

4. Ewens のサンプリング公式、ポアソン・ディリクレ分布とポリアの壺様モデル

4. 1 Ewens のサンプリング公式

生物集団においては一つの種内においても多くの遺伝的多様性が存在する。古くはエンドウ豆の色、シヨウジョウバエの羽の変異、染色体の逆位、血液型など表現形として明確に識別できるものに限られていたが、近年、分子レベルの解析の進歩とともに、タンパク質のアミノ酸配列の多型、遺伝子の DNA 塩基配列の個人差まで識別可能となってきた。生物集団内の遺伝的多様性の推定は集団遺伝学の主要な問題の一つである。遺伝的多様性の源は種々の突然変異であり、DNA 塩基配列がヒトでは 30 億塩基対ということを考えてみると、生じる突然変異はほとんど常に新しいタイプと考えられる。このように突然変異は常に集団中にこれまで存在しなかった新しいタイプと仮定するモデルを無限対立遺伝子モデルという。Ewens(1972)は中立な無限対立遺伝子モデルの下で、集団からランダムに取り出した n 個の遺伝子サンプルの中に存在する異なる対立遺伝子 (アレル) の数、およびサンプル内アレル構成の分布について「Ewens のサンプリング公式」と呼ばれる分布を求めた。対立遺伝子 (アレル) は自然選択に対して中立と仮定する。個々のアレルタイプに個性はないので、我々は集団から取り出した n 個のサンプルのアレルタイプによる分割構造に興味がある。そこでサンプル中に i 個含まれているアレルタイプの数を $a_i (i=1,2,\dots,n)$ としよう。例えば 10 個の遺伝子をサンプルしたとき、アレル A が 3 個、B が 1 個、C が 1 個、D が 2 個、E が 3 個含まれているとすると $a_1 = 2, a_2 = 1, a_3 = 2$ 、他は全て $a_i = 0$ である。 $a = (a_1, a_2, \dots, a_n)$ は遺伝子サンプル中のアレル構成による分割構造を表現している。全てが同じアレルタイプのときは $a_n = 1, a_1 = \dots = a_{n-1} = 0$ である。 $\sum_{i=1}^n ia_i = n$ (サン

ル数) であり、 $|a| = \sum_{i=1}^n a_i$ はサンプル中の異なるアレルタイプ数を表す。

定理 4. 1

集団からランダムに取り出した遺伝子サンプルのアレル分割構造が $a = (a_1, a_2, \dots, a_n)$ である確率は次の式で与えられる。

$$P(a_1, a_2, \dots, a_n) = \frac{n!}{(\mathcal{G})_n} \prod_{j=1}^n \frac{\mathcal{G}^{a_j}}{j^{a_j} a_j!} \quad (\text{Ewens のサンプリング公式}) \quad (4.1)$$

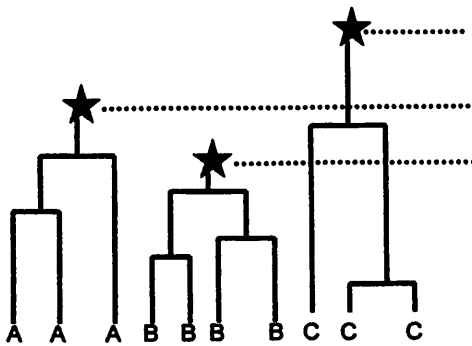
$$\text{またアレルタイプの数については } P(|a| = k) = \frac{\mathcal{G}^k S(n, k)}{\mathcal{G}(\mathcal{G}+1)\dots(\mathcal{G}+n-1)} \quad (4.2)$$

ここで、 $S(n, k)$ は第 1 種スターリング数 (付録(C.1)参照) である。

(証明)

これらの結果は Ewens(1972), Karlin&MacGregor(1972), Watterson(1984), Kingman (1982) などにより求められているが、ここでは第 3 章で紹介した突然変異存在下での遺伝子系図を用いて求めてみよう。まず p 種のアレルを考え、集団から n 個のサンプルを取り出

し、サンプルには1番からn番までナンバーをつける。このサンプル中に A_1, A_2, \dots, A_p の p 種のアレルを順に、それぞれ k_1 個、 k_2 個、 \dots 、 k_p 個ずつ取り出す確率を求める (ただし $k_1 + \dots + k_p = n$)。そのためには、n 個のサンプルは A_1, A_2, \dots, A_p の p 種のアレルに対応して k_1 個、 k_2 個、 \dots 、 k_p 個の p 個のクラスに分けられ、各々のクラスは同じタイプのアレルで一つの共通な mutant を祖先とする一つの血族同値類に属することになる。



上記の条件を満たす系図を描いたとき、i 番目のクラスの k_i 個の遺伝子の系図を遡ると $k_i - 1$ 回の coalesce と最後の突然変異の合計 k_i 回の事象、従って全部で $k_1 + \dots + k_p = n$ 回の事象が存在するはずである。系図上で祖先の数 j のとき、生じる事象がその時点で k 個の祖先を持つアレルクラスで起こる合祖である確率は $\frac{k(k-1)/2}{\{j\theta + j(j-1)\}/2}$ であり、突然変異で

ある確率は $\frac{\theta/2}{\{j\theta + j(j-1)\}/2}$ である。さらに、条件を満たす可能な系図は各 i 番目のクラ

スの遺伝子について合祖(coalesce)と突然変異(mutation)合わせて k_i 個の event があり、全部で合計 n 回の事象を順に並べる方法の数だけあるので $\frac{n!}{\prod_{i=1}^p k_i!}$ 通りとなる。よって1から

$$\prod_{i=1}^p k_i!$$

n までナンバーのついた遺伝子中に A_1, A_2, \dots, A_p の p 種のアレルをそれぞれ k_1 個、

k_2 個、 \dots 、 k_p 個ずつ含まれる確率は次式で与えられる。

$$P(k_1, k_2, \dots, k_p) = \left[\prod_{j=n}^1 \frac{1}{\{j\theta + j(j-1)\}/2} \right] \left\{ \prod_{i=1}^p \left(\prod_{j=k_i}^2 \frac{j(j-1)}{2} \right) \right\} \left(\frac{\theta}{2} \right)^p \times \frac{n!}{\prod_{i=1}^p k_i!}$$

$$= \frac{\mathcal{G}^p}{(\mathcal{G})_n} \prod_{i=1}^p (k_i - 1)! \quad \text{ただし } (\mathcal{G})_n = \mathcal{G}(\mathcal{G}+1)(\mathcal{G}+2)\dots(\mathcal{G}+n-1) \quad (4.3)$$

中立な遺伝子の場合、その個々の遺伝子のアレルタイプに興味は無く、サンプル中のアレルタイプ数とその構成に興味がある。そこで a_i ($1 \leq i \leq n$) を i 個の遺伝子がサンプル中に存在するアレルタイプの数とすると、 $a = (a_1, a_2, \dots, a_n)$, $\sum_{i=1}^n i a_i = n$, $\sum_{i=1}^n a_i = p$ を満たす。

例えば、 $(k_1, k_2, k_3, k_4, k_5) = (5, 5, 3, 3, 1)$ と $(k_1, k_2, k_3, k_4, k_5) = (3, 5, 3, 5, 1)$ は同じアレル分割 $a = (a_i; 1 \leq i \leq 17, a_1 = 1, a_3 = 2, a_5 = 2, a_i = 0 (i \neq 1, 3, 5))$ を持つ。 n 個のサンプル遺伝子に対してその中に、分割構成 (a_1, \dots, a_n) を持つような対立遺伝子を含んでいる確率 $P(a_1, a_2, \dots, a_n)$ を求めてみよう。ある一つの分割構造 (a_1, \dots, a_n) に対して、 $\sum_{i=1}^n a_i = p$ のとき、 i 個のサンプルを含むアレルタイプが a_i 種あるので、ラベルの与え方は $\prod_{i=1}^p a_i!$ 通りある。

その各々のアレルタイプ内で i 個の遺伝子について順番の付け方は $i!$ 通り。 n 個のサンプル遺伝子について、アレルタイプと、そのタイプ内での順番によって $\frac{n!}{\prod_{j=1}^n \{(j!)^{a_j} a_j!\}}$ 通りのラ

ベリングの方法がある。この因子を(4.3)式に掛けて次のサンプリング公式(4.1)を得る。

$$P(a_1, a_2, \dots, a_n) = \frac{n!}{(\mathcal{G})_n} \prod_{j=1}^n \frac{\mathcal{G}^{a_j}}{j^{a_j} a_j!} \quad \text{また付録(C.1)の定理 C2、及び } \sum_{k=1}^n P(|a|=k) = 1 \text{ より}$$

$$P(|a|=k) = \sum_{\substack{a \\ |a|=k}} P(a_1, a_2, \dots, a_n) = \frac{\mathcal{G}^k}{(\mathcal{G})_n} \sum_{\substack{a \\ |a|=k}} \frac{n!}{j^{a_j} a_j!} = \frac{\mathcal{G}^k S(n, k)}{\mathcal{G}(\mathcal{G}+1)\dots(\mathcal{G}+n-1)} \quad \text{を得る。}$$

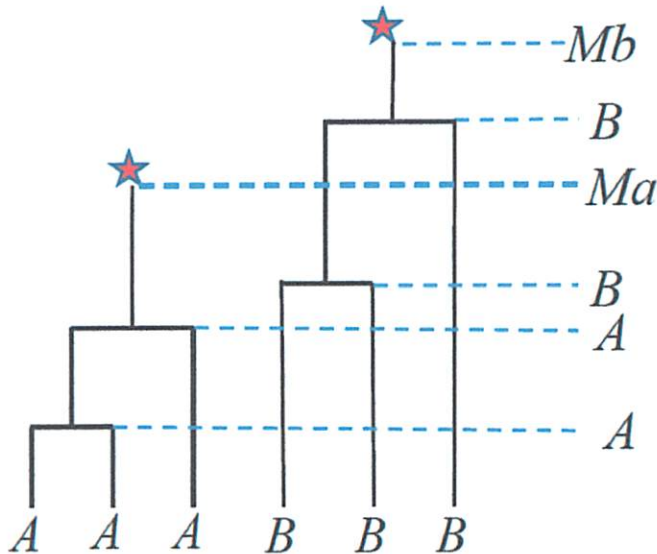
4. 2 アレルの年齢

集団からランダムに取り出した遺伝子の系図を遡るとそのアレルタイプの共通祖先、すなわちそのタイプの起源となる突然変異を起こした祖先遺伝子に到達する。サンプル中の各アレルタイプの起源に到達するまでの時間をそのアレルタイプの年齢という。アレルの年齢については古くは Kimura&Ohta(1973), Maruyama(1974)など多くの研究がある。ここでは遺伝子系図の結果を使ってアレルの年齢を求める。 n 個のサンプルを取り出したとき、その中に A_1, A_2, \dots, A_p の p 種のアレルがそれぞれ k_1 個、 k_2 個、 \dots k_p 個ずつ含まれ、

かつ A_1 アレルが最も若く、次に A_2 アレル、 \dots 、最も古いのが A_p アレルとなる確率

$P(A_1 \triangleright A_2 \triangleright \dots \triangleright A_p; (k_1, k_2, \dots, k_p))$ を求めよう。ここで記号 $A_1 \triangleright A_2$ はアレル A_1 が A_2

よりも若いことを意味する。そのための条件は起源突然変異 A_1, A_2, \dots, A_p がその年齢の順に系図上で生じることである。たとえば2種のアレル A, B が n 個のサンプル中にそれぞれ $n-k$ 個と k 個含まれ、かつ $A \triangleright B$ となる確率を考えてみよう。



上図は $n=6, k=3$ の場合の図である。アレル A での合祖を A 、アレル B での合祖を B 、 A, B の起源となる突然変異を Ma, Mb で表し、上の系図を古い事象から順に表すと

(Mb, B, Ma, B, A, A) と表示される。一般に A, B が n 個のサンプル中にそれぞれ $n-k$ 個と k 個含まれるとき、 $n-k-1$ 個の A と $k-1$ 個の B 、および Ma, Mb の合計 n 個の記号を左端が Mb で始まり、 B の位置は残り $n-1$ 個の場所から $k-1$ 個選ぶ方法だけある。 Ma は残り $n-k$ 個の場所の最も左の位置に置き、他は全て A となる。条件を満たす可能な系図の総数は $\frac{(n-1)!}{k!(n-1-k)!}$ 通りとなる。同様に、 A_1, A_2, \dots, A_p の p 種のアレルがそれぞれ k_1 個、 k_2

個、 \dots 、 k_p 個ずつ含まれる場合には、合祖と突然変異を合わせた計 n 回の事象に対応して

n 個の席を横一列に準備し、左から右へ事象の古い順に並べる。最も古いアレル A_p からそ

の起源となる突然変異の事象は左端の位置に置かれ、アレル A_p で生じた $k_p - 1$ 回の合祖事

象を残り $n-1$ 個の場所に入れるので $\binom{n-1}{k_p-1}$ 通りある。次のアレル A_{p-1} について、起源の

突然変異事象は残りの中の最も左端の位置に置かれ、 $k_{p-1} - 1$ 回と合祖事象は残り

$n - k_p - 1$ 個の位置から選ぶ。これを繰り返して次式を得る。

$$\begin{aligned}
& P(A_1 \triangleright A_2 \triangleright \dots \triangleright A_p; (k_1, k_2, \dots, k_p)) \\
&= \frac{g^p}{(g)_n n!} \left[\prod_{i=1}^p \{k_i!(k_i-1)!\} \right] \times \binom{n-1}{k_p-1} \binom{n-k_p-1}{k_{p-1}-1} \dots \binom{k_1-1}{k_1-1} \\
&= \frac{g^p}{(g)_n n!} \left\{ \prod_{i=1}^p k_i! \prod_{i=1}^p (k_i-1)! \right\} \frac{(n-1)!}{(k_p-1)!(n-k_p)!} \times \frac{(n-k_p-1)!}{(k_{p-1}-1)!(n-k_p-k_{p-1})!} \times \dots \\
&= \frac{g^p}{(g)_n} \times \frac{\prod_{i=1}^p k_i!}{n(n-k_p)(n-k_p-k_{p-1}) \dots k_1} = \frac{g^p}{(g)_n} \times \frac{\prod_{i=1}^p k_i!}{k_1(k_1+k_2) \dots (k_1+k_2+\dots+k_p)} \\
& \hspace{15em} \text{(Donnelly\&Tavare(1986)の(4.1)式)} \quad (4.4)
\end{aligned}$$

n 個のサンプル遺伝子において、どの遺伝子がどのクラス(アレル)に属するか区別を無視するために $\frac{n!}{\prod_{i=1}^p k_i!}$ を掛けると、単に p 種のアレルがそれぞれ k_1 個、 k_2 個、 \dots 、 k_p 個ずつ

サンプルされているとき、 k_1 個のアレルが最も若く、 \dots 、 k_p 個のアレルが最も古い確率 $P(k_1 \triangleright k_2 \triangleright \dots \triangleright k_p)$ について次の定理を得る。

定理 4. 2

$$P(k_1 \triangleright k_2 \triangleright \dots \triangleright k_p) = \frac{g^p n!}{(g)_n k_1(k_1+k_2) \dots (k_1+k_2+\dots+k_p)} \quad (4.5)$$

これは Donnelly\&Tavare(1986)の式(4.2)と一致する。

最も古いアレル A_p の年齢分布は指数分布の畳み込み $\sum_{j=1}^n \text{Exp}\left(\frac{2}{jg+j(j-1)}\right)$ に従い、期待値 $= \sum_{j=1}^n \frac{2}{jg+j(j-1)}$ である。

4. 3 固有変異の多型 (Unique Event Polymorphism)

ヒトでは DNA 塩基の数が 10 の 9 乗のオーダーということを考えると生物のゲノム上に生じる各突然変異は進化の過程でただ 1 回起こる事象と考えられる。このように各突然変異をユニークな事象と見るとき、サンプル中に見られる多型の起源となる固有な突然変異の年齢を考える。その例として n 個のサンプル中に、ある UEP のサンプルが k 個 ($1 \leq k \leq n-1$) 含まれているとき、その UEP の年齢を推定しよう。

T_Δ を突然変異が起きた時刻、 $\frac{\rho}{2}$ を突然変異率、サンプル遺伝この系図過程を α_i で表し、

$|\alpha_i| = i$ の状態の時に突然変異が生じたと仮定しよう ($2 \leq i \leq n-k+1$)。

突然変異型 : $t \leq T_\Delta$ では $k-1$ 回の合祖と最後に突然変異が生じる。 $t > T_\Delta$ では突然変異タイプの祖先も含めた i 個の野生型祖先が $i-1$ 回合祖して最終的な MRCA に到達する。

野生型 : $t \leq T_\Delta$ では $n-k$ 個の系統が $i-1$ 個の系統になるまで合祖する。
 $t > T_\Delta$ では突然変異型と同じ。

事象 U : n 個のサンプル中に $(k, n-k)$ の割合で固有変異の多型が観察される事象

事象 U_i : 事象 U の中で $|\alpha_i| = i$ の状態の時に突然変異が生じた事象

$$P_\rho(U_i) = \left(\frac{\rho/2}{\{\rho i + i(i-1)\}/2} \right) \prod_{j=2}^i \left(\frac{j(j-1)/2}{\{\rho j + j(j-1)\}/2} \right) \times \left\{ \prod_{j=i+1}^n \frac{1}{\{\rho j + j(j-1)\}/2} \right\} \\ \times \left(\prod_{j=1}^{n-k} \frac{j(j-1)}{2} \right) \left(\prod_{j=2}^k \frac{j(j-1)}{2} \right) \times_{n-i} C_{k-1}$$

$$P_\rho(U) = \sum_{i=2}^{n-k+1} P(U_i)$$

事象 U_i の条件下で突然変異の年齢 (T_Δ) の分布 $P_\rho(T_\Delta = t | U_i)$ は

$$P_\rho(T_\Delta = t | U_i) = \prod_{j=n}^i \text{Exp}\left(\frac{2}{j(\rho + j - 1)}\right) \quad (\text{指数分布の畳み込み})、$$

事象 U の条件下で突然変異の年齢 (T_Δ) の分布は

$$P_\rho(T_\Delta = t | U) = \frac{1}{P_\rho(U)} \sum_{i=2}^{n-k+1} P_\rho(T_\Delta = t, U_i) = \sum_{i=2}^{n-k+1} \frac{P_\rho(U_i)}{P_\rho(U)} P_\rho(T_\Delta = t | U_i)$$

固有変異は突然変異率が非常に小さい事象と考えられるので $\rho \rightarrow 0$ の極限を考える。

$$\lim_{\rho \rightarrow 0} \frac{P_\rho(U_i)}{P_\rho(U)} = \frac{n-i C_{k-1}}{\sum_{i=2}^{n-k+1} n-i C_{k-1}} = \frac{n-i C_{k-1}}{n-1 C_k}, \quad \lim_{\rho \rightarrow 0} P_\rho(T_\Delta = t | U_i) = \prod_{j=n}^i \text{Exp}\left(\frac{2}{j(j-1)}\right)$$

定理 4. 3 固有変異の年齢の分布

$$P(T_\Delta = t | U) = \sum_{i=2}^{n-k+1} \frac{n-i C_{k-1}}{n-1 C_k} \left\{ \prod_{j=n}^i \text{Exp}\left(\frac{2}{j(j-1)}\right) \right\}$$

平均は $E[T_\Delta | U] = \sum_{i=2}^{n-k+1} \frac{n-i C_{k-1}}{n-1 C_k} \times \frac{2(n-i+1)}{n(i-1)}$ となる。この結果は Tavaré(2004)(p116,p117)

と一致する。特に $\frac{k}{n} = x$ (一定) として、サンプル数 n を無限に大きくすると

$$\lim_{n \rightarrow \infty} E[T_\Delta | U] = \lim_{n \rightarrow \infty} \sum_{i=2}^{n-k+1} \frac{{}_{n-i}C_{k-1}}{{}_{n-1}C_k} \times \frac{2(n-i+1)}{n(i-1)} = \frac{2x}{1-x} \sum_{i=2}^{\infty} \frac{(1-x)^i}{i-1} = -\frac{2x}{1-x} \log x$$

これは拡散過程を用いた Kimura&Ohta(1973)の結果と一致する。

(cf. Wiuf and Donnelly(1999))

○固有変異多型 (UEP) としての SNP (単一塩基多型) の分布

第2節の分離サイトの数の分布で使った $\vartheta = 4Nu$ は全ゲノム上に発生する SNP に関する全突然変異率である。しかし特定の領域あるいはサイトに着目するとそこで発生する SNP は非常に稀にしか起きない UEP と考えられる。現実のサンプルで検出される SNP はある領域 (例えばある遺伝子のイントロン領域) での SNP でありこのデータから系図を考察するのは UEP として扱うのが適切である。第2節の分離サイトの結果から $\vartheta \rightarrow 0$ の極限を考える。定理 3.6 より以下の結果を得る。

$$\lim_{\vartheta \rightarrow 0} P(T=t | S=k) = \sum_{|r|=k} \frac{\prod_{j=2}^n \left(\frac{1}{j-1}\right)^{r_j}}{P(S=k)} \left\{ \prod_{j=2}^n \text{Gamma}(r_j + 1, \frac{2}{j(j-1)}) \right\}$$

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

$$E[T | S=k] = \sum_{|r|=k} A(r) \left\{ \sum_{j=2}^n \frac{2(r_j + 1)}{j(j-1)} \right\} \quad \text{ただし} \quad A(r) = \prod_{j=2}^n \left(\frac{1}{j-1}\right)^{r_j} / P(S=k)$$

4. 4 ディリクレ分布とポリアの壺モデル

4. 4. 1 ベータ分布、ガンマ分布、ディリクレ分布

相互に関連の深い次の三つの確率分布の説明から始めることにしよう。

(1) ベータ分布 $Beta(\alpha, \beta)$

$$f(x; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} \quad (0 < x < 1)$$

$$\text{ここで } B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \quad ; \quad \Gamma(\cdot) \text{ はガンマ関数。}$$

(2) ガンマ分布 $\Gamma(\alpha, \lambda)$

$$g(x; \alpha, \lambda) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} \quad (0 < x < \infty), \quad \alpha > 0, \quad \lambda > 0$$

α を shape parameter, λ を scale parameter という。

特性関数は以下の様に求められる。

$$\phi(t) = \int_0^{\infty} e^{itx} g(x; \alpha, \lambda) dx = \frac{\lambda^\alpha}{\Gamma(\alpha)} \int_0^{\infty} x^{\alpha-1} \exp[(it - \lambda)x] dx = \left(1 - \frac{it}{\lambda}\right)^{-\alpha}$$

これより、確率変数 X_1, \dots, X_n が独立で、各 X_i が $\Gamma(\alpha_i, \lambda)$ に従うとき (λ は共通)、和 $X_1 + \dots + X_n$ はガンマ分布 $\Gamma(\alpha_1 + \dots + \alpha_n, \lambda)$ に従う。

(3) ディリクレ分布 $D(\alpha_1, \alpha_2, \dots, \alpha_{k+1})$ (ただし $\alpha_i > 0, i = 1, \dots, k+1$)

領域 $\Delta = \{\bar{x} = (x_1, \dots, x_k); 0 \leq x_i \leq 1, \sum_{i=1}^k x_i \leq 1\}$ 上の分布

$$h(\bar{x}; \bar{\alpha}) = \frac{\Gamma(\alpha_1 + \dots + \alpha_{k+1})}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_{k+1})} \prod_{i=1}^{k+1} x_i^{\alpha_i-1}, \text{ ただし } \bar{\alpha} = (\alpha_1, \dots, \alpha_{k+1}), x_{k+1} = 1 - \sum_{i=1}^k x_i.$$

$k=1$ のときはベータ分布となり、ベータ分布の多次元版と考えられる。

補題 4. 4

ディリクレ分布 $D(\alpha_1, \alpha_2, \dots, \alpha_{k+1})$ を変数 x_k について $0 \leq x_k \leq 1 - \sum_{i=1}^{k-1} x_i$ の区間で積分した (x_1, \dots, x_{k-1}) の周辺分布密度はディリクレ分布 $D(\alpha_1, \dots, \alpha_{k-1}, \alpha_k + \alpha_{k+1})$ となる。

(証明)

$$\frac{\Gamma(\alpha_1 + \dots + \alpha_{k+1})}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_{k+1})} \prod_{i=1}^{k-1} x_i^{\alpha_i-1} \int_0^{1 - \sum_{i=1}^{k-1} x_i} x_k^{\alpha_k-1} \left\{1 - \sum_{i=1}^{k-1} x_i - x_k\right\}^{\alpha_{k+1}-1} dx_k \quad (\text{ただし } t = 1 - \sum_{i=1}^{k-1} x_i)$$

$x_k = (1 - \sum_{i=1}^{k-1} x_i)y = ty$ と変数変換すると

$$= \frac{\Gamma(\alpha_1 + \dots + \alpha_{k+1})}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_{k+1})} \prod_{j=1}^{k-1} x_j^{\alpha_j-1} \left(1 - \sum_{i=1}^{k-1} x_i\right)^{\alpha_k + \alpha_{k+1} - 1} \int_0^1 y^{\alpha_k-1} (1-y)^{\alpha_{k+1}-1} dy$$

$$= \frac{\Gamma(\alpha_1 + \dots + \alpha_{k+1})}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_{k+1})} \times B(\alpha_k, \alpha_{k+1}) \left(\prod_{j=1}^{k-1} x_j^{\alpha_j-1}\right) \left(1 - \sum_{i=1}^{k-1} x_i\right)^{\alpha_k + \alpha_{k+1} - 1}$$

$$= \frac{\Gamma(\alpha_1 + \dots + \alpha_{k+1})}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_{k-1}) \Gamma(\alpha_k + \alpha_{k+1})} \left(\prod_{j=1}^{k-1} x_j^{\alpha_j-1}\right) x_k^{\alpha_k + \alpha_{k+1} - 1} \quad \text{ただし } x_k = 1 - \sum_{i=1}^{k-1} x_i$$

次にガンマ分布とディリクレ分布の関係について次の補題を証明する。

補題 4. 5 (Mosimann(1962))

$Y_j (j=1,2,\dots,k)$ をそれぞれガンマ分布 $\Gamma(\alpha_j, \lambda)$ に従う独立な確率変数とする。

このとき、 $P_j = \frac{Y_j}{\sum_{i=1}^k Y_i} (j=1,2,\dots,k)$ とすると、 (P_1, P_2, \dots, P_k) はディリクレ分布

$D(\alpha_1, \dots, \alpha_k)$ に従う。

(証明)

(Y_1, \dots, Y_k) の同時分布密度を $g(y_1, \dots, y_k)$ とすると

$$g(y_1, \dots, y_k) = \left(\prod_{j=1}^k \frac{\lambda^{\alpha_j}}{\Gamma(\alpha_j)} y_j^{\alpha_j-1} \right) \exp \left[-\lambda \left(\sum_{i=1}^k y_i \right) \right] \quad (4.6)$$

$P_j = Y_j / \left(\sum_{i=1}^k Y_i \right), j=1,2,\dots,k, \quad W = \sum_{i=1}^k Y_i$ とする。これより、

$$Y_j = P_j W \quad (j=1,2,\dots,k-1), \quad Y_k = \left(1 - \sum_{i=1}^{k-1} P_i \right) W.$$

(Y_1, \dots, Y_k) の分布から (P_1, \dots, P_{k-1}, W) の分布を求める。変数変換のヤコビアンは

$$\frac{\partial(y_1, y_2, \dots, y_k)}{\partial(p_1, \dots, p_{k-1}, w)} = \begin{vmatrix} \frac{\partial y_1}{\partial p_1} & \frac{\partial y_1}{\partial p_2} & \dots & \frac{\partial y_1}{\partial p_{k-1}} & \frac{\partial y_1}{\partial w} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y_k}{\partial p_1} & \frac{\partial y_k}{\partial p_2} & \dots & \frac{\partial y_k}{\partial p_{k-1}} & \frac{\partial y_k}{\partial w} \end{vmatrix} = \begin{vmatrix} w & 0 & \dots & p_1 \\ 0 & w & \dots & p_2 \\ \vdots & \vdots & \vdots & \vdots \\ -w & -w & \dots & 1 - \sum_{j=1}^{k-1} p_j \end{vmatrix}$$

第 1 行から第 $k-1$ 行までの和を第 k 行に加えると

$$= \begin{vmatrix} w & 0 & \dots & 0 & p_1 \\ 0 & w & \dots & 0 & p_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & w & p_{k-1} \\ 0 & 0 & \dots & 0 & 1 \end{vmatrix} = w^{k-1}.$$

これより (P_1, P_2, \dots, P_k) の分布密度関数は(4.6)を変数変換し、 w について積分すると

$$\begin{aligned} f(p_1, p_2, \dots, p_{k-1}) &= \left\{ \prod_{j=1}^k \frac{\lambda^{\alpha_j}}{\Gamma(\alpha_j)} p_j^{\alpha_j-1} \right\} \int_0^\infty w^{\alpha-1} e^{-\lambda w} dw, \quad \alpha = \sum_{i=1}^k \alpha_i \\ &= \frac{\Gamma(\alpha)}{\prod_{j=1}^k \Gamma(\alpha_j)} \prod_{j=1}^k p_j^{\alpha_j-1}. \end{aligned}$$

補題 4. 6 (Lukacs(1955))

X, Y を互いに独立で退化しない正値の確率変数とする。このとき $P = \frac{Y}{X+Y}$ と

$W = X+Y$ が独立となるための必要十分条件は X, Y が同じスケールパラメーターのガンマ分布に従うことである。

(証明) まず必要性を証明しよう。

X, Y, P, W の分布関数を $F(x), G(y), H_1(p), H_2(w)$ とする。 X, Y の特性関数を

$$f(v) = \int_0^\infty e^{vx} dF(x), \quad g(v) = \int_0^\infty e^{vy} dG(y) \quad (v = s + it) \quad (4.7)$$

とする。ラプラス変換の性質より $t > 0$ で $f(v), g(v)$ は解析的である。さらに、次の連続性を仮定する。

$$\lim_{t \downarrow 0} f(s + it) = f(s), \quad \lim_{t \downarrow 0} g(s + it) = g(s) \quad (4.8)$$

v について微分すると、

$$f'(v) = i \int_0^\infty x e^{vx} dF(x), \quad f''(v) = - \int_0^\infty x^2 e^{vx} dF(x) \quad (4.9a)$$

$$g'(v) = i \int_0^\infty y e^{vy} dG(y), \quad g''(v) = - \int_0^\infty y^2 e^{vy} dG(y) \quad (4.9b)$$

$0 \leq P = \frac{Y}{X+Y} \leq 1$ より確率変数 P のモーメントは常に存在する。

$$\theta_1 = E[P] = E\left[\frac{Y}{X+Y}\right], \quad \theta_2 = E[P^2] = E\left[\left(\frac{Y}{X+Y}\right)^2\right] \quad (4.10)$$

とすると、明らかに $0 < \theta_1^2 \leq \theta_2 < 1$ 。 W と P は独立なので、 (W, P) の特性関数は

$E[\exp(izP + ivW)] = E[\exp(izP)]E[\exp(ivW)]$ 、すなわち

$$E\left[\exp\left\{iz\left(\frac{Y}{X+Y}\right) + ivW\right\}\right] = E\left[\exp\left\{iz\left(\frac{Y}{X+Y}\right)\right\}\right]E[\exp\{iv(X+Y)\}] \quad (4.11)$$

両辺を v および z で微分した後、 $z = 0$ と置くと

$E[Y \exp\{iv(X+Y)\}] = \theta_1 E[(X+Y) \exp\{iv(X+Y)\}]$ 、ただし(4.10)を用いた。

X, Y の独立性より

$E[Y \exp(ivY)]E[\exp(ivX)] = \theta_1 \{E[X \exp(ivX)]E[\exp(ivY)] + E[\exp(ivX)]E[Y \exp(ivY)]\}$

(4.9a, b)より $(-i)g'(v)f(v) = \theta_1 \{(-i)f'(v)g(v) + (-i)g'(v)g(v)\}$ 、これより

$$(1 - \theta_1)g'(v)f(v) = \theta_1 f'(v)g(v) \quad \text{ただし } \text{Im}(v) > 0 \quad (4.12)$$

また $f(0) = g(0) = 1$, と $f(v), g(v)$ の連続性より原点の近傍 (ただし $\text{Im}(v) > 0$) で、
 $f(v) \neq 0, g(v) \neq 0$ が成り立つ。特性関数の定義より明らかに $f(it) > 0, g(it) > 0$ ただし
 $t > 0$ の実数。 $f(v), g(v)$ の連続性とハイネ=ボレルの被覆定理より $0 < \text{Im}(v) \leq 1$ を含む
ある単連結領域で $f(v) \neq 0, g(v) \neq 0$ が成り立つ。以後この領域を D とする。

$v \in D$ のとき、(4.12)の両辺を $f(v)g(v)$ で割ると

$$(1 - \theta_1) \frac{g'(v)}{g(v)} = \theta_1 \frac{f'(v)}{f(v)} \quad (4.13)$$

初期条件 $f(0) = g(0) = 1$ より、積分して次の関係式を得る。

$$\{g(v)\}^{1-\theta_1} = \{f(v)\}^{\theta_1} \quad \text{ただし } \text{Im}(v) > 0 \quad (4.14)$$

次に、(4.11)の両辺を v, z でそれぞれ2回微分し $z = 0$ と置くと、

$$E[Y^2 \exp\{i(X+Y)v\}] = E\left[\left(\frac{Y}{X+Y}\right)^2\right] E[(X+Y)^2 \exp\{iv(X+Y)\}]$$

X, Y の独立性に注意し、(4.7), (4.9a, b) 及び(4.10)を使うと

$g''(v)f(v) = \theta_2 \{f''(v)g(v) + 2f'(v)g'(v) + g''(v)f(v)\}$, $v \in D$ のとき、両辺を $f(v)g(v)$
で割ると

$$\frac{g''(v)}{g(v)} = \theta_2 \left\{ \frac{f''(v)}{f(v)} + 2 \frac{f'(v)}{f(v)} \times \frac{g'(v)}{g(v)} + \frac{g''(v)}{g(v)} \right\} \quad (4.15)$$

$\phi(v) = \log f(v), \psi(v) = \log g(v)$ とすると

$$\frac{f'(v)}{f(v)} = \phi'(v), \quad \frac{f''(v)}{fv} = \phi''(v) + \{\phi'(v)\}^2 \quad (4.16)$$

$\psi(v)$ についても同様の式が成り立つ。(4.13)を $\phi(v), \psi(v)$ で表すと

$$(1 - \theta_1)\psi'(v) = \theta_1\phi'(v) \quad \text{これより } (1 - \theta_1)\psi''(v) = \theta_1\phi''(v) \quad (4.17)$$

(4.15)を(4.16)(4.17)を用いて $\phi(v)$ のみで表すと

$$(1 - \theta_1)(\theta_1 - \theta_2)\phi''(v) = (\theta_2 - \theta_1^2)\{\phi'(v)\}^2 \quad (4.18)$$

(1) $0 < \theta_1^2 < \theta_2 < \theta_1 < 1$ (等号が付かないとき)

$\rho = \frac{(1 - \theta_1)(\theta_1 - \theta_2)}{(\theta_2 - \theta_1^2)}$ と置くと $\rho > 0$ である。

$$(4.18) \text{より } v \in D, \text{Im}(v) > 0 \text{ のとき} \quad \frac{\phi''(v)}{\{\phi'(v)\}^2} = \frac{1}{\rho} \quad (4.19)$$

$$k_1 = E[\exp(-X)], \quad k_2 = E[X \exp(-X)], \quad \lambda = (k_1 \rho - k_2) / k_2 \text{ とする。} \quad (4.20)$$

初期条件 $\phi'(i) = \frac{f'(i)}{f(i)} = \frac{ik_2}{k_1}$, ($i = \sqrt{-1}$) より $C = \frac{\lambda}{i\rho}$ 。故に $\phi'(v) = \frac{i\rho}{\lambda - iv}$ 。積分して

$$f(v) = A \left(1 - \frac{iv}{\lambda}\right)^{-\rho}, \quad (A \text{ は定数})。 \text{連続性の条件より } \lim_{i \downarrow 0} f(it) = f(0) = 1, \text{ よって } A=1。$$

$$(4.14) \text{ より } f(v) = \left(1 - \frac{iv}{\lambda}\right)^{-\rho}, \quad g(v) = \left(1 - \frac{iv}{\lambda}\right)^{-\sigma} \quad (\text{ただし } \sigma = \frac{\theta_1 \rho}{1 - \theta_1}) \quad (4.21)$$

(4.21)は領域 D で成り立つが、 $f(v)$ および $\left(1 - \frac{iv}{\lambda}\right)^{-\rho}$ は複素平面の上半部 $\text{Im}(v) > 0$ で解析

的であつ領域 D で一致する。故に解析関数の一致の定理より $\text{Im}(v) > 0$ で(4.21)が成り立つ。 $g(v)$ についても同様である。連続性の仮定(4.8)より実数 v についても成り立つ。

X, Y の特性関数は(4.21)で与えられ、ともに同じスケールパラメーター λ のガンマ分布に従う。

(2) $\theta_1 = \theta_2$ または $\theta_1^2 = \theta_2$ のとき

$\theta_1 = \theta_2$ のとき、 $\theta_2 = \theta_1 > \theta_1^2$ なので(4.18)より $\phi'(v) = 0$ 、よつて $\phi(v) = C$ 。また $\theta_1^2 = \theta_2$ ならば $\theta_1 > \theta_1^2 = \theta_2$ なので(4.18)より $\phi''(v) = 0$ すなわち $\phi(v) = Av + C$ 。いずれにしても $f(v) = \exp(Av + C)$ となりラプラス逆変換により X は退化した分布になり定理の条件を満たさない。 Y についても同様である。故に P と W が独立ならば X, Y は同じスケールパラメーターの Γ 分布に従う。

十分性の証明： X, Y を独立でそれぞれガンマ分布 $\Gamma(\alpha, \lambda), \Gamma(\beta, \lambda)$ に従う確率変数とする。

このとき $W = X + Y$ と $P = \frac{Y}{X + Y}$ の結合分布の特性関数は

$$\begin{aligned} & E \left[\exp \left\{ iv(X + Y) + iz \left(\frac{Y}{X + Y} \right) \right\} \right] \\ &= \frac{\lambda^{\alpha + \beta}}{\Gamma(\alpha) \Gamma(\beta)} \int_0^\infty \int_0^\infty \exp \left[-(\lambda - iv)(x + y) + iz \left(\frac{y}{x + y} \right) \right] \times x^{\alpha - 1} y^{\beta - 1} dx dy \end{aligned}$$

$\xi = (\lambda - iv)x$, $\eta = (\lambda - iv)y$ とし (x, y) から (ξ, η) に変数変換すると

$$= \frac{\lambda^{\alpha+\beta}}{\Gamma(\alpha)\Gamma(\beta)} \left\{ \frac{1}{(\lambda - iv)^{\alpha+\beta}} \right\} \int_0^\infty \int_0^\infty \exp \left[-(\xi + \eta) + iz \left(\frac{\eta}{\xi + \eta} \right) \right] \times \xi^{\alpha-1} \eta^{\beta-1} d\xi d\eta$$

すなわち特性関数が v の関数と z の関数の積として表されているので $W = X + Y$ と $P = \frac{Y}{X+Y}$ は独立である。 (証明終わり)

証明は省略するが、補題 4. 6 を一般化した次の補題が知られている。

補題 4. 7

(1) (Mosiman(1962)) $\{Y_i; i = 1, 2, \dots, k\}$ を k 個の独立で非退化、正値の確率変数とする。

$k-1$ 個の確率変数 $P_j = Y_j / \left(\sum_{i=1}^k Y_i \right)$, ($j = 1, 2, \dots, k-1$) が $W = \sum_{i=1}^k Y_i$ と独立となるの

は、全ての Y_j ($j = 1, \dots, k$) が同じスケールパラメータ λ の Γ 分布に従う時に限る。

(2) $\{Y_i; i = 1, 2, \dots, k\}$ を k 個の独立で非退化、正値の確率変数とする。 $P_j = Y_j / \left(\sum_{i=1}^k Y_i \right)$,

($j = 1, 2, \dots, k-1$) とすると、各 P_j が和 $W = \sum_{i=1}^k Y_i$ と独立ならば (P_1, \dots, P_{k-1}) はディリクレ分布に従う。

4. 4. 2 完全中立性とディリクレ分布の一般化

P_1, P_2, \dots, P_k を非負で $\sum_{i=1}^k P_i = 1$ を満たす確率変数とする。また $S_0 = 0$, $j \geq 1$ のとき

$$S_j = \sum_{i=1}^j P_i, \quad Z_j = \frac{P_j}{1 - S_{j-1}} = \frac{P_j}{P_j + P_{j+1} + \dots + P_k}$$

を定義する。明らかに $Z_1 = P_1$, $Z_k = 1$ である。これより、逆に $P_j = Z_j \prod_{i=1}^{j-1} (1 - Z_i)$ と表わ

される。 Z_j は全体から P_1, P_2, \dots, P_{j-1} を取り除いた残りの中での P_j の相対的な割合

$P_j / (1 - S_{j-1})$ を表わしている。確率ベクトル $\vec{P} = (P_1, \dots, P_k)$ に対して $\vec{P}_j^1 = (P_1, \dots, P_j)$,

$\vec{P}_j^2 = (P_{j+1}, \dots, P_k)$; ($j < k$) とし、 $\vec{W}_j = \frac{\vec{P}_j^2}{1 - S_j} = \frac{(P_{j+1}, P_{j+2}, \dots, P_k)}{P_{j+1} + \dots + P_k}$ と定義する。

この小節では Connor and Mosimann(1969)に従って、確率ベクトルの完全中立性 (completely neutral) と一般化されたディリクレ分布について紹介する。

定義：上で定義した確率ベクトル $\vec{P} = (P_1, \dots, P_k)$ に対して、 P_1 が確率ベクトル

$$\vec{W}_1 = \left(\frac{P_2}{1-P_1}, \frac{P_3}{1-P_1}, \dots, \frac{P_k}{1-P_1} \right) \text{ と独立のとき } P_1 \text{ は中立(neutral)であると言う。}$$

定理 4. 8

確率ベクトル $\vec{P} = (P_1, \dots, P_k)$ に対して P_1 が中立ならば、 P_1 は \vec{W}_j ($j = 1, 2, \dots, k-1$) と独立である。

(証明) $X_2 = \frac{P_2}{1-P_1}, X_3 = \frac{P_3}{1-P_1}, \dots$ とすると条件より P_1 は (X_2, X_3, \dots, X_k) と独立である。

$P_2 = (1-P_1)X_2, P_3 = (1-P_1)X_3, \dots, P_k = (1-P_1)X_k$ より

$$\frac{P_3}{1-P_1-P_2} = \frac{(1-P_1)X_3}{1-P_1-(1-P_1)X_2} = \frac{X_3}{1-X_2}, \text{ 同様に } \frac{P_j}{1-P_1-P_2} = \frac{X_j}{1-X_2}; (3 \leq j \leq k).$$

故に P_1 は (X_2, \dots, X_k) と独立なので $(X_3, X_4, \dots, X_k)/(1-X_2) = \vec{W}_2$ と独立である。

同様にして P_1 は \vec{W}_j ($j = 1, 2, \dots, k-1$) と独立である。

中立性の概念をベクトルまで拡張しよう。

定義：確率ベクトル $\vec{P} = (\vec{P}_j^1, \vec{P}_j^2), j = 1, 2, \dots, k-1$ に対して、 \vec{P}_j^1 が \vec{W}_j と独立であるとき

確率ベクトル \vec{P}_j^1 は中立であるという。また全ての j について \vec{P}_j^1 が中立のとき、確

率ベクトル $\vec{P} = (P_1, \dots, P_k)$ は完全中立という。

すなわち完全中立とは任意の j に対して最初の j 個の要素の確率 P_1, \dots, P_j を与えたとき

それを除いた残りの確率を他の $k-j$ 個の要素で分割する方法 \vec{W}_j が P_1, \dots, P_j と独立である

ことを意味する。ただしベクトル \vec{P} の成分の順序に依存することを注意しなければならない。すなわち (P_1, P_2, P_3, P_4) の P_1 は中立であるが、順序を入れ替えて (P_2, P_1, P_3, P_4) とする

と P_2 は中立でないこともある。確率ベクトルの中立性と確率変数系 $\{Z_j; j=1,2,\dots,k\}$ の

独立性に関して次の定理が成り立つ。

定理 4. 9

$j=1,2,\dots,r$ について、 \bar{P}_j^1 が中立ならば Z_1, Z_2, \dots, Z_r は互いに独立である。

(証明)

Z_1, Z_2, \dots, Z_r の同時分布 $F_r(z_1, \dots, z_r)$ が各 Z_j の分布 $f_j(z_j)$ の積として

$F_r(z_1, \dots, z_r) = \prod_{j=1}^r f_j(z_j)$ となることを示せばよい。まず、 P_1 が中立であることより、

$Z_1 = P_1$ と $Z_2 = P_2 / (1 - P_1)$ が独立。故に、 $F_2(z_1, z_2) = f_1(z_1)f_2(z_2)$ が成り立つ。次に

$\bar{P}_2^1 = (P_1, P_2)$ について \bar{P}_2^1 が中立より (P_1, P_2) は $\left(\frac{P_3}{1-S_2}, \frac{P_4}{1-S_2}, \dots \right) = \left(Z_3, \frac{P_4}{1-S_2}, \dots \right)$

と独立である。故に (Z_1, Z_2) は (P_1, P_2) の関数なので (Z_1, Z_2) は Z_3 と独立である。これより $F_3(z_1, z_2, z_3) = F_2(z_1, z_2)f_3(z_3) = f_1(z_1)f_2(z_2)f_3(z_3)$ となり、 z_1, z_2, z_3 は独立である。以下同様に Z_1, Z_2, \dots, Z_r が独立であることが示される。

この定理より $\bar{P} = (P_1, \dots, P_k)$ が完全中立ならば $\{Z_1, Z_2, \dots, Z_{k-1}\}$ は互いに独立であるが、この条件が必要十分であることが示される。

定理 4. 10

$\bar{P} = (P_1, \dots, P_k)$ が完全中立であるための必要十分条件は $\{Z_1, Z_2, \dots, Z_{k-1}\}$ が互いに独立であることである。ただし、 $Z_k = 1$ とする。

(証明)

必要性は定理 4.9 で証明されたので十分性を示す。 $\{Z_1, Z_2, \dots, Z_{k-1}\}$ が互いに独立のとき

\bar{P}_j^1 が \bar{W}_j ($j=1,2,\dots,k-1$) と独立であることを示せば十分である。 Z_j の定義より

$\prod_{i=1}^j (1 - Z_i) = 1 - S_j$, ($1 \leq j \leq k-1$)。これより $P_j = (1 - S_j)Z_j = Z_j \prod_{i=1}^{j-1} (1 - Z_i)$ 、また

$$\frac{P_{j+r}}{1-S_j} = \frac{Z_{j+r} \prod_{i=1}^{j+r-1} (1-Z_i)}{\prod_{i=1}^j (1-Z_i)} = Z_{j+r} \prod_{i=j+1}^{j+r-1} (1-Z_i), \quad (1 \leq j \leq k-j).$$

すなわち \bar{P}_j^1 は $\{Z_1, \dots, Z_j\}$

で表され、 \bar{W}_j ($j=1, 2, \dots, k-1$) は (Z_{j+1}, \dots, Z_r) のみに依存する。従って \bar{P}_j^1 と \bar{W}_j は独立

である。この二つの定理より $1 \leq j \leq k-2$ に対して \bar{P}_j^1 が Z_{j+1} と独立ならば (Z_1, \dots, Z_j) は

Z_{j+1} と独立になるので $\{Z_1, Z_2, \dots, Z_{k-1}\}$ は独立である。故に、定理 4. 10 より

$\bar{P} = (P_1, \dots, P_k)$ は完全中立である。すなわち次の定理が成り立つ。

定理 4. 11

$\bar{P} = (P_1, \dots, P_k)$ が完全中立であるための必要十分条件は \bar{P}_j^1 と Z_{j+1} が全て独立となることである。

この完全中立性を利用して一般化されたディリクレ分布を定義し、ディリクレ分布と完全中立性との関係を調べる。 $\{Z_1, Z_2, \dots, Z_{k-1}\}$ を独立な確率変数系で各 Z_j はベータ分布

$Beta(\alpha_j, \beta_j)$ に従うものとする。すなわち Z_j の分布密度は

$$f_j(z_j) = \frac{1}{B(\alpha_j, \beta_j)} z_j^{\alpha_j-1} (1-z_j)^{\beta_j-1}, \quad (0 < z_j < 1).$$

また $P_j = Z_j \left\{ \prod_{i=1}^{j-1} (1-Z_i) \right\}$ によって $\bar{P} = (P_1, \dots, P_k)$ を定義する。

$$Z_j = \frac{P_j}{1-S_{j-1}} = \frac{P_j}{1-(P_1 + \dots + P_{j-1})} \quad \text{より、} \quad j > i \text{ のとき } \frac{\partial Z_i}{\partial P_j} = 0, \quad \frac{\partial Z_j}{\partial P_j} = \frac{1}{1-S_{j-1}},$$

$j < i$ のとき $\frac{\partial Z_i}{\partial P_j} = \frac{P_i P_j}{(1-S_{j-1})^2}$ となる。これより、変数変換のヤコビアンは

$$J = \frac{\partial(Z_1, \dots, Z_{k-1})}{\partial(P_1, \dots, P_{k-1})} = \begin{vmatrix} \frac{\partial Z_1}{\partial P_1} & \dots & \frac{\partial Z_{k-1}}{\partial P_{k-1}} \\ \vdots & & \vdots \\ \frac{\partial Z_{k-1}}{\partial P_1} & \dots & \frac{\partial Z_{k-1}}{\partial P_{k-1}} \end{vmatrix} = \frac{1}{\prod_{i=1}^{k-1} (1-S_{i-1})}. \quad (4.22)$$

この変数変換により $\vec{Z} = (Z_1, \dots, Z_{k-1})$ から $\vec{P} = (P_1, \dots, P_k)$ の分布を求める。

$$\begin{aligned} & \prod_{j=1}^{k-1} \left\{ B(\alpha_j, \beta_j) \right\}^{-1} z_j^{\alpha_j-1} (1-z_j)^{\beta_j-1} \Big\} \\ &= \left\{ \prod_{j=1}^{k-1} B(\alpha_j, \beta_j) \right\}^{-1} \left\{ \prod_{j=1}^{k-1} \left(\frac{P_j}{1-S_{j-1}} \right)^{\alpha_j-1} \right\} \left\{ \prod_{j=1}^{k-1} \left(1 - \frac{P_j}{1-S_{j-1}} \right)^{\beta_j-1} \right\} \times J \\ &= \left\{ \prod_{j=1}^{k-1} B(\alpha_j, \beta_j) \right\}^{-1} \left\{ \prod_{j=1}^{k-1} \left(\frac{P_j}{1-S_{j-1}} \right)^{\alpha_j-1} \right\} \left\{ \prod_{j=1}^{k-1} \left(\frac{1-S_j}{1-S_{j-1}} \right)^{\beta_j-1} \right\} \times \frac{1}{\prod_{j=1}^{k-1} (1-S_{j-1})} \\ &= \left\{ \prod_{j=1}^{k-1} B(\alpha_j, \beta_j) \right\}^{-1} \left\{ \prod_{j=1}^{k-1} P_j^{\alpha_j-1} \right\} \left\{ \prod_{j=1}^{k-1} (1-S_{j-1})^{\beta_{j-1}-(\alpha_j+\beta_j)} \right\} (1-S_{k-1})^{\beta_{k-1}-1} \end{aligned}$$

$$1-S_{j-1} = \sum_{i=j}^k P_i \quad (\text{ただし } \sum_{j=1}^k P_j = 1) \quad \text{および } 1-S_{k-1} = P_k \text{ より}$$

$$= \left\{ \prod_{j=1}^{k-1} B(\alpha_j, \beta_j) \right\}^{-1} P_k^{\beta_{k-1}-1} \left\{ \prod_{j=1}^{k-1} \left(P_j^{\alpha_j-1} \left(\sum_{i=j}^k P_i \right)^{\beta_{j-1}-(\alpha_j+\beta_j)} \right) \right\} \quad (4.23)$$

ここで β_0 は任意の実数としてよい。また $\beta_{i-1} = \alpha_i + \beta_i$ ($i = 2, 3, \dots, k-1$) のとき(4.23)は

$$\left\{ \prod_{j=1}^{k-1} B(\alpha_j, \beta_j) \right\}^{-1} \left\{ \prod_{j=1}^{k-1} P_j^{\alpha_j-1} \right\} P_k^{\beta_{k-1}-1} \quad \text{すなわちディリクレ分布 } D(\alpha_1, \dots, \alpha_{k-1}, \beta_{k-1}) \text{ となる}$$

る。従って、ディリクレ分布に従う確率変数 (P_1, P_2, \dots, P_k) は完全中立であり、ディリクレ

分布は任意の変数の置換をしてもディリクレ分布なので $\vec{P} = (P_1, P_2, \dots, P_k)$ は任意の置換

をしても完全中立な確率ベクトルである。それに対して(4.23)で定義される一般化されたディリクレ分布では一般に (P_1, P_2, \dots, P_k) の順序において完全中立となる。

特に対称ディリクレ分布 $D(\alpha, \dots, \alpha)$ となるのは(4.23)において $\beta_{i-1} = \alpha_i + \beta_i$ かつ

$$\alpha_1 = \alpha_2 = \dots = \alpha_{k-1} = \beta_{k-1} = \alpha \text{ のときなので、 } \alpha_j = \alpha, \beta_j = (k-j)\alpha \quad (1 \leq j \leq k-1) \text{ よ}$$

り Z_j はベータ分布 $Beta(\alpha, (k-j)\alpha)$ に従うことを注意しておこう。

$(P_1, P_2, \dots, P_k), \sum_{i=1}^k P_i = 1$ がパラメーター α の対称ディリクレ分布に従うとき単に $D(\alpha; k)$ と書くことにする。完全中立性と関連して対称ディリクレ分布の **size-biased permutation** について触れておこう。

$\bar{P} = (P_1, P_2, \dots, P_k), P_i \geq 0, \sum_{i=1}^k P_i = 1$ を $\{1, 2, \dots, k\}$ 上のランダム確率測度とする。 \bar{P} の **size-biased permutation** $\bar{P}^S = (P_1^S, P_2^S, \dots, P_k^S)$ とは次のように定義される。

定義 : **Size-biased permutation**

任意の $\{1, 2, \dots, k\}$ の置換 σ に対して、 $\bar{P}^S = (P_1^S, P_2^S, \dots, P_k^S) = (P_{\sigma(1)}, P_{\sigma(2)}, \dots, P_{\sigma(k)})$

となる確率が次式で与えられるとき、 \bar{P}^S は \bar{P} の **size-biased permutation** という。

$$P(P_j^S = P_{\sigma(j)}, 1 \leq j \leq i | \bar{P}) = P_{\sigma(i)} \prod_{j=2}^i \left(\frac{P_{\sigma(j)}}{1 - \sum_{r=1}^{j-1} P_{\sigma(r)}} \right), \quad 1 \leq i \leq k \quad (4.24)$$

すなわち $\{1, 2, \dots, k\}$ の任意の置換 σ および \bar{P} に対して、集団からランダムに一つの個体をサンプルし、その個体のタイプの頻度を P_1^S とする。その確率は $P(P_1^S = P_{\sigma(1)}) = P_{\sigma(1)}$ である。この個体と同じタイプの個体を集団からすべて除き、残りの集団からまたランダムに 1 個体サンプルする。2 番目にサンプルしたタイプの元の全集団での遺伝子頻度を P_2^S とす

る。その確率は $P(P_2^S = P_{\sigma(2)} | P_1^S = P_{\sigma(1)}) = \frac{P_{\sigma(2)}}{1 - P_{\sigma(1)}}$ 。この操作を全てのタイプが無くなるま

で続け、順に $P_1^S, P_2^S, \dots, P_k^S$ とするとき $\bar{P}^S = (P_1^S, P_2^S, \dots, P_k^S) = (P_{\sigma(1)}, P_{\sigma(2)}, \dots, P_{\sigma(k)})$ となる確率を(4.24)式右辺は示している。

集団の遺伝子頻度 $\bar{P} = (P_1, P_2, \dots, P_k)$ が対称ディリクレ分布 $D(\alpha; k)$ に従うとき、その

Size-biased permutation \bar{P}^S の分布を求めてみよう。

$$\begin{aligned}
P(P^S = (x_1, \dots, x_k)) &= \sum_{\sigma \in \mathfrak{S}(k)} P(P_j^S = x_{\sigma(j)}, 1 \leq j \leq k | \bar{P} = (x_1, \dots, x_k)) \times P(\bar{P} = (x_1, \dots, x_k)) \\
&= \sum_{\sigma \in \mathfrak{S}(k)} \left\{ x_{\sigma(1)} \prod_{j=2}^k \left(\frac{x_{\sigma(j)}}{1 - \sum_{i=1}^{j-1} x_{\sigma(i)}} \right) \right\} \times \frac{\Gamma(n\alpha)}{\{\Gamma(\alpha)\}^n} \prod_{j=1}^k x_j^{\alpha-1} \\
&= \sum_{\sigma^{-1} \in \mathfrak{S}(k)} \left\{ x_1 \prod_{j=2}^k \left(\frac{x_j}{1 - \sum_{i=1}^{j-1} x_i} \right) \right\} \frac{\Gamma(n\alpha)}{\{\Gamma(\alpha)\}^n} \left(\prod_{j=1}^k x_{\sigma^{-1}(j)} \right)^{\alpha-1} \\
&= \frac{k! \Gamma(k\alpha)}{\{\Gamma(\alpha)\}^k} \left\{ \prod_{j=1}^k x_j \right\}^{\alpha-1} \left\{ x_1 \prod_{j=2}^k \frac{x_j}{1 - \sum_{i=1}^{j-1} x_i} \right\}
\end{aligned}$$

最後の式は $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$ に注意して整理すると

$$\left\{ \prod_{j=1}^{k-1} B(1 + \alpha, (k-i)\alpha) \right\}^{-1} x_k^{\alpha-1} \left\{ \prod_{j=1}^{k-1} \left(x_j^\alpha \left(\sum_{i=j}^k x_i \right)^{-1} \right) \right\} \quad \text{となる。} \quad (4.25)$$

これは(4.23)で $\alpha_i = 1 + \alpha$, $\beta_i = (k-i)\alpha$ と置いた一般化されたディリクレ分布の形である。

すなわち、対称ディリクレ分布 $D(\alpha; k)$ の size-biased permutation は Z_j ($1 \leq j \leq k-1$)

がベータ分布 $Beta(1 + \alpha, (k-j)\alpha)$ に従う一般化されたディリクレ分布である。そして

$$\bar{P}^S = (P_1^S, P_2^S, \dots, P_k^S), \quad P_1^S = Z_1, P_2^S = (1-Z_1)Z_2, P_3^S = (1-Z_1)(1-Z_2)Z_3, \dots,$$

$$P_k^S = (1-Z_1) \cdots (1-Z_{k-1})Z_k \quad (\text{ただし } Z_k = 1) \quad \text{と表される。}$$

4. 4. 3 ポリアの壺とディリクレ過程

ポリアの壺とディリクレ分布の間の興味ある関係について紹介する。

ポリアの壺とは、壺の中に例えば b 個の黒球と r 個の赤球が入っているとす。まず 1 個の球をランダムに壺の中から取り出し、この球と同色の球を 1 個加えて壺に戻す。この操作を繰り返すとき、壺の中の 2 色の球の割合を問う確率論の問題を、発案者の名よりポリアの壺と呼んでいる。

n 回の試行の後、壺の中の黒球の割合を Y_n とする。 $Y_0 = \frac{b}{b+r}$ であ

る。 Y_n は試行の度に確率的に変化するので確率過程である。 n 回の試行後の 2 色の球の個数を b_n, r_n とすると $Y_n = \frac{b_n}{b_n + r_n}$, $b_n + r_n = b + r + n$ である。このとき、

$$E[Y_{n+1}|Y_n] = \frac{b_n}{b_n + r_n} \times \frac{b_n + 1}{b_n + r_n + 1} + \frac{r_n}{b_n + r_n} \times \frac{b_n}{b_n + r_n + 1} = \frac{b_n}{b_n + r_n} = Y_n, \text{ よって } \{Y_n\} \text{ は有界な}$$

マルチンゲールである。マルチンゲールの収束定理より、極限 $\lim_{n \rightarrow \infty} Y_n = Y$ が存在する。

$$\text{最初の } n \text{ 回が全て黒球である確率を } \mu_n \text{ とすると } \mu_n = \frac{b(b+1)\cdots(b+n-1)}{(b+r)(b+r+1)\cdots(b+r+n-1)},$$

確率変数 X_n を n 回目に取り出した球が黒球のとき $X_n = 1$ 、赤球のとき $X_n = 0$ で定義すると、 (X_1, X_2, \dots, X_n) は簡単な計算により交換可能である。すなわち n 次元同時分布が変数の順序を変えても同じ分布に従う。 $\mu_n = P(X_1 = 1, X_2 = 1, \dots, X_n = 1)$ に注意すると、

ド・フィネッティの定理より $\mu_n = \int_0^1 x^n F(dx)$ となる $[0, 1]$ 上の確率速度 $F(x)$ が存在する

(参考: Feller 「確率論とその応用 II」 VII 章, 4 節)。 Γ 関数の性質を利用すると、

$$\mu_n = \frac{\Gamma(b+n)/\Gamma(b)}{\Gamma(b+r+n)/\Gamma(b+r)} = \frac{\Gamma(b+n)\Gamma(b+r)}{\Gamma(b+r+n)\Gamma(b)}$$
 と表せる。これはパラメーターが (b, r) の

ベータ分布 $f_{(b,r)}(x) = \frac{\Gamma(b+r)}{\Gamma(b)\Gamma(r)} x^{b-1}(1-x)^{r-1}$ の n 次のモーメントに等しい。すなわち、

$$F(x) = f_{(b,r)}(x) \text{ となる。 } S_n = \sum_{i=1}^n X_i \text{ とすると、 } P(S_n = k) = \binom{n}{k} \int_0^1 x^k (1-x)^{n-k} f_{(b,r)}(x) dx$$

すなわち、 S_n の分布は二項分布のランダム化で表される。 Y のラプラス変換を考えると

$$\begin{aligned} E[\exp(-\lambda Y)] &= \lim_{n \rightarrow \infty} E[\exp(-\lambda Y_n)] = \lim_{n \rightarrow \infty} E[\exp(-\lambda S_n / n)] = \lim_{n \rightarrow \infty} \int_0^1 (xe^{-\lambda/n} + 1 - x)^n f_{(b,r)}(x) dx \\ &= \int_0^1 e^{-\lambda x} f_{(b,r)}(x) dx \end{aligned}$$

これより極限 $\lim_{n \rightarrow \infty} Y_n = Y$ はベータ分布 $f_{(b,r)}(x)$ に従うことが分かる。

さらに一般化して、壺の中に k 種の色の球がそれぞれ、 α_1 個、 α_2 個、 \dots 、 α_k 個ずつ入っているとする。壺から 1 個ランダムに球を取り出し、同じ色の球を 1 個付け加えて壺に戻す。

このとき、 n 回の試行の後、各色の球の割合を $Y_n = (Y_n^{(1)}, \dots, Y_n^{(k)})$ とすると、極限

$\lim_{n \rightarrow \infty} Y_n = Y$ が存在し、 Y はディリクレ分布 $D(\alpha_1, \dots, \alpha_k)$ に従うことが示される。

Blackwell&MacQueen(1973)はポリアの壺とディリクレ分布の関係をより一般化した定理を示した。

定義：ディリクレ過程

完備で可分な距離空間 X 上の有界な正值測度 μ とランダム確率測度 μ^* があり、 X の任意の可測な有限分割 (B_1, B_2, \dots, B_r) に対して、 r 次元ベクトル $(\mu^*(B_1), \dots, \mu^*(B_r))$ がパラメーター $(\mu(B_1), \dots, \mu(B_r))$ のディリクレ分布に従うとき μ^* はパラメーター μ のディリクレ過程と呼ばれる。ただし、 $\mu(B_i) = 0$ ならば $\mu^*(B_i) = 0$ とする。詳しくは Ferguson(1973) を参照されたい。

定義：ポリア系列

X 値の確率変数系 $\{X_1, X_2, X_3, \dots\}$ が次の条件を満たすとき、パラメーター μ のポリア系列という。任意の可測部分集合 $B \subset X$ に対して

$$(1) P(X_1 \in B) = \frac{\mu(B)}{\mu(X)}, \quad j=1, 2, \dots, r$$

$$(2) P(X_{n+1} \in B | X_1, X_2, \dots, X_n) = \frac{\mu_n(B)}{\mu_n(X)}, \quad \text{ただし } \mu_n(X) = \mu(X) + n$$

$$\mu_n(B) = \mu(B) + \sum_{i=1}^n \delta_{X_i}(B), \quad \delta_{X_i}(B) = \begin{cases} 1 & \text{if } X_i \in B \\ 0 & \text{if } X_i \notin B \end{cases}$$

これは X が有限集合 $X = \{1, 2, \dots, r\}$ のとき、最初 r 色の球がそれぞれ $\mu_1, \mu_2, \dots, \mu_r$ 個ずつ入っている壺から 1 個球を取り出し、その球と同色の球を 1 個付け加えて壺に戻すという試行を繰り返す、ポリアの壺を一般化したものである。この時、次の定理が成り立つ。

定理 4. 1 2 Blackwell&MacQueen(1973)

$\{X_n; n=1, 2, 3, \dots\}$ をパラメーター μ のポリア系列とする。このとき、次が成り立つ。

$$(1) m_n = \frac{\mu_n}{\mu_n(X)} \text{ は } n \rightarrow \infty \text{ のとき、} X \text{ 上のランダム確率測度 } \mu^* \text{ に確率 1 で収束する。}$$

(2) μ^* はパラメーター μ の Ferguson 分布を持つ。

(3) μ^* を 1 つ与えたとき (μ^*_0 とする)、 X_1, X_2, X_3, \dots は分布 μ^*_0 に従う独立な確率変数系である。

(証明)

ここでは X が有限集合 $X = \{1, 2, \dots, r\}$ の場合のみ示す。一般の場合は Blackwell and MacQueen(1973)を参照されたい。

μ^* を X 上のパラメーター μ のディリクレ過程とする。すなわち、 $(\mu^*(1), \mu^*(2), \dots, \mu^*(r))$

はパラメーター $(\mu(1), \dots, \mu(r))$ のディリクレ分布に従い、 $\sum_{j=1}^r \mu^*(j) = 1$ 。また μ_ω^* を与えた

とき、 X_1, X_2, X_3, \dots を分布 μ_ω^* に従う独立な確率変数系とする。 π_n を X_1, X_2, \dots, X_n の経

験分布、 $\pi_n(j) = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}(j)$ とする。大数の強法則より $n \rightarrow \infty$ のとき π_n は μ_ω^* に確率1で

収束する。 $m_n = \frac{\mu + n\pi_n}{\mu(X) + n}$ より、 $n \rightarrow \infty$ のとき、確率1で μ^* に収束する。従って最後に

$\{X_1, X_2, \dots\}$ がパラメーター μ のポリア系列であることを証明すればよい。その為には
 $A = \{X_1 = x_1, X_2 = x_2, \dots, X_n = x_n\}$, $x_i \in \{1, 2, \dots, r\}$ とすると、その確率が

$$P(A) = \frac{\prod_{j=1}^r (\mu(j))_{(n(j))}}{\mu(X)_{(n)}} \text{、ただし } n(j) = \#\{i; x_i = j\}, a_{(n)} = a(a+1)\cdots(a+n-1) \text{ で与えら}$$

れることを示せば十分である。 μ_ω^* を与えたとき、 X_1, X_2, X_3, \dots は独立で μ_ω^* に従うので

$$P(A|\mu_\omega^*) = \prod_{i=1}^n P(X_i = x_i | \mu_\omega^*) = \prod_{j=1}^r (\mu_\omega^*(j))^{n(j)} \text{、 } \mu_\omega^* \text{ の分布で平均し}$$

$$P(A) = E \left[\prod_{j=1}^r (\mu^*(j))^{n(j)} \right] \quad \text{とすると} \quad (4.26)$$

$(\mu^*(1), \mu^*(2), \dots, \mu^*(r))$ はパラメーター $(\mu(1), \dots, \mu(r))$ のディリクレ分布に従うので、

$\mu(i) = \alpha_i$, $n(j) = n_j$ と書くと(4.26)より

$$\begin{aligned} P(A) &= \frac{\Gamma(\alpha_1 + \dots + \alpha_r)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_r)} \int \cdots \int \left(\prod_{j=1}^r x_j^{\alpha_j + n_j - 1} \right) dx_1 \cdots dx_{r-1} \\ &= \frac{\Gamma(\alpha_1 + \dots + \alpha_r)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_r)} \times \frac{\Gamma(\alpha_1 + n_1) \Gamma(\alpha_2 + n_2) \cdots \Gamma(\alpha_r + n_r)}{\Gamma(\alpha_1 + \dots + \alpha_r + n_1 + \dots + n_r)} \\ &= \frac{\Gamma(\alpha_1 + \dots + \alpha_r)}{\Gamma(n + \alpha_1 + \dots + \alpha_r)} \times \frac{\Gamma(\alpha_1 + n_1)}{\Gamma(\alpha_1)} \times \frac{\Gamma(\alpha_2 + n_2)}{\Gamma(\alpha_2)} \times \cdots \times \frac{\Gamma(\alpha_r + n_r)}{\Gamma(\alpha_r)} \\ &= \frac{\alpha_{1(n_1)} \alpha_{2(n_2)} \cdots \alpha_{r(n_r)}}{(\alpha_1 + \dots + \alpha_r)_{(n)}} = \frac{\prod_{j=1}^r (\mu(j))_{(n(j))}}{(\mu(X))_{(n)}} \end{aligned}$$

よって $\{X_1, X_2, X_3, \dots\}$ はポリア系列である。

(証明終わり)

4. 5 ポアソン・ディリクレ分布

4. 5. 1 ポアソン・ディリクレ分布

領域 $\Delta_n = \{\bar{x} = (x_1, \dots, x_n); x_i \geq 0, \sum_{i=1}^n x_i \leq 1\}$ 上の対称ディリクレ分布 $D(\alpha, n+1)$ の分

$$\text{布密度は } f_n(\bar{x}; \alpha) = \frac{\Gamma((n+1)\alpha)}{\{\Gamma(\alpha)\}^n} \left(\prod_{i=1}^{n+1} x_i \right)^{\alpha-1} \quad \text{ただし } x_{n+1} = 1 - \sum_{i=1}^n x_i \quad (4.27)$$

明らかに、 $E[x_i] = \frac{1}{n+1}$ で $\lim_{n \rightarrow \infty} E[x_i] = 0$ より $n \rightarrow \infty$ における $\bar{x} = (x_1, \dots, x_n)$ の極限分布を

うまく構成できない。Kingman(1975)は (x_1, x_2, \dots, x_n) の順序統計量 $x_{(1)} \geq x_{(2)} \geq \dots \geq x_{(n)}$

の分布を考え、 $n \rightarrow \infty$ で極限分布 (ポアソン・ディリクレ分布) が存在することを示した。

$$\nabla_\infty = \left\{ \bar{x} = (x_1, x_2, x_3, \dots; x_1 \geq x_2 \geq x_3 \geq \dots, \sum_{i=1}^{\infty} x_i = 1 \right\} \text{ とする。}$$

定理 4. 13 (Kingman(1975))

$0 < \theta < \infty$ に対して、次の性質を持つ ∇_∞ 上の測度 P_θ が存在する。

(x_1, x_2, \dots, x_n) が対称ディリクレ分布に従うとき、任意の k に対してその順序統計量

$x_{(1)} \geq x_{(2)} \geq \dots \geq x_{(k)}$ の分布は $n\alpha \rightarrow \theta, n \rightarrow \infty$ の極限で P_θ の k 次元周辺分布に収束する。

(証明)

$\{\xi(t); t \geq 0\}$ を Γ 過程とする。すなわち、 $\xi(0) = 0$, $\xi(t_2) - \xi(t_1)$ は shape parameter $\alpha = t_2 - t_1$, scale parameter $\lambda = 1$ の Γ 分布に従う。 Γ 過程は独立増分を持つ加法過程でジャンプのみで増加する。各 n について

$$q_j = q_j(n) = \{\xi(j\alpha) - \xi((j-1)\alpha)\} / \xi(n\alpha) \quad (4.27)$$

とすると、定理 4. 5 より $q = (q_1, q_2, \dots, q_n)$ は対称ディリクレ分布 $D(\alpha; n+1)$ に従う。

$q = (q_1, q_2, \dots, q_n)$ の順序統計量を $(q_{(1)}, q_{(2)}, \dots, q_{(n)})$ とする。(4.27) より $n\alpha \rightarrow \theta, n \rightarrow \infty$

の極限を取ると、後述するように極限が存在し

$$\lim_{n \rightarrow \infty} q_{(j)} = \frac{\delta \xi_{(j)}}{\xi(\theta)} \quad (4.28)$$

ここで $\delta \xi_{(j)}$ は Γ 過程 $\xi(t)$ の区間 $(0, \theta)$ における j 番目に大きなジャンプを表す。

$\xi(\theta) = \sum_{j=1}^{\infty} \delta \xi_{(j)}$ より $\left(\frac{\delta \xi_{(1)}}{\xi(\theta)}, \frac{\delta \xi_{(2)}}{\xi(\theta)}, \dots \right) \in \nabla_\infty$ であり、極限分布 P_θ の存在が示される。

そこで最後に(4.28)を証明する。Γ過程の一つのサンプルに対して N_0 を十分大きく取り、 $N \geq N_0$ のとき大きな方から k 個のジャンプ $\{\delta_{\xi_{(j)}}; j=1,2,\dots,k\}$ が $(0, N\alpha_N)$ を N 等分した異なる区間 $((i-1)\alpha_N, i\alpha_N)$ に生じているようにしておく。(4.27)より

$$\xi(N\alpha_N)q_{(j)} \geq \delta_{\xi_{(j)}} \quad , \quad j=1,2,\dots,k \quad (N \geq N_0), \quad \text{故に} \quad \liminf_{N \rightarrow \infty} q_{(j)} \geq \frac{\delta_{\xi_{(j)}}}{\xi(\theta)} \quad (4.29)$$

k は任意なので全ての $j=1,2,3,\dots$ について(4.29)は成り立つ。

$$\text{また} \quad \limsup_{N \rightarrow \infty} q_{(j)} = \limsup_{N \rightarrow \infty} \left(1 - \sum_{i \neq j} q_{(i)} \right) \leq 1 - \sum_{i \neq j} \liminf_{N \rightarrow \infty} q_{(i)} \leq 1 - \sum_{i \neq j} \frac{\delta_{\xi_{(i)}}}{\xi(\theta)} = \frac{\delta_{\xi_{(j)}}}{\xi(\theta)} \quad (4.30)$$

(4.29)(4.30)より(4.28)が成り立つ。 (証明終わり)

より一般に非対称ディリクレ分布 $D(\alpha_1, \alpha_2, \dots, \alpha_{n+1})$ の場合はΓ過程 $\xi(t)$ に対して、

$$q_{(j)}(n) = \frac{\xi(\alpha_1 + \dots + \alpha_j) - \xi(\alpha_1 + \dots + \alpha_{j-1})}{\xi(\alpha_1 + \dots + \alpha_{n+1})} \quad (4.31)$$

と置き、 $n \rightarrow \infty$, $\alpha_1 + \dots + \alpha_{n+1} \rightarrow \theta$ かつ $\max(\alpha_1, \dots, \alpha_{n+1}) \rightarrow 0$ となる様に極限をとると

$q_{(j)}(n)$ が収束することが同様に示される。この極限分布をポアソン・ディリクレ分布と呼び、 $PD(\theta)$ で表す。 $PD(\theta)$ の具体的な形は扱い易い形ではない。Watterson(1976)は対称

ディリクレ分布の順序統計量 $(x_{(1)}, x_{(2)}, \dots, x_{(k)})$ の分布から求めているので、結果だけを紹

介しよう。 $PD(\theta)$ の k 次元周辺分布を $f_\theta(x_1, x_2, \dots, x_k)$ とすると

$$f_\theta(x_1, x_2, \dots, x_k) = \theta^k \Gamma(\theta) e^{\gamma\theta} g\left(\frac{1-x_1-\dots-x_k}{x_k}\right) \left(\prod_{i=1}^{k-1} x_i\right)^{-1} x_k^{\theta-2} \quad (4.32)$$

ここで $1 \geq x_1 \geq x_2 \geq \dots \geq x_k \geq 0$, $\sum_{i=1}^k x_i \leq 1$ $k=1,2,3,\dots$

$g(\cdot)$ はラプラス変換が次の式で与えられる無限分解可能な分布密度関数である。

$$E[\exp(-tz)] = \int_0^\infty e^{-tz} g(z) dz = \exp(-\gamma\theta) t^{-\theta} \exp[-\theta E_1(t)]$$

ただし $E_1(t) = \int_t^\infty e^{-x} x^{-1} dx = -\gamma - \log t - \sum_{k=1}^\infty \frac{(-t)^k}{k \cdot k!}$, $\gamma = 0.57721$ (オイラーの定数)。

4. 5. 2 ポアソン・ディリクレ分布の size-biased permutation と Ewens のサンプリング公式

前節で述べたようにディリクレ分布から順序統計量の分布の極限としてえられたポアソン・ディリクレ分布は扱い易い形とは言い難い。しかし、以下で示すように 4. 4. 2

節で導入した size-biased permutation の分布を考えると非常に扱いやすい表現が得られる。

すなわち、対称ディレクレ分布 $D(\alpha; n)$ の size-biased permutation は Z_j ($1 \leq j \leq n-1$)

がベータ分布 $Beta(1+\alpha, (n-j)\alpha)$ に従う一般化されたディレクレ分布であり、

$$\bar{P}^S = (P_1^S, P_2^S, \dots, P_k^S), \quad P_1^S = Z_1, P_2^S = (1-Z_1)Z_2, P_3^S = (1-Z_1)(1-Z_2)Z_3, \dots,$$

$$P_n^S = (1-Z_1)\dots(1-Z_{n-1})Z_n \quad (\text{ただし } Z_n = 1) \text{ と表される。これより、 } n\alpha \rightarrow \theta,$$

$n \rightarrow \infty$ の極限をとると、 Z_j ($j=1,2,3,\dots$) はベータ分布 $Beta(1, \theta)$ に従う独立な確率変数で

$$\text{あり、 } P_1^S = Z_1, P_2^S = Z_1(1-Z_2), \dots, P_n^S = Z_1(1-Z_2)\dots(1-Z_{n-1})Z_n, \dots \text{ と表される。}$$

この $P^S = (P_1^S, P_2^S, P_3^S, \dots)$ がポアソン・ディレクレ分布の size-biased permutation であ

る。ベータ分布 $Beta(1, \theta)$ に従う独立な確率変数系 $\{Z_j (j=1,2,3,\dots)\}$ に対して上で定義した

確率変数系 $P_1^S, P_2^S, \dots, P_n^S, \dots$ は G.E.M.(Griffiths-Engen-McCloskey)分布と呼ばれている。

4. 6 Ewens のサンプリング公式とポリアの壺様モデル

定理 4. 1 で Ewens のサンプリング公式を遺伝子系図を用いて導いたが、ポリアの壺モデルとの関係が Hoppe(1984)によって示された。このポリアの壺モデルを用いたサンプリング公式に関連した種々の性質の証明と遺伝子系図との関係などについて Hoppe(1984, 1987)および Donnelly(1986)に従って紹介する。

4. 6. 1 ポリアの壺様モデル

1 個の黒球と他の種々の色の球が入った壺を考える。黒球の重さを θ 、他の色の球は全て重さ 1 とする。壺からそれぞれの球の重さに比例した確率でランダムに 1 個の球を取り出す。取り出した球が黒以外の球であれば、その球と同じ色の球を 1 個付け加えて壺に戻す。取り出した球が黒球のときは壺の中にはない新しい色の球を付け加えて黒球と一緒に壺に戻す。最初 ($t=0$)、壺の中には黒球 1 個のみ入っているとす。時刻 $t=1$ に黒球を取り出し、最初に黒球と一緒に壺に入れる球のラベル (色) を 1 とする。以後は新たに壺に付加する色を順に 2, 3, ... と自然数でラベルする。 n 回目の試行で壺に戻す球のラベルを確率変数 X_n で現わす。例えば $X_1 = 1, X_2 = 1, X_3 = 2, X_4 = 1, X_5 = 3$ などとなる。この 5 回の試行の結果、壺の中にはラベル 1 の球が 3 個とラベル 2, 3 の球がそれぞれ 1 個、黒球が 1 個の計 6 個の球が入っている。 n 回の試行の後、壺の中の球のラベルの種類を確率変数 K で表す。この壺モデルは黒球を入れることにより新しい色の球が付加されて行くのでポリ

アの壺様モデルと呼ぶことにする。ラベル $i (1 \leq i \leq K)$ の球の数を n_i とすると $n_1 + n_2 + \dots + n_K = n$ であり、この占有数 $\{n_1, n_2, \dots, n_K\}$ によって次の様に n の分割 (partition) が決まる。 $1 \leq i \leq n$ に対して、 i 個の球を含むラベルの種類を a_i とする。すなわち、 $a_i = \#\{j; n_j = i\}$, $\bar{a} = (a_1, \dots, a_n)$ とする。明らかに $\sum_{i=1}^n ia_i = n$ である。 n 個の球が全て同じラベル (同色) ならば $\bar{a} = (0, \dots, 0, 1)$ であり、全て異なるならば $\bar{a} = (n, 0, \dots, 0)$ である。 $\bar{a} = (a_1, a_2, \dots, a_n)$ を n のアレル分割 (allelic partition) と呼ぶ。アレル分割と呼ぶ理由は n 回の試行によって得られる確率変数系 $\{X_1, X_2, \dots, X_n\}$ によって決まる n のアレル分割を Π_n とするとき、 Π_n の分布が Ewens のサンプリング公式 (定理 4. 1) と同じ分布になることによる。すなわち次の定理が成り立つ。

定理 4. 1 4 (Hoppe(1984))

Π_n は次の分布に従うマルコフ過程である。

$$P(\Pi_n = \bar{a}) = \frac{n!}{g_{(n)}} \prod_{i=1}^n \frac{g^a}{i^{a_i} a_i!}, \quad g_{(n)} = g(g+1)\dots(g+n-1), \quad \bar{a} = (a_1, \dots, a_n), \quad \sum_{i=1}^n ia_i = n$$

(証明)

ポリアの壺様モデル $\{X_1, X_2, \dots, X_n\}$ から決まる占有数を $\{n_1, n_2, \dots, n_K\}$ 、さらにこの占有数から得られるアレル分割を $\bar{a} = (a_1, a_2, \dots, a_n)$ とする。占有数が $\{n_1, n_2, \dots, n_K\}$ となる一つのサンプルを $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ とすると、その確率は次式で与えられる。

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = \frac{g^K}{g_{(n)}} \prod_{i=1}^K (n_i - 1)! \quad (4.33)$$

なぜなら 1 回球を取り出すごとに壺の中の球の数が一つずつ増えることより、分母は $g_{(n)}$ となる。分子の g^K は K 種類のラベルの球があることより、黒球が K 回取り出され、そのたびに新たなラベルの球が付加されたことに対応する。さらにその各々に対してラベル $j (1 \leq j \leq K)$ の球が 1 個ずつ $n_j - 1$ 回付加されるので、その場合の数が $\prod_{i=1}^K (n_i - 1)!$ となる。以上より (4.33) が得られる。従って、アレル分割 $\bar{a} = (a_1, a_2, \dots, a_n)$ となる異なるサンプル

$$\{X_1, X_2, \dots, X_n\} \text{ が } \frac{n!}{\left(\prod_{i=1}^n i^{a_i} a_i!\right) \prod_{j=1}^K (n_j - 1)!} = \frac{n!}{\left(\prod_{i=1}^n a_i!\right) \left(\prod_{j=1}^K n_j!\right)} \text{ 通りあることを示せば、こ}$$

の因子を (4.33) に掛けて証明が終わる。 i 個のサンプルを含むアレルタイプの数が a_i 種あり、

それらのラベルの与え方はアレル分割 $\vec{a} = (a_1, a_2, \dots, a_n)$ に対して $\prod_{i=1}^n a_i!$ 通りある。

また、 n_j 個は同じアレルタイプなので、その並べ方は全部で $\prod_{j=1}^K n_j!$ 通りある。ポリアの壺様モデルでは出現した順にラベルを自然数で与えるので、全て異なるときの順列 $n!$ をアレルのラベルと同じアレルの順列から求めた上記の場合の数の積 $\prod_{i=1}^n a_i! \prod_{j=1}^K n_j!$ で割るとアレル分割 $\vec{a} = (a_1, a_2, \dots, a_n)$ に対応するポリアの壺様モデルのサンプル $\{X_1, X_2, \dots, X_n\}$ の場合の数となる。

定理 4. 1 2 でポリアの壺モデルとディリクレ分布の関係について述べた。この定理の応用として、修正されたポリアの壺モデル（ポリアの壺様モデル）とポアソン・ディリクレ分布について関係が導かれる。 $i = 1, 2, \dots, K$ に対して、ラベル i で重さ 1 の球が α_i 個ずつ入っている壺を考える。この壺から毎時刻 1 個の球をランダムに取り出し、その球と同じラベルの球を付け加えて戻すポリアの壺モデルを考える。ただし、 α_i が非整数の場合は α_i に比例した確率でラベル i の球は取り出されたとする。 n 回目の試行で壺に付加した球のラベルを X_n , $\sum_{i=1}^K \alpha_i = \tilde{g}$ とする。

このとき、明らかに $P(X_{n+1} = i | X_1, X_2, \dots, X_n) = \frac{\alpha_i + \nu_i(n)}{\tilde{g} + n}$ が成り立つ。ここで、 $\nu_i(n)$ は n 回の試行の後、壺の中のラベル i の球の数を表す。定理 4. 1 2 より

$\lim_{n \rightarrow \infty} \frac{\nu_i(n)}{n} = p_i$, $\sum_{i=1}^K p_i = 1$ (a.s.) であり、 $P = (p_1, p_2, \dots, p_K)$ はディリクレ分布 $D(\alpha_1, \dots, \alpha_K)$ に従う。以上を準備として次の定理を証明する。

定理 4. 1 5 (Hoppe(1987))

ポリアの壺様モデルのサンプル $\{X_1, X_2, \dots, X_n\}$ において n 回の試行の後の壺の中のラベル i の球の数を $S_i(n)$ とすると $\lim_{n \rightarrow \infty} \frac{S_i(n)}{n} = P_i$, $\sum_{i=1}^{\infty} P_i = 1$ (a.s.) が成り立つ。

(証明)

$\{\lambda_i; i = 1, 2, \dots\}$ を分布 $P(\lambda_i = 1) = \frac{g}{g+i-1}$, $P(\lambda_i = 0) = \frac{i-1}{g+i-1}$ に従う独立な確率変数系とする。すなわち、 $\lambda_i = 1$ のとき黒球が取り出され、新しいラベルの球が壺に入れられる。マルコフ時刻 τ_i ($i = 1, 2, \dots$) を $\tau_1 = 1$, $i \geq 2$ に対して $\tau_i = \min\{j > \tau_{i-1}; \lambda_j = 1\}$ で定義する。 τ_i はラベル i の球が初めて壺に入れられた時刻を表す。整数 $k (\geq 1)$ を固定し、 $n = 1, 2, \dots$

に対して、球が初めて壺に入れられた時刻によって球の再ラベリングをしたマルコフ連鎖 $\{Y_n^{(k)}; n=1,2,3,\dots\}$ を次のように定義する。

$$Y_n^{(k)} = \begin{cases} j & (1 \leq j \leq k, \lambda_j = 1, X_{n+k} = X_j \text{ のとき}) \\ 0 & (\text{その他}) \end{cases} \circ$$

最初のポリアの壺様モデル $\{X_1, X_2, \dots, X_n\}$ では出現する順に 1, 2, 3, \dots と球にラベリングしたが、 $\{Y_n^{(k)}; n=1,2,\dots\}$ は k 回の試行 $\{X_1, X_2, \dots, X_k\}$ の条件下での黒球を含めて最大 $k+1$ 色のポリアの壺モデルである。そして初期条件としての k 回の試行 $\{X_1, X_2, \dots, X_k\}$ の中でその色の球が最初に付加された時刻をその色のラベルとし、黒球はラベルを 0 とした再ラベリングである。そして、 $\alpha_i^{(k)} = \begin{cases} \lambda_i S_{X_i}(k) & (1 \leq i \leq k) \\ \vartheta & (i=0) \end{cases}$ によって

$(\alpha_0^{(k)}, \alpha_1^{(k)}, \dots, \alpha_k^{(k)})$ を定義する。例えば $k=7$ で、サンプルが

$(X_1, X_2, X_3, X_4, X_5, X_6, X_7) = (1, 1, 2, 1, 3, 2, 3)$ とすると、

$(\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5, \lambda_6, \lambda_7) = (1, 0, 1, 0, 1, 0, 0)$ であり、 $\alpha_1^{(7)} = \lambda_1 S_{X_1}(7) = S_1(7) = 3$,

$\alpha_2^{(7)} = \lambda_2 S_{X_2}(7) = 0$, $\alpha_3^{(7)} = \lambda_3 S_{X_3}(7) = S_2(7) = 2$, $\alpha_4^{(7)} = \lambda_4 S_{X_4}(7) = 0$,

$\alpha_5^{(7)} = \lambda_5 S_{X_5}(7) = S_3(7) = 2$, $\alpha_6^{(7)} = \lambda_6 S_{X_6}(7) = 0$, $\alpha_7^{(7)} = \lambda_7 S_{X_7}(7) = 0$

すなわち、 $(\alpha_0^{(7)}, \alpha_1^{(7)}, \alpha_2^{(7)}, \alpha_3^{(7)}, \alpha_4^{(7)}, \alpha_5^{(7)}, \alpha_6^{(7)}, \alpha_7^{(7)}) = (\vartheta, 3, 0, 2, 0, 2, 0, 0)$ となる。

$\{Y_n^{(k)}; n=1,2,\dots\}$ は初期条件 $(\alpha_0^{(k)}, \alpha_1^{(k)}, \dots, \alpha_k^{(k)})$ のポリアの壺モデルになる。ポリアの壺

モデル $\{Y_n^{(k)}; n=1,2,3,\dots\}$ において n 回の試行でラベル i 、すなわち時刻 i に最初に壺に付加

された球が $k+1$ 回目以降、 n 回の試行で壺に付加された回数を

$S_i(n, k) = \#\{j; Y_j^{(k)} = i, 1 \leq j \leq n\}$ ($i=0,1,2,\dots,k$) とする。定理 4. 1 2 より

$$P\left(\lim_{n \rightarrow \infty} \frac{S_i(n, k) + \alpha_i^{(k)}}{n + k + \vartheta} = \gamma_i^{(k)} \mid X_1, X_2, \dots, X_k\right) = 1, \quad \sum_{i=0}^k \gamma_i^{(k)} = 1. \quad \gamma_i^{(k)} \text{ は } (\alpha_0^{(k)}, \dots, \alpha_k^{(k)})$$

に依存して決まる確率変数である。全てのサンプル $\{X_1, X_2, \dots, X_k\}$ について、期待値を取

ると、確率 1 で極限 $\lim_{n \rightarrow \infty} \frac{S_i(n, k) + \alpha_i^{(k)}}{n + k + \vartheta}$ が存在することが分かる。他方、 $S_i(n, k)$ の定義よ

り $S_i(n, k) + \alpha_i^{(k)} = \alpha_i^{(n+k)}$ なので、 $\lim_{n \rightarrow \infty} \frac{S_i(n, k) + \alpha_i^{(k)}}{n + k + \vartheta} = \lim_{n \rightarrow \infty} \frac{\alpha_i^{(n+k)}}{n + k + \vartheta} = \gamma_i^{(k)}$ であり極限值は k によらない。この極限を $\lim_{n \rightarrow \infty} \frac{\alpha_i^{(n+k)}}{n + k + \vartheta} = \lim_{n \rightarrow \infty} \frac{\alpha_i^{(n)}}{n} = \gamma_i$ ($1 \leq i < \infty$) (a.s) とする。

条件 $\{X_1, X_2, \dots, X_k\}$ の下で、定理 4. 1 2 より $(\gamma_0^{(k)}, \gamma_1^{(k)}, \dots, \gamma_k^{(k)})$ はディリクレ分布

$D(\alpha_0^{(k)}, \alpha_1^{(k)}, \dots, \alpha_k^{(k)})$ (ただし $\alpha_0^{(k)} = \theta$) に従う。その周辺分布として $\gamma_0^{(k)}$ はベータ分布

$Beta(\theta, \sum_{i=1}^k \alpha_i^{(k)}) = Beta(\vartheta, k)$ に従うので、 $E[\gamma_0^{(k)} | X_1, \dots, X_k] = \frac{\vartheta}{\vartheta + k}$ となる。

全ての初期条件 $\{X_1, X_2, \dots, X_k\}$ に関して期待値を取ると $E[\gamma_0] = E\left[1 - \sum_{i=1}^k \gamma_i\right] = \frac{\vartheta}{\vartheta + k}$ 。

これより、 $E\left[1 - \sum_{i=1}^{\infty} \gamma_i\right] = E\left[1 - \lim_{k \rightarrow \infty} \sum_{i=1}^k \gamma_i\right] = \lim_{k \rightarrow \infty} E\left[1 - \sum_{i=1}^k \gamma_i\right] = \lim_{k \rightarrow \infty} \frac{\vartheta}{\vartheta + k} = 0$ 。

$1 - \sum_{i=1}^{\infty} \gamma_i \geq 0$ (a.s.) なので、 $\sum_{i=1}^{\infty} \gamma_i = 1$ である。最初のポリアの壺様モデル $\{X_1, X_2, \dots, X_n\}$

に戻って、ラベル i の球が最初に付加される時刻を τ_i とすると、 $n \geq \tau_i$ のとき、 $S_i(n) = \alpha_{\tau_i}^{(n)}$

である。これより $\lim_{n \rightarrow \infty} \frac{S_i(n)}{n} = \gamma_{\tau_i} \equiv P_i$ (a.s.) と置くことにする。 $(\gamma_1(\omega), \gamma_2(\omega), \gamma_3(\omega), \dots)$

の中で正の値は $(\gamma_{\tau_1}(\omega), \gamma_{\tau_2}(\omega), \gamma_{\tau_3}(\omega), \dots)$ であり、ポリアの壺様モデルにおいて有限個の

色 (ラベル) しか現れない確率は

$$P(\text{有限個の } i \text{ のみ } \tau_i < \infty) = \lim_{n \rightarrow \infty} P(\lambda_i = 0; \text{全ての } i \geq n) = \lim_{n \rightarrow \infty} \prod_{i=n}^{\infty} \frac{i-1}{\vartheta+i-1} = 0。$$

よって、確率 1 で無限個のラベルが現れるので $\sum_{i=1}^{\infty} P_i = \sum_{i=1}^{\infty} \gamma_{\tau_i} = 1$ である。

定理 4. 1 6

$P_1, P_2, \dots, P_n \dots$ を定理 4. 1 5 で定義される確率変数とすると、

(1) $Z_n = \frac{P_n}{\sum_{i=n}^{\infty} P_i}$, ($n = 1, 2, 3, \dots$) とすると、 $\{Z_1, Z_2, Z_3, \dots\}$ はベータ分布 $Beta(1, \vartheta)$ に従う

独立な確率変数系である。

(2) $P_1 = Z_1, P_2 = Z_2(1 - Z_1), \dots$ 、一般に $P_n = Z_n \prod_{i=1}^{n-1} (1 - Z_i)$ ($n \geq 2$)。

すなわち $\{P_1, P_2, \dots, P_n \dots\}$ はポアソン・ディリクレ分布の size-biased permutation で

ある。

(3) $\lambda_j = 1$ のとき、 γ_j は $Beta(1, \vartheta + j - 1)$ に従う。

(証明)

(1) 第4. 4. 2節の式(4.23)で述べたように、 Y_1, Y_2, \dots, Y_{n+1} がディリクレ分布

$D(\alpha_1, \alpha_2, \dots, \alpha_{n+1})$ に従うとき、 $\vec{Y} = (Y_1, \dots, Y_n)$ は完全中立な確率ベクトルであり、

$$U_j = Y_j / \left(1 - \sum_{i=1}^{j-1} Y_i\right) = Y_j / \left(\sum_{i=j}^{n+1} Y_i\right) \quad (j=1, 2, \dots, n)$$

は独立でベータ分布 $Beta(\alpha_j, \sum_{i=j+1}^{n+1} \alpha_i)$ に従う。ポリアの壺様モデル $\{X_i; i \geq 1\}$ において、ラベル n の球が最初に壺に入れられた時刻を τ_n とする。時刻 τ_n でのラベル k ($1 \leq k \leq n$) の球の数を $\alpha_k = S_k(\tau_n)$ とする。明らかに $\alpha_n = 1$ であり、壺の中には黒球が 1 個あるので $\alpha_{n+1} = \vartheta$ とする。初期条件

$(\alpha_1, \dots, \alpha_n, \alpha_{n+1})$ の下で、ポリアの壺モデルを $\{Y_k^{(n)}; k=1, 2, 3, \dots\}$ とする。ここでラベル $n+1$ は黒色であり、黒色の球が取り出されたときは、黒球を 1 個付加して壺に戻す。条件 $(X_1, X_2, \dots, X_{\tau_n})$ の下で $(P_1, P_2, \dots, P_{n+1})$ はディリクレ分布 $D(\alpha_1, \alpha_2, \dots, \alpha_{n+1})$ に従う。故

に、 $(X_1, X_2, \dots, X_{\tau_n})$ の下で $Z_i = P_i / \sum_{j=i}^{n+1} P_j$ ($1 \leq i \leq n$) はベータ分布 $Beta(\alpha_i, \sum_{j=i+1}^{n+1} \alpha_j)$ に従

う独立な確率変数である。特に $\alpha_n = 1, \alpha_{n+1} = \vartheta$ より条件 $(\alpha_1, \dots, \alpha_n, \alpha_{n+1})$ の下で Z_1, Z_2, \dots, Z_n の同時分布は

$$\left[\prod_{i=1}^{n-1} \frac{\Gamma(\alpha_i + \beta_i)}{\Gamma(\alpha_i)\Gamma(\beta_i)} z_i^{\alpha_i-1} (1-z_i)^{\beta_i-1} \right] \times \vartheta (1-z_n)^{\vartheta-1} \quad \text{ただし } \beta_i = \sum_{j=i+1}^{n+1} \alpha_j.$$

初期条件 $(\alpha_1, \dots, \alpha_n, \alpha_{n+1})$ の下で、 (Z_1, \dots, Z_{n-1}) と Z_n は独立である。

初期条件 $(\alpha_1, \dots, \alpha_{n-1}, \alpha_n, \alpha_{n+1}) = (\alpha_1, \dots, \alpha_{n-1}, 1, \vartheta)$ の確率分布は $(\alpha_1, \dots, \alpha_{n-1})$ に関する確率分布なのでこれを上式に掛けると、 Z_1, Z_2, \dots, Z_n の同時分布は適当な関数 Φ を用いて

$\Phi(z_1, \dots, z_{n-1}) \vartheta (1-z_n)^{\vartheta-1}$ と書ける。これより (Z_1, \dots, Z_{n-1}) と Z_n は独立で Z_n の分布は

$\vartheta (1-z_n)^{\vartheta-1}$ となることが分かる。 n は任意なので、全ての自然数について成り立つ。

$$(2) P_n = Z_n \left(\sum_{i=n}^{\infty} P_i \right) = Z_n (1 - P_1 - \dots - P_{n-1}) = Z_n \left(\frac{1 - P_1 - \dots - P_{n-1}}{1 - P_1 - \dots - P_{n-2}} \right) (1 - P_1 - \dots - P_{n-2})$$

$$\begin{aligned}
&= Z_n \left(1 - \frac{P_{n-1}}{1 - P_1 - \dots - P_{n-2}} \right) (1 - P_1 - \dots - P_{n-2}) = Z_n (1 - Z_{n-1}) (1 - P_1 - \dots - P_{n-2}) \\
&= Z_n \prod_{i=1}^{n-1} (1 - Z_i) \quad (n \geq 2)
\end{aligned}$$

(3) $\lambda_j = 1$ のとき、この球のラベルを k とする。条件 (X_1, X_2, \dots, X_j) の下で

$(P_1, P_2, \dots, P_k, P_{k+1})$ はディリクレ分布 $D(\alpha_1, \alpha_2, \dots, \alpha_{k+1})$ に従う、ただし $\alpha_k = 1$ 、 $\alpha_{k+1} = \mathcal{G}$
 $\sum_{i=1}^{k+1} \alpha_i = j + \mathcal{G}$ である。これより、 P_k の分布はディリクレ分布の周辺分布として、

補題 4. 4 よりベータ分布 $Beta(1, \sum_{i \neq k} \alpha_i) = Beta(1, \mathcal{G} + j - 1)$ に従う。

4. 6. 2 ポアソン・ディリクレ集団からのサンプルとポリアの壺様モデル

4. 5. 2 節で述べたようにポアソン・ディリクレ分布からのサンプルのアレル分割分布は Ewens のサンプリング公式に従う。Watterson(1976)は K 個のアレルを含む定常な集団モデルからのサンプルの分布から極限操作により Ewens のサンプリング公式が導かれることを示した。 K 個のアレルを $1, 2, \dots, K$ でラベルし、その集団内頻度を X_1, X_2, \dots, X_K とする。集団からランダムに n 個の遺伝子をサンプルしたとき、各々のアレルが n_1 個、 n_2 個 \dots, n_k 個含まれる確率は

$$P_n(n_1, \dots, n_K | X) = \frac{n!}{\prod_{i=1}^K n_i!} \left(\prod_{j=1}^K X_j^{n_j} \right), \quad X = (X_1, \dots, X_K), \quad (n_1 + \dots + n_K = n).$$

$X = (X_1, \dots, X_K)$ がディリクレ分布 $D(\varepsilon, K)$ に従うとし、上式の期待値を取ると

$$P_n(n_1, \dots, n_K) = \frac{n!}{\prod_{i=1}^K n_i!} \int_{\Delta_K} \prod_{j=1}^K x_j^{n_j} \times \left\{ \frac{\Gamma(K\varepsilon)}{(\Gamma(\varepsilon))^K} \left(\prod_{j=1}^K x_j^{\varepsilon-1} \right) \right\} dx_1 \dots dx_{K-1}$$

ただし、 $\Delta_K = \{x = (x_1, \dots, x_{K-1}); x_j \geq 0, \sum_{j=1}^{K-1} x_j \leq 1\}$ 。ディリクレ分布の積分より

$$P_n(n_1, \dots, n_K) = \frac{n!}{\prod_{i=1}^K n_i!} \times \frac{\Gamma(K\varepsilon) \prod_{j=1}^K \Gamma(\varepsilon + n_j)}{(\Gamma(\varepsilon))^K \Gamma(K\varepsilon + n)} \quad \text{となる。}$$

n_1, n_2, \dots, n_K の中で $n_j \geq 1$ であるアレルの数を k とし、大きさの順に $n_{(1)} \geq n_{(2)} \geq \dots \geq n_{(k)}$

と並べる。 k と $n_{(1)} \geq n_{(2)} \geq \dots \geq n_{(k)}$ を与えたとき、確率 $P(n_{(1)}, n_{(2)}, \dots, n_{(k)}; k)$ は

$$P(n_{(1)}, n_{(2)}, \dots, n_{(k)}; k) = \frac{n!}{\prod_{i=1}^k n_{(i)}!} \times \frac{\Gamma(K\varepsilon) \prod_{i=1}^k \Gamma(\varepsilon + n_{(i)})}{(\Gamma(\varepsilon))^k \Gamma(K\varepsilon + n)} \times M$$

ここで、Mは異なるK種のアレルを $n_{(1)} \geq n_{(2)} \geq \dots \geq n_{(k)}$ となるように割り当てる場合の数

であり、 $n_{(1)}, n_{(2)}, \dots, n_{(k)}, 0, \dots, 0$ のK個の数を異なるK個の部屋に割り当てる場合の数に等

しい。 $\alpha_j = \# \{i; n_{(i)} = j, 1 \leq i \leq k\}$ とすると、同じものを含む順列になるので

$$M = \frac{K!}{\alpha_1! \alpha_2! \dots \alpha_n! (K-k)!}$$

である。ここで、アレルの数を無限に増やし $K \rightarrow \infty, \varepsilon \rightarrow 0$,

そして $K\varepsilon \rightarrow \vartheta$ となるように極限を取ると

$$\begin{aligned} \lim_{\substack{K \rightarrow \infty \\ K\varepsilon \rightarrow \vartheta}} \frac{M}{(\Gamma(\varepsilon))^k} &= \lim_{K \rightarrow \infty} \left\{ \frac{1}{\alpha_1! \dots \alpha_n!} \times \frac{K!}{(\Gamma(\varepsilon))^k (K-k)!} \right\} \\ &= \frac{1}{\alpha_1! \alpha_2! \dots \alpha_n!} \lim_{K \rightarrow \infty} \left\{ \frac{K}{\Gamma(\varepsilon)} \right\}^k \left(1 - \frac{1}{K} \right) \left(1 - \frac{2}{K} \right) \dots \left(1 - \frac{k-1}{K} \right) \end{aligned}$$

ここで、 $\lim_{\varepsilon \rightarrow 0} \varepsilon \Gamma(\varepsilon) = \lim_{\varepsilon \rightarrow 0} \int_0^\infty \varepsilon t^{\varepsilon-1} e^{-t} dt = \lim_{\varepsilon \rightarrow 0} \Gamma(\varepsilon + 1) = \Gamma(1) = 1$ より

$$\lim_{\substack{K \rightarrow \infty \\ K\varepsilon \rightarrow \vartheta}} \frac{K}{\Gamma(\varepsilon)} = \lim_{\varepsilon \rightarrow 0} \frac{\vartheta}{\varepsilon \Gamma(\varepsilon)} = \vartheta \text{ なので } \lim_{K \rightarrow \infty} \frac{M}{(\Gamma(\varepsilon))^k} = \frac{\vartheta^k}{\alpha_1! \alpha_2! \dots \alpha_n!}.$$

故に $K \rightarrow \infty, K\varepsilon \rightarrow \vartheta$ の

極限を取ると次の Ewens のサンプリング公式を得る。

$$\lim_{\substack{K \rightarrow \infty \\ K\varepsilon \rightarrow \vartheta}} P(n_{(1)}, n_{(2)}, \dots, n_{(k)}; k) = \frac{n! \vartheta^k \Gamma(\vartheta)}{n_{(1)}! n_{(2)}! \dots n_{(k)}! \alpha_1! \alpha_2! \dots \alpha_n! \Gamma(\vartheta + n)}.$$

N個の半数体個体から成る有限集団を Q^N とする。可算無限個のアレルタイプを仮定する。集団 Q^N から非復元抽出により n 個の遺伝子サンプルを取り出したとき、そのサンプル構成

が $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$, $\sum_{i=1}^n i \alpha_i = n$ となる確率を $P_N(\alpha)$ とする。 $N \rightarrow \infty$ のとき、この分布

が Ewens のサンプリング公式(4.1)に収束するとき、この集団の系列 Q^N は Ewens のサンプリング性を持つことにしよう。Kingman(1977)は有限集団の系列 Q^N が Ewens のサンプリング性を持つための必要十分条件はこの集団がポアソン・ディリクレ極限を持つことであることを示した。

4. 6. 3 ポリアの壺様モデルおよび逆進プロセスの推移確率

ポリアの壺様モデルによるアレル分割 Π_n のプロセスは壺の中に黒球 1 個の初期状態 (これを $\Pi_0 = \phi$ で表す) から出発するマルコフ連鎖である。より一般に任意の初期状態 Π_0 から出発するアレル分割 Π_n のマルコフプロセスの 1 ステップ推移確率を

$P(\vec{a}, \vec{b}) = P(\Pi_{n+1} = \vec{b} | \Pi_n = \vec{a})$, $\vec{a} = (a_1, \dots, a_n)$, $\vec{b} = (b_1, \dots, b_{n+1})$ とすると、明らかに

$$P(\vec{a}, \vec{b}) = \begin{cases} \frac{g}{g+n} & \text{if } \vec{b} = (a_1 + 1, a_2, \dots, a_n, 0) \\ \frac{ia_i}{g+n} & \text{if } \vec{b} = (a_1, \dots, a_i - 1, a_{i+1} + 1, \dots, a_n, 0) \\ \frac{na_n}{g+n} & \text{if } \vec{b} = (a_1, \dots, a_n - 1, 1) \end{cases} \quad (4.34)$$

これより時間を逆向きにした逆進プロセスの推移確率を考えると

$$P(\Pi_n = \vec{a} | \Pi_{n+1} = \vec{b}) = \frac{P(\Pi_{n+1} = \vec{b} | \Pi_n = \vec{a})P(\Pi_n = \vec{a})}{P(\Pi_{n+1} = \vec{b})}。 \text{これは初期分布 } \Pi_0 \text{ に依存するの}$$

でポリアの壺様モデルの初期条件 $\Pi_0 = \phi$ とすると、 $P(\Pi_n = \vec{a})$ は定理 4. 14 で与えられる。よって

$$P(\Pi_n = \vec{a} | \Pi_{n+1} = \vec{b}) = \begin{cases} \frac{a_1 + 1}{n+1} & \text{if } \vec{b} = (a_1 + 1, a_2, \dots, a_n, 0) \\ \frac{(i+1)(a_{i+1} + 1)}{n+1} & \text{if } \vec{b} = (a_1, \dots, a_i - 1, a_{i+1} + 1, \dots, a_n, 0) \\ 1 & \text{if } \vec{b} = (a_1, \dots, a_n - 1, 1) \end{cases} \quad (4.35)$$

この逆進プロセスの推移確率行列を $M_{n,n+1}$ で現わす。

$$P_n(\vec{a}) = P(\Pi_n = \vec{a}) = P(\Pi_n = \vec{a} | \Pi_0 = \phi) \text{ とすると}$$

$$P_n(\vec{a}) = \sum_{\vec{b}} P(\Pi_n = \vec{a} | \Pi_{n+1} = \vec{b}) P(\Pi_{n+1} = \vec{b}) \text{ より次の漸化式を得る。}$$

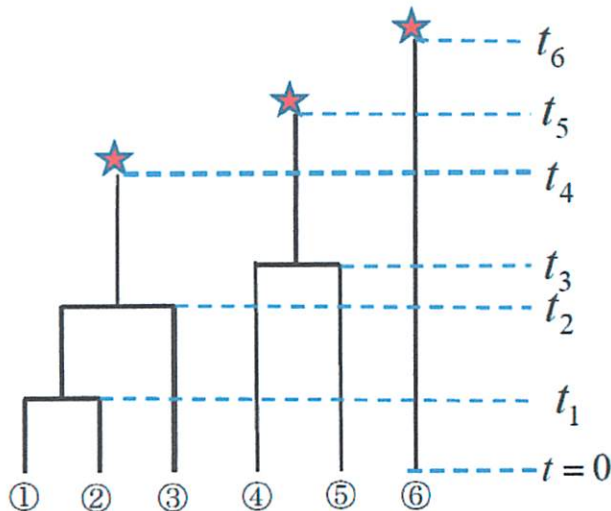
$$\begin{aligned} P_n(a_1, \dots, a_n) &= \frac{a_1 + 1}{n+1} P_{n+1}(a_1 + 1, a_2, \dots, a_n, 0) \\ &\quad + \sum_{r=2}^n \frac{r(a_r + 1)}{n+1} P_{n+1}(a_1, \dots, a_{r-1} - 1, a_r + 1, \dots, a_n, 0) + P_{n+1}(a_1, \dots, a_n - 1, 1) \end{aligned} \quad (4.36)$$

(4.36)は Kingman(1978)によっても Ewens のサンプリングの分割構造(partition structure)が満たす式として導出されている。同様にして任意の自然数 $k < m < n$ に対して行列 M_{mn} をその成分が $M_{mn}(\vec{a}, \vec{b}) = P(\Pi_m = \vec{a} | \Pi_n = \vec{b})$ で定義される行列、 P_m, P_n をその成分が $P_m(\vec{a}), P_n(\vec{b})$ のベクトルとすると、線形変換 $P_m = M_{mn}P_n$ および $M_{kn} = M_{km}M_{mn}$ が成り立つ。これは逆進プロセスの Chapman-Kolmogorov 方程式である。

4. 6. 4 合祖モデルとポリアの壺様モデル

第3章では(3.32)式で Ewens のサンプリング公式が Moranモデルを用いて遺伝子系図過程から導かれることを示した。また、第4章の定理4. 1では推移確率が(3.8)式で与えられる突然変異を含む遺伝子系図過程を用いて同様に Ewens のサンプリング公式を導いた。この節では突然変異を含む遺伝子系図、サンプリング公式およびポリアの様壺モデルの三者の間に深い関連があることを示す。

推移確率(3.8)の連続時間マルコフ連鎖を \mathfrak{R}_t とする。 \mathfrak{R}_t は合祖と突然変異によって1つずつ減少する死滅過程であるが、 \mathfrak{R}_t (ただし $\mathfrak{R}_0 = n$) のサンプルパスとして1つの系図が与えられたとき、次のような同値関係を定義する。 $\mathfrak{R}_t = k$ のとき k 個のメンバーの2つの遺伝子についてその2つの遺伝子が突然変異を経ることなく、共通祖先に到達するとき、この2つの遺伝子は同値とする。すなわちこの2つの遺伝子はある突然変異を祖先として共有していることを意味する。下の図で①と②は同値であるが②と⑤は同値ではない。



また、合祖をしたとき、その合祖によって現れた共通祖先には2つの遺伝子の若い番号をラベルとして与えることにする。上図で①と②が合祖した共通祖先は①というラベルを持つ。 \mathfrak{R}_t の同値関係のプロセスを $\tilde{\mathfrak{R}}_t$ とする。上の系図では $\tilde{\mathfrak{R}}_0 = \{(1,2,3), (4,5), (6)\}$ 、

$$\tilde{\mathfrak{R}}_1 = \{(1,3),(4,5),(6)\}, \tilde{\mathfrak{R}}_2 = \{(1),(4,5),(6)\}, \tilde{\mathfrak{R}}_3 = \{(1)(4)(6)\}, \tilde{\mathfrak{R}}_4 = \{(4),(6)\}$$

$\tilde{\mathfrak{R}}_5 = \{(6)\}, \tilde{\mathfrak{R}}_6 = \phi$ となる。突然変異を起こしていない個体について、その同値関係を時

間を遡って追跡したプロセスである。同値関係のプロセス $\tilde{\mathfrak{R}}_i$ はまた、次のように分割

(partition)のプロセス Ω_i を定義する。上の例では初期時刻にサンプル数がそれぞれ1, 2, 3の3つの部分に分割されているので $\Omega_0 = (a_1, a_2, a_3) = (1,1,1)$ 、以下同様にして

$$\Omega_1 = (1,2,0), \Omega_2 = (2,1,0), \Omega_3 = (3,0,0), \Omega_4 = (2,0,0), \Omega_5 = (1,0,0), \Omega_6 = (0,0,0)$$

となる。 Ω_i のジャンププロセスを $\Omega_k^* (k = n, n-1, \dots, 0)$ とする。上の例では $\Omega_6^* = (1,1,1)$ 、

$$\Omega_5^* = (1,2,0), \Omega_4^* = (2,1,0), \Omega_3^* = (3,0,0), \Omega_2^* = (2,0,0), \Omega_1^* = (1,0,0), \Omega_0^* = (0,0,0)$$

となる。ジャンププロセス $\Omega_k^* (k = n, n-1, \dots, 0)$ の逆進過程を $\Omega_0^*, \Omega_1^*, \Omega_2^*, \dots, \Omega_n^*$ とすると、

次の定理が成り立つ。

定理4. 17

逆進プロセス $\Omega_0^*, \Omega_1^*, \Omega_2^*, \dots, \Omega_n^*$ はポリアの壺様モデルから定義されるアレル分割のプロセス $\Pi_0, \Pi_1, \Pi_2, \dots, \Pi_n$ と確率論的に同値である。

(証明)

$\Omega_k^* (k = 0, 1, 2, \dots, n)$ と $\Pi_k (k = 0, 1, 2, \dots, n)$ の推移確率が一致することを示す。分割

$$\bar{a} = (a_1, a_2, \dots, a_n), \sum_{i=1}^n ia_i = n-1; \bar{b} = (b_1, b_2, \dots, b_n), \sum_{i=1}^n ib_i = n$$

に対して、推移確率を $P(\Omega_n^* = \bar{b} | \Omega_{n-1}^* = \bar{a})$ とする。

(1) $a_{r-1} = b_{r-1} + 1, a_r = b_r - 1$ 、かつ $j \neq r, r-1 (j \leq n-1)$ のとき $a_j = b_j$ の場合

時刻 T_{n-1} に分割 \bar{b} の中の同値な r 個の遺伝子を含む同値類クラスにおいて合祖が起き、そのクラスはサイズが $r-1$ のクラスとして a_{r-1} 個のクラスの1つとなることを意味する。

そこで、結合確率 $P(\Omega_n^* = \bar{b}, \Omega_{n-1}^* = \bar{a}) = P(\Omega_n^* = \bar{b} | \Omega_{n-1}^* = \bar{a})P(\Omega_{n-1}^* = \bar{a})$ を求めよう。

$$P(B) = P(\text{時刻 } T_{n-1} \text{ に起きた事象が合祖である} | \Omega_{n-1}^* = \bar{a})$$

分割 \bar{a} のサイズ $r-1$ のクラスを C_1, C_2, \dots, C_m ; ($m = a_{r-1}$) とする。

$$P(C_i) = P(\text{合祖によってクラス } C_i \text{ を生じる} | \text{時刻 } T_{n-1} \text{ で合祖, } \Omega_{n-1}^* = \bar{a}) \text{ とすると}$$

$$P(\Omega_n^* = \bar{b} | \Omega_{n-1}^* = \bar{a}) = \sum_{i=1}^m P(B)P(C_i) \text{ ただし } m = a_{r-1}. \text{ ポアソン過程の性質より、}$$

$$P(B) = \frac{n(n-1)}{2} / \frac{n(n-1+\theta)}{2} = \frac{n-1}{n-1+\vartheta}.$$

合祖が起きるといふ条件の下で、時間を逆向きにして分割 \bar{a} の a_{r-1} 個のサイズ $r-1$ のある

クラス C_i が合祖の結果生じる確率は $n-1$ 個体間で等確率なので、 $P(C_i) = \frac{r-1}{n-1}$ となる。

$$\text{よって、} P(\Omega_n^* = \bar{b} | \Omega_{n-1}^* = \bar{a}) = \sum_{i=1}^m P(B)P(C_i) = \frac{n-1}{n-1+\vartheta} \times \frac{r-1}{n-1} \times m = \frac{(r-1)a_{r-1}}{n-1+\vartheta}.$$

(2) $a_1 = b_1 - 1, a_j = b_j$ ($j \leq n-1$) の場合

時刻 T_{n-1} に分割 \bar{b} の中の 1 個の遺伝子を含む同値類クラスにおいて突然変異が起き、

除かれることを意味する。従って

$$P(\Omega_n^* = \bar{b} | \Omega_{n-1}^* = \bar{a}) = P(\text{時刻 } T_{n-1} \text{ に突然変異} | \Omega_{n-1}^* = \bar{a}) = \frac{\vartheta}{n-1+\vartheta}$$

$$\text{以上をまとめると} \quad P(\Omega_n^* = \bar{b} | \Omega_{n-1}^* = \bar{a}) = \begin{cases} \frac{(r-1)a_{r-1}}{n-1+\vartheta} & (1) \\ \frac{\vartheta}{n-1+\vartheta} & (2) \end{cases} \quad (4.37)$$

これはポリアの壺様モデルの推移確率(4.34)と同じである。初期条件は Ω_1^*, Π_1 いずれも 1 個の状態から始まるので、2つのプロセスは確率的に同値である。

系図過程 \mathfrak{R}_t のジャンプが起きる時刻を $T_i = \inf\{t: \mathfrak{R}_t = i\}, T_n = 0$ とする。 $T_{i-1} - T_i$ は状態 i への滞在時間で指数分布 $P(T_{i-1} - T_i > t) = \exp\left(-\frac{i(i-1+\vartheta)}{2}t\right)$ に従う。

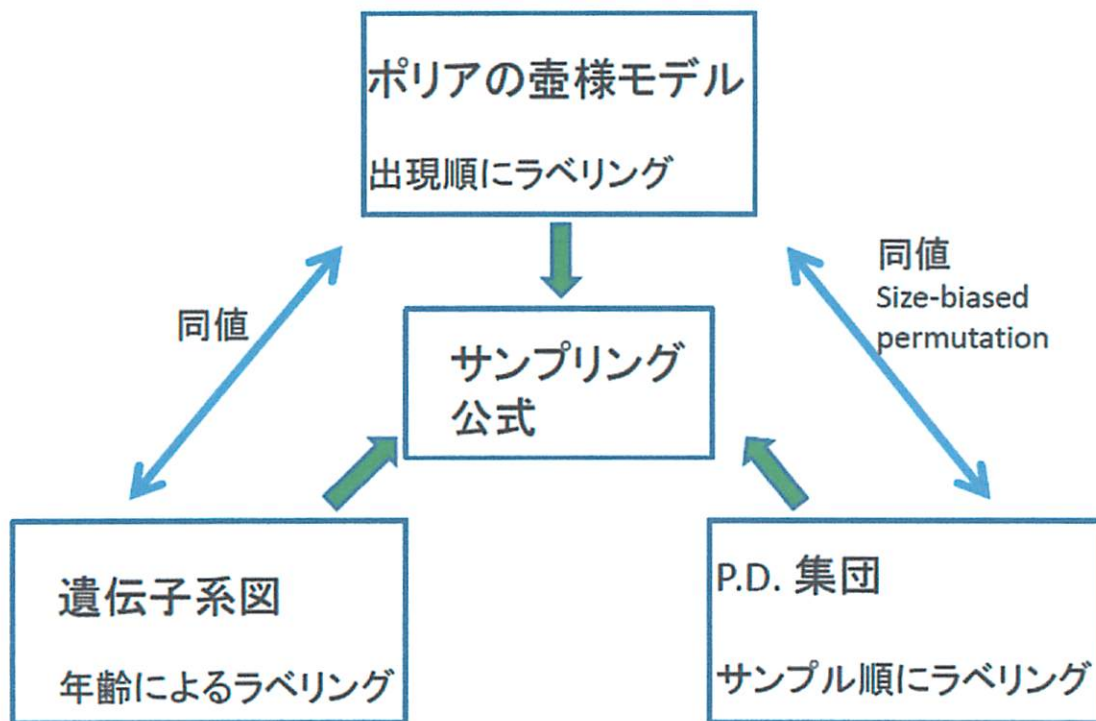
この時刻 $T_n = 0 < T_{n-1} < T_{n-2} < \dots < T_2 < T_1 < T_0$ に突然変異または合祖が起こっている。

従って、逆進過程 Ω_k^* ($k = 0, 1, 2, \dots, n$) は突然変異によって新しい遺伝子タイプが生じ、合祖においては、時間を逆さに見ると、分枝によって同じタイプが追加されるプロセスである。逆進プロセスにおいて古い突然変異から順にラベルをつけ、合祖では同じラベルを与

え、ラベルを生じた順に Y_1, Y_2, \dots, Y_n とすると、定理4. 17よりポリアの壺のラベルのプロセス X_1, X_2, \dots, X_n と同値であることが分かる。以上のことから、ポリアの壺様モデルにおいて現れるボールのラベリングは突然変異を含む合祖過程を時間について古い方から現れる順にラベリングしたものに一致する。すなわち、サンプル中に現れる順に遺伝子にラベリングすることは系図上で古い方から順にラベリングすることに確率論的に同値である。

系4. 18

年齢によるラベリングされたサンプル中のアレル数の分布はサンプル中に現れた順によるラベリングされたアレル数の分布と確率的に等しい。



以上をまとめると、上図のようになる。すなわち、ポリアの壺様モデルのサンプルのラベリング、ポアソン・ディリクレ(P.D.)集団からのサンプルの出現順のラベリングおよび遺伝子系図で古いアレルから年齢によるラベリング、この3者が確率論的に同値であるということである。これらの事実より、集団遺伝学における種々の結果について、この3つの視点から統一的な解釈を与えることができる。次の節で紹介しよう。

4. 6. 5 アレルの年齢

ポリアの壺様モデルからアレルの年齢に関する種々の公式を導くことができる。 $\{X_1, X_2, X_3, \dots\}$ をポアソン・ディリクレ集団 P の観測系列とし $\{X_1, X_2, \dots, X_n\}$ を観測サンプル、 $\{X_{n+1}, X_{n+2}, \dots\}$ で集団を現わす。

ε を n 個サンプルからランダムに取り出した一つのアレルすなわち $P(\varepsilon = X_j) = \frac{1}{n}$ 、

$S \equiv \#\{1 \leq j \leq n; X_j = \varepsilon\}$ を ε と同じタイプのアレルの数とする。さらに、サンプル中で

i 番目のサンプルと同じタイプのアレルの数を $S_i = \#\{1 \leq j \leq n; X_j = X_i\}$ で表す。

1. (Watterson and Guess(1977)) :

(a) 集団内で頻度が Y のアレルが集団内で最も古いアレルである確率は Y である。

(b) もし或るアレルの頻度が Y であるならば、それが標本中で最初にサンプルされる確率は Y である。

(証明) 頻度 Y のアレルが最初にサンプルされる確率は Y であり、それは同時に最も古いアレルである確率でもある。

2. (a) Watterson and Guess(1977), Kelly(1977) : サイズ n のサンプル中に或るアレルが

i 個含まれているとき、そのアレルが集団中で最も古いアレルである確率は $\frac{i}{g+n}$ 。

$$(b) P(X_{n+1} = \varepsilon | S = i) = \frac{i}{g+n}$$

(証明) $P(X_{n+1} = X_j | S_j = i) = \frac{i}{g+n}$ より、

$$\begin{aligned} P(X_{n+1} = \varepsilon | S = i) &= \frac{P(X_{n+1} = \varepsilon, S = i)}{P(S = i)} = \frac{1}{P(S = i)} \sum_{j=1}^n P(\varepsilon = X_j, X_{n+1} = X_j, S_j = i) \\ &= \sum_{j=1}^n P(\varepsilon = X_j) \frac{P(X_{n+1} = X_j, S_j = i)}{P(S = i)} = \sum_{j=1}^n \frac{1}{n} P(X_{n+1} = X_j | S_j = i) \frac{P(S_j = i)}{P(S = i)} = \frac{i}{g+n} \end{aligned}$$

3. (a) Watterson and Guess(1977), Kelly(1977) : サイズ n のサンプル中に或るアレルが i 個含まれているとき、それがサンプル中で最も古いアレルである確率は $\frac{i}{n}$ である。

$$(b) P(X_1 = \varepsilon | S = i) = \frac{i}{n}$$

(証明) (b) n 個のサンプル中にあるアレルが i 個含まれるとき、それが最初のサンプルである確率、すなわちサンプル中で最も古い確率は $\frac{i}{n}$ である。

4. (a) Kelly(1977) : サイズ n のサンプル中に、最も古いアレルが i 個含まれる確率は

$$\frac{\mathcal{G}\binom{n}{i}}{\binom{\mathcal{G}+n-1}{i}}.$$

$$(b) P(S_1 = i) = \frac{\mathcal{G}\binom{n}{i}}{\binom{\mathcal{G}+n-1}{i}}.$$

(証明) ポリアの壺様モデルより、最初にサンプルするアレルがラベル1となる。そのアレルがその後の $n-1$ 個のサンプルに何個含まれるかは二色の球を含むポリアの壺モデルで表現される。ラベル1のアレル頻度はベータ分布 $Beta(1, \mathcal{G})$ に従い、ラベル1の頻度が x のとき、サンプル中に含まれるラベル1の個数は二項分布 $B(n-1, x)$ に従う。よって、 $P(S_1 = i)$ は次の式で表される。

$$\begin{aligned} P(S_1 = i) &= \int_0^1 \binom{n-1}{i-1} x^{i-1} (1-x)^{n-i} \times \mathcal{G}(1-x)^{\mathcal{G}-1} dx = \mathcal{G} \binom{n-1}{i-1} \frac{\Gamma(i)\Gamma(n+\mathcal{G}-i)}{\Gamma(n+\mathcal{G})} \\ &= \frac{\mathcal{G}\binom{n}{i}}{\binom{\mathcal{G}+n-1}{i}} \end{aligned}$$

5. (a) Saunders, Tavare and Watterson(1984) : サイズ n のサンプル中で最も古いアレルより古い集団中のアレルの数を N_n とすると

$$P(N_n = k) = \frac{n}{\mathcal{G}+n} \left(\frac{\mathcal{G}}{\mathcal{G}+n} \right)^k, \quad k = 0, 1, 2, \dots$$

(b) サンプル $\{X_1, X_2, \dots, X_n\}$ の中の最も古いアレルよりも古いアレルが集団 $\{X_{n+1}, X_{n+2}, \dots\}$ において見出されるアレルタイプの数 ν_n とすると

$$P(\nu_n = k) = \frac{n}{n+\mathcal{G}} \left(\frac{\mathcal{G}}{\mathcal{G}+n} \right)^k, \quad k = 0, 1, 2, \dots$$

(証明) 事象 $\{\nu_n = k\}$ は集団で最も古いアレル、二番目に古いアレル、...、 k 番目に古いアレルがサンプル $\{X_1, X_2, \dots, X_n\}$ の中に現れず、 $k+1$ 番目に古いアレルはサンプル中に現れる事象である。また最も古いアレルが n 個のサンプル中に現れない確率はポアソン・ディリクレ分布の size-biased permutation より

$$E[(1-Z_1)^n] = \int_0^1 (1-z)^n \mathcal{G}(1-z)^{\mathcal{G}-1} dz = \mathcal{G} \int_0^1 (1-z)^{n+\mathcal{G}-1} dz = \frac{\mathcal{G}}{n+\mathcal{G}}$$

以下同様にして、 Z_1, Z_2, Z_3, \dots はベータ分布 $Beta(1, \theta)$ に従う独立な確率変数なので

$$\begin{aligned} P(\nu_n = k) &= E[(1-Z_1)^n (1-Z_2)^n \dots (1-Z_k)^n \{1 - (1-Z_{k+1})^n\}] \\ &= \prod_{i=1}^k E[(1-Z_i)^n] \times E[\{1 - (1-Z_{k+1})^n\}] = \left(\frac{\mathcal{G}}{n+\mathcal{G}} \right)^k \frac{n}{n+\mathcal{G}} \end{aligned}$$

6. Watterson(1974) : サイズ n のサンプル中に i 個含まれるアレルの種類数を a_i とする

$$\text{と } E[a_i] = \frac{\vartheta \binom{n}{i}}{\binom{\vartheta+n-1}{i}}$$

(証明) $P(S_1 = i | a_1, \dots, a_n) = \frac{ia_i}{n}$ 、より、 (a_1, \dots, a_n) の全てについて期待値を取ると 4.

$$\text{の計算より } E[a_i] = \frac{n}{i} P(S_1 = i) = \frac{n}{i} \times \frac{\vartheta \binom{n}{i}}{\binom{\vartheta+n-1}{i}} = \frac{\vartheta \binom{n}{i}}{\binom{\vartheta+n-1}{i}}.$$

これは、サンプル中で最も古いアレルがサンプル中に i 個含まれる確率 $P(S_1 = i)$ が

$P(S_1 = i) = \frac{i}{n} E[a_i]$ であることを意味する。

7. Ewens(1972) : 集団からランダムに取り出した二つの遺伝子が同じアレルタイプである

$$\text{確率は } E[\sum P_i^2] = \frac{1}{1+\vartheta}.$$

(証明) ポリアの壺様モデルで 2 番目の球が 1 番目の球である確率 $P(X_1 = X_2) = \frac{1}{1+\vartheta}$ 。

8. Watterson and Guess(1977) : 最大の頻度を $P_{(1)} = \text{Max}_i P_i$ とする、このとき

$0.5 \leq x \leq 1$ に対して、 $P_{(1)}$ の分布密度は $\vartheta x^{-1}(1-x)^{\vartheta-1}$ である。

(証明) λ_j を定理 4. 15 で導入した確率変数とする、

$$\begin{aligned} P(P_{(1)} > x) &= \sum_{j=1}^{\infty} P(P_{(1)} = P_j, P_j > x) = \sum_{j=1}^{\infty} P(\gamma_j > x \text{ and } P_{(1)} = \gamma_j) \\ &= \sum_{j=1}^{\infty} P(\gamma_j > x) = \sum_{j=1}^{\infty} P(\gamma_j > x | \lambda_j = 1) P(\lambda_j = 1) \end{aligned}$$

定理 4. 16 (3) より $\lambda_j = 1$ のとき、 γ_j は $\text{Beta}(1, \vartheta + j - 1)$ に従うので、

$$P(\lambda_j = 1) = \frac{\vartheta}{\vartheta + j - 1} \text{ より}$$

$$\begin{aligned} P(P_{(1)} > x) &= \sum_{j=1}^{\infty} \int_x^1 \frac{\Gamma(\vartheta + j)}{\Gamma(\vartheta + j - 1)} (1-t)^{\vartheta+j-2} dt \times \frac{\vartheta}{\vartheta + j - 1} \\ &= \vartheta \int_x^1 \sum_{j=1}^{\infty} (1-t)^{\vartheta+j-2} dt = \vartheta \int_x^1 t^{-1} (1-t)^{\vartheta-1} dt \end{aligned}$$

9. Ewens(1972) : サイズ n のサンプル中での異なるアレルタイプ数を K_n とする。

$$P(K_n = k) = \frac{\mathcal{G}^k S(n, k)}{\mathcal{G}_{(n)}}, \text{ただし } S(n, k) \text{ は第 1 種スターリング数 (付録(C.1)参照)}.$$

(証明) $K_n = \lambda_1 + \lambda_2 + \dots + \lambda_n$ なので、 K_n の母関数は $\lambda_1, \lambda_2, \dots, \lambda_n$ が独立より

$$E[s^{K_n}] = E[s^{\lambda_1 + \lambda_2 + \dots + \lambda_n}] = \prod_{i=1}^n E[s^{\lambda_i}] = \prod_{i=1}^n \left\{ 1 \times \frac{i-1}{\mathcal{G} + i - 1} + s \times \frac{\mathcal{G}}{\mathcal{G} + i - 1} \right\} = \prod_{i=1}^n \left(\frac{i-1 + \mathcal{G}s}{i + \mathcal{G} - 1} \right)$$

s^k の係数に注意すると第 1 種スターリング数及び $P(K_n = k)$ が得られる。

10. Ewens(1972) : ポアソン・ディリクレ集団 \underline{P} の頻度スペクトルを $\phi(x)$ とすると

$$\phi(x) = \mathcal{G}x^{-1}(1-x)^{\mathcal{G}-1} \quad (1 < x < 1).$$

(証明) $f(x)$ を $[0, 1]$ 上の有界可測関数とする。 $Q_1 = P_1^S$ をポアソン・ディリクレ \underline{P} の

size-biased Permutation の最初の頻度とすると、 $\sum_i f(P_i)P_i = E[f(Q_1)|\underline{P}]$ 。 Q_1 は

$Beta(1, \theta)$ に従うので、上式の期待値を取ると

$$E\left[\sum_i f(P_i)P_i\right] = E[f(Q_1)] = \int_0^1 f(x)\mathcal{G}(1-x)^{\mathcal{G}-1} dx.$$

ここで $f(x) = \begin{cases} 1/x & (x > t \text{ のとき}) \\ 0 & (x \leq t \text{ のとき}) \end{cases}$ とすると

$$E\left[\sum_i f(P_i)P_i\right] = E\left[\sum_i 1(P_i > t)\right] = \int_t^1 \mathcal{G}x^{-1}(1-x)^{\mathcal{G}-1} dx,$$

ただし $1(P_i > t) = \begin{cases} 1 & (P_i > t \text{ のとき}) \\ 0 & (P_i \leq t \text{ のとき}) \end{cases}$ とする。これより頻度スペクトル $\phi(x)$ は

$\phi(x) = \mathcal{G}x^{-1}(1-x)^{\mathcal{G}-1} \quad (1 < x < 1)$ 。 また $g(x) = xf(x)$ とすると

$$E\left[\sum_i g(P_i)\right] = E\left[\sum_i \frac{g(P_i)}{P_i} P_i\right] = E\left[\frac{g(Q_1)}{Q_1}\right] = \int_0^1 \frac{g(x)}{x} \mathcal{G}(1-x)^{\mathcal{G}-1} dx = \int_0^1 g(x)\phi(x) dx.$$