

種類數估計量之一致性 — 紀念陳文成教授去世12年

林惠賢
清華大學
統計所

趙蓮菊
清華大學
統計所

李桑銘
逢甲大學
統計系

摘 要

本文利用陳文成教授著作中的一個定理之推廣以證明 Chao et al. (1993)對估計種類數之鞅估計函數 (martingale estimating function)所得之估計量的一致性。

關鍵詞：多項式模式，差異性，樣本涵蓋率，零均值鞅 (ZMM)。
美國數學會分類索引：主要 62P10，次要 62F10。

1. 前言

1.1 問題描述

母體種類數的估計在生物學的應用上是一個古典的問題。Bunge and Fitzpatrick (1993)回顧關於這個問題的不同模式和估計方法，他們也收集了與這個題目相關的參考論文超過 550篇。

本文主要討論最常用的多項式模式(multinomial model)：考慮一個母體，其組成的個體可分類為 N 個未知的不同種類。我們每次選取一個個體，記下它的種類別，再把它放回母體。 n 次個體選取完成後，稱此選取為 n 階段選取(n -stage selection)。假設這些種類以任意次序標上 $1, 2, \dots, N$ 。令 p_i 為被選取者屬於第 i 種類的機率，且令 X_{ik} 為 k 階段選取屬於第 i 種類的個體數，則對所有的 $k = 1, 2, \dots, n$ ， $(X_{1k}, X_{2k}, \dots, X_{Nk})$ 服從多項式分配。我們的目的是在 n 次選取完成後估計種類數 N 。

1.2 文獻回顧及研究動機

Bunge et al. (1992)在其論文中提及過去處理本問題在頻率論(frequentist)觀點下有三個主要的估計法：

- (1) 假定多項式分配中的種類機率有某一函數形式(functional form) (例如：McNeil, 1973)，亦即存在某一函數 f 使 $p_i = f(i), i = 1, \dots, N$ 。
- (2) 以一有母數(parametric)機率密度函數(p.d.f)描述或近似種類機率的分配(例如：Sichel, 1986)。
- (3) 無母數樣本涵蓋率(sample coverage)方法：藉由樣本涵蓋率的估計來估計種類數(例如：Chao and Lee, 1992)。在樣本涵蓋率法中種類機率的變異係數(coefficient of variation)於測度種類機率之差異性(heterogeneity)時扮演重要的角色。

最近Chao et al. (1993)提出另一個無母數方法—使用鞅估計函數(martingale estimating function)。這個方法提供了另一類估計量—鞅估計函數計量(或簡稱鞅估計量)。此類估計量包括種類機率相等時的最大概似估計量(maximum likelihood estimator)(Darroch, 1958)和種類機率不等時的無母數樣本涵蓋率估計量(Chao and Lee, 1992)。Chao and Lee (1992)並列有樣本涵蓋率應用在種類數和母體個體總數估計的歷史可供參考；亦可參看Becker (1984)，Becker and Heyde (1990)，Yip (1989,1991)和Yip et al. (1993)關於重複捕取(capture-recapture)模式下鞅估計函數的應用。然而，前述的鞅估計函數法只處理種類機率均等的情況。而Chao, Yip and Lin (1993)可同時處理種類機率不等的情況。為了整篇文章之完整性，我們將在第二節簡略複習樣本涵蓋率估計量及鞅估計函數估計量

。而本文主要將在第三節證明在某些條件下，鞅估計函數估計量具一致性 (consistency)。

2. 樣本涵蓋率估計量及鞅估計函數估計量

2.1 樣本涵蓋率估計量及其變異數之估計 (Chao and Lee, 1992)

令 f_{ik} 為在 k 階段選取恰好出現 i 次的種類數，亦即 $f_{ik} = \sum_{j=1}^N I[X_{jk} = i]$ ，其中 $I(\cdot)$ 是常用的指示函數 (indicator function)。令 D_k 為 k 階段選取的不同種類數， $D_k = \sum_{i=1}^k f_{ik} = \sum_{j=1}^N I[X_{jk} > 0]$ 。定義 k 階段選取的樣本涵蓋率 C_k 為

$$C_k = \sum_{i=1}^N p_i I[\text{第 } i \text{ 類在 } k \text{ 階段選取曾出現過}]。 \quad (2.1)$$

則估計種類數之樣本涵蓋率估計量 (Chao and Lee, 1992) 為

$$\hat{N}_s = \frac{D_n}{\hat{C}_n} + \frac{f_{1n}}{\hat{C}_n} \hat{\gamma}_n^2, \quad (2.2)$$

其中為 $\hat{C}_n = 1 - f_{1n}/n$ 為 C_n 之估計量 (Good, 1953; 或 Robbins, 1968)， $\hat{\gamma}_n$ 為種類機率的變異係數之估計量 (定義及估計量見下節及 (2.20) 式)。而此估計量近似變異數的估計值為

$$\hat{v}ar(\hat{N}_s) = \sum_{i=1}^n \sum_{j=1}^n \frac{\partial \hat{N}_s}{\partial f_{in}} \frac{\partial \hat{N}_s}{\partial f_{jn}} cov(f_{in}, f_{jn}),$$

其中

$$cov(f_{in}, f_{jn}) = \begin{cases} f_{in}(1 - f_{in}/\hat{N}_s) & \text{if } i = j \\ -f_{in}f_{jn}/\hat{N}_s & \text{if } i \neq j \end{cases}。$$

2.2 鞅估計函數估計量及其變異數之估計 (Chao et al. 1993)

令 F_k 表由 k 階段選取的選取過程所生成的 σ -域 (σ -field)，亦即 $F_k = \sigma\{X_{1i}, X_{2i}, \dots, X_{Ni}, i = 1, 2, \dots, k\}$ 。令 $m_k = I[\text{第 } k \text{ 次選取為一曾出現過的種類}]$ 及 $u_k = I[\text{第 } k \text{ 次選取為一不曾出現過的種類}]$ ，則對所有 $k, u_k + m_k = 1$ ，且

$$E(u_k | F_{k-1}) = 1 - C_{k-1}。 \quad (2.3)$$

同理，我們可得

$$E(m_k | F_{k-1}) = C_{k-1} \circ \quad (2.4)$$

我們進一步可得到變異數

$$\text{var}(u_k | F_{k-1}) = C_{k-1}(1 - C_{k-1}), \quad (2.5)$$

$$\text{var}(m_k | F_{k-1}) = C_{k-1}(1 - C_{k-1}), \quad (2.6)$$

及共變異數

$$\text{cov}(u_k, m_k | F_{k-1}) = -C_{k-1}(1 - C_{k-1}) \circ \quad (2.7)$$

根據(2.3)和(2.4)我們考慮下列鞅差：

$$\begin{aligned} D_k &= NC_{k-1}[u_k - (1 - C_{k-1})] - N(1 - C_{k-1})[m_k - C_{k-1}] \\ &= NC_{k-1}u_k - N(1 - C_{k-1})m_k \circ \end{aligned} \quad (2.8)$$

若 NC_{k-1} 和 $N(1 - C_{k-1})$ 是 F_{k-1} 可測 (F_{k-1} -measurable)，則過程 $M = \{M_k : k = 1, 2, \dots\}$ ，其中 $M_k = \sum_1^k D_i$ ，是一個 ZMM (zero mean martingale)。爲了得到一個 N 的一般化的鞅估計函數，我們可取一個有界且相對於 M 是可預測 (predictable) 的加權函數 w_{k-1} ，使估計函數成爲

$$\begin{aligned} M_n^* &= \sum_{k=1}^n w_{k-1} D_k \\ &= \sum_{k=1}^n w_{k-1} [NC_{k-1}u_k - N(1 - C_{k-1})m_k] \end{aligned} \quad (2.9)$$

$$= \sum_{k=1}^n w_{k-1} [(NC_{k-1}) - N m_k] \circ \quad (2.10)$$

在本文中，我們討論兩種可能的加權函數：對所有的 $k = 1, 2, \dots, N$, $w_{k-1} = 1$ ，和 Godambe (1985) 所建議的最佳加權函數。他證明對一已知鞅差函數 D_k 之最佳加權函數是

$$w_{k-1} = \frac{E[\partial D_k / \partial N | F_{k-1}]}{E[D_k^2 | F_{k-1}]} \circ \quad (2.11)$$

在此“最佳”意指可造一個 N 的最短信賴區間。

對所有 p_i 都相等的特別情況，亦即 $p_1 = p_2 = \dots = p_N = 1/N$ ，我們可得 $NC_{k-1} = D_{k-1}$ ($D_0 \equiv 0$) 和 $N(1 - C_{k-1}) = N - D_{k-1}$ 。所以估計函數 (2.9) 式變成

$$\sum_{k=1}^n w_{k-1} [D_{k-1} u_k - (N - D_{k-1}) m_k] \circ \quad (2.12)$$

這個估計函數類似於 Yip (1989, 1991), Yip et al. (1993), Becker (1984) 和 Becker and Heyde (1990) 對重複捕取模式所考慮的。令 (2.12) 等於其均值零，求解得 N 的估計量為

$$\sum_{k=1}^n w_{k-1} D_{k-1} / \sum_{k=1}^n w_{k-1} m_k \circ \quad (2.13)$$

假如 $w_{k-1} = 1$ ，上式簡化成

$$\hat{N}_0 = \sum_{k=1}^n D_{k-1} / \sum_{k=1}^n m_k \circ \quad (2.14)$$

這和重複捕取模式中的 Schnabel 估計量完全相同（假設將每次選取看成一個陷阱 (trapping) 樣本）（Schnabel, 1938；或 Seber, 1982）。注意 Schnabel 估計量和選取的次序有關。換言之，由於我們的估計量是根據前面的歷史而定，所以 Schnabel 估計量當選取的次序不同時會有不同的結果。

當所有 p_i 都相等時，最大概似估計量 (MLE) $\hat{N}_{0,mle}$ 是下列方程式 (Darroch, 1958) 的解

$$\sum_{i=0}^{D_n-1} (N - i)^{-1} = n/N, \quad (2.15)$$

其漸近變異數 (asymptotic variance) 為

$$\text{var}(\hat{N}_{0,mle}) = N / [\exp(n/N) - (n/N) - 1] \circ$$

注意此最大概似估計量和選取的次序無關。

由 (2.6), (2.11) 和 (2.12) 可導出此情況下的最佳加權函數為 $w_{k-1}^* = 1/(N - D_{k-1})$ 。則對應於此最佳加權數的最佳估計量為下列方程式之解

$$\sum_{k=1}^n [D_{k-1} - N m_k] / (N - D_{k-1}) = 0 \circ \quad (2.16)$$

我們可證明上方程式和 (2.15) 是完全相同的。類似 Yip (1991) 的導法， \hat{N}_0 和 $\hat{N}_{0,mle}$ 的漸近標準分別為下式中的 w_{k-1} 以 1 和 $1/(N - D_{k-1})$ 替代而得：

$$\left\{ \sum_{k=1}^n w_{k-1}^2 D_{k-1} (N - D_{k-1}) \right\}^{1/2} / \sum_{k=1}^n w_{k-1} m_k。$$

在大部分的實際應用上，種類機率均等的假設是不成立的。以前處理有差異情況的方法多數都對種類機率的分配做有參數的假設。例如：有些作者假設 p_1, p_2, \dots, p_N 是一組來自某一分配的隨機取樣 (random sample) (Fisher et al. 1943 是其中一例)。但在我們的方法中，我們把 p_1, p_2, \dots, p_N 當做固定的參數 (fixed parameters)，但在導變異數估計量時例外。注意此地雖然允許 N (N 可能很大) 個不同之參數，但並無辨認性 (identifiability) 之問題，因為在本論文中將展示種類之估計只與最重要之參數平均值 $\bar{p} = \sum_{i=1}^N p_i / N = 1/N$ 和異係數 (CV) $\gamma = [\sum_{i=1}^N (p_i - \bar{p})^2 / N]^{1/2} / \bar{p}$ 有關。只要平均值和變異係數相同，實際上 p_i 各自的大小並不影響估計結果。故在整個討論中只看成二個參數。此處基本的動機如下：我們不直接估計 N ，而藉由 (NC_{k-1}) 的估計來估計 N 。從 (2.10) 式，假如可找到 NC_{k-1} 的估計量 \widehat{NC}_{k-1} ，則我們有下列 N 的估計量

$$\sum_{k=1}^n w_{k-1} (\widehat{NC}_{k-1}) / \sum_{k=1}^n w_{k-1} m_k。 \quad (2.17)$$

當估計量 \widehat{NC}_{k-1} 是 F_{k-1} 可測時，(2.17) 之漸近標近標準差可以證明為

$$N \left\{ \sum_{k=1}^n w_{k-1}^2 [C_{k-1}(1 - C_{k-1})] \right\}^{1/2} / \sum_{k=1}^n w_{k-1} m_k。 \quad (2.18)$$

到此我們的過程就是要找出 $E(NC_{k-1})$ 的估計量，既而用它當做 (NC_{k-1}) 的估計量。我們可導出

$$E(NC_{k-1}) \approx E(D_{k-1}) + \gamma^2 \cdot E(f_{1,k-1}), \quad (2.19)$$

其中 $f_{10} \equiv 0$ 。若每得一觀測值，即作一次 CV 之修正估計，即對每一 m ，求一 $(CV)^2$ 估計如下：

$$\hat{\gamma}_m^2 = \max \left\{ \frac{D_m \sum_{i=1}^m i(i-1) f_{im}}{m(m-1)[1 - f_{1m}/m]} - 1, 0 \right\}。 \quad (2.20)$$

從(2.17), (2.19)和(2.20), 我們提出的估計量在 $w_{k-1} = 1$ 時為

$$\hat{N} = \sum_{k=k_0}^n [D_{k-1} + \hat{\gamma}_{k-1}^2 f_{1,k-1}] / \sum_{k=k_0}^n m_k \quad (2.21)$$

其中所有累加從起始值 k_0 開始是因為我們需要足夠的觀測值去得到一個CV的穩定(stable)估計量。在Chao et al. (1993)中, 我們選擇 $k_0 = n/2$ 。CV之估計實為此過程中最困難之處。另一個替代的方法是我們等所有選取都做完後再估計此參數。換言之, 我們考慮下列含修正的CV估計之估計量

$$\begin{aligned} \hat{N} &= \sum_{k=1}^n [D_{k-1} + \hat{\gamma}_n^2 f_{1,k-1}] / \sum_{k=1}^n m_k \\ &= \hat{N}_0 + \hat{\gamma}_n^2 \cdot \sum_{k=1}^n f_{1,k-1} / \sum_{k=1}^n m_k \end{aligned} \quad (2.22)$$

此時最佳加權函數可以證明為

$$w_{k-1}^* = 1/(1 - C_{k-1}) \quad (2.23)$$

所以藉由使用 C_{k-1} 之估計量 $\hat{C}_{k-1} = 1 - f_{1,k-1}/(k-1)$ 可以得到 w_{k-1}^* 之估計量 \hat{w}_{k-1}^* 。因此對應於此最佳加權函數估計量, 我們有兩個估計量

$$\hat{N}_w = \sum_{k=k_0}^n \hat{w}_{k-1}^* [D_{k-1} + \hat{\gamma}_{k-1}^2 \cdot f_{1,k-1}] / \sum_{k=k_0}^n \hat{w}_{k-1}^* m_k \quad (2.24)$$

和

$$\bar{N}_w = \sum_{k=1}^n \hat{w}_{k-1}^* [D_{k-1} + \hat{\gamma}_n^2 \cdot f_{1,k-1}] / \sum_{k=1}^n \hat{w}_{k-1}^* m_k \quad (2.25)$$

有關各種估計量之比較可參Chao et al. (1993)。

在Chao et al. (1993)中採用重抽法(bootstrap resampling method)來估計變異數。作法如下: 首先假設 p_1, p_2, \dots, p_N 是一組來自分配 $F(p)$ 的隨機取樣。然後無條件地(unconditionally), $(f_{0n}, f_{1n}, \dots, f_{nn})$ 服從含參數 N 和種類機率 $\binom{n}{k} \int p^k (1-p)^{n-k} dF(p)$, $k = 0, 1, \dots, n$ 之多項式分配。由含參數 \hat{N} 和估計的種類機率 f_{kn}/\hat{N} , $k = 1, 2, \dots, n$; $f_{0n} = \hat{N} - \sum_{k=1}^n f_{kn}$ 之多項式分配重複生成 $(f_{0n}^*, f_{1n}^*, \dots, f_{nn}^*)$ 。然而我們仍需要製造每一種類的出現歷史, 以便計算鞅估計量。對出現 i 次的每一種類, 我們從 $1, 2, \dots, n^*$ 中,

$n^* = \sum_{i=1}^n i f_{in}^*$ ，隨機選取 i 個觀測值（但被選取過的時間除外），令此種類出現在所選取的時間且不出現在其他時間。

對每一組資料，我們重複生成 B 次，因此得到 B 個重抽法估計值 $\hat{N}_i^*, i = 1, 2, \dots, B$ 。 \hat{N} 的重抽法變異數就是此 B 個 \hat{N}_i^* 的樣本變異數，即

$$\text{var}(\hat{N}) = \left[\sum_{i=1}^B (\hat{N}_i^*)^2 - \left(\sum_{i=1}^B \hat{N}_i^* \right)^2 / B \right] / (B-1)。$$

3. 估計量之一致性

3.1 樣本涵蓋率估計量之一致性

Lee and Chao (1993) 根據 Chen (1980, 1981a, 1981b) 的結果已證明樣本涵蓋率估計量的一致性。現將它轉述如下：

定理 3.1: 在多項式模式下，假如種類機率 (p_1, p_2, \dots, p_N) 服從一含參數 α 且對稱的 Dirichlet 分配，則其 $(CV^2)\gamma^2 \approx 1/\alpha$ 。當 $\lim_{N \rightarrow \infty} \frac{n}{N} = \lambda > 0$ 時，假如 CV 已知，則樣本涵蓋率估計量具一致性，即我們有

$$N^{-1} \left[\frac{D_n}{\hat{C}_n} + \frac{f_{1n}}{\hat{C}_n} \gamma^2 \right] \xrightarrow{P} 1。$$

在證明此定理時，主要用到 Chen (1980, 1981a, 1981b) 的下列定理：即在定理 3.1 之條件下，對 $i = 0, 1, \dots$

$$\frac{f_{in}}{N} \xrightarrow{P} \frac{\Gamma(i+\alpha)}{i! \Gamma(\alpha)} \left(\frac{\lambda}{\lambda+\alpha} \right)^i \left(\frac{\alpha}{\lambda+\alpha} \right)^\alpha。 \quad (3.1)$$

3.2 缺估計量之一致性

本文主要利用 (3.1) 之推廣以證明缺估計量亦具有一致性。

定理 3.2: 在多項式模式下，假如種類機率 (p_1, p_2, \dots, p_N) 服從一含參數 α 且對稱的 Dirichlet 分配，則其 $(CV^2)\gamma^2 \approx 1/\alpha$ 。當 $\lim_{N \rightarrow \infty} \frac{n}{N} = \lambda > 0$ 時，假如 CV 已知，則缺估計量具一致性。換言之，對任意起始值 k_0 ，我們可以證明

$$N^{-1} \left\{ \frac{\sum_{k=k_0}^n [D_{k-1} + \gamma^2 f_{1,k-1}]}{\sum_{k=k_0}^n m_k} \right\} \xrightarrow{P} 1。$$

假如 CV 已知，則鞅估計量具一致性。換言之，對任意起始值 k_0 ，我們可以證明

$$N^{-1} \left\{ \frac{\sum_{k=k_0}^n [D_{k-1} + \gamma^2 f_{1,k-1}]}{\sum_{k=k_0}^n m_k} \right\} \xrightarrow{P} 1。$$

在證明定理 3.2 之前，先證明幾個需要的引理。在引理 3.1, 3.3 和 3.4 中，我們令

$$f_N(x) = \left(1 - \frac{1}{N} \cdot \frac{\alpha + 1}{\alpha + \lambda x} \right)^{N(\alpha + \lambda x) - \alpha - \frac{3}{2}}。$$

引理 3.1：當 $N \rightarrow \infty$ 時，在 $[0, 1]$ 上 $f_N(x) \xrightarrow{u} e^{-(\alpha+1)}$ ，其中 \xrightarrow{u} 表均勻收斂。

引理 3.2：若 φ_N 和 φ 在 $[0, 1]$ 上可積，且當 $N \rightarrow \infty$ 時，在 $[0, 1]$ 上 $\varphi_N \xrightarrow{u} \varphi$ ，則對任意極限值為無窮大之序列 $\{a_N\}_{N=1}^\infty$ ，我們有：當 $N \rightarrow \infty$ 時，

$$\frac{1}{a_N} \sum_{k=1}^{a_N} \varphi_N\left(\frac{k}{a_N}\right) \rightarrow \int_0^1 \varphi(x) dx。$$

引理 3.3：當 $\lim_{N \rightarrow \infty} \frac{n}{N} = \lambda$ 時，

$$\begin{aligned} T(N) &= \frac{1}{n} \left| \sum_{k=1}^n \frac{k}{N} \left(\frac{\alpha}{\alpha + \frac{k}{N}} \right)^{\alpha+1} \left(1 - \frac{1}{N} \cdot \frac{\alpha+1}{\alpha + \frac{k}{N}} \right)^{N(\alpha + \frac{k}{N}) - \alpha - \frac{3}{2}} \right. \\ &\quad \left. - \sum_{k=1}^n \frac{k}{n} \lambda \left(\frac{\alpha}{\alpha + \frac{k}{n} \lambda} \right)^{\alpha+1} f_N\left(\frac{k}{n}\right) \right| \rightarrow 0。 \end{aligned}$$

證明：

令 $T_0(k) = \frac{k}{N} \left(\frac{\alpha}{\alpha + \frac{k}{N}} \right)^{\alpha+1} \left(1 - \frac{1}{N} \frac{\alpha+1}{\alpha + \frac{k}{N}} \right)^{N(\alpha + \frac{k}{N}) - \alpha - \frac{3}{2}}$ ，且令 T_i 為將 T_{i-1} 中第一個 $\frac{k}{N}$ 以 $\frac{k}{n} \lambda$ 替代而成之函數，則

$$T(N) \leq \frac{1}{n} \sum_{k=1}^n \sum_{i=0}^3 |T_i(k) - T_{i+1}(k)|。 \quad (3.2)$$

因為對 $k \geq 1$ ，我們都有 $\left(\frac{\alpha}{\alpha + \frac{k}{N}} \right)^{\alpha+1} \leq 1$ ， $\left(1 - \frac{\alpha+1}{N\alpha+k} \right)^{N\alpha+k} \leq e^{-(\alpha+1)}$ 和 $\left(1 - \frac{\alpha+1}{N\alpha+k} \right)^{-\alpha-\frac{3}{2}} \leq \left(1 - \frac{\alpha+1}{N\alpha} \right)^{-\alpha-\frac{3}{2}}$ ，所以

$$\begin{aligned}
& \frac{1}{n} \sum_{k=1}^n |T_0(k) - T_1(k)| \\
& \leq 1 \cdot e^{-(\alpha+1)} \left(1 - \frac{\alpha+1}{N\alpha}\right)^{\alpha \frac{3}{2}} \cdot \frac{1}{n} \sum_{k=1}^n \left| \frac{k}{N} - \frac{k}{n} \lambda \right| \\
& = e^{-(\alpha+1)} \left(1 - \frac{\alpha+1}{N\alpha}\right)^{\alpha \frac{3}{2}} \left(\frac{1}{n^2} \sum_{k=1}^n k \right) \cdot \left| \frac{n}{N} - \lambda \right| \\
& \rightarrow 0
\end{aligned}$$

同理，當 $1 \leq i \leq 3$ 時，我們有

$$\frac{1}{n} \sum_{k=1}^n |T_i(k) - T_{i+1}(k)| \rightarrow 0 \circ$$

故由(3.2)，可得證此引理。 ■

引理 3.4：當 $\lim_{N \rightarrow \infty} \frac{n}{N} = \lambda$ 時，

$$\begin{aligned}
& \frac{1}{n} \sum_{k=1}^n \frac{k}{N} \left(\frac{\alpha}{\alpha + \frac{k}{N}} \right)^{\alpha+1} \left(1 - \frac{\alpha+1}{N\alpha+k} \right)^{(N-1)\alpha+k \frac{3}{2}} \\
& \rightarrow \int_0^1 \lambda x \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} e^{-(\alpha+1)} dx \circ
\end{aligned}$$

證明：因為當 $x \in [0, 1]$ 時，我們有

$$\begin{aligned}
& \left| \lambda x \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} f_N(x) - \lambda x \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} e^{-(\alpha+1)} \right| \\
& \leq \lambda |f_N(x) - e^{-(\alpha+1)}|,
\end{aligned}$$

且由引理 3.1 可得當 $N \rightarrow \infty$ 時，在 $[0, 1]$ 上

$$\lambda x \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} f_N(x) \xrightarrow{u} \lambda x \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} e^{-(\alpha+1)} \circ$$

再由引理 3.2 我們可得：當 $\lim_{N \rightarrow \infty} \frac{n}{N} = \lambda$ 時，

$$\begin{aligned} & \frac{1}{n} \sum_{k=1}^n \lambda \frac{k}{n} \left(\frac{\alpha}{\alpha + \lambda \frac{k}{n}} \right)^{\alpha+1} f_N\left(\frac{k}{n}\right) \\ & \rightarrow \int_0^1 \lambda x \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} e^{-(\alpha+1)x} dx. \end{aligned}$$

最後由引理 3.3 即得證此引理。 ■

現在開始證明定理 3.2 如下：不失一般性，可假設 $k_0 = 1$ 。我們首先推廣 (3.1) 至以下形式：

$$\begin{aligned} \text{(a)} \quad & \frac{1}{n} \sum_{k=1}^n \frac{f_{1,k-1}}{N} \xrightarrow{P} \int_0^1 \lambda x \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} dx, \\ \text{(b)} \quad & \frac{1}{n} \sum_{k=1}^n \frac{D_{k-1}}{N} \xrightarrow{P} 1 - \int_0^1 \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha} dx \text{ 及} \\ \text{(c)} \quad & \frac{1}{n} \sum_{k=1}^n m_k \xrightarrow{P} 1 - \int_0^1 \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} dx. \end{aligned}$$

首先讓 k 固定。因為

$$\begin{aligned} E \left(\frac{f_{1k}}{N} \right) &= Pr(X_{1k} = 1) \\ &= \int_0^1 kp(1-p)^{k-1} p^{\alpha-1} (1-p)^{(N-1)\alpha-1} \frac{\Gamma(N\alpha)}{\Gamma(\alpha)\Gamma((N-1)\alpha)} dp \\ &= \frac{k \cdot \Gamma(N\alpha)}{\Gamma(\alpha)\Gamma((N-1)\alpha)} \int_0^1 p^{\alpha}(1-p)^{(N-1)\alpha+k-2} dp \\ &= \frac{k \cdot \Gamma(N\alpha)}{\Gamma(\alpha)\Gamma((N-1)\alpha)} \cdot \frac{\Gamma(\alpha+1)\Gamma((N-1)\alpha+k-1)}{\Gamma(N\alpha+k)} \\ &= k\alpha \cdot \frac{\Gamma(N\alpha)\Gamma[(N-1)\alpha+k-1]}{\Gamma[(N-1)\alpha]\Gamma(N\alpha+k)}, \end{aligned}$$

所以

$$E \left(\frac{f_{1k}}{N} \right) = \frac{k}{N} \left\{ N\alpha \cdot \frac{\Gamma(N\alpha)\Gamma[(N-1)\alpha+k-1]}{\Gamma[(N-1)\alpha]\Gamma(N\alpha+k)} \right\}.$$

令

$$h(N, k) = N\alpha \frac{\Gamma(N\alpha)\Gamma[(N-1)\alpha+k-1]}{\Gamma[(N-1)\alpha]\Gamma(N\alpha+k)}.$$

利用 $\log\Gamma(z) = (z-1/2)\log(z) - z + \frac{1}{2}\log(2\pi) + O(\frac{1}{z})$ (參看 Whittaker and Watson, 1927)，我們得到

$$\begin{aligned} \log h(N, k) &= (\alpha + 1) \log(N\alpha) - (\alpha + 1) \log(N\alpha + k) \\ &\quad + 1 - [(N - 1)\alpha - 1/2] \log\left(1 - \frac{1}{N}\right) \\ &\quad + [(N - 1)\alpha + k - 1 - 1/2] \log\left(1 - \frac{\alpha + 1}{N\alpha + k}\right) \\ &\quad + O\left(\frac{1}{N}\right) + O\left(\frac{1}{N\alpha + k}\right) \circ \end{aligned}$$

既而令

$$h(N, k) = \left(\frac{\alpha}{\alpha + \frac{k}{N}}\right)^{\alpha+1} \cdot R(N, k),$$

其中

$$\begin{aligned} R(N, k) &= e \cdot \left(1 - \frac{1}{N}\right)^{-[(N-1)\alpha-1/2]} \cdot \left(1 - \frac{\alpha+1}{N\alpha+k}\right)^{(N-1)\alpha+k-3/2} \\ &\quad \cdot e^{O(\frac{1}{N})} e^{O(\frac{1}{N\alpha+k})} \circ \end{aligned}$$

故

$$\begin{aligned} &\frac{1}{n} \sum_{k=1}^n E\left(\frac{f_{1k}}{N}\right) \\ &= \frac{1}{n} \sum_{k=1}^n \frac{k}{N} \left(\frac{\alpha}{\alpha + \frac{k}{N}}\right)^{\alpha+1} \cdot R(N, k) \\ &= e \cdot \left(1 - \frac{1}{N}\right)^{-[(N-1)\alpha-1/2]} \cdot e^{O(1/N)} \cdot \frac{1}{n} \\ &\quad \cdot \sum_{k=1}^n \left\{ \frac{k}{N} \left(\frac{\alpha}{\alpha + \frac{k}{N}}\right)^{\alpha+1} \left(1 - \frac{\alpha+1}{N\alpha+k}\right)^{(N-1)\alpha+k-3/2} \right\} \circ \end{aligned}$$

再由引理 3.4，我們可得

$$\begin{aligned} &\frac{1}{n} \sum_{k=1}^n E\left(\frac{f_{1,k-1}}{N}\right) \\ &= \frac{n-1}{n} \cdot \frac{1}{n-1} \sum_{k=1}^{n-1} E\left(\frac{f_{1k}}{N}\right) \\ &\rightarrow \int_0^1 \lambda x \left(\frac{\alpha}{\alpha + \lambda x}\right)^{\alpha+1} dx \circ \end{aligned}$$

又

$$\text{var} \left(\frac{1}{n} \sum_{k=1}^n \frac{f_{1k}}{N} \right) = \frac{1}{n^2} \sum_{k=1}^n \frac{\text{var}(f_{1k})}{N^2} + \frac{1}{n^2} \sum_{k \neq l} \frac{\text{cov}(f_{1k}, f_{1l})}{N^2},$$

由繁瑣的計算可得

$$\text{var}(f_{1k}) = O(N), \text{cov}(f_{1k}, f_{1l}) = O(N),$$

所以

$$\text{var} \left(\frac{1}{n} \sum_{k=1}^n \frac{f_{1,k-1}}{N} \right) = \frac{n \cdot O(N)}{n^2 N^2} + \frac{(n^2 - n) \cdot O(N)}{n^2 N^2} = O(1/N),$$

即 $\text{var} \left(\frac{1}{n} \sum_{k=1}^n \frac{f_{1,k-1}}{N} \right) \rightarrow 0$, 故得證(a)。

同理, 因

$$\begin{aligned} \frac{1}{n} \sum_{k=1}^n E \left(\frac{D_{k-1}}{N} \right) &= \frac{1}{n} \sum_{k=1}^n E \left(\frac{N - f_{0,k-1}}{N} \right) \\ &= 1 - \frac{1}{n} \sum_{k=1}^n E \left(\frac{f_{0,k-1}}{N} \right) \rightarrow 1 - \int_0^1 \left(\frac{\alpha}{\alpha + \lambda x} \right)^\alpha dx, \end{aligned}$$

且

$$\text{var} \left(\frac{1}{n} \sum_{k=1}^n \frac{D_{k-1}}{N} \right) = \text{var} \left(\frac{1}{n} \sum_{k=1}^n \frac{f_{0,k-1}}{N} \right) \rightarrow 0,$$

故得證(b)。而(c)的證明如下：

因

$$\begin{aligned} \frac{1}{n} \sum_{k=1}^n E(m_k) &= \frac{1}{n} \sum_{k=1}^n E \left(1 - \frac{f_{1k}}{k} \right) \\ &= 1 - \frac{1}{n} \sum_{k=1}^n \frac{1}{k} E(f_{1k}) \\ &\rightarrow 1 - \int_0^1 \left(\frac{\alpha}{\alpha + \lambda x} \right)^{\alpha+1} dx, \end{aligned}$$

且

$$\begin{aligned}
 \text{var} \left(\frac{1}{n} \sum_{k=1}^n m_k \right) &= \frac{1}{n^2} \text{var}(n - D_n) \\
 &= \frac{1}{n^2} \text{var}(D_n) \\
 &= \frac{1}{n^2} \text{var}(N - f_{0n}) \\
 &\rightarrow 0,
 \end{aligned}$$

故得證(c)。最後由(a), (b), (c)且因

$$\gamma^2 = \frac{1}{\alpha} \left(1 - \frac{\alpha+1}{N\alpha+1} \right),$$

故

$$\begin{aligned}
 &N^{-1} \left\{ \sum_{k=1}^n [D_{k-1} + \gamma^2 f_{1,k-1}] / \sum_{k=1}^n m_k \right\} \\
 &= \frac{\frac{1}{Nn} \sum_{k=1}^n [D_{k-1} + \gamma^2 f_{1,k-1}]}{\frac{1}{n} \sum_{k=1}^n m_k} \\
 &\xrightarrow{P} \frac{1 - \int_0^1 \left(\frac{\alpha}{\alpha+\lambda x} \right)^\alpha dx + \frac{1}{\alpha} \int_0^1 \lambda x \left(\frac{\alpha}{\alpha+\lambda x} \right)^{\alpha+1} dx}{1 - \int_0^1 \left(\frac{\alpha}{\alpha+\lambda x} \right)^{\alpha+1} dx} \\
 &= 1. \quad \blacksquare
 \end{aligned}$$

因此，很明顯地 *CV* 在證明一致性時扮演最重要的角色。假如 *CV* 已知，我們可以一致地 (consistently) 估計看不到的種類數。否則，一致性就變成一個達不到的理想而已。誠如 Bunge and Fitzpatrick (1993) 中引述 I. J. Good 之語 “總是會有很多極稀少而看不到的種類”。

誌謝詞：阮感謝台美文化基金會 kap 陳文憲教授 (Miami 大學) 提供陳文成教授生前所有 e 著作。阮 beh 用本文來紀念陳文成教授去世 12 年，並衷心願望伊生前 e 理想卡早實現。(1993 年 7 月 1 日)

參考文獻

- Becker, N. G. (1984). Estimating Population Size from Capture-Recapture Experiments in Continuous Time. *Australian Journal of Statistics*. 26, pp1-7.
- Becker, N. G. and Heyde, C. C. (1990). Estimating Population Size from Capture-Recapture Experiments. *Stochastic Processes and Their Applications*. 36, pp77-83.
- Bunge, J. and Fitzpatrick, M. (1993). Estimating the Number of Species: Recent Developments. *Journal of American Statistical Association*. 88, pp364-373.
- Bunge, J., Fitzpatrick, M. and Handley, J. (1992). Simulation Comparison of Three Estimators of the Number of Species. Cornell University Statistics Center, Technical Report #92.25.
- Chao, A. and Lee, S.-M. (1992). Estimating the Number of Classes via Sample Coverage. *Journal of American Statistical Association*. 87, pp210-217.
- Chao, A., Yip, P. and Lin, H.-S. (1993). Estimating the Number of Species via a Martingale Estimating Function. Technical Report, Institute of Statistics, Tsing-Hua University, Taiwan. Feb. 1993. (Submitted).
- Chen, W.-C (1980). On the Weak Form of Zipf's Law. *Journal of Applied Probability*. 18, pp611-622.
- Chen, W.-C (1981a). Limit Theorems for General Size Distributions. *Journal of Applied Probability*. 18, pp139-147.
- Chen, W.-C (1981b). Some Local Limit Theorems in the Symmetric Dirichlet-Multinomial Urn Models. *Annals of Institute Statistical Mathematics*. 33A, pp405-415.

- Darroch, J. N. (1958). The Multiple-Recapture Census I: Estimation of a Closed Population. *Biometrika*. 45, pp343-359.
- Fisher, R. A., Corbet, A. S., and Williams, C. B. (1943). The Relation Between the Number of Species and Species and the Number of Individuals in a Random Sample of an Animal Population. *Journal of Animal Ecology*. 12, pp42-58.
- Godambe, V. P. (1985). The Foundations of Finite-Sample Estimation in Stochastic Processes. *Biometrika*. 72, pp419-428.
- Good, I. J. (1953). The Population Frequencies of Species and the Estimation of Population Parameters. *Biometrika*. 40, pp237-264.
- Lee, S.-M. and Chao, A. (1993). Consistency of Sample Coverage Estimator of the Number of Classes Under a Dirichlet-Multinomial Model. Technical Report, Institute of Statistics, Tsing-Hua University, Taiwan, Feb. 1993.
- McNeil, D. (1973). Estimating an Author's Vocabulary. *Journal of American Statistical Association*. 68, pp92-96.
- Robbins, H. (1968). Estimating the Total Probability of the Unobserved Outcomes of an Experiment. *Annals of Mathematical Statistics*. 39, pp256-257.
- Sicel, H. (1986). Parameter Estimation for a Word Frequency Distribution Based on Occupancy Theory. *Communications in Statistics. Part A-Theory and Methods*. 15, pp935-949.
- Schnabel, Z. E. (1938). The Estimation of the Total Fish Population of a Lake. *American Mathematical Monthly*. 45, pp348-352.
- Seber, G. A. F. (1982). *The Estimation of Animal Abundance*. 2nd edition. London: Griffin.
- Whittaker, E. T. and Watson, G. N. (1927). *A Course of Modern Analysis*. Cambridge University Press.

- Yip, P. (1989). An Inference Procedure for a Capture and Recapture Experiment with Time-Dependent Capture Probabilities. *Biometrics*. 45, pp471-479.
- Yip, P. (1991). A Martingale Estimating Equation for a Capture-Recapture Experiment. *Biometrics*. 47, pp1081-1088.
- Yip, P., Fong, D.Y.T. and Wilson, K. (1993). Estimating Population Size by Recapturing Sampling via Estimating Function. *Stochastic Models*. 9, pp179-194.

[民國 82 年 6 月 26 日收稿]

Consistency of an Estimator of the Number of Species ... In Memory of Professor Wen-Chen Chen

H-S Lin

Institute of Statistics

National Tsing Hua University

Anne Chao

S.-M. Lee

Institute of Statistics

Department of Statistics

National Tsing Hua University

Feng-Chia University

ABSTRACT

A theorem in Chen (1980,1981a,1981b) is extended to prove the consistency of an estimator of the number of species derived from a martingale estimating function (Chao et al. 1993).

Key words and phrases. number of classes, multinomial, heterogeneity, sample coverage, zero mean martingale.

AMS 1991 subject classifications. Primary 62P10, Secondary 62F10.