)$ がジャンプした時点でのサブサンプルの過程 $A_2(\cdot)$ の状態のみに依存しているので、連続時間モデルにおいてもそのジャンプ過程に着目すればそのまま成り立つ。定理3.13で $n \rightarrow \infty$ とすると

$$P(A_2(t) = j | A_1(t) = k, A_1(0) = \infty, A_2(0) = m) \\ = \lim_{n \rightarrow \infty} \frac{(n-m)!(n-k)!m!(m-1)!k!(k-1)!}{(m-j)!(k-j)!n!(n-1)!j!(j-1)!} \times \frac{(n+j-1)!}{(k+m-1)!(n+j-k-m)!} \\ = \left\{ \binom{k}{j} \binom{m-1}{j-1} \right\} / \binom{k+m-1}{k-1} \quad j = 1, 2, \dots, \min(m, k) \quad (3.51)$$

これは全集団が時間 t 遡った祖先集団の k 個の祖先に由来するという条件の下で集団からランダムに取り出された m 個の遺伝子が同じ祖先集団の j 個の祖先に由来する確率を表わしている。(3.46)の分母は j 個の個体を k 個の祖先に割り当てる重複組み合わせの数 $\binom{k+m-1}{k-1}$ に等しく、分子は k 個の祖先から j 個の祖先を選ぶ組み合わせ $\binom{k}{j}$ と m 個の個体を j 個の祖先に空部屋が無いように割り当てる方法 $\binom{m-1}{j-1}$ の積となっている。すなわちランダムな組み合わせから期待される確率であり、中立遺伝子ということから帰結される当然の結果でもある。

次に突然変異を仮定すると、離散モデルから $N \rightarrow \infty$ として連続時間近似を行ない同様に

$\lim_{N \rightarrow \infty} (A_1^N([\frac{N^2 t}{2}], A_2^N([\frac{N^2 t}{2}])) = (A_1(t), A_2(t))$ とすると $(A_1(t), A_2(t))$ は空間 \mathfrak{S} の連続時間マルコフ連鎖で、その生成作用素は次式のようになる。

$$Q((n, m), (k, j)) = \begin{cases} -\frac{n(n+g-1)}{2} & (k=n, j=m) \\ \frac{(n-m)(n+m+g-1)}{2} & (k=n-1, j=m) \\ \frac{m(m+g-1)}{2} & (k=n-1, j=m-1) \end{cases} \quad (3.52)$$

3. 5 集団サイズの変動を含む遺伝子系図

これまでの議論は全て集団の大きさが一定という前提で考えてきた。しかし、現実の集団は常に様々な環境変動等の要因による個体数の変動がある。遺伝子系図におけるその影響を考えてみよう。離散世代モデルで集団のサイズが世代 τ の関数として $\{N(\tau); \tau=0, 1, 2, \dots\}$ 、ただし $N(0) = N$ と与えられると仮定する。この集団から任意に二つの遺伝子を取り出したとき、共通な祖先に達するまでの時間（世代数）を $\tau_2(N)$ とすると $[Nt]$ 世代以上合祖が起きない確率は $P(\tau_2(N) > [Nt]) = \prod_{\tau=1}^{[Nt]} \left(1 - \frac{1}{N(\tau)}\right)$ となる。両辺の対

数を取ると $\log P(\tau_2(N) > [Nt]) = \sum_{\tau=1}^{[Nt]} \log \left(1 - \frac{1}{N(\tau)}\right)$

集団サイズの相対比を表す関数として $f_N(t) = \frac{N([Nt])}{N} = \frac{N(\tau)}{N}$ ($\frac{\tau-1}{N} < t \leq \frac{\tau}{N}$ のとき)

とする。 $f_N(t)$ は t の階段関数である。さらに $f(t) = \lim_{N \rightarrow \infty} f_N(t)$ とする。

$x > 1$ のとき不等式 $x \leq -\log(1-x) \leq \frac{x}{1-x}$ より $x = \frac{1}{N(\tau)} = \frac{1}{N f_N(\tau)}$ とすると、

$$\frac{1}{N} \sum_{\tau=1}^{[Nt]} \frac{1}{f_N(\tau)} \leq -\sum_{\tau=1}^{[Nt]} \log \left(1 - \frac{1}{N(\tau)}\right) \leq \frac{1}{N} \sum_{\tau=1}^{[Nt]} \frac{1}{f_N(\tau) - (1/N)}$$

従って $N \rightarrow \infty$ の極限を取り、 $\tau_2 = \lim_{N \rightarrow \infty} T_2(N)/N$ とすると、区分別積分法により

$$\log P(\tau_2 > t) = \lim_{N \rightarrow \infty} \log P(\tau_2(N) > [Nt]) = \lim_{N \rightarrow \infty} \sum_{\tau=1}^{[Nt]} \log \left(1 - \frac{1}{N(\tau)}\right) = -\int_0^t \frac{1}{f(s)} ds$$

すなわち $P(\tau_2 > t) = \exp[-\Lambda(t)]$ 、ただし $\Lambda(t) = \int_0^t \frac{1}{f(s)} ds$ 。集団サイズが一定のときは

$f(t) = 1$ として $P(\tau_2 > t) = \exp[-t]$ となる。十分時間を遡ると任意の二つの遺伝子は確率

1 で共通祖先を持つので $\lim_{t \rightarrow \infty} P(\tau_2 > t) = 0$ より $\lim_{t \rightarrow \infty} \Lambda(t) = \infty$ を仮定する。 τ_2 の分布密度は

$$\frac{d}{dt} P(\tau_2 \leq t) = \frac{1}{f(t)} \exp[-\Lambda(t)], \text{ 平均は } E[\tau_2] = \int_0^{\infty} \exp[-\Lambda(t)] dt \text{ である。}$$

一般に遺伝子の数が n 個のとき、最初の合祖までの待ち時間を τ_n とすると、同様の議論により $N \rightarrow \infty$ の極限を取ると、 $P(\tau_n > t) = \exp[-\frac{n(n-1)}{2} \Lambda(t)]$ となる。集団サイズが変動

する場合も 3. 1 節と同様に n -合祖過程 α_t^ν を導くことができるが、これは時間的に非一様な純粋死滅過程であり、上添え字 ν は変動サイズモデルを表わす。集団サイズ変動モデルでの合祖過程 α_t の祖先の数を $A^\nu(t) = |\alpha_t^\nu|$ とすると、微小時間 h での推移確率は

$$P(A^\nu(t+h) = j | A^\nu(t) = k) = \begin{cases} \frac{k(k-1)}{2f(t)} h + o(h) & (j = k-1) \\ 1 - \frac{k(k-1)}{2f(t)} h + o(h) & (j = k) \\ o(h) & (\text{その他}) \end{cases} \quad (3.53)$$

各状態 $A^\nu(\cdot) = n, n-1, \dots, 2$ での滞在時間を $\tau_n, \tau_{n-1}, \dots, \tau_2$ とすると、その同時分布密度は

$$P(\tau_n = t_n, \tau_{n-1} = t_{n-1}, \dots, \tau_2 = t_2) = \prod_{j=2}^n \left\{ \frac{j(j-1)}{2f(s_j)} \exp[-\frac{j(j-1)}{2} \{\Lambda(s_j) - \Lambda(s_{j+1})\}] \right\} \quad (3.54)$$

ここで $0 \leq t_n, \dots, t_2 < \infty$, $s_{n+1} = 0$, $s_n = t_n$, $s_j = t_j + t_{j-1} + \dots + t_n$ ($j = 2, \dots, n-1$)。

$S_j = \tau_j + \tau_{j+1} + \dots + \tau_n$ とすると S_2 は一つの共通祖先までの合祖時間であり、条件

$$S_{j+1} = s \text{ の下で } \tau_j \text{ の分布は } P(\tau_j > t | S_{j+1} = s) = \exp[-\frac{j(j-1)}{2} \{\Lambda(t+s) - \Lambda(s)\}]$$

祖先数の過程 $A^\nu(t)$ は時間的に非一様なマルコフ連鎖である。集団サイズ一定の Kingman の n -合祖過程の祖先数を $A(t) = |\alpha_t|$ とすると、集団サイズ変動における過程 $A^\nu(t)$ は推移

確率 $P(A^\nu(t+h) = j | A^\nu(t) = k)$ が (3.53) 式で与えられることから

$A^\nu(t) = A(\Lambda(t))$, ($A^\nu(0) = A(0) = n$, $t \geq 0$) と表現される。すなわち、集団サイズ変動モ

デルの合祖過程 $A^\nu(t)$ は集団サイズ一定の n -合祖過程の関数 $\Lambda(t)$ による決定論的時間変更によって得られることが分かる。

Kingman の n -合祖過程の推移確率を $g_{n,k}(t) = P(A(t) = k | A(0) = n)$ とすると

過程 $A^\nu(t)$ の推移確率は次式で与えられる。

$$P(A^\nu(t) = k | A^\nu(0) = n) = g_{n,k}(\Lambda(t)) = \sum_{j=k}^n \frac{(2j-1)(-1)^{j-k} k_{(j-1)} n_{[j]}}{k!(j-k)! n_{(j)}} \exp\left[-\frac{j(j-1)}{2} \Lambda(t)\right] \quad (3.55)$$

(Tavare(2004))

この章で紹介したモデルは全て半数体生物、すなわち雌雄の区別がない生物という条件のモデルであるが、ヒトをはじめとして、有性生殖を行う生物についての合祖過程の導出は厳密には証明が必要である。この問題については、Möhle(1998c)等を参照されたい。また、合祖過程が導かれるためには、定理 3. 2 で述べた条件 $\lim_{N \rightarrow \infty} \text{Var}(\nu_1) = \sigma^2 > 0$ 、すべての

$p \geq 1$ に対して $\text{Sup}_N E[\nu_1^p] < \infty$ を満たさなければならない。これは、個々の個体が次世代に

残す子供の数の分散が $N \rightarrow \infty$ においても有限という条件であるが、例えば 1 個体の子供が全集団のある一定の割合 (例えば 10% など) を占めるような確率が正であるような場合、あるいは親世代の各個体は繁殖能力に違いはないが、ランダムに選ばれた数個体のみが次世代に子供を残せると言う様な場合この条件は満たされない。このようなとき、サンプルした 3 個体以上の個体が同時に合祖するという事態が生じる。このように多数の個体が同時に合祖を起こすようなプロセスは Λ -Coalescent と呼ばれているが、この話題については、Sagitov(1999), Pitman(1999), Möhle and Sagitov(1999, 2003)などを参照されたい。