

α -エントロピーを用いた多様性推定について

種市信裕

(教養課程統計学教室)

1989年10月31日受理

On the estimation of diversity using α -entropy

Nobuhiro TANEICHI

摘要

従来多様性の指標として用いられてきた Shannon-エントロピー関数や Gini-index を、それぞれ $\alpha = 1$, $\alpha = 2$ の場合として含む Havrda-Charvat の α -エントロピー関数を多様性の指標として用いる場合における、多様性の指標を推定するうえでのバイアス修正、及び推定量の分散に基づく指標選択の基準について考察する。

1. 多様性の指標

m 個の種とそれに属する個体数 N_1, \dots, N_m が与えられた時、種の数 m (richness) と、それぞれの種に属する個体数の均等性 (evenness) を同時に測る尺度として、多様性の指標は、生物学、特に生態学、遺伝学等でよく用いられている。多様性の指標として有名な尺度としては、Shannon-エントロピー関数 $H_{k,1}$ 及び、Gini-index $H_{k,2}$ がある。Shannon-エントロピー関数 $H_{k,1}$ は、

$$S_k = \{ \mathbf{p} = (p_1, \dots, p_k) \mid \sum_{i=1}^k p_i = 1, p_i \geq 0 \quad (i=1, \dots, k) \}$$

上で定義される関数として、

$$H_{k,1}(p) = - \sum_{i=1}^k p_i \log_e p_i$$

と表される。Gini-index $H_{k,2}$ は、 S_k 上で定義される関数として、

$$H_{k,2}(p) = 1 - \sum_{i=1}^k p_i^2$$

と表される。 $H_{1,k}, H_{2,k}$ には、次のような性質がある。

1. $H_k(p) = 0 \Leftrightarrow \exists i \quad p_i = 1$
2. $\forall p \in S_k$ に対して $H_k(p) \leq H_k(1/k, 1/k, \dots, 1/k)$
3. $H_k(1/k, \dots, 1/k) \leq H_{k+1}(1/(k+1), \dots, 1/(k+1))$

1, 2 の性質により、これらの尺度は均等性 (evenness) を、3 の性質により、種の数 (richness) を示す尺度として用いられる。ところが 1, 2, 3 の性質は、 $\alpha = 1$ の場合として Shannon-エントロピー、 $\alpha = 2$ の場合として、Gini-index となる α -エントロピー関数

$$H_{k,\alpha}(p) = \begin{cases} (\alpha - 1)^{-1} (1 - \sum_{i=1}^k p_i^\alpha) & (\alpha \neq 1, \alpha > 0) \\ - \sum_{i=1}^k p_i \log_e p_i & (\alpha = 1) \end{cases}$$

(Havrda and Charvat, 1967)

全体が満たす性質である。このことからオーダー α のエントロピー関数全体は、多様性の指標としての条件を満たしていることがわかる。Rao, C.R. (1982), Nayak, T.K. (1984) らによって、多様性の指標としてオーダー α のエントロピーの関数の研究が行なわれているが、どの α を用いることが好ましいかについては考察されていない。

m 個の母集団があり、それぞれの母集団は k 個のカテゴリーがあるとする。 m 個の母集団からのそれぞれの標本数を N_1, \dots, N_m として、 N_i ($i = 1, \dots, m$) 個の標本の中で j 番目 ($j = 1, \dots, k$) のカテゴリーに反応している個体を n_{ij} とする。 $X_i = (n_{i1}, \dots, n_{ik})$, ($i = 1, \dots, m$) とおき、 $X_i \sim \text{Mult}(N_i, p_i)$ ただし、 $p_i = (p_{i1}, \dots, p_{ik})$ 、すると、それぞれの母

集団における多様性の指標の M.L.E. は、

$$\widehat{H_k, \alpha(p_i)} = (\alpha - 1)^{-1} (1 - \sum_{j=1}^k \widehat{p}_{ij}^{\alpha-1}), \quad (i = 1, \dots, m)$$

ただし、 $\widehat{p}_{ij} = (n_{ij} / N_i)$ となる。実際の問題ではこの推定値を用いて、m 個の母集団の多様性の指標の値を比較している。

2. 推定量のバイアス修正

m 個の母集団における推定量の期待値 $E(\widehat{H_k, \alpha(p_i)})$ は次のように展開される。

$$\begin{aligned} E(\widehat{H_k, \alpha(p_i)}) &= \\ H_k, \alpha(p_i) &+ (1/N_i) A_{1i}(p_i) + (1/N_i^2) A_{2i}(p_i) \\ &+ (1/N_i^3) A_{3i}(p_i) + O(1/N_i^4), \quad (i = 1, \dots, m) \end{aligned}$$

ただし、

$$A_{1i}(p_i) = (\alpha/2) \left\{ \sum_{j=1}^k p_{ij}^{\alpha-1} (p_{ij}-1) \right\}$$

$$\begin{aligned} A_{2i}(p_i) &= \{ \alpha(\alpha-2)/24 \} \left\{ \sum_{j=1}^k p_{ij}^{\alpha-2} (p_{ij}-1) \times \right. \\ &\quad \left. [-5 + 3\alpha + p_{ij}(1-3\alpha)] \right\} \end{aligned}$$

$$\begin{aligned} A_{3i}(p_i) &= \{ \alpha(\alpha-2)(\alpha-3)/48 \} \left[\sum_{j=1}^k p_{ij}^{\alpha-3} (p_{ij}-1) \times \right. \\ &\quad \left. [(\alpha-2)(\alpha-3) + 4(\alpha-1)p_{ij} + \alpha(\alpha-1)p_{ij}^2] \right] \end{aligned}$$

$p_i = (p_{i1}, \dots, p_{ik})$ は、各母集団においてそれぞれカテゴリーに反応する真の確率。

ここで、 $H_k, \alpha(p_i)$ の推定量を

$$\begin{aligned} \theta_{ia} &= H_k, \alpha(p_i) - (1/N_i) A_{1i}(p_i) - (1/N_i^2) A_{2i}(p_i) \\ &- (1/N_i^3) A_{3i}(p_i) \end{aligned}$$

とおくと、

$$E(\theta_{i\alpha}) = H_k, \alpha(p_i) - O(1/N_i^4)$$

$\theta_{i\alpha}$ の近似として

$$\begin{aligned}\hat{\theta}_{i\alpha} &= H_k, \alpha(\hat{p}_i) - (1/N_i) A_{1i}(\hat{p}_i) - (1/N_i^2) A_{2i}(\hat{p}_i) \\ &\quad - (1/N_i^3) A_{3i}(\hat{p}_i)\end{aligned}$$

とおき、これを $H_k, \alpha(p_i)$ の推定量とする。（この推定量を用いると少なくとも $1/N_i$ のオーダーまでバイアスは修正されている。）

3. 多様性の指數としてどの α を用いれば良いか決める基準

本報告では α を決める基準として、多様性の指數を推定するときの推定量の分散を用いることを考える。 $H_k, \alpha(p_i)$ の推定量 $\hat{\theta}_{i\alpha}$ の分散の近似として

$$\text{Var}(\theta_{i\alpha}) \approx \text{Var}(\hat{\theta}_{i\alpha}) = E\{\theta_{i\alpha} - E(\theta_{i\alpha})\}^2$$

とおくと、

$$\begin{aligned}N_i E\{\theta_{i\alpha} - E(\theta_{i\alpha})\}^2 \\ = B_{1i}(p_i) + (1/N_i) B_{2i}(p_i) + (1/N_i^2) B_{3i}(p_i) \\ + O(1/N_i^3) \quad (i = 1, \dots, m)\end{aligned}$$

ただし、

$$B_{1i}(p_i) = \{\alpha/(\alpha-1)\}^2 \left\{ \sum_j \sum_{j=1}^k p_{ij}^{2\alpha-1} - (\sum_j \sum_{j=1}^k p_{ij}^\alpha)^2 \right\}$$

$$\begin{aligned}B_{2i}(p_i) &= \{\alpha^2/2(\alpha-1)\} \left\{ 3(\alpha-1) \sum_j \sum_{j=1}^k p_{ij}^{2\alpha-2} + 2(1-2\alpha) \times \right. \\ &\quad \left. \sum_j \sum_{j=1}^k p_{ij}^{2\alpha-1} - 2(\alpha-1) \sum_i \sum_{i=1}^k \sum_j \sum_{j=1}^k p_{ij}^\alpha p_{il} \alpha^{-2} + \right. \\ &\quad \left. (3\alpha-1) (\sum_j \sum_{j=1}^k p_{ij}^\alpha)^2 \right\}\end{aligned}$$

α -エントロピーを用いた多様性推定について

$$\begin{aligned}
 B_{3i}(\mathbf{p}_i) = & \{\alpha^3/24(\alpha-1)\} \{2(14\alpha^3-60\alpha^2+87\alpha-43) \sum_{j=1}^k p_{ij}^2 \alpha^{-3} \\
 & + \{-36(\alpha-1)^2(2\alpha-3)\} \sum_{j=1}^k p_{ij}^2 \alpha^{-2} \\
 & + 4(12\alpha^3-32\alpha^2+25\alpha-6) \sum_{j=1}^k p_{ij}^2 \alpha^{-1} \\
 & + \{-2(\alpha-2)^2(3\alpha-5)\} \sum_{j=1}^k \sum_{l=1}^k p_{ij}\alpha p_{il}\alpha^{-1} \\
 & + 12(\alpha-1)^2(3\alpha-4) \sum_{j=1}^k \sum_{l=1}^k p_{ij}\alpha p_{il}\alpha^{-1} \\
 & + 4(-7\alpha^3+16\alpha^2-9\alpha+1) (\sum_{j=1}^k p_{ij}\alpha)^2 \\
 & + \{-6(\alpha-1)^3\} (\sum_{j=1}^k p_{ij}\alpha^{-1})^2\},
 \end{aligned}$$

$\mathbf{p}_i = (p_{i1}, \dots, p_{ik})$ は、各母集団においてそれぞれのカテゴリーに反応する真の確率。

各母集団からの標本数が比較的大きいときには B_{1i} 、比較的小さいときには B_{2i} 、
 B_{3i} まで考慮にいれることにする。ここでエントロピーの推定量の大きさで
 $\text{Var}(\theta_{i\alpha})$ を標準化したものを、

$$C_{i\alpha}(\mathbf{p}_i) = \text{Var}(\theta_{i\alpha}) / (\hat{\theta}_{i\alpha})^2, \quad (i = 1, \dots, m)$$

とおき、

$$C_{i\alpha}(\hat{\mathbf{p}}_i) = \text{Var}(\theta_{i\alpha}) / (\hat{\theta}_{i\alpha})^2 \mid \mathbf{p}_i = \hat{\mathbf{p}}_i$$

として、

$$\min_{\alpha} \{ \max \{ C_{1\alpha}(\hat{\mathbf{p}}_1), \dots, C_{m\alpha}(\hat{\mathbf{p}}_m) \} \} \quad (3.1)$$

を α の値を決める基準とする。つまり、各母集団 ($1 \sim m$) における多様性の指標の推定量の分散の近似を推定量の大きさで標準化し、その最大値が最小となるように α の値を定めることとする。

4. 数値例

表1で与えられる、生態学における生活形スペクトル（緯度系列）を例にとり多様性の指標 $H_{k,1}(p)$ (Shannon-エントロピー), $H_{k,2}(p)$ (Gini-index), $H_{k,3}(p)$ を比較する。

表1. 生活形スペクトル（緯度系列）

	地上植物	地表植物	半地表植物	地中植物	一年生植物
亜熱帯林	0.65	0.17	0.02	0.05	0.10
温暖帯林	0.54	0.09	0.24	0.09	0.04
冷温帯林	0.10	0.17	0.54	0.12	0.07
ツンドラ	0.01	0.22	0.60	0.15	0.02

いま、標本数は $N_1 = N_2 = N_3 = N_4 = 100$ とした場合を考える、それぞれの地域の推定確率ベクトルは、

$$p_1 = (0.65, 0.17, 0.20, 0.05, 0.1),$$

$$p_2 = (0.54, 0.09, 0.24, 0.09, 0.04),$$

$$p_3 = (0.1, 0.17, 0.54, 0.12, 0.07),$$

$$p_4 = (0.01, 0.02, 0.6, 0.15, 0.02),$$

であり、一般に用いられている Shannon-エントロピー関数 $H_{5,1}$ を多様性の指標とすると、バイアス修正後の指定値は、 $\hat{\theta}_{11} = 1.053$, $\hat{\theta}_{21} = 1.2379$, $\hat{\theta}_{31} = 1.3052$, $\hat{\theta}_{41} = 1.0509$ となり、多様性の指標の大小関係は $\hat{\theta}_{41} < \hat{\theta}_{11} < \hat{\theta}_{21} < \hat{\theta}_{31}$ となる、つまり多様性の大小関係は、ツンドラ < 亜熱帯林 < 温暖帯林 < 冷温帯林、となる。一方、標準化した推定分散を用いた α を決定する基準 (3.1) で $\alpha = 1, 2, 3$ の場合を比較すると、分散として漸近分散 $B_{1,i}$ のみを考えた場合、 $1/N_i^2$ のオーダーまで考えた場合、 $1/N_i^3$ のオーダーまで考えた場合、いずれも $\alpha = 3$ の場合が好ましいという結果である。そこで、この $H_{k,3}$ のバイアス修正後の推定値を求めると、 $\hat{\theta}_{13} = 0.3595$, $\hat{\theta}_{23} = 0.4136$

$\hat{\theta}_{33} = 0.4172$, $\hat{\theta}_{43} = 0.3849$ となり, 多様性の指標の大小関係は $\hat{\theta}_{13} < \hat{\theta}_{43} < \hat{\theta}_{23} < \hat{\theta}_{33}$ であることにより, 多様性の大小関係は, 亜熱帯林 < ツンドラ < 溫暖帯林 < 冷温帯林となり, 推定分散の基準の意味では, 従来の Shannon-エントロピーによる多様性の関係よりもこちらの関係の方が適切と思われる。

[参考文献]

- 1) Havrda, J. and Charvat, F. (1967). Quantification method in classification processes; concept of α -entropy, Kybernetika, 33, 30-35.
- 2) Nayak, T. K. (1984). On diversity measure based on entropy functions, Tech. Rept. No 84-16. Univ. of Pittsburgh.
- 3) Pielou, E. C. (1969). An introduction to mathematical ecology, John Wiley and Sons.
- 4) Rao, C. R. (1982). Diversity; its measurement, decomposition, apportionment and analysis, Sankya A, 44, 1-22.