

經營	105
測定	34

森林測定研究資料

9

(標本調査關係論文)

昭和37年3月

農林省林業試驗場經營部

目 次

1	有限母集団よりの抽出理論	1
	(<i>M. M. Hansen, W. N. Hurwitz</i>)	
2	系統的抽出の理論について (I)	42
	(<i>W. G. Madoc, L. H. Madoc</i>)	
3	ある種の母集団に対する系統的標本と無作為標本の相対的 精度	75
	(<i>W. G. Cochran</i>)	
4	系統的抽出の理論について	95
	(<i>W. G. Madoc</i>)	
5	二次元のサンプリングについて	125
	(<i>M. H. Quenouille</i>)	
6	二次元の系統的抽出およびこれと関連する層化および無作 為抽出について	154
	(<i>A. C. Das</i>)	
7	種々の形式の二重抽出法による確定値の誤差について	175
	(<i>K. C. Seal</i>)	
8	プロットの大きさにとりなう収量の分散の変動について	193
	(<i>P. Whittle</i>)	
9	系統的抽出研究の展望	210
	(<i>W. R. Buckland</i>)	
10	系統的抽出の理論 (II)	228
	(<i>W. G. Madoc</i>)	
11	平面における確率過程について	230
	(<i>P. Whittle</i>)	
12	系統的標本の平均値の分散	264
	(<i>R. M. Williams</i>)	

1 有限母集団よりの抽出理論

On the theory of sampling from populations

1 近代的標本抽出理論の歴史的基礎

抽出単位と分析の単位が一致している母集団から、要素を無作為に抽出する理論は200年以上の昔 *Bernoulli* によって展開された。それより一世紀遅れて、サンプリングに層化を導入することによってえられる利益を評価する理論が *Poisson* によって示された。引続いて *Lexis* がそれまでの研究を体系づけ、要素の集塊抽出 (*sampling cluster of elements*) に対する理論的な基礎づけを行った⁽¹⁾。 *Bernoulli* および *Poisson* の有限母集団からの標本抽出理論に対する修正は *Poisson* の研究から約1世紀遅れて1926年に *Bowley* がとりまとめた。

Pearson, *Fisher* および他の人々の寄与に引続く標本調査理論の急速な進歩は、*Neyman* が1934年にその論文で *representative method* の二つの異なる側面について発表した[8]ことに利戟されたものである。その論文の中で彼は、標本内の全抽出単位数は一定であるという条件のもとにおいて、色々な層に対する抽出単位の最適割当の概念を含む調査資源の最適利用の基準を導入したのである。

もし標本が抽出されても、費用が標本に含まれる要素の数に全く支配されるのであれば、層内の要素の独立な無作為抽出を取扱

(1) 要素の集塊の調査とは、1つ以上の要素を含む単位の抽出を意味する。集塊抽出の例としては、個々の人または個々の世帯から構成される母集団特性を知ることを目的とする調査で *city block* または *county* を抽出単位とする場合である。これらの例で *city block* または *county* は要素の集塊と考えられ、個人または世帯は要素と考えられる。

った Bernoulli および Poisson の古典的理論以上の、有限母集団理論の拡張、あるいは抽出単位の最適割当に対する拡張などで補正された理論を用いる必要性は少ないのである。しかし統計的研究で母集団要素の独立な無作為抽出を行なうことは、不可能でないまでも非常に屢々多額の費用を要するのである。実際の場合そのような標本調査では母集団の全要素を確認することのできるリストが必要であるが、この記録は存在しないことが多いか、または入手に非常に費用がかかる。そのような記録が利用できても、標本が非常に広範囲にわたっていると、要素を数えるための費用が極めて大となる。同様に標本計画 (sample design) に他の制約があることも多い。例えば計算者の作業を限られた人数の監督者の厳重な管理下におくとか、限定された行政中心地で野外作業を行なわなければならない場合などである。集落抽出 [2, 3, 4, 5, 6, 7, 8, 10]、副次抽出 (subsampling)、二重抽出 (double sampling) [9] のような抽出法は、利用できる資料の最も有効な使用を目的として、一方では設けられた管理的制限の範囲内で止まりながら他方ではこれらの資料および制限内で可能な最大限の情報を入手するために発展してきたものである。この点については Neyman [8], Yates および Zaccagny [10], Cochran [2], Mahalanobis [7] および他の人々が重要な寄与を行っている。

我々は上に引用した多くの進歩を簡単ではあるが純粋に一貫的な副次抽出計画として説明することができる。この計画は層化母集団の要素の集落抽出および抽出された各集落要素の副次抽出を含むものである。この場合層内の一次抽出単位 (primary sampling unit) のそれぞれにおける要素の数は同一である。

我々は母集団を L ケの層に分ち、 i 番目の層はそれぞれ N_i の要素を有する M_i ケの一次抽出単位を含むものと仮定する。 X_{ijk} は i 層の j 一次抽出単位の k 番目のある性値を示す値とする。さうして推定すべき性値は、

$$(1) \bar{X} = \frac{\sum_i^L \sum_j^{M_i} \sum_k^{N_i} X_{ijk}}{\sum_i^L M_i N_i}$$

であると仮定する。例えば \bar{X} は与えられた街 (city) の一世帯当りの平均収入である。 X_{ijk} は i 区 (ward) の j 通り (city block)、 k 世帯の収入である。ここで世帯は副次抽出単位であり、通りは一次抽出単位、かつ層化は区によるものである。更に我々は i 層から m_i ケの一次単位を抽出し、この一次抽出単位の n_i から n_i ケの単位を副次抽出するものとする。

標本からの \bar{X} の "最良線型不偏推定値 (best linear unbiased estimate)" [8] は

$$(2) \bar{X}' = \frac{\sum_i^L \frac{M_i N_i}{m_i n_i} \sum_j^{m_i} \sum_k^{n_i} X_{ijk}}{\sum_i^L M_i N_i}$$

でありまた \bar{X}' の分散は

$$(3) \sigma_{\bar{X}'}^2 = \sum_i^L \frac{M_i^2 N_i^2}{M_i - 1} \left\{ \frac{M_i - m_i}{M_i - 1} \frac{\sum_j^{m_i} (\bar{X}_{ij} - \bar{X}_i)^2}{M_i m_i} + \frac{N_i n_i}{N_i - 1} \frac{\sum_j^{m_i} \sum_k^{n_i} (X_{ijk} - \bar{X}_{ij})^2}{M_i N_i m_i n_i} \right\} / \left(\sum_i^L M_i N_i \right)^2$$

$$\text{ここで } \bar{X}_{ij} = \sum_k^{n_i} X_{ijk} / n_i, \quad \bar{X}_i = \sum_j^{m_i} \sum_k^{n_i} X_{ijk} / M_i N_i$$

である。

これらの公式は更に費用の区別を考慮に入れなければ実際の抽出計画には役立たない。時には費用の関係を m_i および n_i の函数として正確に表わしうることがあるが、これは多くの場合、計画に含まれる種々の可能性の中から直観または理論を通じて妥当な次第を導くのに十分な近似を与えるものである。

もし我々に費用函数がわかっているならば、定められた費用に対し、また加えられるかもしれない他の制限のもとにおいて $\sigma_{\bar{X}'}$ を最小ならしめる m_i および n_i の値を定めることができる。この理論は色々な層ならびに層内の一次および二次抽出単位の抽出比の最適配分を決定するための基礎となるものである。

しかしこのような発展は標本計画の才ノ段階としてのみ考えらるべきである。我々は最適抽出計画が母数と費用因子のある種の函数であることを知っただけでは前に進むことはできない。それと同時に種々のサンプリングまたは推定に附随する費用の知識と考察している特定の母集団内のある母数の相対的な大きさについての知識が必要である。

このように最近の多数の論文の取扱っているものは種々なタイプの抽出単位の単位同および単位内の分散および共分散の相対的な大きさ、ならびにそれに作用する費用函数のタイプについての研究である。この方面の研究は農業的な対象については *Department of Agriculture* が行なっており、また *Bureau of the Census* あるいは他の場所においても研究されている。

II 最近の研究動向

上に指摘した標本抽出法は、才ノ段階としての、例えば抽出法は某簿抽出法、二重抽出法または副次抽出法を含むかどうかというように標本抽出方式 (*system of sampling*) の定義、もしくはこれとともに層化および抽出単位の定義を含むものである。才ノ段階は抽出単位の配分と抽出法を決定することである。

才ノ段階の抽出方式の定義は管理上の実行可能性と効率の面から行なわれ、しかし非実際的なリストの作成、またはすべての可能な方法を考え、その中からある基準で最良のものを一つ選ぶというようなことをする以外にはいつでも同じ方式を選び出せる簡単な方法は存在しない。しかし推定さるべき母集団特性の与えられた定義に対し、我々が最良線型不偏推定値のような「最良」についてのある種の基準を受け入れるならば、任意の費用あるいは規定された制限内で才ノ段階に対する単一の解を与える簡単な方法がある。このような方法は、推定法および定義された抽出方式内におけるサンプリングの配分法の双方を我々に与えるものである。

「最良」の抽出方式を示すと同時に「最良」の推定法および抽出配分を示すための実際に応用できる理論はこれまで展開されていないが、改善された抽出方式と推定方式の選択については若干の進歩があった。我々は次の二つの方向の進歩が特に当を得たもののように思われる。

- 1 あるタイプの標本抽出方式で、より確実な結果を与えるための、全く一般に受け入れられる良好な推定値というものについての基準の修正(これらのいくつかは *Sec. III* で述べる)
- 2 抽出単位、層および抽出方式の他の面の決定を改善しうるある種の原則が現われ、このような原則を構成しようとする試みについては *Sec. IV, V* および *VI* で述べる。

我々は主としてセンサスにおける最近の研究をまとめることにする。したがって密接な関係をもつ他の研究についても論及するであろう。我々が要約しようとする研究の大部分は、抽出単位が要素の集落で、その大きさが種々変化するという問題に關するものである。

III 良い推定値の基準に対する修正

Sec. I において提起した副次抽出の一般問題で与えられた推定値は「最良線型不偏推定値 (*best linear unbiased estimate*)」の基準を踏たすものである。我々の堅執によれば、これは層同の各抽出単位中の要素の数々等しい母集団では、層々最も有効 (*efficient*) なものである。しかし抽出単位同の要素数が違っている場合でも、偏ってはいるが大体において最良線型不偏推定値よりも平均平方誤差 (*mean square error*) の小さい一致推定

3) $E(x - \hat{x})$ なるときの $E(x - \hat{x})^2$ を指すためこの論文では「平均平方誤差」と「分散」を同じ意味に用いる。 $E(x)$ は推定さるべき母集団特性である。しかし $\sum x^2 / N$ のときは $E(x - \hat{x})^2$ はただ「平均平方誤差」だけを指すものとする。後者の場合は $E(x - \hat{x})^2 = E(x - EX)^2 + E(x - \hat{x})^2$ となって平均平方誤差は x の分散 + 偏りによる変動となる。

量 (consistent estimate) が存在することが多い。

例えば集落が抽出単位の時、ある特性の要素当りの平均値 $\bar{X} = \sum_{i=1}^M X_i / \sum_{i=1}^M N_i$ を推定したいとする。ここで M は母集団の抽出単位の数で、 X_i は i 番目の集落の全要素に対する特定の性質の合計値 (aggregate value) である。 N_i はこの集落の要素の数、 X_i と N_i の同時分布は未知であるが、 $\sum_{i=1}^M N_i = N$ はわかっている。このような条件のもとで m 々の集落を含む標本からの "最良線型不偏推定値" は $\frac{M}{m} \sum_{i=1}^m X_i / N$ なることかわかる。しかし標本から $\sum_{i=1}^m X_i / \sum_{i=1}^m N_i$ なる比推定を作れば多くの場合更に小さい平均平方誤差がえられる。この推定値の偏りは通例無視できるものであり、一致性をも有するが、偏りのあることと非線型という理由で "最良不偏推定値" の基準にはあてはまらない。 \bar{X} の最良線型不偏推定値には N がわかっているなければならないから、この標本比は、 N が未知の場合でも使えるという点で優れている。

最近の Cochran [3] の論文は、 X_i の N_i への線型回帰の最小二乗推定を用いて、一致性はあるが、偏りのある \bar{X} の多くの推定値を与えている。

これらの推定値の平均平方誤差は一般に上にあげた最良線型不偏推定値あるいは単純比推定値のそれよりも小さい。しかしこれらに対しては最良線型推定値と同じく N の知識が必要であり、その上推定の手続の一部として相当な計算を要する。

上に述べた偏りのある推定値はいずれも一致推定量であって、普通抽出単位の大さが異なるような抽出方式では、最良線型不偏推定値よりも平均平方誤差が小さい。したがって、"最良線型不偏推定値" の基準に、非線型、一致推定量ではあるが平均平方誤差が最良線型不偏推定値よりも小さいような推定値を含むことにすれば、改善された標本推定値がえられる。

IV 抽出方式の指定に関する改善

母集団構成に関する知識が少なく、その特定の母集団母数に

いての知識が皆無の場合でさえ、抽出方式の指定の改善を通じて、抽出計画を大きく改善することが出来る (Sec. VII をみよ)。

1 抽出単位の大さ

最近の研究の多くは、圍場調査の抽出単位に要素の集落を用いるとき、考えられている費用のもとで集落の大きさをできるだけ小さくすることが望ましいことを示している [2, 5, 6, 7, 8]。しかしこの原則は副次抽出方式では必ずしもあてはまらないことに注意しなければならない。副次抽出を含む抽出方式では、一次抽出単位に大きい集落を用いると副次抽出を用いないで小さい集落を用いた場合よりもっと別の利点がある。その上大規模調査において屢々生ずる問題の一つは、調査を行なう場所数に管理上の嚴重な制約があつて、そのような条件のもとで調査計画をたてなければならないことである。このような制限がある場合によく用いられる方法は、限られた数の county などの行政単位を一次抽出単位とし、副次抽出単位には block や小さい地方的地域あるいは世帯をとるものである。このような事情のもとで標本内に含まれる一次抽出単位の数を一定に保つとすれば、実際に存在する行政的単位より大きい一次抽出単位を用いることによって抽出分散を減少させることが出来る。

副次抽出で大きい抽出単位を用いることの利点は、例えば、 m 回数の要素を含む元の単位 (original unit) を、それぞれ $\frac{1}{2}$ の大きさをもつ $\frac{1}{2}m$ の数の一次抽出単位に結びつけるような簡単な場合を考えれば明らかである。拡大された一次抽出単位間の分散は $\sigma_{ic}^2 = \frac{1}{2} \sigma_{ib}^2 (1 + \rho)$ となる。ここで σ_{ib}^2 は元の一次抽出単位間の分散で、 ρ は対にした単位間の相関係数である。この相関係数は対が確率的に作られている場合には殆んど 0 に近い (正確には $1/(M-1)$ に等しい)。ここで M は元の一次抽出単位の数である。したがって county 間の分散は少くとも略々半分になる。対にする単位は隣りあつていなければ

ばならないという条件があれば、普通 f の値は ρ より大とならう。しかしなるべく近いの大きい隣接単位を選んでこのような単位の結合をつくれれば、 f は可能な限り小さくなり、ある場合にはこの極小値が負になることさえある。何れの場合でも、 f の値が小さい程拡大単位を用いて生ずる抽出単位間の抽出分散 (*sampling variance*) の減少は大きくなる。一方、一次抽出単位間の抽出分散はこのような結合を行なうことによって増大するが、その増加は僅少だから、全抽出分散は常に減少する (Appendix Sec. 1)。一次抽出単位の結合を大きくする場合の制約は、地域が増すにつれて副次抽出の費用が増大することによって生ずる。このような費用の増加は分散の減少と逆の働きをする。もし費用の制限がそうきびしくなければ、結合は一次抽出を全然用いずに副次抽出単位を各層から独立に抽出する方向に進められるであろう。

2 一次抽出単位の大きさが異なる場合の副次抽出。副次抽出で大きさに比例する確率を用いること。

一次抽出単位の大きさが変化すると否とにかかわらず層々用いられる副次抽出組織では、それぞれの層から層内の一次抽出単位の抽出される確率が同じになるようにして、 1 つまたはそれ以上の一次単位を抽出し、その中から定められた割合の副次抽出単位を抽出しなければならない。一次抽出単位の大きさが異なる場合には、抽出された一次抽出単位内の要素数につれて標本中に含まるべき副次抽出単位の数が変化するため、この抽出方式には若干の管理的な不利益が生ずる (この論文において用いる "抽出単位の大きさ (*size*)" とは、抽出単位内に含まれる要素の数を指すものである)。上記の抽出方式の不利益はたとえ一次抽出単位が等確率で抽出されても、副次抽出が一定比率というよりはむしろ一定数であるように指定された才の副次抽出方式の場合に生じたものであった。上記の二つの抽出方式より更に推奨できる才の抽出方式は、一次抽出単位に

の大きさに相応する抽出確率を与え、したかつて一定数の副次抽出単位を副次抽出するものである。我々は3つの抽出方式の何れについても、各層からは唯一つの一次抽出単位を抽出するものと仮定する。この程度まで層化を行なうと、層化が不完全な場合よりはるかに抽出分散が小さくなる。比較のための単純化として、我々は更に副次抽出単位が分析単位であり、用いられる標本推定値が $\bar{X}' = \sum N_k \bar{X}_k / \sum N_k$ なる形のものとする。ここで \bar{X}_k は推定している特性の各層に対する標本平均であり、 N_k はその層の大きさである。層々用いられる推定値は、最初の二つの抽出方式では偏りがあるが、推奨された抽出方式では不偏である。しかし不偏推定値すなわち最初の二つの抽出方式に対する "最良" 線型不偏推定値の平均平方誤差は、一般にこれらの比較で用いられる偏りのある推定値のものよりずっと大きいから、次の比較では考慮しなかつた (Sec. III, 脚注 A を見よ)。

前にあげた最初の二つの副次抽出方式は、何れも適当な大きさの一次抽出単位から副次抽出を行なえば、その効率は大體同じである。しかしどちらも推奨した抽出方式より普通平均平方誤差が大きい。最初の二つの抽出組織と推奨した計画との平均平方誤差の差は近似的に何れも

$$(4) \frac{1}{N^2} \sum_k Q_k \bar{N}_k \sigma_k^2 \left[\sum_j \beta_{kj} \bar{N}_k - \sum_j \beta_{kj} N_{kj} \right]$$

で表わされる。ここで N_{kj} は k 番目の層内の j 番目の一次抽出単位内の要素数であり、 \bar{N}_k は一次抽出単位の平均の大きさである。 Q_k は一次抽出単位の数、 β_{kj} は j 番目の単位内要素間の級内相関係数 (*intra-class correlation*)、 σ_k^2 はこの層内の個々の要素間の分散で、 L は層の数である (Appendix Sec. 2 にあるこの差の展開をみよ)。

もし N_{kj} と β_{kj} の相関が負なら、 N_{kj} と β_{kj} 間の平均の共分散 (*covariance*) であるところのこの差は正になり、そうして

これは我々が社会統計および経済統計で屢々遭遇する大部分の実際的問題においてみられるものである(*Sec. VI* をみよ)。

推奨される計画において平均平方誤差が減少するのはその大抽出単位と小抽出単位間の抽出単位の配分が、他の二つのものより最適に近いことによるものである。もちろん他のものと同様に、大きさを一次抽出単位を層化し、最適配分を考慮して色々な層に抽出をふりわけるともできる。しかしこれによって何か他の、しかもおそらく更に重要な層化の方法が犠牲になり、加えて大部分の実際問題で大抽出単位と小抽出単位間の最適配分は推測による他はないであろう。その上、普通層内の単位の大きさが同じ場合には大きさにによる層化は不可能である。

推奨された抽出方式からえられる推定値は他の二つの推定値が通例通りをもち、ある場合には非常な偏倚を示すにもかかわらず不偏である(このことの証明は *Appendix Sec 1* の数値例については *Sec. VII* をみよ)。

大きさに比例した確率を用いても一次抽出単位間の抽出変動が減少するだけであるから、層内の抽出変動を減少させる上では効果が少ない。したがって一次抽出単位間の平均平方誤差の全平均平方誤差に対する寄与が大きい場合には推奨された計画は他の二つのものより非常に有利である。

通例一次抽出単位の実際の大きさは不明であるが、これと密接な関係をもつ単位の数はわかるものである。たとえば普通 *block*, *city*, あるいは *county* 等の母集団は、標本がとられたときには我々にはわからないものである。しかし以前のセンサスにおけるこれらの母集団については知ることができ、これらの状況のもとで、一次抽出単位はすでに知られていく(あるいは推定された)大きさに比例した確率で抽出される。しかしこのようなことが行なわれれば、二つの調査時点間の抽出単位の大きさの変化を考慮するため、この副次抽出は修正されなければならない。

もし実際の大きさが既知なら、 k -層で抽出された一次抽出単位の中からとる副次抽出単位の数は、

$$n_k = t_k N_k$$

である。そこで t_k はこの層に対して定められた抽出比であり、 N_k は層内の全要素の数である。したがって抽出された一次抽出単位内の副次抽出率は $t_k N_k / N_{kj}$ である。 N_{kj} は抽出された単位内の要素数である。一方もし単位の実際の大きさ N_{kj} と密接な関係にある大きさの尺度 P_{kj} だけしか利用できなくて、一次抽出単位の抽出確率を P_{kj} に比例するようにとった場合には、抽出された一次抽出単位内の副次抽出率は $t_k P_k / P_{kj}$ に等しくなる。ここで P_k は層全体に対するこの尺度であり、 P_{kj} は抽出された一次抽出単位の大きさの尺度である。大きさの尺度を用いたときの標本推定値の分散はこの論文のあとで述べる(*Eq. (9)* をみよ)。

3. 副次抽出方式の場合一次層内で地域的副次層化を用いること。
もし地域にわたる母集団から比較的小さい標本をとり、その中心を小数の中心地で行なわれなければならないような場合には、一次層内(*primary strata*)で地域的副次層化(*area sub-stratification*)とよばれる他の修正を行なうことが特に有用である。地域的副次層化について説明する前に少し準備が必要である。

地域的副次層化には、

- (a) 標本調査を行なつた全母集団を一次抽出単位となる地域に分割する。
 - (b) これらの単位を更に細分して幾つかの副次区域(*sub-area*)に分割する。そうして
 - (c) 標本を抽出する前に各副次区域についてある種のまとまった統計的知識が利用できる。
- ことが必要である。この副次的な地域について知らねばならない知識にはそれらの大きさと身分全母集団、全世界(*dwelling*)

unit) あるいは全農家) に対する妥当かつ良好な尺度およびこの地域の特性を示すた(たとえば農業地域 (predominantly farm) であるか、非農業地域 (predominantly non-farm) であるか、あるいは白人の多い地域 (predominantly white) か有色人種の多い地域 (predominantly colore) であるかなどの他の知識も含まれる。齊一なクラスに群わけした場合、副次区域は次に述べる副次層 (sub-strata) を決定するためだけに用いられるもので、普通それと独立に定められる副次抽出単位を与えらるものではない。前に述べたように層内の一次抽出単位から選ばれるだけ齊一になるように、また層間で可能な限り代表的になるように一次抽出単位の定義およびその層への群わけを行なう。この場合各層からは一次単位を一つだけ抽出し、その層の j 一次単位の選ばれる確率は P_{ji} に比例させるということも仮定される。ここで P_{ji} は一次単位の大きさの尺度で、これはその中に含まれる副次区域の大きさの尺度の和に等しい。また各層内で用いる全抽出率 w_i は全部の層に対して最適割合の考えをもととして決定されているということも仮定される。

そうすると地域的副次層化は次のようにして一次層内に導入できる。

- (a) それぞれの一次層内の副次区域をそれらの特性をもととして副次層に群わけする(たとえば農業地域と非農業地域という二つの副次区域に群わけし、次にこれらの農家の平均的な大きさまたは世帯 (dwelling unit) の平均地代によって更に分類する。そのようなとき、農業地域にあり、かつ地代の平均が指定された幅の中にあるような一次層内の副次区域は一つの副次層を構成する)。
- (b) 一次層の各々から選ばれた一次単位内の副次区域は同じ副次層内に群わけする。
- (c) 副次抽出単位は選ばれた一次単位内の各副次層内で定められる。 j 一次抽出単位内に含まれる i 番目の副次層の中

められる副次抽出単位の個数は M_{ij} で表わされる(種々の副次抽出単位、たとえば個人、農場、世帯またはその構成員、あるいは非常に小さい地域などとして定義することができ、副次抽出単位は選ばれた一次抽出単位の中でのみ定義しなければならない)。

- (d) 選ばれた j 番目の一次抽出単位内の i 番目の副次層からとった標本に含まれる副次抽出単位の数は、

$$m_{ij} = M_{ij} \cdot p_{ji} / P_{ji}$$

である。ここで P_{ji} は j 番目の一次単位内にある i 番目の副次層の大きさを表わす。 $P_{ji} = \sum P_{kij}$ は j 番目の一次層 (Primary stratum) の i 番目の副次層に含まれる副次区域の大きさの和である。

このような配分法によって副次抽出を行なえば、選ばれた一次単位からとった標本は、その層から標本内にたまたま選出された特定の一次単位を代表するというよりはむしろ層全体を可能な限り代表することになる。説明のため、1940年のセンサスからえられた副次区域内の人数をそれらの大きさを表わす尺度にとり、これらの副次区域を1940年の Decennial Census of Population で示された1940年におけるそれらの特性にもとずいて副次層に分類するものとしよう。そうすると上述の副次抽出の配分から、農業地域内の副次区域に居住する人口比率が30%なら、標本内で予想される1940年の人口の30%が農業地域内の副次区域からのものとなるように標本を抽出する。このとき選ばれた一次抽出単位内では1940年の人口中そのような地域に住んでいたのはおそらく15%にしかならないかも知れないがそれは問題外である。

- (e) 推定すべき母集団特性は、

$$(6) \quad X = \frac{1}{N} \sum_i \sum_j \sum_k \sum_h M_{ijkh} X_{ijkh}$$

である。ここで X_{ijk} は j 番目の一次単位、 i 番目の副次層を k 番目の副次抽出単位に含まれる全要素についてある特性の値を加え上げたものである。また S_k は副次層の個数、 Q_k は k 番目の一次層内の一次単位の数であり L は一次層の数である (X は米国の全労働者数でもよいし、また全農業労働者数等でもよい)。標本からの推定値は、

$$(7) \quad \bar{X} = \sum_k \frac{1}{L} \frac{S_k}{Q_k} \sum_{i=1}^{M_{jk}} X_{ijk}$$

である。 j に対する和は不要である。何故なら k 番目の層から抽出される一次単位は L だけだからである。これは一次層の段階だけで重みをつけた和を含む非常に単純な推定値である。もし Q_k がすべて L に等しい、すなわち標本の抽出比を一定とすれば、推定値は標本内のある特性を有する全要素の数に抽出比の逆数 $1/Q_k$ をかけただけのものとなる。

上に述べた副次抽出の概念と違った方法をとることもできるが、その場合でも標本推定値に正しい修正を加えれば地域的副次層化の拘束が維持される。このときには一次層ばかりでなく副次層の段階にも別の重みを導入しなければならない。

異なる一次抽出単位の定義やそれらを正しく層化すること、あるいは一次単位を選ぶ際に大きさに比例した確率を用いることなどは、地域的副次層化を利用するときには特に望ましい。もしこれらを導入しないときは、地域的副次層化を行なっても実質的な利益のえられる可能性は少なくなるであろう。一次層の定義は副次層の定義と関連させて行なうべきであり、また各一次単位も一次層内で定められるそれぞれの副次層を十分代表するものでなければならない。ここに見られた制約のために、定めらるべき層の数は一層の数より一次層の異質性によって制限を受けらるであろう。したがって農業地域と非農業地域に副次層化する場合には農業地域と非農業地域の両方が各単位内で代表されるように一次抽出単位を定めなければならない。この方法は副次層化を一層効果的に

するばかりでなく、母集団のそのようなクラスの推定を別々に行なう場合標本の効率を改善するものである。実際の場合にこの手続を正確にあてはわれなければ、ある副次層を代表しないような抽出単位が標本内に入ってくることが多い。このような場合にとるべき方法の一つとしては特定の副次層を合併することであり、もう一つは標本からそのような一次単位を取除くことである。

一次層の数は抽出すべき一次単位の数によって制約を受けらるから、地域的副次層化で十分抑制できる変動因について一次の段階 (primary level) で層を作ることは不経済である。例えば副次層で農業地域と非農業地域を区別すべきものとするれば、一次単位を農家比率 (percent farm, 全母集団中の農業による生計をたてている一次単位の比率) で多数の層に分類してしまつてはならない。それは農家比率の変動をコントロールするのは副次層化たからである。一次の段階で農家比率のクラスの数を制限することによって、農家の型や非農家母集団の生産的特性、あるいは他の関連な基準をコントロールできる別な形の層化を利用することができる。

地域的副次層化は、要素の色々なクラスの各々から、標本内に含まれる要素数を定めるのにごく普通に用いられている方法——一次層の全部または選ばれた一次抽出単位の特定の特性と標本とを一致させるために、そのような制当が固定されていてもいなくても——と区別しなければならない。

与えられた数の、指定された色々な特性をもつ要素 (人、居住単位、農場あるいは有権者等) を与えるために、制当を定めたり、面接者または調査者に指令を与えたりする方法には根本的な欠点があつて、一次層内の地域的副次層化では避けなければならない。普通そのような制当は既往の情報や大雑把な推定をもととしなければならないことから、母集団の变化した特性を正確に示すことはできない。一方地域的副次層化は既往の情報を利用して、色々なタイプの地域が標本内に正しく代表されることを保証するもの

である。指定された種々の特性についてとられる要素の数は過去の時点ではなく現在の母集団から決定される。変化が急激な時期に過去の情報に基づいて一定の割当を用いるとますます著しい偏りを生ずることになる。地域を層化する際、利用できる過去の情報を用いることによってえられる利益は、地域の時期毎の特性同に数年にわたって高い相関があるという事実から生ずるものである。ある時点において農家比率の高かった地域(農業地域)は数年後でも農家比率は高いのが普通である。同じように、母集団に事実上非常に大きい変化が生じても、ある時点における一組の地域内の人口は普通数年後の人口と非常に高い相関をもっている。しかし地域的副次層化は変化が生じないという事実に頼るものでなく、変化が生じたならそれを評価するものである。変化が大部分の小地域の特性を全く変えてしまう程大きい場合でも、この方法は母集団特性の変化を明らかにする推定値を与えるであろうが、しかしそのような事情のもとではこの方法の効率は減少する。

V 上述の原則を具体化する副次抽出方式の期待値と分散

一次抽出単位の拡大、大きさの尺度に比例した確率をもって行なう一次抽出単位の抽出、および地域的副次層化の原則を具体化する抽出方式を以下において一層詳細に論ずることとする。これを便宜上指定された副次抽出方式とよぶことにする。

1 指定された副次抽出方式に対する総計の推定値の期待値

以下の公式中の和は特に断らない限りすべて母集団全体にわたるものとする。

(7)式で定義された X' の期待値は、

$$EX' = \sum_k \sum_j \sum_i \sum_r \left(\frac{1}{t_k} \right) \left(\frac{P_{kj}}{P_k} \right) \left(\frac{m_{kij}}{M_{kij}} \right) X_{kijr}$$

である。(5)から $t_k = m_{kij} P_{kij} / M_{kij} P_{ki}$ だから

$$EX' = \sum_k \sum_j \sum_i \sum_r \left(\frac{P_{ki}}{P_{kij}} \right) \left(\frac{P_{kj}}{P_k} \right) X_{kijr}$$

$$= \sum_k P_k \sum_j \sum_r \left(\frac{P_{ki}}{P_k} \right) \left(\frac{P_{kj}}{P_k} \right) \left(\frac{X_{kijr}}{P_{kij}} \right) = \sum_k P_k R_{k(CA)}$$

ただし、

$$P_k = \sum_i P_{ki} = \sum_j P_{kj} \quad ; \quad R_{k(CA)} = \sum_j \left(\frac{P_{kj}}{P_k} \right) R_{kj(CA)} \quad ;$$

$$R_{kj(CA)} = \sum_i \left(\frac{P_{ki}}{P_k} \right) R_{kij} \quad ;$$

$$R_{kij} = \sum_r X_{kijr} / P_{kij} = X_{kijr} / P_{kij}$$

この $R_{k(CA)}$ は j -一次単位に対する補正された比と考えられる。これは j 単位内における R_{kij} の重み付き平均値である。ここでは層内のそれぞれの一次単位の R_{kij} に対し同一組の重みが適用される。 $R_{k(CA)}$ は調整された比の j -層内の平均である。よって

$$(8) \quad EX' = X + \sum_k P_k (R_{k(CA)} - R_k)$$

ここで $R_k = X_k / P_k$, $X_k = \sum_j \sum_i X_{kij}$

は、 j 層内要素の指定された特性値の総計と層の大きさとの比で、(8)で推定している母集団特性は $X = \sum X_k = \sum P_k R_k$ に等しい。(8)からわかるように、通常それは実際には僅かだが X' は X の偏りのある推定値である。この偏り $\sum P_k (R_{k(CA)} - R_k)$ は色々な一次層に対する偏りの和である。現実には多くの場合これらのあるものは僅かに正となり、あるものは僅かに負となるだろうから、全体の偏りは比較的小となるであろう。この偏りは地域的副次層化を用いないか、または標本推定値の式を正しく修正すれば生じなかったであろうか。この場合もまた Sec. III で述べた不偏推定値の代わりに偏りのある推定値を代入した場合と同じく、僅かの偏りを導入することによって分散を事実上減少させうるのである。

地域的副次層化を行なったときの標本推定値(7)が不偏であるための十分条件(必要条件ではない)は、比 P_{kij} / P_{ki} が各副

次層内の R_{hij} と無相関であることである。これらの条件のもとでは、

$$\sum_j \frac{P_{hj}}{P_h} \frac{P_{hij}}{P_{hj}} R_{hij} = \sum_j \frac{P_{hj}}{P_h} \frac{P_{hij}}{P_{hj}} \sum_j \frac{P_{hj}}{P_h} R_{hij} = \frac{P_{hi}}{P_h} \sum_j \frac{P_{hj}}{P_h} R_{hij}$$

だから

$$R_h = \sum_i \sum_j \frac{R_{hij}}{P_h} \frac{P_{hij}}{P_{hj}} R_{hij} = \sum_i \sum_j \frac{P_{hi}}{P_h} \frac{P_{hij}}{P_h} P_{hij} = R_{h(CA)}$$

例をあげれば、1940年の人口を大きさの尺度とすると、標本推定値は、種々の副次層内の一次抽出単位に対する1940年の人口比率が対応する R_{hij} と無相関なら不偏である。前に注意したように、これらの条件は多くの実際問題において、特に才一次層化 (primary stratification) が効果的に行なわれた場合には近似的に満足される。更にこの条件が近似的に満たされなくとも、導かれる偏りは非常に小さいものである (数値例は sec. 7 をみよ)。

2 指定された副次抽出方式に対する総計の推定値の平均平方誤差

指定された副次抽出方式に対する X' の平均平方誤差の展開は附録 sec. 2 をみよ。それによれば X' の平均平方誤差は、

$$\begin{aligned} \sigma_{X'}^2 &= \sum_h \sum_i \sum_j P_{hi}^2 \frac{P_{hij}}{P_h} \frac{M_{hij} - m_{hij}}{M_{hij} - 1} \frac{\sigma_{hij}^2}{m_{hij} P_{hij}} \\ (9) \quad &+ \sum_h P_h^2 \sum_j \frac{P_{hj}}{P_h} (R_{hij(CA)} - R_{h(CA)})^2 + \left[\sum_h P_h (R_{h(CA)} - R_h) \right]^2 \end{aligned}$$

ただし $\sigma_{hij}^2 = \sum_k (X_{hijk} - \bar{X}_{hij})^2 / M_{hij}$ は副次抽出単位に対する指定された特性値総計の副次層内における副次抽出単位間の分散であり、 $\bar{P}_{hij} = P_{hij} / M_{hij}$ は才 h の i - j 番目の地域内における副次抽出単位の平均の大きさである。

(9) の才1項は副次抽出単位間の分散の寄与を表わすもので、これは副次抽出単位を正しく定義することにより、また勿論副次抽出比を増すことによって小さくできる。(9) の才2項は層

内の一次抽出単位間の分散の寄与を示す。才3項は偏りの寄与であるが、前に指摘したように通例無視できる大きさであるから、この平均平方誤差と分散は近似的に等しくなる。

多くの副次抽出で全分散に最も大きく寄与するのは一次抽出単位間の分散である。そうしてこの論文で提案した修正が効果をあげるのは主としてこの寄与に対してである。地域的副次層化の効果は、上に与えた一次単位間の分散と、地域的副次層化を用いず、また計画の他の部面を変更しなかったときにえられる分散を比較すればわかる。この場合一次単位間の分散には補正された比 $R_{hij(CA)}$ の分散でなく、比 $R_{hij} = \sum_k X_{hijk} / R_h = X_{hij} / P_{hj}$ の分散が含まれる。

R_{hij} の分散と才1項一次層内の $R_{hij(CA)}$ の分散の間の関係は、

$$(10) \quad \sigma_{R_{hij}}^2 = \sigma_{R_{hij(CA)}}^2 + \sigma_{R_{hij} - R_{hij(CA)}}^2 + 2\rho \sigma_{R_{hij(CA)}} \sigma_{R_{hij} - R_{hij(CA)}}$$

で与えられる。ここで $\sigma_{R_{hij} - R_{hij(CA)}}$ は補正した比と補正しない比の差の分散であり、 ρ は補正された比と補正量との相関である。したがってもしこの相関が0に近いか、あるいは正なら副次層化を導入することによって利益が生れるが、しかしこの相関が負の大きい値をとるときには損失が生ずるのである。本質的に、 ρ が0に等しいかまたは0に近い値となるための条件は、標本推定値が不偏となるための条件と同じである。すなわち P_{hij} / P_{hj} が各層内の R_{hij} と無相関であるか、または僅かの相関しかもたないことである。

X' の分散の中に現われるのは R_{hij} のものよりはむしろ $R_{hij(CA)}$ の分散である。なぜなら、副次抽出単位の数はたまたま標本内ほどの一次抽出単位が選ばれたかには関係なく、 R_{hi} に比例し

(11) 実際には ρ が0に等しいための十分条件 (必要条件ではない) は P_{hij} / P_{hj} がすべての層の対に対して、比 R_{hij} および cross-product $R_{hij} P_{hij}$ の両方と無相関なことである。

て配分されているからである。 $R_{ij(A)}$ と同じく比 R_{ij} は R_{ij} の重みつき平均値と考えられるが、この重みは P_{ij} でなく P_{ij}^* に等しい。したがって一次単位ごとに変化する、ゆえに上に与えた分散の関係から、もし副次層が効果的に設けられ、また P_{ij}^* が副次層の実際の大きさと高い相関をもつならば、すべての一次単位において一定の重みを用いた重みつき平均値は重みを变化させた場合のものより分散が相当小さくなる等である。これは多くの実際の状況のもとにおいて明らかとなるものであるが、その説明はあとで (Sec. Ⅶ をみよ) 与えることにする。

3. 指定された抽出方式に対する比推定の平均平方誤差

標本から比を推定する必要が生ずるのは二つの場合である。その一つはこの比が推定したい母集団特性である場合であり、その二は希望する総計についての改善された推定値を求めるために、利用できる追加情報を用いて既知の総計に標本から求めた比を適用する場合である。

比の推定値はたとえば、時点間の特性の変化を考えているようなときは最終結果として要求される。したがってもしある日時における農業労働者の総所得の推定値を Y' とし、才二の日時における対応する推定値を X' とすると、 $Y' = X'/Y'$ はこの期間内の農業労働者の総所得についての相対的変化の一つの推定値である。同じように失業率のようなパーセンテージの推定値は標本から求めた二つの確率変数の比からなっている。標本からえた比の推定値は、それが分子または分母の何れの推定値より信頼性が高いとき(このようなことは屢々ある)特に有用となる。

比の推定値は、 X と高い相関をもつ才二の特性の総計 Y が独立な出所から正確にわかっている X と Y の推定値 X' と Y' が標本からえられるとき、指定された特性値総計の推定値を求める手段として利用できる。すなわち

$$(11) \quad X'' = (X'/Y') Y = r' Y$$

は指定された特性総計の推定値である。もし連続してとられた標本で、 X' と Y' の相関が十分高ければ、この比の推定値は前に (7) で与えた総計の単純推定値 X' より一層有効な X の推定を与えるであろう。しかし相関が低い場合には X' の方より信頼できる推定値であることが証明される⁽⁴⁾。故に X' と Y' との間の相関が十分高い場合には、 X' を推定するよりも X'' を求めることにより、 X の推定に利用できる情報を一層適切に利用しようである。

指定された副次抽出方式に対する比の推定値の適用を以下において考察しよう。

(a) 比の推定値とその平均平方誤差

母集団の比 $r = X/Y$ の推定値は

$$(12) \quad Y' = \frac{X'}{Y'} = \frac{\frac{1}{N} \sum_A \frac{1}{P_A} \sum_i \sum_j \sum_k \frac{M_{ijk}}{M_{ij}} X_{ijk}}{\frac{1}{N} \sum_A \frac{1}{P_A} \sum_i \sum_j \sum_k \frac{M_{ijk}}{M_{ij}} Y_{ijk}}$$

である。ここで X' は上の (7) で与えられまた Y' は才二の特性値総計の同様な推定値である。 r' の平均平方誤差は近似的に

$$\begin{aligned} \sigma_{r'}^2 = & \frac{1}{Y'^2} \left\{ \sum_A \sum_i \sum_j \sum_k \frac{P_{ij}^2}{P_A} \frac{P_{ij}}{P_A} \frac{M_{ij} - m_{ij}}{M_{ij} - 1} \frac{\sum_k Y_{ijk} (r_{ijk} - r_{ij})^2}{m_{ij} M_{ij} \bar{P}_{ij}^2} \right. \\ & + \sum_A \sum_i \sum_j \sum_k \frac{P_{ij}^2}{P_A} \frac{P_{ij}}{P_A} \frac{M_{ij} - m_{ij}}{M_{ij} - 1} \sigma_{r_{ijk}:Y}^2 \frac{(r_{ijk} - r)^2}{m_{ij} \bar{P}_{ij}^2} \\ & + \sum_A \frac{P_A^2}{P_A} \sum_j \frac{P_{ij}}{P_A} R_{A(j):Y}^2 (\bar{r}_{A(j)} - \bar{r}_{A})^2 \\ & \left. + \sum_A \frac{P_A^2}{P_A} (\bar{r}_{A(j)} - r)^2 \sum_j \frac{P_{ij}}{P_A} (R_{A(j):Y} - R_{A(j):Y})^2 \right\} \end{aligned}$$

(4) $r = X/Y$ なる形の確率変数の比の分布は、近似的に $\sigma_{r'}^2 = r^2 (V_X^2 + V_Y^2 - 2r_{XY} V_X V_Y)$ で表わされる。ここで V は添字で示された変数の変動係数で、 r_{XY} は相関係数である。したがってもし r_{XY} が十分大きければ $V_{r'}$ は V_X より小さい。 r_{XY} の大きさは X と Y の変動係数の相対的な大きさに依存せねばならない。

用いられる際でも、 X_{ijk} が少なくとも Y_{ijk} と平均値にかなり高い相関関係をもつことを確かめるのが無難である。ここで考えている相関は一次抽出単位内の副次層内のものである。この相関が低く、また副次抽出単位の大きさが相当大きく変動するならば、比推定の効率は単純推定の推定値よりもかなり低下するであろう。一方もし種々の副次層の大きさの尺度および一次抽出単位の大きさが実際の大きさの尺度に近いことが明らかであり、また大きさがあまり大きく変化しないように副次抽出単位を注意深く定義してあげればこの二つの推定値は大體同じ効率を与えるものと思われる。

VI この論文で推奨した抽出原理の基礎となる母集団に属々みられる自然的性質

現実の母集団の多くは次のような自然的な性質で特徴づけられる。

- (i) 集落内の要素が指定された特性と正の相関をもつ
- (ii) 多数の要素を含む集落は少数の要素を含む集落より内部の異質性が大きい
- (iii) 集落の大きさが増すと相関をもった要素が導びかれる(たとえば人口または農業調査において、大きい集落は隣接地域の世帯または農場を含めることによって作られる)

この最初の性質は、抽出単位に大きい集落を用いたため効率に損失を生じたことを例証している多くの文献中で暗黙のうちに認められている。我々の経験では、才2、才3の性質は現実の母集団で全く普通に認められるもので、これは通例才1の性質の成立する母集団と同じものである。

層内でこれらの自然的な性質が組合わせて現われると、この論文でこれまで使ってきた次のような数学的な関係が導びかれる。

- (a) 一次抽出単位の大きさ N_{ij} は単位内の級内相関 ρ_{ij} と負の相関をもつ

- (b) N_{ij} と $N_{ij} \rho_{ij}$ は正の相関をもつ

- (c) N_{ij} と σ_{ij}^2 は正の相関をもつ

- (d) N_{ij} と σ_{ij}^2/N_{ij} は負の相関をもつ

この論文ではこれらの関係を用いていままで色々な方法間の選択を行なってきた。これらの関係はもちろん常に成立するものではないが、それらの(5)に示されている。これらの性質によって特徴づけられる母集団が属々出現するということは、成立すると思われるこれらの性質や他の性質を一層効果的に利用する更に進んだ研究を正当化するものである。

VII この論文において記述した原理の現実の抽出問題に対する適用について

以下に大要を示す分析は、労働力および他の特性に対する国家月例標本 (*monthly national sample*) の改訂において、色々な抽出手続から選択を行なうために実行されたものである。予算と管理との制限から、外票は国内に散在する限られた数の行政中心地で行なう必要があった。そうしてこれらから合衆国人口の 1/1000 以下の標本を抽出すべきことが要求された。

(改訂すべき)もとの標本は一次抽出単位と1/2カウンティを用い、副次抽出単位として世帯または世帯の小集落を用いる普通の副次抽出計画であった。改訂された標本においては、それが管理上可能な限り、個々のカウンティより異質な一次単位を作るため隣接するカウンティを合併した。合衆国の3,000のカウンティから大體2,000個の一次抽出単位が作られた。カウンティの組合せ、一次層化、地域的副次層化、および大きさの尺度は1940年の *Decennial Census* のデータと利用することのできた最近のデータをもとにして決定された⁽⁴⁾。

- (5) 用いた層化基準の概要を含む計画された改訂標本の完全な記述については(11)をみよ。この論文は指定された副次抽出方式の簡単な説明としても有用である。

この論文で提示した種々の原則の適用は、1940年と更に最近のデータにもとづいて層化された標本から1930年センサスにおける労働力の標識 (characteristics) を推定することによって評価された。これは1930年から1940年にいたる10年間に起った大層な変化のため幾つかの方法に特に厳格な検討になっている。

この節で大要を述べた分析は主として三つの抽出原則、すなわち、

- (1) 拡大された一次単位 (Sec. IV-1 をみよ)
- (2) 大きさの尺度に比例した確率を用いる一次単位の抽出 (Sec. IV-2 をみよ)
- (3) 地域的副次層化 (Sec. IV-3 をみよ)

の導入に好適な状況のもとで獲得しうる利益を取扱うものである。

これとともに、他の標本推定公式を用いたときの効果を示すため幾つかの比較をおこなった。計算は最近労働月報に含まれているところの主要項目について行なわれた。すなわち男子および女子労働者の総数、農業労働者の男子と女子の総数、および非農業労働者の男子と女子の総数である。他の色々な抽出方式との比較は一次層化の基準と標本内に抽出する人員の期待数の両方を一定にしておいて行なった。

下に与えた利得のパーセントは平均平方誤差に対する一次単位同の寄与 (これには偏りの寄与を含む) の減少を示したものである。⁽⁴⁾ 特に断らない限り使用した標本推定値は (A) によるものである。

1 拡大した一次単位を導入することによってえられる利得

拡大された一次単位の使用による利得は、個々のカウンティを一次単位とする抽出計画の平均平方誤差と、カウンティの

(6) 一次単位内分散の全平均平方誤差に対する寄与はすべての例において比較的小であり、また事実上色々な原則を導入しても変化しない。

合せを一次単位とする計画の平均平方誤差を比較して計算される。何れの計画においても一次単位は等確率で抽出し、地域的副次層化は用いなかった。この比較では、限られた数の層および以上に与えた労働力の二つの項目、すなわち男子および女子労働者の総数についてのみ予備計算を完了した。拡大した一次単位を導入することによってえられた抽出誤差の減少は男女労働者総数に対しては 48%、女子労働者総数については 26% と推定された。

2 大きさの尺度に比例した確率を導入することによってえられる利得

確率を大きさの尺度に比例させる抽出原理の採用によってえられる利益は、単位を等確率で抽出する計画でえられた平均平方誤差と、その大きさに比例した確率で単位を抽出する計画でのものとを比較することによって計算される。どちらの計画でも一次単位はカウンティの組合わせて、いずれも地域的副次抽出は用いなかった。推定された利得の%は次の通りである。

全労働者		農業労働者		非農業労働者	
男	女	男	女	男	女
50	8	27	6	19	21

この利益は抽出分散の減少と、抽出単位が等確率で抽出される場合に生ずる偏りが除去されることの両方を反映している。⁽⁷⁾

3 地域的副次層化を導入することによってえられる利益

地域的副次層化の原理を用いることによってえられる利益は、

(7) 前に述べたように上記比較を行なった二つの計画では推定値 (A) が用いられた。この推定値は一次単位をその大きさに比例した確率で抽出する計画では不偏であるが等確率で抽出する計画では偏りをもつ。しかし標本の計画 (A) は通例最良標型不偏推定値よりも偏りのある推定値の方がより有効である。6つの労働力の項目の場合、最良標型不偏推定値は偏りのある推定値の平均平方誤差の数倍の大きさの分散をもつ。

地域的副次層化を用いない計画でえられる平均平方誤差と、地域的副次層化を導入した計画のそれと比較すれば計算できる。これらのいずれの計画においても一次単位はカウンティの組合せであり、これらはその大きさの尺度に比例した確率で抽出された。推定された利益の%は次の通りである。

全労働者		農業労働者		非農業労働者	
男	女	男	女	男	女
6	31	46	51	32	22

4 上記の原理を単一の副次抽方式(指定された副次抽方式)にすべて組入れた場合にえられる利益

三つの原理をすべて用いることによってえられる利益は、指定された副次抽方式(これには三つの原理がすべて用いられている)に対する平均平方誤差と、これらの原理を全く用いない抽方式の平均平方誤差を比較することによって計算される。指定された副次抽方式においては、一次抽出単位はカウンティの組合せであり、これはその大きさの尺度に比例した確率をもって抽出され、また地域的副次層化が用いられた。もう一つの抽方式では、一次単位は個々のカウンティで、抽出は等確率でありまた地域的副次層化は用いていない。この比較のための予備的な計算は6つの労働力の項目中の二つ、すなわち男子労働者の総数、女子労働者の総数のみについて行なうことができた。推定された利益は、男子労働者については76%、女子労働者については53%であった。

最後の二つの原理、すなわち大きさに比例した確率での抽出と、地域的副次層化を同時に用いてえられる利益を知るための計算は6項目のすべてについて適用できる。これらの利益の大きさを知る方法は、両方の計画の一次単位がカウンティであることを除いては上述のものと同じである。推定された利益の%は次の通りである。

全労働者		農業労働者		非農業労働者	
男	女	男	女	男	女
54	37	88	54	45	39

指定された副次抽方式と、いま比較を行なった一方の抽方式はともに偏りをもつ計画であるが、しかし指定された方式での偏りは明らかに後者の偏りより小さい。たとえば全男子労働者数の推定に指定された抽方式を用いたときの偏りは、真の全男子労働者数の1/2%より小さかったのに対して、同じ母集団特性の推定にもう一つの計画を用いたときの偏りは1/2%以上であった。

5 指定された抽方式で用いる推定値の選択

指定された抽方式に対して与えられた単純推定値(ア)は、同様の手法によって改善することができ(see II, 参照)、しかしそのような手法を用いるには多くの計算が必要だから、実際には利用できないことが多い。しかし Sec. V の最後に述べたように、もし全母集団の知識のようなある独立な情報を利用できれば、(2)の形の簡単な比の推定値を用いることによって(ア)以上の利益の望みかれることが多い。一次抽出単位の大きさの尺度と実際の大きさとの相関がそう高くなく、これとともに実際の大きさを推定している特性の値と高度の相関をもつときは比の推定値の利用が特に望ましいものとなろう。1930年の全男子労働者数を推定するために示された労働力の項目について行なった小規模の抽出試験では、比の推定値(2)に対する一次単位間の分散、一次単位内の分散は何れも単純推定値(ア)のもの約1/2であった。労働力の残りの5つの特性の推定では、比推定を用いても僅かの効果しかえられなかった。男子の全雇傭人口の分散が減少したのは、1930年以降の報告が1930年と1940年の大きさ(size)の間の相関を強め、更に男子の労働者数が全母集団と高度の相関をもっていたことから生じたものである。他の5項目の分散では、他の項目に対する実際の

大きさと相関がそれ程高くなかったから、このような減少はみられなかった。

6 最後にあたっての注意

いまえられた利益は、上に列挙した抽出原理を適用することによってえられたものである。これらの原理を適用できるような状況は好ましいものであるが、しかしこれは実際に屢々実現するものである。この原理の効率は、研究している母集団のもつ特定の属性によって左右される。拡大された一単位が小さい単位より更に内部的に異質なときは、常に拡大単位を用いることが望ましい。Sec. II で述べた一般的な母集団で単位の大きさが相当大きく変動しているときは、一次単位の抽出を大きさに比例した確率で行なうことが望ましい。地域的副次層化の利用は大きい一次抽出単位の用いられる抽出の場合に限られる。副次抽出を用い、一次単位は大きいが、大きさが変動しており、また標本内に含まれる一次単位の数が費用や管理上の条件で制限される場合には、上の三つの原理を同時に適用することによって最も大きい利益がえられる。Sec. III で説明したタイプの推定値は、この論文で指摘した以外の多数の具体的条件のもとで有効であろう。

謝 辞

多くの有益な注意を与えられ、また Sec. III に要約した分析の計画と指導に当られた Harold Nisselson 氏に深く感謝するものである。また我々は原稿を通読され、多くの有益な注意を与えられた W. G. Cochran, W. Edwards Deming, L. R. Frankel, William G. Madow, および Frederick F. Stephan 氏に厚くお礼を申し上げる。この研究は J. C. Capt. Director and Philip M. Hauser, Assistant Director の一般的な指示のもとに、種々の標本調査の展開に関して Bureau of the Census (国勢調査局) において実行されたものである。

付 録

1 一次単位の合併と抽出分散に与える影響標本平均を

$$\bar{X}_i = \frac{\sum_j \sum_k X_{jki}}{g \cdot n}$$

とおく、ただし一次単位はもとの単位で X_{jki} は j 一次単位の i 要素の値、また g は標本内の一次単位の数、 n は g 個の一次単位の各々から抽出された要素の数である。

\bar{X}_i の分散は、

$$(14) \sigma_{\bar{X}_i}^2 = \frac{N-n}{(N-1)ng} \sigma_{iw}^2 + \frac{Q-g}{(Q-1)g} \sigma_{ie}^2$$

である。ここで Q は母集団におけるもとの一次単位の数、 N はもとの一次単位内の要素の数で、 $\sigma_{iw}^2 = \sum \sum (X_{jki} - \bar{X}_i)^2 / QN$ は $\bar{X}_i = \sum_k X_{jki} / N$ なるもとの一次単位内の分散である。また $\sigma_{ie}^2 = \sum (\bar{X}_j - \bar{X})^2 / Q$ は $\bar{X} = \sum \bar{X}_j / Q$ なるもとの一次単位間の分散である。

$$(15) \sigma^2 = \sum \sum (X_{jki} - \bar{X})^2 / QN = \sigma_{iw}^2 + \sigma_{ie}^2$$

であるから、

$$(16) \sigma_{ie}^2 = \sigma^2 (1 + \rho_1 (N-1) / N)$$

ここで $\rho_1 = [\sigma_{ie}^2 - \frac{\sigma_{iw}^2}{N-1}] \frac{1}{\sigma^2}$ はもとの単位内の要素間の概内相関^(*)

である。

(15) と (16) から

$$(17) \sigma_{iw}^2 = \frac{N-1}{N} \sigma^2 (1 - \rho_1)$$

ゆえに

$$(18) \sigma_{\bar{X}_i}^2 = \frac{N-n}{n} \cdot \frac{\sigma^2}{ng} (1 - \rho_1) + \frac{Q-g}{(Q-1)g} \cdot \frac{\sigma^2}{N} [1 + \rho_1 (N-1)]$$

同様に \bar{X}_i の分散は

$$(19) \sigma_{\bar{X}_i}^2 = \frac{CN-n}{CN} \frac{\sigma^2}{ng} (1 - \rho_2) + \frac{Q-gC}{(Q-C)g} \frac{\sigma^2}{CN} [1 + \rho_2 (CN-1)]$$

(*) 概内相関の定義および性質は R. A. Fisher の研究者のための統計的方法、および (5) をみよ。

ここで \bar{x}_2 は拡大された一次単位に対する平均値で、 f_2 は拡大された一次単位内の要素同の級内相関で、 C は各拡大単位を作るために配合されたもとの単位数である。そうすると、

$$(20) \sigma_{\bar{x}_1}^2 - \sigma_{\bar{x}_2}^2 = \frac{\sigma^2}{gN} \left\{ \frac{(g-1)(C-1)}{(Q-1)(Q-C)} + f_1 a_1 - f_2 a_2 \right\}$$

ただし

$$a_1 = \frac{(Q-g)(N-1)}{Q-1} - \frac{N-n}{n}$$

および

$$a_2 = \frac{(Q-Cg)(CN-1)}{(Q-C)C} - \frac{CN-n}{Cn}$$

$$a_1 - a_2 = \frac{(C-1)(g-1)(QN-1)}{(Q-1)(Q-C)} \geq 0$$

かつ

$$\frac{(g-1)(C-1)}{(Q-1)(Q-C)} \geq 0$$

だから $f_1 > f_2$ なる限り拡大された一次単位を用いることにより利益がえられる。

ただし f_1 および f_2 はいずれも正である。

2 一次単位の大きさが等しくない場合の他の副次抽出方式の分散の比較

sec IV-2 で比較を行なった標本推定値の分散同の差の公式 (19) の誘導を述べる。

我々は問題を限定して、標本に各層からただ一つの抽出単位が選ばれているという簡単な場合を考えることにする。

$$(21) \bar{x}' = \sum N_k \bar{X}_k / N$$

は比較すべき三つの計画の各々について用いられる標本推定値とする。ここで

$$\bar{x}'_k = \bar{X}_k = \sum_{j=1}^{n_k} X_{kj} / n_k$$

であり、 X_{kj} は k 層の j 一次単位の k 要素の値である。し

は層の数、 n_{kj} は k 層の j 一次単位から標本内に抽出される要素数で N_{kj} はこれに対応する全要素数である。

$$N_k = \sum_j N_{kj} \text{ で } Q_k \text{ は } k \text{ 層内の一次単位の数、また } N = \sum_k N_k$$

である。層内の副次抽出が前に述べた最初の副次抽出方式のように一定の割合 C なら、上の推定値の n_{kj} は $C N_{kj}$ に等しい。もし推奨された抽出方式および前に述べた k 層の抽出方式の場合のように、層内の副次抽出が一定数であれば、 n_{kj} は

$$n_{kj} = C \sum_j N_{kj} / Q_k = C \bar{N}_k$$

に等しい。

我々は k 層の計画に対する標本推定値を \bar{x}'_k で、また k 層の計画を \bar{x}'_2 で、また推奨された計画のものを \bar{x}'_1 と書くことにする。最初の二つの計画の標本推定値の期待値 $E \bar{x}'_1$ と $E \bar{x}'_2$ は等しく、

$$E \bar{x}'_1 = E \bar{x}'_2 = \bar{x} = \frac{1}{N} \sum_k \frac{N_k}{Q_k} \sum_j \frac{n_{kj}}{N_{kj}} \sum_{i=1}^{N_{kj}} \frac{X_{kij}}{n_{kj}} \\ = \frac{1}{N} \sum_k \frac{N_k}{Q_k} \sum_j \bar{X}_{kj}$$

である。ここで $\bar{X}_{kj} = \sum_i X_{kij} / N_{kj}$ したがって一般に \bar{x} は $\sum_{k,j} X_{kij} / \sum_{k,j} N_{kj} = \bar{x}$ だから \bar{x}'_1 および \bar{x}'_2 は \bar{x} の偏りのある推定値である。

一次単位がその大きさに比例した確率で抽出され、また層内から一定個数の単位をとる推奨された計画では、標本推定値の期待値は、

$$(22) E \bar{x}'_1 = \frac{1}{N} \sum_k \sum_j N_k \frac{N_{kj}}{N_k} \frac{n_{kj}}{N_{kj}} \frac{X_{kj}}{n_{kj}} = \frac{1}{N} \sum_k \sum_j X_{kj} = \bar{x}$$

だから、推奨された計画に対する推定値は不偏である。

\bar{x}'_1 の平均平方誤差は、

$$(23) \sigma_{\bar{x}'_1}^2 = \frac{1}{N^2} \sum_k \frac{N_k^2}{Q_k} \left[\sum_j \frac{N_{kj} - n_{kj}}{(N_{kj} - 1) n_{kj}} \sigma_{kj}^2 + \sum_j (\bar{X}_{kj} - \bar{x}_k)^2 \right] \\ + (\bar{x} - \bar{x})^2 - \frac{1}{N^2} \sum_k N_k^2 (\bar{x}_k - \bar{x})^2$$

である。ここで $\sigma_{h_j}^2 = \sum_k (X_{h_j k} - \bar{X}_{h_j})^2 / N_{h_j}$ は才A層の、才j-次抽出単位内の要素間分散である。

(22) の大括弧内の才/項は一次単位内分散の寄与であり、才2項は近似的に一次単位間の平均平方誤差を表わし、残りの項のこの近似の誤差を与える。 \bar{X}_2 の平均平方誤差もこれと同じ式で表わされるが、ただそのときには n_{h_j} を n_h で置きかえねばならない。

$\sigma_{\bar{X}_1}^2$ と $\sigma_{\bar{X}_2}^2$ の差は、

$$(24) \quad \sigma_{\bar{X}_1}^2 - \sigma_{\bar{X}_2}^2 = \frac{1}{CN^2} \sum_k \frac{N_k^2}{Q_k} \sum_j \sigma_{h_j}^2 \frac{N_{h_j}}{N_{h_j} - 1} \left(\frac{1}{N_{h_j}} - \frac{1}{N_h} \right)$$

である。

これは実際の場合には殆んどいつもそうであるように、 $\sigma_{h_j}^2 / N_{h_j}$ が N_{h_j} と負の相関をもつときには正となる (Sec. III をみよ)。したがって普通 $\sigma_{\bar{X}_1}^2$ は $\sigma_{\bar{X}_2}^2$ より大きいから、推奨された抽出方式が上述の最初の二つの計画のいずれより有効 (efficient) であることを示すには $\sigma_{\bar{X}_1}^2$ と $\sigma_{\bar{X}_2}^2$ を比較すれば十分である。

推奨された計画の分散は

$$(25) \quad \sigma_{\bar{X}_2}^2 = \frac{1}{N^2} \sum_k N_k^2 \left[\sum_j \frac{N_{h_j}}{N_h} \frac{N_{h_j} - \bar{n}_h}{N_{h_j} - 1} \frac{\sigma_{h_j}^2}{\bar{n}_h} + \sum_j \frac{N_{h_j}}{N_h} (\bar{X}_{h_j} - \bar{X}_h)^2 \right]$$

である。 \bar{X}_2 の平均平方誤差を \bar{X}_2 の分散と比較するため、才j-一次単位内の要素間の級内相関を

$$\rho_{h_j} = \frac{1}{\sigma_h^2} \left[(\bar{X}_{h_j} - \bar{X}_h)^2 - \frac{\sigma_{h_j}^2}{N_{h_j} - 1} \right]$$

と定義する。ここで σ_h^2 は才A層内のすべての要素間の分散である。(25) の大括弧の外は平均平方誤差への寄与が正か、または無視できる位の大きさなのでこの比較では考慮に入れなかった。このとき、

$$(26) \quad \sigma_{\bar{X}_1}^2 - \sigma_{\bar{X}_2}^2 = \frac{1}{N^2} \sum_k \frac{N_k^2}{Q_k} \left\{ \sum_j \frac{N_{h_j}}{N_{h_j} - 1} \frac{\sigma_{h_j}^2}{\bar{n}_h} \left(1 - \frac{N_{h_j}}{N_h} \right) + \sigma_h^2 \sum_j \rho_{h_j} \left(1 - \frac{N_{h_j}}{N_h} \right) \right\}$$

この差の才2項は Sec. IV-2 で、才/項を無視したときの近似的な差として与えたものである。この項の相対的な大きさを調べるため、

$$(27) \quad \frac{N_{h_j}}{N_{h_j} - 1} \sigma_{h_j}^2 = \sigma_h^2 (1 - \delta_{h_j})$$

とおく。このとき、

$$(28) \quad \sigma_{\bar{X}_1}^2 - \sigma_{\bar{X}_2}^2 = \frac{1}{N^2} \sum_k \frac{N_k^2}{Q_k} \sigma_h^2 \left\{ \frac{1}{\bar{n}_h} \sum_j \delta_{h_j} \left(\frac{N_{h_j}}{N_h} - 1 \right) - \sum_j \rho_{h_j} \left(\frac{N_{h_j}}{N_h} - 1 \right) \right\}$$

Sec. II で示した一般的な母集団のクラスでは δ_{h_j} と N_{h_j} との間の共分散および ρ_{h_j} と N_{h_j} との共分散も負になる。更にこのクラスに属する母集団での多くの実際の問題では、この二つの共分散の大きさは近似的に等しくなる。そのような場合、(22) の才/項は才2項の $\frac{1}{\bar{n}_h}$ 倍となるから、 $\bar{n}_h > 1$ なるときは常に才2倍より小となり、また十分大きい \bar{n}_h の値については更に小となる。たとえば、一定の大きさの集落で特定の性質をもつ要素の条件付確率がすべての大きさの集落について等しい場合、このような種々の大きさの集落からなる母集団ではこの二つの共分散の大きさは非常に近い値となる。実際問題は近似的にこのような状況にあることが多い。更に δ_{h_j} と N_{h_j} の共分散が ρ_{h_j} と N_{h_j} とのものの数倍たとえば5倍になっている場合でも、才2項は $\bar{n}_h > 5$ なる値では常に才/項より小さくなる。

推奨された抽出方式の使用によつてえられる利益を Sec. III で幾つかあげておいた。そこで結果を示した若干の項目については、この利益はかなり大きいものであった。

3 分散公式 (17) および (19) の誘導

確率変数の比の平均平方誤差は一般に Taylor 展開で近似される。

X' と Y' が確率変数で Y' > 0、また X'/Y' = r' によつて推定され

る母集団特性を r とする。

このとき

$$(29) E\left\{\frac{X'}{Y} - r\right\}^2 = E\frac{Y^2}{(EY')^2} \left(\frac{X'}{Y} - r\right)^2 + E\left(1 - \frac{Y'}{EY'}\right)^2 \left(\frac{X'}{Y} - r\right)^2$$

が成立する。(29)の右辺の第1項は Taylor 展開でえられる平均平方誤差の第1近似で第2項はこの近似の誤差項である。

等式(29)および特別な場合として(28)は次のようにして導かれる。

$$(30) E(r' - r)^2 = E\left\{\frac{\sum_k \frac{1}{t_k} \sum_i \sum_j \sum_k^{m_{ij}} X_{kij} - r}{\sum_k \frac{1}{t_k} \sum_i \sum_j \sum_k^{m_{ij}} Y_{kij}} - r\right\}^2$$

$$\psi_{kij} = Y_{kij} (r_{kij} - r) \quad \text{および}$$

$$Y' = \sum_k \frac{1}{t_k} \sum_i \sum_j \sum_k^{m_{ij}} Y_{kij}$$

とする。このとき

$$(31) \theta = \frac{\sum_k \frac{1}{t_k} \sum_i \sum_j \sum_k^{m_{ij}} \psi_{kij}}{E\left(\sum_k \frac{1}{t_k} \sum_i \sum_j \sum_k^{m_{ij}} Y_{kij}\right)} = \frac{Y'}{EY'} \left(\frac{X'}{Y} - r\right)$$

とおくと

$$E\theta^2 = EY'^2 (r' - r)^2 / (EY')^2$$

は平均平方誤差の第1近似である。 EY' は $EX'(r)$ と同じ方法で計算できるから、 $E\theta^2$ の分子 $EY'^2 (r' - r)^2$ を求めればよい。いま

$$EY'^2 (r' - r)^2 = E\left[\sum_k \frac{1}{t_k} \sum_i \sum_j \sum_k^{m_{ij}} \psi_{kij}\right]^2 = E\sum_k \frac{1}{t_k^2} \psi_k^2 + E\sum_{\substack{k, l \\ k \neq l}} \frac{\psi_k}{t_k} \frac{\psi_l}{t_l}$$

である。ここで

$$\psi_k = \sum_i \sum_j \sum_k^{m_{ij}} \psi_{kij} = \sum_i \psi_{ki}$$

$$E\sum_k \frac{1}{t_k^2} \psi_k^2 = E\sum_k \sum_i \psi_{ki}^2 / t_k^2 + E\sum_k \sum_{\substack{l, r \\ l \neq r}} \psi_{kl} \psi_{kr} / t_k^2$$

であるから、

$$(32) EY'^2 (r' - r)^2 = E\sum_k \frac{1}{t_k^2} \sum_i \psi_{ki}^2 + E\sum_k \frac{1}{t_k^2} \sum_{\substack{l, r \\ l \neq r}} \psi_{kl} \psi_{kr} + E\sum_{\substack{k, l \\ k \neq l}} \frac{\psi_k}{t_k} \frac{\psi_l}{t_l}$$

である。

(32)の右辺の第1項は、

$$(33) E\sum_k \frac{1}{t_k^2} \sum_i \psi_{ki}^2 = \sum_{k, i, j} \frac{1}{t_k^2} \frac{P_{ij}}{P_k} \frac{m_{kij}}{M_{kij}} \frac{M_{kij} - m_{kij}}{M_{kij} - 1} \sum_k \psi_{kij}^2 + \sum_{k, i, j} \frac{1}{t_k^2} \frac{P_{ij}}{P_k} \frac{m_{kij}}{M_{kij}} \frac{m_{kij} - 1}{M_{kij} - 1} \left(\sum_k \psi_{kij}\right)^2$$

(32)の第2項は、

$$(34) E\sum_k \frac{1}{t_k^2} \sum_{\substack{l, r \\ l \neq r}} \psi_{kl} \psi_{kr} = \sum_{k, l} \frac{1}{t_k^2} \left(\sum_i \frac{m_{kij}}{M_{kij}} \psi_{kij}\right)^2 - \sum_{k, l} \frac{1}{t_k^2} \frac{P_{ij}}{P_k} \sum_i \frac{m_{kij}}{M_{kij}} \psi_{kij}^2$$

である。ただし $\psi_{kij} = \sum_k \psi_{kij}$ である。また(32)の第3項は

$$(35) E\sum_{\substack{k, l \\ k \neq l}} \frac{\psi_k}{t_k} \frac{\psi_l}{t_l} = \left[\sum_{k, l} \frac{1}{t_k} \frac{P_{ij}}{P_k} \frac{m_{kij}}{M_{kij}} \psi_{kij}\right]^2 - \sum_k \frac{1}{t_k^2} \left[\sum_i \frac{P_{ij}}{P_k} \frac{m_{kij}}{M_{kij}} \psi_{kij}\right]^2$$

である。したがって

$$EY'^2 (r' - r)^2 = (33) + (34) + (35)$$

であるから、 ψ_{kij} の代わりに $Y_{kij} (r_{kij} - r)$ を代入すると

$$(36) EY'^2 (r' - r)^2 = \sum_k \frac{1}{t_k^2} \left[\sum_{i, j} \frac{P_{ij}}{P_k} \frac{m_{kij}}{M_{kij}} \frac{M_{kij} - m_{kij}}{M_{kij} - 1} \sum_k Y_{kij}^2 (r_{kij} - r)^2 + \sum_{i, j} \frac{P_{ij}}{P_k} \frac{m_{kij}}{M_{kij}} \frac{m_{kij} - 1}{M_{kij} - 1} \left[\sum_k Y_{kij} (r_{kij} - r)\right]^2 + \sum_j \frac{P_{ij}}{P_k} \left[\sum_i \frac{m_{kij}}{M_{kij}} Y_{kij} (r_{kij} - r)\right]^2 - \sum_{i, j} \frac{P_{ij}}{P_k} \frac{m_{kij}^2}{M_{kij}^2}\right]$$

$$-Y_{kij}^2 (Y_{kij} - r)^2 - \left\{ \sum_{ij} \frac{P_{ij}}{P_A} \frac{M_{kij}}{M_{kij}} (Y_{kij} - r) Y_{kij} \right\}^2$$

$$+ \left[\sum_{kij} \frac{1}{t_A} \frac{P_{ij}}{P_A} \frac{M_{kij}}{M_{kij}} Y_{kij} (Y_{kij} - r) \right]^2$$

となる。

(36)の才1項の $(Y_{kij} - r)^2$ に $(Y_{kij} - Y_{kij} + Y_{kij} - r)^2$ を代入し、才1、才2、才4項の $1/t_A$ を $P_{ki} M_{kij} / P_{kij} M_{kij}$ で置きかえらば、これら三つの項は、

$$(37) \sum_{kij} \frac{P_{ij}}{P_A} \frac{M_{kij}}{M_{kij}} \frac{P_{kij}}{P_{kij}} F_{kij} Y_{kij}^2 (Y_{kij} - r)^2$$

$$+ 2 \sum_{kij} \frac{P_{ij}}{P_A} \frac{M_{kij}}{M_{kij}} \frac{P_{kij}}{P_{kij}} F_{kij} Y_{kij}^2 (Y_{kij} - r)(Y_{kij} - r)$$

$$+ \sum_{kij} \frac{P_{ij}}{P_A} \frac{M_{kij}}{M_{kij}} \frac{P_{kij}}{P_{kij}} F_{kij} (Y_{kij} - r)^2 \left[\sum_{kij} Y_{kij}^2 - \frac{Y_{kij}^2}{M_{kij}} \right]$$

ここで $F_{kij} = (M_{kij} - m_{kij}) / (M_{kij} - 1)$ と $r_{kij} = \sum_{kij} X_{kij} / \sum_{kij} Y_{kij}$ である。

(36)の才3、才5、才6項の $1/t_A$ に適当な値を代入すると、これらの項の和は、

$$(38) \sum_{ij} \frac{P_{ij}}{P_A} \left[\sum_i \frac{P_{ki}}{P_{kij}} Y_{kij} (Y_{kij} - r) \right]^2 - \sum_{ij} \left[\sum_{ij} \frac{P_{ij}}{P_A} \frac{P_{ki}}{P_{kij}} Y_{kij} (Y_{kij} - r) \right]^2$$

$$+ \left[\sum_{kij} \frac{P_{ij}}{P_A} \frac{P_{ki}}{P_{kij}} Y_{kij} (Y_{kij} - r) \right]^2$$

となる。

いま、

$$(39) \sum_i \frac{P_{ki}}{P_{kij}} Y_{kij} (Y_{kij} - r) = \sum_i P_{ki} \left(\frac{X_{kij}}{P_{kij}} - \frac{Y_{kij}}{P_{kij}} r \right)$$

$$= P_A (R_{kij} - r R_{kij} : Y)$$

$$= P_A R_{kij} : Y (\bar{Y}_{kij} - r)$$

ただし、 $\bar{Y}_{kij} = R_{kij} / R_{kij} : Y$

また

$$(40) \sum_{ij} \frac{P_{ij}}{P_A} \frac{P_{ki}}{P_{kij}} Y_{kij} (Y_{kij} - r) = \sum_i P_{ij} (R_{kij} - r R_{kij} : Y)$$

$$= P_A (R_{kij} - r R_{kij} : Y)$$

$$= P_A R_{kij} : Y (\bar{Y}_{kij} - r)$$

ただし、 $\bar{Y}_{kij} = R_{kij} / R_{kij} : Y$ である。

(38)に(39)と(40)を代入すると、

$$(41) \sum_{ij} \left(\frac{P_{ij}}{P_A} \right) P_A^2 R_{kij} : Y (\bar{Y}_{kij} - r)^2 - \sum_{ij} P_A^2 R_{kij} : Y (\bar{Y}_{kij} - r)^2$$

$$+ \left[\sum_{ij} P_A R_{kij} : Y (\bar{Y}_{kij} - r) \right]^2$$

(41)の才1項の $(\bar{Y}_{kij} - r)^2$ の代りに $(\bar{Y}_{kij} - \bar{Y}_{kij} + \bar{Y}_{kij} - r)^2$ を代入して展開すると、(41)は、

$$(42) \sum_{ij} P_A^2 \frac{P_{ij}}{P_A} R_{kij} : Y (\bar{Y}_{kij} - \bar{Y}_{kij})^2$$

$$+ 2 \sum_{ij} P_A^2 \frac{P_{ij}}{P_A} R_{kij} : Y (\bar{Y}_{kij} - \bar{Y}_{kij})(\bar{Y}_{kij} - r)$$

$$+ \sum_{ij} P_A^2 (\bar{Y}_{kij} - r)^2 \left[\sum_{ij} \frac{P_{ij}}{P_A} R_{kij} : Y - R_{kij} : Y \right]$$

$$+ \left[\sum_{ij} P_A R_{kij} : Y (\bar{Y}_{kij} - r) \right]^2$$

となる。

ゆえに $(EY')^2 E\theta^2 = (39) + (42)$ となる。

$$(43) (EY')^2 E\theta^2 = \sum_{ij} F_{kij}^2 \frac{P_{ij}}{P_A} \frac{M_{kij} - m_{kij}}{M_{kij} - 1} \frac{\sum_{kij} Y_{kij}^2 (Y_{kij} - r)^2}{M_{kij} M_{kij} P_{kij}^2}$$

$$+ 2 \sum_{ij} P_{kij}^2 \frac{P_{ij}}{P_A} \frac{M_{kij} - m_{kij}}{M_{kij} - 1} \frac{\sum_{kij} Y_{kij}^2 (Y_{kij} - r)(Y_{kij} - r)}{M_{kij} M_{kij} P_{kij}^2}$$

$$+ \sum_{ij} P_{kij}^2 \frac{P_{ij}}{P_A} \frac{M_{kij} - m_{kij}}{M_{kij} - 1} \sigma_{kij}^2 \frac{(Y_{kij} - r)^2}{M_{kij} P_{kij}^2}$$

$$+ \sum_{ij} P_A^2 \left(\frac{P_{ij}}{P_A} \right) R_{kij} : Y (\bar{Y}_{kij} - \bar{Y}_{kij})^2$$

$$+ 2 \sum_{ij} P_A^2 \left(\frac{P_{ij}}{P_A} \right) R_{kij} : Y (\bar{Y}_{kij} - \bar{Y}_{kij})(\bar{Y}_{kij} - r)$$

$$+ \sum_{ij} P_A^2 (\bar{Y}_{kij} - r)^2 \left(\frac{P_{ij}}{P_A} \right) (R_{kij} : Y - R_{kij} : Y)^2$$

$$+ \sum_k P_k R_{k(A):Y} (\bar{Y}_{k(A)} - r)^2$$

ここで $\sigma_{k(A):Y}^2 = \frac{M_{kij}}{n} (Y_{kij} - \bar{Y}_{kij})^2 / M_{kij}$ かつ

$$\bar{Y}_{kij} = \frac{M_{kij}}{n} Y_{kij} / M_{kij} = Y_{kij} / M_{kij}$$

$E(r' - r)$ に対する近似式は (43) を $(EY')^2$ で割ってえられる。才2, 才5, 才7項は殆んど母集団で無視できる大きさだから、これを無視すると、(13) がえられる。

X' の分散は (43) で単に Y_{kij} のところに \bar{P}_{kij}/P を代入すればえられる。これは下に述べる考察から導かれるものである。

$r' = X'/Y'$ で、 X' は r' の分子だから、分母 Y' が、繰返し抽出 (repeated sampling) で恒等的に $1/n$ に等しい場合、 $\sigma_{X'}^2$ は $\sigma_{Y'}^2$ で与えられる。

(5) から

$$\frac{1}{t_k} = \frac{M_{kij} P_{kij}}{m_{kij} P_{kij}} = \frac{P_{kij}}{m_{kij} P_{kij}}$$

だから $\frac{1}{n} \sum_k \frac{1}{t_k} \sum_{i,j} \frac{P_{kij}}{m_{kij} P_{kij}} Y_{kij}$ に等しい r' の分母は、 Y_{kij} が \bar{P}_{kij}/P に等しくとられている場合、繰返し抽出 (repeated sampling) において恒等的に $1/n$ に等しい。ただし、 $P = \sum_k P_k$ である。

X' の平均平方誤差に対する公式 (9) はもちろん誤差項が

$$E\{Y'^2 / (EY')^2\} \{r' - r\}^2 = 0$$

だから正確なものである。 X' は確率変数の比から推定されるものではないから、 $\sigma_{X'}^2$ は (29) を使わずに直接もっと簡単な方法で求めることができる。

(29) から $E(r' - r)^2$ に対する近似式の誤差項 (43) / $(EY')^2$ は $E\{1 - \frac{Y'^2}{(EY')^2}\} \times \{r' + r\}^2$ で与えられる。これは個々の観測値の簡単な函数としては表わせないが、しかし求める最大および最小値を与えるものとして有用である。

r' の分散の上限、下限を求める方法は、 X' と Y' の同時分布 (

joint distribution) とかかわりなく成立する次の不等式から容易にえられる。

$$(44) \quad \frac{EY'^2}{Y_{max}^2} (r' - r)^2 \leq E(r' - r)^2 \leq \frac{EY'^2}{Y_{min}^2} (r' - r)^2$$

ここで Y_{max} は簡単に各層の最大の Y_k を選ぶか、あるいは推定するかしてえられる。 Y_{min} (Y' の最小値) についても同様である。(44) を計算すると、

$$(45) \quad \frac{(EY')^2 E\theta^2}{Y_{max}^2} \leq E(r' - r)^2 \leq \frac{(EY')^2 E\theta^2}{Y_{min}^2}$$

となる。ここで $(EY')^2 E\theta^2$ は (43) で与えられる。(45) は層内の Y の変動性が限られているような抽出方式では、 $E\theta^2$ の正確さの測定すべき指標となるであろう。しかし層化を用いず、 Y の変動性が制限されていないような他の計画では、(45) の与える限界ではなすぎた役にたたないであろう。

M. H. Hansen and W. N. Hurwitz

(A.M.S. Vol. 14, 1943 pp 333~362)

2 系統的抽出の理論について (I)

On the theory of Systematic Sampling

1 緒言

標本設計理論の必要性は説明するまでもない。政府ならびに私企業のいづれにおいても、政策ないし行動 (*operating*) の決定は標本の知識に頼っている。政府や企業においては益々標本調査の理論を利用する傾向がある⁽¹⁾。

不幸なことに標本調査理論と実際との間にはまだかなりの隔りがある。これらの原因は、一方では標本調査理論の与えうる実際の寄与に關係のある行政官 (管理者) の無知があり、また一方では有用なサンプリング設計の価値を認めさせるような抽出理論が欠けているからである。

今迄にもまた現在でも、理論と実際との一致を計る多くの努力が払われている⁽²⁾。管理者と抽出者 (*Sampler*) は互いにうまく訓練し合っている。しかし経験によってその理論が発展すれば今までに理論の組立てられているどのデザインよりも優れていることが証明されると思われるにもかかわらず、理論の解明されていないデザインが依然として存在するのである。おそらく今日における抽出理論の主要な欠陥は、 N 個の要素の母集団から、 n 要素からなる完全に無作為な標本をとった方が良いか、あるいは、 i 要素から始めて、 $i, i+k, \dots, i+(n-1)k$ なる要素を含む系統的標本 (出発点 i は無作為に選び、 N は大体 k になるよう

(1) 標本抽出理論の知識を有する統計家の必要性は現在幾つかの大学において与えられているサンプリングの課程を生み出すことになった。

(2) 応用の分野での研究だけでなく、標本抽出技術について助言を与えるということをも含めた責任の、近年における立場の變化を力えるだけで十分である。

にする)をとる方がよいかをきめる統計的手段のないことである⁽³⁾。理論のない系統的抽出法と系統的抽出より悪い結果の出やすい無作為抽出法との二者択一を迫られていることによって統計家のおらいつていゝジレンマは系統的標本をとるか無作為標本をとるかというこの問題に關係するものである。

この論文の目的は系統的抽出の満足な理論を与えることによつてこの不一致を解決しようとするものである。次節において系統的抽出理論の我々の研究の最初の部分を提出する。この研究が、単一要素の抽出理論と、要素の集落 (*cluster*) の抽出理論の双方にわたるものであるとしても、この論文での単位は要素の集落でなく単一の要素でなければならぬ。後者の問題はあとの論文で取扱う。我々は層化しない母集団と層化した母集団の双方からえられた系統的抽出の理論を与えることにする。推定値の平均値および分散に対する公式を導いた。無作為抽出計画と層化無作為抽出計画との比較を行なった。更に分散の推定値と標本の "最適" な大きさおよび配分 (*allocation*) を導いた。分析の基礎的な部分は、母集団の分散や、系列相関あるいは *serial variance* の知識から系統的標本にもとづく推定値の分散が推定できるというこ

(3) この論文では $N = kn$ と仮定する。この仮定をはずしても一般性という点ではあまり損るところはないが、その場合にはかなり細かい論議が必要となるであろう。 N が丁度 kn でない場合 k の出発点がすべて等確率で選ばれらるような系統的抽出の手続きは偏りをもっている。しかしこの偏りは普通問題にする必要のない程小さい。もし N が既知なら系統的標本の可能な大きさに比例させた抽出を行なうことによって、この偏りを除去できることが知られている。

(4) 我々が系統的抽出手続を定義するとき、ある系統的抽出手続は N 個の中から n 個をとる C_N^n 個の選び方の中の多数を除外するような無作為抽出手続である。

との説明にある⁽⁵⁾ 基本的な結果は、

- a. 系列相関が正の和 (positive sum) なら系統的標本は無作為標本より悪い。
- b. 系列相関の和が近似的に0なら系統的標本は大体無作為標本と同等で
- c. 系列相関が負の和 (negative sum) のときには系統的標本は無作為標本より良い

2 有限母集団の利用

この論文では期待値の計算の際、極限分布が使えらる程母集団が大きい場合でも有限母集団から抽出を行なったと仮定する。これは数学的には屡々述べた問題である。同様な結果は正しく定義された多次元正規分布を仮定し、条件付確率を用いても行うことができる。しかし現実的な観点からは、有限母集団の利用を導く幾つかの因子が存在する。我々が最も普通に抽出を行なうのは、変換の法則 (law of transformation) が知られていないか、あるいは数学的に表示されないような実在する (existing) 母集団である⁽⁶⁾。結局我々が抽出を行なう正規とか、その他の特定の分布の概念および条件付確率の利用ということは、実際問題に關する我々の思考の対象外のものである。

一方我々が母集団は有限でなければならぬと考えて、有限母集団から乱数表を用いて標本を抽出するならば、我々は現実的な問題において漠然とした形の数学を用いているにすぎないことを見る。天に乱数を選ぶことによって繰返し可能な実験を行なうことができるが、これは統計的管理の状態にあるということを我々は

(5) 特に限定しないで“母集団の分散”というときは、母集団の N 要素からなる無作為標本の分散を意味するものとする。

(6) 言い換えては、我々の母集団は時間の上で統計的なコントロールを及べていない状態にある。

知っている。

普通の無作為抽出の理論では、色々な標本平均値を求めたための可能な標本の数は十分大きいから、十分大きい母集団および標本での標本平均値は、近似的に正規分布すると仮定できる。しかし系統的抽出においては、可能な標本平均値の数は普通非常に小さく、母集団や標本が大きい場合でもそれが正規分布すると仮定することは難かしい。結局系統的抽出の平均値および分散に対する我々の解釈においては、母集団要素は確率変数の単一観測値の結果であり、その分布は要素毎に変動するという考えをとらざるを得ない。したがって我々が次に行なう解釈は条件付確率についてのものであるが、もし母集団と標本が十分大きいなら、我々は可能な各標本平均値の算術平均は正規分布に従うという仮定を設けることができる。

適当な正規多次元母集団の仮定のもとにおける系統的抽出の理論は後に述べるであらう。

3 定義

N 個の要素 z_1, \dots, z_N からなる有限母集団から標本抽出を行なうと仮定する。標本設計とは、これらの N 個の要素を k (重なることもあり重ならないこともある) 個のクラスに分割する方法とする。各クラスが定められた確率で抽出されるようにして、 k 個の中から n 個のクラスを選び出す方法を組合せた手続を意味する。与えられた標本設計 (sample design) と組合せた抽出手続は、標本設計において定められた方法に従って、 k 個のクラスから n 個を選び出す作業である。標本は抽出手続によって与えられた特定のクラスである。

無作為抽出手続とは、標本設計で k 個のクラスが与えられたら、これらのどの n 個を抽出する確率も $1/k^n$ となるようにする抽出方法である。

無作為抽出手続と組合せられる任意の標本設計は無作為抽出設計

計 (random sampling design) である。用いられている無作為でない抽出手続の1つは、その中のクラスに大きさ (size) と数を結びつけ、与えられたクラスの標本となる確率がその大きさに等しくするような手続である。(9) 勿論その他の無作為でない抽出手続も用いられている。

N個の要素からk要素を選ぶ単純(無制限)無作為抽出は、Nからkを選ぶ可能な C_k^N 通りのクラスがあって、その各々の標本となる確率が $1/C_k^N$ になっているような抽出設計である。これと結合される抽出手続は、各クラスを番号し、 $i=1, \dots, C_k^N$ と同一視し、乱数表を用いて番号iを選ぶものである。同様に無作為抽出手続はN個の要素を $j=1, \dots, N$ の番号と同じと考えて、乱数表から番号jを選び次に乱数表から別の番号j'を選ぶ。このようにして1~Nの中からk個の数が繰返えしなしに乱数表から抽かれるまでこの手続を続けることでもある。これらの整数と結合された要素は無作為標本である。この二つの手続が同等であることは容易にわかる。無制限でない無作為抽出手続は制限付 (restricted) 抽出計画とよばれる。制限付無作為抽出計画には色々な種類があるが、その中で我々が系統的抽出とよぶものは唯一つしかない。

系統的抽出計画とはN個の要素をk個のクラス、 S_1, S_2, \dots, S_k に分割し、 S_i の1つを選ぶのに無作為抽出手続を用いる計画である。ここで S_i は $x_i, x_{i+k}, \dots, x_{i+(N-k)}$ である。したがって系統的抽出は一種の集落抽出計画 (cluster sampling design) であることが容易にわかる。すぐわかるように、系統的抽出計画において導入された集落抽出の新らしい側面は、級内相関係数の値および、標本の大きさが変化するときのこの値の変化を求めたため

(9) この問題の論議については "On the theory of sampling from finite populations" 及び題名の (M. H. Hansen and W. N. Hurwitz, A. M. S., Vol. 14, 1943) をみよ。

に母集団における要素の順序の知識が用いられるということである。

無作為抽出と無作為でない抽出手続の場合と同じように、抽出計画には無作為と系統的抽出計画の組合せたものがあることでもある。

これらの標本を抽出した母集団は層化されていることもあり層化されていないこともある。また抽出単位も単一の要素のときもあるし要素の集落 (cluster) であることもある。

4 抽出計画を選ぶ際の基礎

望ましい推定値をうるために構成できる多くの抽出計画の中から、人は管理上の(行政的な)考慮、費用および抽出誤差等をもととしてどれを用いるかを定めるであろう。抽出誤差の尺度としては極限分布の理論や最良線型不偏推定値の理論をもとにして、推定された特性に対する標本推定量の標準偏差を用いるのが普通になっている。この論文では以後このやり方を用いるのであるが、ここで注意すべきことは多くの抽出計画を類立てた場合、そのうちの幾つかについては極限分布の理論が成立しないので、標準誤差を用いること自体が、分析の結果以上の問題となるおそれがあるということである。この危険は、系統的抽出計画で生ずるものであるがこれについては後に研究することにする。

偏り (bias)、一貫性 (consistency) および効率 (efficiency) などは用いられた抽出計画および推定函数の性質であって、実際にはえられた特定の標本の性質でないということに注意しなければならない。標本にもとづくいかなる推定値も、おそらく推定された特性とは違ってくるであろう。このような食い違いの大きさなど位々ということを示すのが統計解析の役目である。

記号

記号 P に適当な添字をつけて母集団および層などの副母集団

を表わすことにする。

層の数を L と書き、また i 層内の要素の数を N_i で表わす。標本の大きさは n に適当な添字をつけて表わす。

抽出計画によって定義されたときの、母集団のある特別なクラスは S に添字をつけて示す。

6 層化しない系統的抽出 (抽出単位は一つの要素からなる)

この節で用いた添字のとら値の範囲については附録 A に示してある。

母集団 P は N 個の要素 x_1, \dots, x_N からなっているとす。

P の算術平均 \bar{x} の推定値を求めたい。

$N = kn$ とし、クラス S_i が n 個の要素 $x_i, x_{i+k}, \dots, x_{i+(n-1)k}$ からなるものとする。したがって、大きさ n の標本から \bar{x} を推定するための系統的抽出計画には n 個のクラス S_1, S_2, \dots, S_n があり、また S_i がこの抽出手続で選ばれる確率が $1/n$ とならなければならぬという要求がある。

\bar{x}_i は S_i の要素の算術平均値、すなわち $n\bar{x}_i = x_i + x_{i+k} + \dots + x_{i+(n-1)k}$ とし、また \bar{x} は標本平均値、すなわちもし S_i が抽出手続で抽出されるなら

$$\bar{x} = \bar{x}_i$$

である。

系統的抽出を取扱う場合、我々は屢々系列相関および附随する *serial variance* の定義として *circular* および非-*circular* な定義の両方を用いることにする。

我々は、 $k > kn$ なら $x_k = x_{k-kn}$ と仮定する。これは *circular* な定義に用いる。

$$kn \sigma^2 = \sum_{v=1}^n (x_v - \bar{x})^2$$

および

$$kn C_{k\mu} = \sum_{v=1}^n (x_v - \bar{x})(x_{v+k\mu} - \bar{x})$$

とする。

このとき系列相関係数 $\rho_{k\mu}$ の *circular* な定義は

$$\sigma^2 \rho_{k\mu} = C_{k\mu}$$

である。この式は σ^2 の公式を簡単にするため $\rho_{k\mu/2}$ を

$$2\sigma^2 \rho_{k\mu/2} = C_{k\mu}$$

で定義したときは n が奇数のときだけ用いる。同様に *serial variance* $s_{k\mu}$ を

$$kn s_{k\mu} = \sum_{v=1}^n (x_v - x_{v+k\mu})^2$$

で定義すれば、*serial variance* の *circular* な定義を用いたことになる。

中に *serial variance* の比 $v_{k\mu}$ の *circular* な定義は

$$\sigma^2 v_{k\mu} = s_{k\mu}$$

である。これは $v_{k\mu/2}$ を

$$2\sigma^2 v_{k\mu/2} = s_{k\mu}$$

で定義した場合は n が奇数のときだけに用いることにする。

系列相関および *serial variance* の非-*circular* な定義は、

$$(1) \quad k(n-k) C'_{k\delta} = \sum_{j=1}^n (x_j - \bar{x})(x_{j+k\delta} - \bar{x})$$

$$\sigma^2 \rho'_{k\delta} = C'_{k\delta}$$

$$k(n-k) s'_{k\delta} = \sum_{j=1}^n (x_j - x_{j+k\delta})^2$$

および

$$\sigma^2 v'_{k\delta} = s'_{k\delta}$$

を与えられる。

級内相関係数 \bar{r}_k は式

$$\sigma^2 \bar{r}_k = \bar{C} (x_u - \bar{x})(x_v - \bar{x})$$

で定義される。

この場合無作為化の手続 (*random process*) は、まず S_i を無作為に選び、次にこの選ばれた S_i から無作為に 2 つの x を選ぶことによって実行される。そうすると、

$$k \sigma_{\bar{x}}^2 = \sum_{i=1}^n (\bar{x}_i - \bar{x})^2$$

で、また

$$(2) \sigma^2 \bar{P}_k = (n/n-1) \sigma_{\bar{x}}^2 - (1/n-1) \sigma^2$$

であるから

$$(3) \sigma_{\bar{x}}^2 = \frac{1}{n} (1 + (n-1) \bar{P}_k)$$

が成立する。

(1) から級内相関係数は

$$\begin{aligned} \bar{P}_k &= \frac{2}{n(n-1)} \sum_{\delta} (n-\delta) P'_{k\delta} \\ &= \frac{2}{n-1} \sum_{\mu} P_{k\mu} \end{aligned}$$

で、結局、 n が奇数のとき \bar{P}_k は $P_{k\mu}$ の算術平均となるが、 n が偶数のときは \bar{P}_k は $P_{k\mu}$ の $n/n(n-1)$ 倍の算術平均値となることわかる。

定理⁽³⁾: 系統的抽出計画による推定値 \bar{x} は \bar{x} の不偏推定量で、その分散 $\sigma_{\bar{x}}^2$ は

$$\begin{aligned} (4) \sigma_{\bar{x}}^2 &= \sigma^2 \left\{ 1 - \frac{1}{n^2} \sum_{\delta} (n-\delta) v_{k\delta} \right\} \\ &= \sigma^2 \left(1 - \frac{1}{n} \sum_{\mu} v_{k\mu} \right) \\ &= \frac{\sigma^2}{n} \left(1 + 2 \sum_{\mu} \bar{P}_{k\mu} \right) \\ &= \frac{\sigma^2}{n} \{ 1 + (n-1) \bar{P}_k \} \end{aligned}$$

である。

証明、期待値、 \bar{x}_i, \bar{x} および系統的抽出計画の定義から、 \bar{x} はそのとりうる値が $\bar{x}_1, \dots, \bar{x}_k$ で $\bar{x} = \bar{x}_i$ となる確率は $\frac{1}{k}$ となっているような変量である。ゆえに、

(8) 定理1の幾分わかりやすい証明は、(2), (3) を用いて \bar{P}_k を代入することによってえられるが、筆者の考えではこれは以下に述べるもの程有益なものでない。附録BのLemmaはそれ自身有限母集団よりの抽出理論において重要であることは勿論である。

$$(5) \text{を } E \bar{x} = \bar{x}_1 + \dots + \bar{x}_k$$

で、 \bar{x}_i の値を (5) に代入すると、 $E \bar{x} = \bar{x}$ すなわち \bar{x} が \bar{x} の不偏推定量であることがわかる。

$E \bar{x}$ が計算されれば $\sigma_{\bar{x}}^2$ を評価するため $E \bar{x}^2$ を計算しなければならない。期待値の定義から

$$(6) \text{を } E \bar{x}^2 = \bar{x}_1^2 + \dots + \bar{x}_k^2$$

\bar{x}_i の値を (6) に代入すると

$$(7) n^2 E \bar{x}^2 = \sum_{i, j} x_{i+(n-1)k} x_{j+(n-1)k}$$

ゆえに附録BのLemma 6 の $f(u)$ を u とおくと、分散の定義から

$$\begin{aligned} \sigma_{\bar{x}}^2 &= \left(\frac{1}{kn} \right) \sum_{\nu} (x_{\nu} - \bar{x})^2 - \left(\frac{2}{kn^2} \right) \sum_{\delta} (x_j - x_{j+k\delta})^2 \\ &= \sigma^2 - \frac{1}{n^2} \sum_{\delta} (n-\delta) c_{k\delta} \end{aligned}$$

また Lemma 8 の $f(u)$ を u でおきかえると

$$\begin{aligned} \sigma_{\bar{x}}^2 &= \left(\frac{1}{kn^2} \right) \sum_{\nu} (x_{\nu} - \bar{x})^2 + \left(\frac{2}{kn^2} \right) \sum_{\delta} (x_j - \bar{x})(x_{j+k\delta} - \bar{x}) \\ &= \left(\frac{1}{n} \right) \sigma^2 + \frac{2}{n^2} \sum_{\delta} (n-\delta) c_{k\delta} \end{aligned}$$

いま附録BのLemma 9 で、 $f(x_j, x_{j+k\delta})$ を $(x_j - x_{j+k\delta})^2$ でおきかえると

$$\sigma_{\bar{x}}^2 = \sigma^2 - \left(\frac{1}{n} \right) \sum_{\mu} v_{k\mu}$$

で、 $f(x_j, x_{j+k\delta})$ を $(x_j - \bar{x})(x_{j+k\delta} - \bar{x})$ でおきかえれば

$$\sigma_{\bar{x}}^2 = \frac{1}{n} \left(\sigma^2 + 2 \sum_{\mu} c_{k\mu} \right)$$

よって結局

$$\sigma_{\bar{x}}^2 = \sigma^2 \left(1 - \frac{1}{n} \sum_{\mu} v_{k\mu} \right)$$

で、また

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \left(1 + 2 \sum_{\mu} P_{k\mu} \right)$$

をうる。

7 $\rho_{k\delta}^i$, $\rho_{k\delta}$ および $\sigma_{\bar{x}}^2$ のとりうる値

ここで我々は変動の型が色々違った場合に、 $\rho_{k\delta}^i$ および $\sigma_{\bar{x}}^2$ の値がどのようになるかを簡単に調べてみよう。

いま、

$$\sigma^2 \rho_{k\delta}^i = \frac{1}{k(n-\delta)} \sum_i (x_i - \bar{x})(x_{i+k\delta} - \bar{x})$$

である、 $x_i = x_{i+k\delta}$, $\delta = 1, \dots, n-1$, $i = 1, \dots, k$ とすると、

$$\sum_i (x_i - \bar{x})^2 = n \sum_i (x_i - \bar{x})^2$$

および

$$\sum_i (x_i - \bar{x})(x_{i+k\delta} - \bar{x}) = (n-\delta) \sum_i (x_i - \bar{x})^2$$

である。これを代入すると、

$$\rho_{k\delta}^i = 1, \quad \sigma_{\bar{x}}^2 = \sigma^2$$

がえられる。

$\sigma_{\bar{x}}^2$ に対するこの結果は直観的に明らかである。なぜなら変動はすべてとりうる標本同のものであるから、任意の特定の系統的標本は単一観測値と同等になるからである。

一方 $x_{i\delta+\alpha} = x_{i\delta+\beta}$, $\alpha, \beta = 1, \dots, k$; $\delta = 1, \dots, n-1$ とおいてみよう。このとき任意の i , $i = 1, \dots, k$ に対し

$$\sum_i (x_i - \bar{x})^2 = k \sum_i (x_{i+\alpha-\delta k} - \bar{x})^2$$

および

$$\sum_i (x_i - \bar{x})(x_{i+k\delta} - \bar{x}) = k \sum_i (x_{i+\alpha-\delta k} - \bar{x})(x_{i+\alpha+\delta-\delta k} - \bar{x})$$

が成立する。更に、

$$0 = \left[\sum_i (x_{i+\alpha-\delta k} - \bar{x}) \right]^2 = \sum_i (x_{i+\alpha-\delta k} - \bar{x})^2 + 2 \sum_i (x_{i+\alpha-\delta k} - \bar{x})(x_{i+\alpha+\delta-\delta k} - \bar{x})$$

である。よって

$$2 \sum_i \frac{n-\delta}{n} \rho_{k\delta}^i = -1, \quad \sigma_{\bar{x}}^2 = 0$$

である。

任意の特定の $\rho_{k\delta}^i = -1$ であるような例を作ることができるが、このような場合には残りの $\rho_{k\delta}^i$ はすべて0となる。よく知られているように \bar{P}_k の最小値は $-1/(n-1)$ である。

最後に x の添字が無作為に定められているときの $\rho_{k\delta}^i$ および $\sigma_{\bar{x}}^2$ の期待値を考えよう。これらの値は、

$$\rho_{k\delta}^i = -1/(nk-1)$$

および

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \left(\frac{nk-n}{nk-1} \right)$$

となる。

系統的抽出を実際に用いる際は、添字が無作為に定められているといえるような分布にしたがう x を取扱うことはまずないといえてよい。一般に、 x にはある種の基本的なトレンド(傾向)があるということを期待できる論理的な理由が存在するのである。したがって、その情報入手できることも多いし、またそうでなければ小標本からそれを求めて、これらの基礎のもとに x の添字が無作為になっていると仮定した場合とは違った何らかのやり方を採用することができるのである。

8 パラメーターの推定

さきでえた系統的標本の平均値の分散の公式は母集団に対する公式である。これらの値は母集団のすべての要素の値に関係する。しかし抽出手続のテストが可能の場合でも、分母集団を研究するための資料を利用できることは殆んどない。結局、問題は標本から母集団分散と系列相関を推定することが可能かどうかを調べることになってくる。単一の S_i から導びかれた分散および相関係数には偏りがあり、一貫性 (consistency) のないことが示される

か、一方 S_i の一つ以上の標本から不偏あるいは一致性のある推定値を構成することが可能である。これらの推定値の抽出分散は将来の研究にまたねばならない。

我々は S_i においてしたように S_i をノつだけ抽出するのではなく、 p 個の S_i を無作為に抽出すると仮定する。このとき標本は S_p 内のすべての要素からなる。標本平均 \bar{x} は、我々の標本クラス C_1, \dots, C_p であるとき、

$$p\bar{x} = \sum_p \bar{x}_p$$

で定義される。そうすると \bar{x} が不偏であることが容易にわかる。更に我々はこの抽出手続が p 個の要素を無作為に抽出するものと考えられることできるから、

$$\sigma_{\bar{x}}^2 = \frac{k-p}{k-1} \frac{1}{p} \sigma_{\bar{x}_i}^2$$

が導かれる。我々は S_i において $\sigma_{\bar{x}_i}^2$ を評価しておいた、

$$k\sigma_{\bar{x}_i}^2 = \sum_i (\bar{x}_i - \bar{x})^2$$

であるから $\sigma_{\bar{x}_i}^2$ を $d_{\bar{x}_i}^2$ で推定することを考えよう。

$$p d_{\bar{x}}^2 = \sum_p (\bar{x}_p - \bar{x})^2$$

である。

$$\text{いま } E \hat{x}^2 = \sigma_{\bar{x}}^2 + \bar{x}^2$$

$$\text{かつ } E \sum_p \bar{x}_p^2 = \frac{p}{k} \sum_i \bar{x}_i^2 = p(\sigma_{\bar{x}_i}^2 + \bar{x}^2)$$

(9) 我々がなぜこれらの S_i を系統的にではなく無作為に抽出したかについては不思議に思うかもしれない。もし S_i を系統的に抽出すれば、これは標本要素間の間隔がより小さい単一の系統的標本をとることと同等となる。その上ここで述べたような抽出分散の不偏推定値を導くことはできない。

であるから

$$E d_{\bar{x}}^2 = \sigma_{\bar{x}}^2$$

で、結局 $d_{\bar{x}}^2$ は $\sigma_{\bar{x}}^2$ の不偏推定値であることがわかる。更に、

$$E d_{\bar{x}}^2 = \frac{p(k-1)}{k-p} \sigma_{\bar{x}_i}^2$$

ここで我々は方向をかえて $P_{k\mu}$ および σ^2 の推定を考えよう。

$$p k \hat{d}_{\bar{x}}^2 = \sum_{p,\mu} (x_{p+(k-1)\mu} - \bar{x})^2$$

および

$$p k \hat{C}_{k\mu} = \sum_{p,\mu} (x_{p+(k-1)\mu} - \bar{x})(x_{p+(k-1)\mu} - \bar{x})$$

と置く。このとき

$$E \hat{d}_{\bar{x}}^2 = \sigma^2 - \sigma_{\bar{x}}^2$$

および

$$E \hat{C}_{k\mu} = C_{k\mu} - \sigma_{\bar{x}}^2$$

を示すことができる。よって

$$E d_{\bar{x}}^2 + d_{\bar{x}}^2 \left(\frac{k-p}{p(k-1)} \right) = \sigma^2$$

および

$$E \hat{C}_{k\mu} + d_{\bar{x}}^2 \left(\frac{k-p}{p(k-1)} \right) = C_{k\mu}$$

よって

$$\hat{C}_{k\mu} + d_{\bar{x}}^2 \left(\frac{k-p}{p(k-1)} \right) = Y_{k\mu} \left[\hat{d}_{\bar{x}}^2 + d_{\bar{x}}^2 \left(\frac{k-p}{p(k-1)} \right) \right]$$

で定義される $P_{k\mu}$ の推定値 $Y_{k\mu}$ には偏りがある。しかし多くの場合この偏りは小であらう。勿論 k が偶数で $\mu = \frac{k}{2}$ なら前に定義したように $Y_{k\mu}$ を推定するには $Y_{k\mu}$ を二倍にする。

ひとつの方法は

$$p k \mu \hat{d}_{\bar{x}}^2 = \sum_{p,\mu} (x_{p+(k-1)\mu} - \bar{x}_p)^2$$

および

$$g n_w \hat{C}_{A\mu g} = \sum_{\beta, \lambda} (x_{\beta+(\alpha-1)k} - \bar{x}_\beta)(x_{\beta+(\alpha+\mu-1)k} - \bar{x}_\beta)$$

を考ふるものである。このとき

$$E_w \hat{S}_g^2 = \sigma^2 - \sigma_{\bar{x}}^2$$

$$E_w \hat{C}_{A\mu g} = C_{A\mu} - \sigma_{\bar{x}}^2$$

および

$$E(w \hat{S}_g^2 + S_g^2) = \sigma^2$$

$$E(w \hat{C}_{A\mu g} + C_{A\mu g}) = C_{A\mu}$$

が導びかれる。よって $\rho_{A\mu}$ のもう一つの推定値は

$$w \hat{C}_{A\mu g} + S_g^2 = w \gamma_{A\mu} (w \hat{S}_g^2 + S_g^2)$$

なる式で定義される。 $g=1$ のときは $S_g^2=0$ だから標本から σ^2 , $C_{A\mu}$ および $\sigma_{\bar{x}}^2$ の不偏推定値を与えることは不可能である。しかし、

$$\frac{1 - \gamma_{A\mu}}{1 - \gamma_{A\mu}'} = \frac{\hat{S}_g^2 - \hat{C}_{A\mu g}}{\hat{S}_g^2 - \hat{C}_{A\mu g}'}$$

$$\text{で } E[\hat{S}_g^2 - \hat{C}_{A\mu g}] = \sigma^2 - C_{A\mu}$$

だから近似的に

$$\frac{1 - \rho_{A\mu}}{1 - \rho_{A\mu}'} = E \frac{1 - \gamma_{A\mu}}{1 - \gamma_{A\mu}'}$$

が導びかれる。

同様な方程式は $w \gamma_{A\mu}$ に対しても成立する。

$\rho_{A\delta}$ を推定すると "クラス内" の定義は簡単にできる。

$$g(n-\delta)_w \hat{C}_{A\delta g} = \sum_{i, \lambda} (x_{i+(\alpha-1)k} - 1\delta \bar{x}_i)(x_{i+(\alpha+\delta-1)k} - 2\delta \bar{x}_i)$$

とおく。

ここで

$$(n-\delta)_{1\delta} \bar{x}_i = \sum_{\lambda} x_{i+(\alpha-1)k}$$

$$(n-\delta)_{2\delta} \bar{x}_i = \sum_{\lambda} x_{i+(\alpha+\delta-1)k}$$

また

$$g(n-\delta)_w \hat{C}_{A\delta g} = \sum_{\beta, \lambda} (x_{\beta+(\alpha-1)k} - 1\delta \bar{x}_\beta)(x_{\beta+(\alpha+\delta-1)k} - 2\delta \bar{x}_\beta)$$

とする。更に

$$g_c \hat{C}_{A\delta g} = \sum_i (1\delta \bar{x}_i - \bar{x})(2\delta \bar{x}_i - \bar{x})$$

および

$$g_c \hat{C}_{A\delta g} = \sum_{\beta} (1\delta \bar{x}_\beta - \bar{x})(2\delta \bar{x}_\beta - \bar{x})$$

とおくと、

$$C_{A\delta} = w \hat{C}_{A\delta g} + c \hat{C}_{A\delta g}'$$

および

$$E(w \hat{C}_{A\delta g} + c \hat{C}_{A\delta g}') = C_{A\delta}$$

である。

よって $\rho_{A\delta}$ の推定値として $\gamma_{A\delta}'$ を用いる。ここで

$$w \hat{C}_{A\delta g} + c \hat{C}_{A\delta g}' = \gamma_{A\delta}' (w \hat{S}_g^2 + S_g^2)$$

\bar{x} が既知のときには

$$E \sum_{\beta, \lambda} (x_{\beta+(\alpha-1)k} - \bar{x})(x_{\beta+(\alpha+\mu-1)k} - \bar{x}) = g n C_{A\mu}$$

$$E \sum_{\beta, \lambda} (x_{\beta+(\alpha-1)k} - \bar{x})(x_{\beta+(\alpha+\delta-1)k} - \bar{x}) = g(n-\delta) C_{A\delta}$$

$$\text{および } E \sum_{\beta, \lambda} (x_{\beta+(\alpha-1)k} - \bar{x})^2 = g n \sigma^2$$

だから $\rho_{A\mu}$, $\rho_{A\delta}$ および σ^2 の単純な推定値が容易に求まる。

ゆえに予備調査で \bar{x} が既知なら単一の標本からでも $\sigma_{\bar{x}}^2$ のパラメータを推定することが可能である。

3 標本の大きさが変わったときの分散の変化

系統的抽出計画の分散を無作為標本の分散および系列相関係数と表わすということの主な理由は

- 1 無作為抽出と他の抽出計画の比較をおこなうため
- 2 系統的と無作為計画の効率の差が何に原因するかを分析し

やすくするため

3 色々な大きさの標本に対する分散推定値を簡単に求めるため

である。この節においてはこれらの理由の三番目を取扱う。\$P_{k\mu}\$ による分析は容易だから、我々は \$P_{k\mu}\$ のみを論ずることとする。

問題は \$k\$ の函数 \$\bar{P}_k\$ を推定することにある。\$k\$ の一つの値について \$P_{k\mu}\$ が計算されているとき、どのようにすれば \$k\$ のすべての値に対する \$\bar{P}_k\$ の推定ができるかを示すため、まず \$\sigma^2\$ が \$k\$ と無関係なことから、我々の考察を \$C_{k\mu}\$ に限ってよいことに注意する。\$S_6\$ において、\$C_{k\mu}\$ を

$$k n C_{k\mu} = \sum_v (x_v - \bar{x})(x_{v+k\mu} - \bar{x})$$

で定義した。

したがって \$k' n' = k n = N\$ なる \$k'\$ について \$C_{k'\mu}\$ を評価したければ \$C_{k'\mu} = C_{k\mu}\$ である。したがって \$k'\mu' = k\mu\$ なら任意の \$k\$ および \$\mu\$ について

$$C_{k'\mu'} = C_{k\mu} / k$$

をうる。この場合 \$\mu\$ を \$k'\mu'/k\$ でおきかえている。

この方法は \$k' < k\$ ならある種の内挿をおこなうことになるが、もし \$P_{k\mu}\$ が \$\mu\$ に対してプロットされれば、これは図上で行なえることが多い。しかし普通 \$k\$ の値が \$k' > k\$ となるように \$k\$ をとるのがよい。

ある場合には相関々数を作れることがある。例えば \$x_v\$ を \$v\$ の多項式で表現できれば、\$P_{k\mu}\$ は \$k\$ の多項式で表わせる。このことからもし \$x_v\$ が滑らかな傾向に従って変るなら、\$P_{k\mu}\$ もまた滑らかな傾向に従って変化する。したがって内挿が可能である。この問題については更に検討が必要である。

10 層化系統抽出

サンプリングを実行する場合には層化母集団を取扱うことがあ

い。層化母集団にもとづく推定値の分散は普通層間の変動性を含まない。結局、母集団がうまく層化されていれば、大きさ \$n\$ の標本から求めた推定値の変動は、層を考えずにとられた大きさ \$n\$ の無作為標本の推定値より非常に変動が小さいのが普通である。ここで我々は層化された母集団からの系統的抽出の理論を考えることにする。

母集団 \$P\$ が \$L\$ 個の層 \$P_1, \dots, P_L\$ からなっていて、その \$a\$ 番目には \$N_a\$ 個の要素 \$x_{a1}, \dots, x_{aN_a}\$ が含まれるものとする。要求されているのは \$P\$ の算術平均値 \$\bar{x}\$ の推定である。\$P_a\$ の算術平均を \$\bar{x}_a\$ と書く。\$N_a = k_a n_a\$ とする。

我々は可能な二つの場合を考える。その一つは、異なる場所でも標本抽出に同じ手順を用いてよいという管理上の簡単さのために屢々用いられるものである。この前の結果から、この方法を用いてよい場合が指示されるであろう。

抽出方式 I — \$k_1 = k_2 = \dots = k_L = k\$ で、抽出方式は \$1, \dots, k\$ の中から \$k\$ の整数を無作為に選ぶものであるとする。これらの各整数の抽出確率は \$1/k\$ になっている。このとき選ばれた整数がたとえば \$i\$ なら \$P_a\$ の標本は \$x_{ai}, x_{a(i+k)}, \dots, x_{a(i+(n/k)-1)k}\$ からなる。したがってこの抽出方式をとることにより、実際の標本となる確率が何れも \$1/k\$ であるような丁度 \$k\$ 個の可能な標本、すなわち \$S_1, \dots, S_k\$ が存在することになる。

抽出方式 II — この抽出方式はそれぞれの \$a\$ について \$1, \dots, k_a\$ の整数の \$k\$ を、各整数の選ばれる確率が \$1/k_a\$ となるようにして無作為に抽出するものである。故にこの抽出方式を適用したとき、実際に標本となる確率がそれぞれ \$1/k_1, \dots, 1/k_L\$ であるような可能な標本が丁度、\$k_1, \dots, k_L\$ 個あることになる。

層化抽出には勿論他の抽出方式もある。しかし上にあげた二つは、系統抽出に含まれるものを除けば実際的問題の殆んど全部を包括している。讀者については後の論議で取扱うことにする。更にこれらの方式について等びかれた結論から、他の層化抽出方式

に対する結論を推測することも可能である。

S_{ai} を要素 $X_{ai}, X_{ai+k}, \dots, X_{ai+(n_a-1)k}$ のクラスとする。まず抽出方式 I を考察する。Pa から大きさ n_a の系統的標本を抽出しなければならない。

可能な標本は S_1, \dots, S_k である。ここで S_i は S_{Li}, \dots, S_{Li} 内のすべての要素からなる。 S_{ai} 内の要素の算術平均値を \bar{x}_{ai} と書く。 Pa からとられた標本の算術平均値を \bar{x}_a 、標本平均値を \bar{x} と書く。ここで

$$N\bar{x} = N_1\bar{x}_1 + \dots + N_k\bar{x}_k$$

である。そうすると、

$$N\bar{x} = \sum_a N_a \bar{x}_a = \sum_a N_a \frac{1}{k} \sum_i \bar{x}_{ai} = N\bar{x}$$

附録 C から

$$\sigma_{\bar{x}}^2 = \frac{1}{N^2} \sum_{a,c} N_a N_c \sigma_{\bar{x}_a \bar{x}_c}$$

がえられる。ここで

$$\begin{aligned} \sigma_{\bar{x}_a \bar{x}_c} &= E(\bar{x}_a - \bar{x}_a)(\bar{x}_c - \bar{x}_c) \\ &= \frac{1}{k^2} \sum_i (\bar{x}_{ai} - \bar{x}_a)(\bar{x}_{ci} - \bar{x}_c) \end{aligned}$$

である。

$\sigma_{\bar{x}_a \bar{x}_c}$ の式はこれ以上簡単にはならないが、重要なことは、異なる層の対応する item 間に正の相関があれば、抽出誤差以外の考慮が重大でない限り、抽出方式 I は用いない方がよい。しかし対応する item 間に負の相関があれば、抽出方式 I の分散は II の分散より小さくなる。

次に抽出方式 II を考えよう。 I と II の違いは、 II では各層内の標本抽出を別々に行なっているから、もし α キムなら $\sigma_{\bar{x}_a \bar{x}_c} = 0$ となることである。故に方式 II の場合の分散は、

$$\sigma_{\bar{x}}^2 = \frac{1}{N^2} \sum_a N^2 \sigma_{\bar{x}_a}^2$$

となり、ここで $\sigma_{\bar{x}_a}^2$ はオノ節で導びかれている。

11 系統的抽出と無作為抽出方式の効率の比較

どんな抽出方法の研究でも他の可能な抽出方法と比較をしなければ完成したとはいえない。よってこの節では、系統的抽出方式と無制限無作為および層化無作為抽出方式との比較を行なうことにする。

無作為および層化無作為抽出方式に対する平均値と分散はそれぞれ (') および (") で区別することにする。

そうすると

$$\sigma_{\bar{x}}^2 / \sigma_{\bar{x}}^2 = (1 + 2 \sum_{\mu} P_{\mu\mu}) \left(\frac{k\alpha - 1}{k\alpha - n} \right)$$

なることがわかるから、略号

$$\sum P_{\mu\mu} < -(\alpha - 1) / 2(k\alpha - 1)$$

なら $\sigma_{\bar{x}}^2 < \sigma_{\bar{x}}^2$ となり、 n が α に比して大きければ $-(\alpha - 1) / 2(k\alpha - 1)$ の近似として $-\frac{1}{2}k$ を使ってもよい。

もっと細かい比較を行なうには、母集団要素 X_v が v のある函数で与えられるものと仮定し、その函数が

$$X_v = A_0 + A_1 v + \dots + A_k v^k$$

あるいは

$$\begin{aligned} X_v &= B_0 + A_1 \sin \frac{2\pi v}{N} + B_1 \cos \frac{2\pi v}{N} \\ &+ \dots \\ &+ A_k \sin \frac{2\pi k v}{N} + B_k \cos \frac{2\pi k v}{N} \end{aligned}$$

のようなものであるとして、種々の可能な抽出方式の効果を、そのように仮定された X_v の分布に基づいて研究するのが有用である。系統的抽出を用いる際には、母集団の要素を論理的なやり方で順序づけ、またこの順序を用いて標本を系統的に抽出できると仮定しているということに注意しなければならない。

我々はここで幾つかの可能性を考える。まず一つの層から / つ

の要素しか抽出しなければ層標本平均値の分散は抽出が無作為でも系統的でも同じだということを注意しておこう。一方才ノ節から、母集団をL個の層に分けて、各層からj番目の要素をとることにより大きさnの系統的標本を作ったとすると、層化無作為標本平均値の分散は、系統的標本内の層標本平均値間の平均的相関が負であるか正であるかによって系統的標本の平均値の分散より大きくなったり小さくなったりすることになる。

ここで周期性をも母集団については系統的標本を用いてはならないという警告がどうして出てきたかを考えてみよう。周期があるなら系統的標本の層平均値間の相関は+1であるから無作為標本の方が優れている。しかし周期が2をなら、系統的標本の方がおそろく分散が小さくなるということを証明しよう。

周期が2をなら、大きさnの隣接する二つの層内では常に

$$x_1 = x_{2k}, x_2 = x_{2k-1}, \dots, x_n = x_{n+1} \quad \text{で}$$

$$x_i - \bar{x} = -(x_{n+1} - \bar{x})$$

であると仮定する。このとき各層から1つの要素を抽出すると系統的標本平均値(この場合には個々の要素)間の相関は層番手の差が奇数なら-1で、偶数なら+1である。

n個の層のそれぞれの中の分散は σ_1^2 である。ここで

$$n \sigma_1^2 = \sum_{i=1}^n (x_i - \bar{x})^2$$

ある。層平均値間の分散は σ^2 である。よって $\sigma^2 = \sigma_1^2$ 、ゆえに $n = L$ なるときの大きさnの無制限無作為標本平均値の分散は

$$\sigma_{\bar{x}}^2 = \frac{N-L}{N-1} \frac{\sigma_1^2}{L}$$

で、層化無作為標本平均値の分散は $\sigma_{\bar{x}}^2$ である。一方系統的抽出平均値の分散は

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{L^2} \sum_{i,j=1}^L (-1)^{i+j} (2 - \delta_{ij})$$

である。ただし $i=j$ なら $\delta_{ij}=1$ 、 $i \neq j$ なら $\delta_{ij}=0$ である。

よってLが偶数なら $\sigma^2=0$ であるがもしLが奇数なら $\sigma^2 = \frac{1}{L} \sigma_1^2$ である。

結局母集団の分布が周期性を有し、層の大きさが周期の半分なら系統的標本平均値の効率は層化無作為標本の効率より大きいということになる。それが周期に等しいかまたは周期の半分に等しい場合に成立するのと同様だが、それが周期の半分の奇数倍または偶数倍に等しいときにもなり立つことに注意しなければならない。

しかし母集団の要素が直線分布に従うと仮定すれば事情は全く変わってくる。一般性を失わずに、直線上の分布は $x_j = j$ で与えられると仮定してよい。そうすると大きさnの無制限無作為標本については、標本平均を \bar{x}' と書くと

$$\bar{x}' = \bar{x} = \frac{1}{2} (n+1)$$

$$\sigma^2 = \frac{n^2 n^2 - 1}{12}$$

および

$$\sigma_{\bar{x}'}^2 = \frac{(n-1)(n+1)}{12}$$

をうる。

層化無作為抽出計画では $N_1 = \dots = N_L = \frac{CN}{n}$ と仮定する。ここでCは整数1, 2, ..., nのどれと同じでもよい。すなわち $L = \frac{n}{C}$ 、 $n_1 = \dots = n_L = C$ とする。このとき

$$\sigma_{\bar{x}'}^2 = \frac{C^2(n-1)}{n^2(Cn-1)} \sum_{\alpha} \sigma_{\alpha}^2$$

ここで \bar{x}' は層化無作為標本の標本平均値である。もし才a層に $x_{(a-1)C+1}, \dots, x_{aC}$ が含まれていれば $\sigma_{\alpha}^2 = \frac{C^2 n^2 - 1}{12}$ および $\sigma_{\bar{x}'}^2 = \frac{C}{N} \cdot \frac{n-1}{Cn-1} \cdot \frac{C^2 n^2 - 1}{12}$ である。

結局

$$\sigma_{\bar{x}'}^2 = \sigma^2 - \frac{1}{n^2} \sum_{\alpha} \sum_{\beta} (x_j - x_{j+k\delta})^2$$

$$= \sigma^2 - \frac{k^2(n^2-1)}{12}$$

$$= \frac{k^2-1}{12}$$

これをまとめると、

$$\sigma_{\bar{x}}^2 = \frac{(k-1)(kn+1)}{12} = \frac{(k-1)(N+1)}{12}$$

$$\sigma_{\bar{x}'}^2 = \frac{c(k-1)(ck+1)}{12} = \frac{(k-1)(\frac{N}{L}+1)}{12L}$$

$$\sigma_{\bar{x}''}^2 = \frac{k^2-1}{12}$$

のようになる。

$\sigma_{\bar{x}}^2$ および $\sigma_{\bar{x}'}^2$ のどちらも $\sigma_{\bar{x}''}^2$ より小さいことは明らかである。しかし

$$\sigma_{\bar{x}}^2 / \sigma_{\bar{x}'}^2 = \frac{L(k+1)}{\frac{N}{L}+1}$$

である。

$kn = N$, $cL = n$ だから $k = \frac{N}{cL}$ 。 よって

$$N > L(L-1) \text{ かつ } c \geq \frac{L}{1 - \frac{L(L-1)}{N}} \text{ なら}$$

$$\sigma_{\bar{x}}^2 < \sigma_{\bar{x}'}^2$$

である。いかなる場合にも $\frac{N}{L} \geq c$ 。 よって c が上記の値となるためには $\frac{N}{L} > \frac{NL}{N-L(L-1)}$ 。 すなわち N は $2L^2 - L$ 以上でなければならぬことがわかる。 よって $N > 2L^2 - L$ および $c \geq \frac{L}{1 - \frac{L(L-1)}{N}}$ なら $\sigma_{\bar{x}}^2 \leq \sigma_{\bar{x}'}^2$ である。 それ以外の場合は $\sigma_{\bar{x}}^2 > \sigma_{\bar{x}'}^2$

である。

この結果は二つの事実によってわかる。

- (1) 各層から1つの要素を抽出してれば、層平均値間の平均的な相関が高いので、層内の分散が等しいにもかかわらず層化無作為標本平均値の効率が、系統的抽出の平均値の効率より大きくなって来る。

(2) 各層から二つ以上の要素が抽出されていると、系統的抽出平均の層内分散は層化無作為抽出平均の層内分散より小であって、層および層からとられた標本の大きさが十分大なら系統的標本の層内分散の減少は層平均値間の相関によるものより更に大となる。もちろん、直線的分布 (*straight line distribution*) のときには、我々が用いたものより更に有効な層化無作為標本の定義の仕方がある。その上抽出手続もここに論じたものより一層有効なものが利用できる。しかしこの例は、実際に生ずる一般的問題のみならず、それらの解決に利用できる方法を指摘する上で有用である。

系統的抽出と層化無作為抽出の他の比較は、 x_v が二つの要素、すなわち傾向函数 (*trend function*) および周期函数 (*periodic function*) からなっていると考えられる。したがって傾向 (*trend*) からの偏差は周期函数となっている。

$x_v = \varphi_1(v) + \varphi_2(v)$ とする。ここで $\varphi_1(v)$ は傾向函数で $\varphi_2(v)$ は周期函数、 $N = 2AQ$ の周期函数である。

$\varphi_2(v) = \sum_j \varphi_{2j}$ とする。このとき $\varphi_{2j} = \varphi_{2At+j} = \dots = \varphi_{2A(c-1)+j}$, $j = 1, 2, \dots, 2A$ および $\varphi_{2At+j} - \bar{\varphi} = -(\varphi_{2A(c-1)+j} - \bar{\varphi})$, $j = 1, \dots, A$, $a = 0, \dots, Q-1$ 。この比較で我々が考える標本の大きさはすべからぬの倍数であるから、必要な情報が全部一度でえられるように、分散、共分散を計算する。

$\varphi_1(v)$ の平均値を $\bar{\varphi}_1$, $\varphi_2(v)$ の平均値を $\bar{\varphi}$ と書く。そうすると $\bar{x} = \bar{\varphi}_1 + \bar{\varphi}$ である。

$$N\sigma^2 = \sum_v (x_v - \bar{x})^2$$

$$= \sum_v [\varphi_1(v) - \bar{\varphi}_1]^2 + \sum_v (\varphi_2(v) - \bar{\varphi})^2 + 2 \sum_v [\varphi_1(v) - \bar{\varphi}_1] (\varphi_2(v) - \bar{\varphi})$$

$$= \sum_{a,i} (\varphi_{1(a-1)+i} - \bar{\varphi}_{1a})^2 + A \sum_a (\bar{\varphi}_{1a} - \bar{\varphi}_1)^2$$

$$+ \sum_{a,i} (\varphi_{2(a-1)+i} - \bar{\varphi}_a)^2 + A \sum_a (\bar{\varphi}_a - \bar{\varphi})^2$$

$$+ \sum_{a,i} [g_{i(a-1)h+i} - \bar{g}_{ia}] (g_{i(a-1)h+i} - \bar{g}_{ia})$$

$$+ h \sum_a (\bar{g}_{ia} - g_i) (\bar{g}_{ia} - \bar{g})$$

ここで $a=1, \dots, 2Q; i=1, \dots, h$. \bar{g}_{ia} は $g_{i(a-1)h+i}, \dots, g_{i(a)h}$ の算術平均値で \bar{g}_a は $g_{1(a-1)h+1}, \dots, g_{ah}$ の算術平均値である。

g_v についての仮定から $\bar{g}_a = \bar{g}$ および $\sum_i (g_i - \bar{g}_a)^2$ は a のそれぞれ別の値について同一で、 $\sum_{i=1}^h (g_i - \bar{g})^2$ に等しいことがわかる。

$g_i = g_{2h+i} = \dots = g_{2h(a-1)+i}$ および $g_{hi} = g_{2h+i} = \dots = g_{2h(a-1)+hi}$ であるから

$$\sum_v [g_i(v) - \bar{g}_i] (g_v - \bar{g}) = \sum_{i=1}^h (g_i - \bar{g}) \sum_{a=1}^Q [g_{i(2a-2)h+i} - \bar{g}_{i(2a-2)}]$$

$$+ \sum_{i=h+1}^{2h} (g_i - \bar{g}) \sum_{a=1}^Q [g_{i(2a-2)h+i} - \bar{g}_{i(2a-2)}]$$

をうる。

また $g_i - \bar{g} = (g_{i+h} - \bar{g})$ だから

$$\sum_v [g_i(v) - \bar{g}_i] (g_v - \bar{g}) = \sum_{i=1}^h (g_i - \bar{g}) \left\{ \sum_{a=1}^Q [g_{i(2a-2)h+i} - \bar{g}_{i(2a-2)}] \right.$$

$$\left. - g_{i(2a-1)h+i} - \bar{g}_{i(2a-1)} + g_{i(2a-1)} \right\}$$

がえられる。これは $g_i(v)$ が直線であるかまたは $g_i(v)$ がそれぞれ長さ $2h$ の直線の系列のときは 0 になる。

ここで $g_i(v) = A + Bv$ とする。そうすると

$$\bar{g}_{ia} = A + B[(a-1)h + \frac{A+1}{2}],$$

$$g_{i(a-1)h+i} - \bar{g}_{ia} = i - \frac{A+1}{2}$$

$$\sum_i [g_{i(a-1)h+i} - \bar{g}_{ia}]^2 = \frac{B^2 A^2 (h^2 - 1)}{12}$$

$$\sum_a (\bar{g}_{ia} - \bar{g}_i)^2 = \frac{B^2 A^2 (4Q^2 - 1)(2Q)}{12}$$

$$\bar{g}_i = A + B \frac{2A+1}{2},$$

$$\sum_v (g_i(v) - \bar{g}_i)^2 = \frac{2hQB^2}{12} [4h^2Q^2 - 1]$$

よって $\sigma^2 = \sigma_{\bar{g}}^2 + \frac{B^2}{12} [4h^2Q^2 - 1]$. ただし $h\sigma_{\bar{g}}^2 = \sum_i (g_i - \bar{g})^2$. また大きさ n の単純無作為標本の平均値の分散は

$$\sigma_{\bar{x}}^2 = \frac{N-n}{N-1} \frac{\sigma^2}{n}$$

いま層の大きさが mh で、 m は $2Q$ の因数すなわち $2Q/m = L_m$ としよう。そうすると L_m 個の各層内の分散は一定値すなわち σ_i^2 である。ただし m が偶数なら $\sigma_i^2 = \sigma_{\bar{g}}^2 + \frac{B^2}{12} [h^2 m^2 - 1]$. m が奇数のときは L_m が偶数で、層の半分では層内分散は $\sigma_i^2 + \frac{1}{4m} \sum_{i=1}^m (g_i - \bar{g})(i - \frac{A+1}{2})$ である。一方他の層における層内分散は

$$\sigma_i^2 - \frac{1}{4m} \sum_{i=1}^m (g_i - \bar{g})(i - \frac{A+1}{2})$$

そうすると L_m 個の層の各々から C 個の要素を無作為に抽出すると

$$\sigma_{\bar{x}}^2 = \frac{1}{L_m} \left(\frac{mh - C}{mh - 1} \right) \frac{\sigma_i^2}{C}$$

$$= \frac{1}{L_m} \left(\frac{mh - C}{mh - 1} \right) \frac{1}{C} \left(\sigma_{\bar{g}}^2 + \frac{B^2}{12} [h^2 m^2 - 1] \right)$$

なることがわかる。

系統的抽出の平均値の分散を求めたため $\sum_i (x_i - x_{i+h})^2 = h(n-h) s_{h\delta}^2$ を計算してみよう。いま x_v を代入すると $h(n-h) s_{h\delta}^2 = \sum_i (g_i - g_{i+h})^2 - 2Bh\delta \sum_i (g_i - g_{i+h}) + h(n-h) B^2 h^2 \delta^2$ なることがわかる。

よって h が h の倍数なら $\sum_i (g_i - g_{i+h}) = 0$ となる。又は h が h の奇数倍ならば $g_i = g_{i+h}$ だから $\sum_i (g_i - g_{i+h})^2 = 0$. 結局 h が h の奇数倍ならば δ は奇数となって $g_i - g_{i+h} = 2(g_i - \bar{g})$. 一方 δ が偶数ならば $g_i - g_{i+h} = 0$ となるから h が h の奇数倍で δ が奇数のときは

$$\begin{aligned} \sum_i (y_i - y_{i+\delta k})^2 &= 4 \sum_i (y_i - \bar{y})^2 \\ &= 4 \frac{k(n-\delta)}{k} \sum_i (y_i - \bar{y})^2 \end{aligned}$$

となる。

したがって n が偶数であることを注意しておく。

$$\sigma_x^2 = \sigma^2 - \frac{1}{n^2} \sum_i (n-\delta) A_{k\delta}^2$$

だから $\sum_i (n-\delta) A_{k\delta}^2$ を評価しなければならない。いま n が k の偶数倍なら $(n-\delta) A_{k\delta}^2 = (n-\delta) B^2 k^2 \delta^2$ および

$$\begin{aligned} \sum_i (n-\delta) A_{k\delta}^2 &= B^2 k^2 \left\{ n \sum_i \delta^2 - \sum_i \delta^2 \right\} \\ &= B^2 k^2 \frac{n^2(n^2-1)}{12} \end{aligned}$$

が成立する。よってもし n が k の偶数倍なら

$$\sigma_x^2 = \sigma^2 - \frac{B^2 k^2 (n^2-1)}{12}$$

なることがわかる。

一方もし n が k の偶数倍で δ が奇数なら

$$(n-\delta) A_{k\delta}^2 = (n-\delta) B^2 k^2 \delta^2 + 4(n-\delta) \sigma_y^2$$

で、 δ が偶数なら

$$(n-\delta) A_{k\delta}^2 = (n-\delta) B^2 k^2 \delta^2$$

となる。

よって

$$\sum_i (n-\delta) A_{k\delta}^2 = \frac{B^2 k^2 n^2 (n^2-1)}{12} + n^2 \sigma_y^2$$

ゆえに n が k の奇数倍なら

$$\sigma_x^2 = \sigma^2 - \frac{B^2 k^2 (n^2-1)}{12} - \sigma^2$$

なることがわかる。

また n が k の偶数倍のときは

$$\sigma_x^2 = \sigma_y^2 + \frac{B^2}{12} (k^2 n^2 - 1) - \frac{B^2}{12} k^2 (n^2 - 1)$$

$$= \sigma_y^2 + \frac{B^2}{12} (k^2 - 1)$$

で、 n が k の奇数倍なら

$$\sigma_x^2 = \frac{B^2}{12} B^2 (k^2 - 1)$$

となる。

このことから

$$C > \frac{L}{1 + \frac{12 \sigma_y^2}{B^2 (A_m)(A_m-1)} - \frac{L(L-1)}{N}}$$

のときには系統的抽出が優れた結果を与える。

$\frac{N}{L} > C$ であるから C の解が存在するためには

$$N > 2L^2 - L - \frac{12 \sigma_y^2}{B^2} \left(\frac{L^2}{N-L} \right)$$

でなければならない。

12 摘要

この論文では、分散の公式、可能な母数の値の検討、母数の推定、標本の大きさを変えたときの影響、および系統的抽出、単純無作為抽出、層化無作為抽出間の比較をも含めて、層化または非層化母集団に対する系統的抽出法の理論的基礎を提示した。この論文には、抽出単位が単一の要素からなっているような場合において必要とされる理論のみならず、更に系統的抽出を採用すべき条件についても若干の分析を行ないまた分散を計算するための公式が含まれている。

このあとの論文においては、抽出単位が要素の集塊である場合の系統的抽出の理論を提示するが、この理論は有限母集団からの抽出でなく、その各要素が正規分布するときの無限母集団からの抽出を仮定するものである。そうして更に系統的抽出理論の色々な部分ならびに実際を与えるであろう。

附録 A

特定の変数の示す値

変数の限界の繰返えしを避けるためこれらの範囲をこの附録に与える。

第 1 表 添字のとる値

文 字	この文字は / から F に示す値までのすべての整数をとるものとする
i	1
λ	$n - \delta$
j	1 から $(n - \delta)$
δ	$n - 1$
v, v'	1 から n
α	n
γ	n
μ, μ'	n が偶数なら $n/2$, n が奇数なら $\frac{n-1}{2}$
a, b	L

文字 β は i_1, i_2, \dots, i_p なる値を仮定する, たゞし i_1, \dots, i_p は互いの整数 $1, \dots, n$ のうちから選んだ p 個の数である。

附録 B

若干の有限和の範囲について

有限和を変換するとき生ずる困難は重積分の変換理論においゝ生ずるもの, すなわち変換する変数の影響または和の範囲についてとる和の順序等に非常によく似ている。この論文で有用なことが証明された幾つかの Lemma を別々に更に一般的に形を与えようとする。

$f(u), f(u, v)$ は u と v のすべての可能な値について有限な u, v の函数と仮定する。

Lemma 1

$$\sum_{i=1}^n f(x_{i+(n-\delta)k}, x_{i+(n-1)k}) = \sum_{i=1}^n f(x_{i+(n-\delta)k}, x_{i+(n-\delta-1)k})$$

証明: $\alpha = \lambda, \gamma = \lambda + \delta$ とする, $1 \leq \alpha < \gamma$ および $\gamma \leq n$ だから δ の可能な値は $1, \dots, (n-1)$ である, $\lambda = \gamma - \delta$ だから δ の各値について λ のとりうる値は $n - \delta$ である, また δ を一定としたとき λ の最大値は $\gamma = n$ なるとき決定される, これらの範囲に対して左辺の式の f の各項は右辺において唯一回だけしか現れない, 更に右辺においてはこれ以外の項は出現しない。

Lemma 2

$$\sum_i \sum_{\lambda} f(x_{i+(n-\delta)k}, x_{i+(n-\delta-1)k}) = \sum_j f(x_j, x_{j+k\delta})$$

証明: $j = i + (n-1)k$ とおく, このとき j は i および λ の単調増加函数である, j の最小値は $i=1$ のときに生ずる, そのとき $j=1$ である, j の最大値は $i=n$, $\lambda=n-\delta$ なるときに現れる, そのとき $j=(n-\delta)k$ である, これらの範囲に対して左辺の式の f の各項は右辺において唯一回しか現れない, 更に右辺においゝはこれ以外の項は出現しない。

Lemma 3

$$\sum_{\substack{i, \lambda \\ i \neq \lambda}} f(x_{i+(n-\delta)k}, x_{i+(n-1)k}) = \sum_{i \neq j} f(x_j, x_{j+k\delta})$$

証明: まず Lemma 1 を $\sum_{i \neq j} f(x_{i+(n-\delta)k}, x_{i+(n-1)k})$ に適用し次に与えられた式に Lemma 2 を適用する。

Lemma 4

$$\sum_{j=1}^m [f(x_j) + f(x_{j+k\delta})] = (n+1) \sum_{j=1}^m f(x_j)$$

証明: $m = j + k\delta$ とおく, このとき δ を一定とすると m の最小値は $j=1$ なるときに生ずる, そのとき $m = k\delta + 1$ である, δ の任意の一定値に対して m の最小値は $j = (n-\delta)$ なるときに生ずる, その場合 $m = nk$ である, 文字 m は $k\delta + 1$ から nk までのすべての整数値をとるものとする, よって

$$\sum_{j=1}^m f(x_j) + \sum_{j=1}^m f(x_{j+k\delta}) = \sum_{j=1}^m f(x_j) + \sum_{j=m}^n f(x_m)$$

$\sum_{j=1}^m f(x_m)$ を 1 から $n-1$ まで $n-1$ から 1 までの δ につい

て加えあげると

$$\sum_{j=1}^{k(n-1)} f(x_j) + \sum_{m=1}^{kn} f(x_m) = \sum_{j=1}^{k(n-1)} f(x_j) + \sum_{m=k(n-1)+1}^{kn} f(x_m) + \dots + \sum_{j=1}^k f(x_j) + \sum_{m=k+1}^{kn} f(x_m)$$

なることがわかる。ここで x_j の和は $\sum_{j=1}^k f(x_j)$ の項で x_m の和は $\sum_{m=1}^{kn} f(x_m)$ の項である。しかるに

$$\sum_{j=1}^{k(n-1)} f(x_j) + \sum_{m=k(n-1)+1}^{kn} f(x_m) = \sum_{v=1}^n f(x_v)$$

であるから Lemma 4 が証明された。

Lemma 5

$$\sum_{i=1}^n f(x_{i+(n-1)k}) f(x_{i+(n-1)k}) = A$$

とする

$$A = n \sum_{v=1}^n [f(x_v)]^2 - \sum_{j=1}^{n-1} [f(x_j) - f(x_{j+k})]^2$$

である。

証明 :

$$A = \sum_{i=1}^n [f(x_{i+(n-1)k})]^2 + 2 \sum_{\substack{i,j \\ i < j}} f(x_{i+(n-1)k}) f(x_{j+(n-1)k})$$

Lemma 3 によって

$$2 \sum_{\substack{i,j \\ i < j}} f(x_{i+(n-1)k}) f(x_{j+(n-1)k}) = 2 \sum_{j=1}^{n-1} f(x_j) f(x_{j+k})$$

また

$$2f(x_j) f(x_{j+k}) = f(x_j)^2 + f(x_{j+k})^2 - [f(x_j) - f(x_{j+k})]^2$$

が成立するから Lemma 4 を用いれば証明が完結する。

Lemma 6 $kn\bar{f} = \sum_{v=1}^n f(x_v)$ とする。このとき

$$A \left(\frac{1}{kn^2} \right) - \bar{f}^2 = \left(\frac{1}{kn} \right) \sum_{v=1}^n [f(x_v) - \bar{f}] - \left(\frac{1}{kn^2} \right) \sum_{j=1}^{n-1} [f(x_j) - f(x_{j+k})]^2$$

である。

証明 : この Lemma は Lemma 5 の直接の結果である。

Lemma 7

$$A = \sum_{v=1}^n [f(x_v) - \bar{f}]^2 + 2 \sum_{j=1}^{n-1} [f(x_j) - \bar{f}][f(x_{j+k}) - \bar{f}] + kn^2 \bar{f}^2$$

が成立する。

証明 : Lemma 4 から

$$A = n \sum_{v=1}^n f(x_v)^2 - \sum_{j=1}^{n-1} \{ [f(x_j) - \bar{f}]^2 + [f(x_{j+k}) - \bar{f}]^2 + 2 \sum_{j=1}^{n-1} [f(x_j) - \bar{f}][f(x_{j+k}) - \bar{f}] \}$$

が導びかれるから。

$$A = n \sum_{v=1}^n [f(x_v) - \bar{f}]^2 + n^2 k \bar{f}^2 - (n-1) \sum_{v=1}^n [f(x_v) - \bar{f}]^2 + 2 \sum_{j=1}^{n-1} [f(x_j) - \bar{f}][f(x_{j+k}) - \bar{f}]$$

なることがわかる。このとき Lemma 4 が証明される。

Lemma 8

$$A \left(\frac{1}{kn^2} \right) - \bar{f}^2 = \left(\frac{1}{kn^2} \right) \sum_{v=1}^n [f(x_v) - \bar{f}]^2 + \left(\frac{2}{kn^2} \right) \sum_{j=1}^{n-1} [f(x_j) - \bar{f}][f(x_{j+k}) - \bar{f}]$$

この Lemma は Lemma 7 の直接の結果である。

Lemma 9 $k > kn$ なら x_k は x_{k-4n} に等しいものとする。

$f(u, v) = f(v, u)$ 。すなわち f は対称とする。

このとき

$$d_{ks} = \sum_j f(x_j, x_{j+ks})$$

とおけば

$$d_{ks} + d_{kn-s} = \sum_v f(x_v, x_{v+ks})$$

が成立する。

証明 : 明らかに

$$\sum_v f(x_v, x_{v+ks}) = d_{ks} + B$$

ここで

$$B = \sum_{j=k(n-s)+1}^n f(x_j, x_{j+ks})$$

いま $k = j - (n-s)k$ とおくと

$$B = \sum_{k=1}^{sk} f(x_{k+(n-s)k}, x_{k+kn})$$

$x_{A+n} = x_A$, $f(x_{A+(n-1)h}, x_A) = f(x_A, x_{A+(n-1)h})$ だから
 $B = d_{h(n-1)}$ となり Lemma が証明される。 $f(u, v)$ の対称性は
 十分条件であるとともに必要条件でもあることに注意せよ。何
 となれば $f(x_v, x_{v+nh}) = x_v - x_{v+nh}$ ならこの定理は成立しな
 い。

附録 C

層化抽出

母集団 P は L 個の層 P_1, \dots, P_L からなっているとす。 P の算術
 平均を \bar{x} とし、 P_a の算術平均を \bar{x}_a とする。 \bar{x}_a の標本推定値を \tilde{x}_a とし、
 $\tilde{x} = \sum_a C_a \tilde{x}_a$ とおく。このとき $E\tilde{x} = \sum_a C_a A_a = A$ 。ここで $E\tilde{x}_a = A_a$
 $\sigma_{\tilde{x}}^2$ は $\sigma_{\tilde{x}}^2 = E(\tilde{x} - \bar{x})^2$ で定義されるものとする。そうすると、
 $\sigma_{\tilde{x}}^2 = E(\tilde{x} - A)^2 + (A - \bar{x})^2$ となって、容易に $\sigma_{\tilde{x}}^2 = \sum_a C_a C_a \sigma_{\tilde{x}_a}^2 + (A - \bar{x})^2$
 なることがわかる。ただし

$$\sigma_{\tilde{x}_a}^2 = E(\tilde{x}_a - A_a)(\tilde{x}_b - A_b)$$

$$(A - \bar{x})^2 = \left[\sum_a (C_a A_a - \frac{N_a}{N} \bar{x}_a) \right]^2$$

で、もし $N C_a = N_a$ なら

$$(A - \bar{x})^2 = \sum_a C_a C_a (A_a - \bar{x}_a)(A_b - \bar{x}_b)$$

かつ $\sigma_{\tilde{x}}^2 = \sum_a C_a C_a \sigma_{\tilde{x}_a}^2$

である。ここで

$$\sigma_{\tilde{x}_a}^2 = E(\tilde{x}_a - \bar{x}_a)(\tilde{x}_b - \bar{x}_b)$$

これらの公式はオシ層の抽出が如何なる方法でなされることと
 成立する。もし \tilde{x} が \bar{x} の不偏推定値で \tilde{x}_a が \bar{x}_a と独立なら普通の公
 式 $\sigma_{\tilde{x}}^2 = \sum_a C_a^2 \sigma_{\tilde{x}_a}^2$ が成立する。勿論 $\sigma_{\tilde{x}_a}^2$ の公式は、無作為、系
 統的、および他の抽出手段のどれを使つたかによって変わら
 ず。

3) ある種の母集団に対する系統的標本と無作為標本
 の相対的精度。

Relative Accuracy of Systematic and Stratified Random Samples for a Certain Class of populations

1) 要約 (summary)

色々なサンプルにおいて層々出現する母集団のタイプは、そ
 の中の要素を群分けしたとき、群が大きくなるにつれて群内分散が
 急速に増加してくるようなものである。この様な母集団の Class
 は、その要素に系統的な関係がある、即ち二つの要素間の相関は正
 で、それらの距離が増加するに従って単調に減少する種数。模型に
 よつて表わすことができる。この型の母集団について、各要素毎の
 系統的標本 (Systematic sample)、層当りノ要素を抽出する層
 化無作為標本 (stratified random sample) 及び単純無作為標
 本 (random sample) の相対的な効率 (efficiency) を比較した。
 平均的には、層化無作為標本は少くとも単純無作為標本より常に精
 度が高く、その相対的効率は標本サイズの増加に伴って増加する単
 調増加種数で表わすことができる。系統的標本の相対的な効率につ
 いては何一つとして普遍性のある結果を与える事は出来ない。実際
 的な抽出率では、系統的標本の方が層化無作為標本より精度の良い
 事もあるが、他の抽出率では逆に悪くなるような母集団の Class
 が存在するのである。しかしもしコレログラム (correlogram)
 が上に凹であれば、任意の標本サイズに対して、平均的には、層化
 無作為標本よりも系統的標本の精度が良い。コレログラムが (i) 線型
 (ii) 指数型 (exponential) の場合について幾つかの数字が得られた。

2) 結言

我々は要素 x_1, x_2, \dots, x_{nh} からなる有限母集団を考へる。こいで

n は整数である。系統的標本を抽出するには、 X_1, \dots, X_n から任意の一要素を選び、それから次々に巻目毎の要素をこつてゆけばよい。即ち、もし最初に選ばれた要素が X_i ならば、求める系統的標本には要素 $X_i, X_{i+n}, \dots, X_{i+(n-1)n}$ が含まれる。このタイプの標本は

(i) 単純無作為標本あるいは層化無作為標本よりも抽出及び取扱いが簡単であること。

(ii) 標本が母集団全体にわたって一様に配置されるという事が直観的に示されること等の理由から、実際的に可成り数多く用いられる。しかし、この系統的標本にも、比較し得る単純無作為もしくは層化無作為標本との相対的精度について研究すべきことが数多く残されている。おそらく最も当を得た比較は系統的標本と層当り/個の要素を抽出する層化無作為標本とを比較することであろう。後者の場合母集団は n 個の層 $\{X_1, \dots, X_n\}, \{X_{n+1}, \dots, X_{2n}\}, \dots$ に分割され、この各層から無作為 (random) が一つ一つの要素を抽出する。この型の標本は多くの点で系統的標本と似通っている。二つとも母集団を n 個の層からなる n 個の層に分割し、各層から一つの要素を抽出する。その上どちらの標本も、少なくともその推定値が要素 X_i のどの標本母集団に対しても不偏 (unbiased) であるという意味においては、平均値のまわりの抽出分散 (sampling variance) の不偏推定法 (unbiased estimate) のデータを手えることはできない。系統的標本の性質に関する最初の完全な研究は W. F. 及び Madocw (1) によって行なわれた。特にこれらの著者は、色々な型の有限母集団に対して、上述の様な系統的標本と層化無作為標本の精度の比較を行なった。母集団の要素が直線 $X_i = c$ の上にある場合層当り/要素の層化無作為標本のほうが系統的標本より精度の高い事が示された。もし母集団が周期分布をすれば、 n が周期の整数倍のときは層化無作為標本のほうがずっと優れているが、 n が周期の2分の1の奇数倍のときは系統的標本のほうが優れている。

著者は、母集団に傾向 (trend) 及び周期函数 (periodic function) が同時に含まれている様なより複雑な場合についても考察した。この論文の目的は、色々なサンプリングで非常に多く見られる他のタイプの母集団について、同様な比較を行なうことである。この母集団は隣接する要素を群にまとめたとき、群の大きさが増すにつれて群内の要素間分散が急激に増加する様なものである。これまで長い間この型の母集団はブロック内のプロット間分散が、ブロックの大きさとともに増大する様な圃場実験の研究に適用できるものと考えられてきた。40 回の斉一性試験 (uniformity trial) から得たデータを要約し Fairfield Smith [2] はこの考え方を証明し、増加率を推定することのできる実験的な関係を導いた。最近の広範な標本調査のいくつかに対する論文にも、同様な母集団が考察されている。このようにして Jessen は農場の母集団 (farm population) のサンプリング法の研究で grid 内の農場間の分散は grid サイズの大きさの単調増加函数になるという法則を見出し、抽出単位 (sampling unit) 中に含まるべき最適な農場数を推定するのみに用いた。Mahalanobis [4] は大規模標本調査 (large scale sample survey) の広範な研究の中で Fairfield Smith と同じ法則を独立に展開した。Hansen と Hurwitz [5] は、多くの現実的母集団の典型的なものでは、集落 (cluster) 内の分散が集落サイズが大きさとともに増大するということを指摘している。その他にも多くの参考文献をあげることができる。

2. 母集団の明確化

群内分散が群の大きさが大きくなるにつれて増加する状態を表わすには色々な数学的模型を設けることができる。例えば c に対して要素 X_i が規則的に変化している様な色々な母集団から抽出された要素 X_i を考えることができる。逆に c は同一母集団に属するが、それ

には系列的相関がある (serial correlated) と考える事も出来る。簡単の爲に、これから \$X_i\$ と \$X_{i+u}\$ の間の系列相関は \$u\$ のみに依存するある量 \$f_u\$ であると仮定する。そうすると \$f_u\$ が正で、\$u\$ の単調減少函数であれば直観的に (後で証明する) 要素 \$X_i, X_{i+1}, \dots, X_{i+k}\$ からなる群内分散は \$k\$ の単調増加函数となる事が予想される。この模型は我々の目的によくあてはまる様に思われる。というのは、多くの着着が分散増大現象の根拠としてあげているのは、\$X\$ の間の相関が常に正となることであるからである。上記の特徴づけはある点に関して批判を受ける。普通の大きさの実際の有限母集団で、\$f_u\$ が厳密に単調であるという仮定は現実的でない様に思われる。一オコレログラムにはっきり下降の傾向 (trend) が現われている場合でも、この傾向のまわりの個々のばらつきによってコレログラムが厳密には単調といえない事がある。更に適切であると思われるのは、有限母集団自体を単調な \$f_u\$ をもつ無限母集団からの標本と考えることである。

私はこのような態度は上に引用した着着達の考えと一致していることを信ずるものである。論文を理解したところによれば、彼等は分散法則 (variance law) が理想化された母集団において成立するものと考えている。このように、系統的標本と層化無作為標本の比較は、一つの有限母集団のみについて行なうのではなく、単調減少の \$f_u\$ を有する無限母集団から抽出された有限母集団の平均について行なわれるのである。個々の有限母集団に対する結果は、それらの中の \$\gamma\$ が、期待値 \$f_u\$ のまわりで変動するが、平均的な結果とはかけ離れたものになる。有限母集団が大きくなるに従い、結果は平均的な結果に次第に接近するようになる。このことから要素 \$X_i\$, \$i = 1, 2, \dots, n\$, は

$$E(X_i) = \mu \quad E(X_i - \mu)^2 = \sigma^2$$

$$E(X_i - \mu)(X_{i+u} - \mu) = f_u \sigma^2$$

ここで \$u < v\$ なるときは常に \$f_u > f_v > 0\$ なる母集団から抽出するものと仮定する。

4. 有用な幾つかの予備公式

特定の有限母集団の平均値を \$\bar{x}\$ とすると、分散分析で屡々有用である次の等式を容易に証明することができる。

$$(1) (kn) \sum_{i=1}^{kn} (X_i - \bar{x})^2 = \sum_{i,j}^{kn} (X_i - X_j)^2$$

(\$X_i, X_j\$) には \$(kn)(kn-1)/2\$ 個の対をとることができるから上式から

$$(2) \sum_{i=1}^{kn} (X_i - \bar{x})^2 = \frac{(kn-1)}{2} E(X_i - X_j)^2 = \frac{(kn-1)}{2} E\{(X_i - \mu) - (X_j - \mu)\}^2$$

ここで \$E\$ は有限母集団の全体にわたってとるものとする。いま二次式を展開してすべての有限母集団について平均する。\$(kn)(kn-1)/2\$ 個の組合せ (combination) 中には、\$j\$ が \$i\$ より 1 だけ多いものが \$(kn-1)\$ 個、\$j\$ が \$i\$ より 2 だけ多いものが \$(kn-2)\$ 個等々となっているから

$$(3) E \sum_{i=1}^{kn} (X_i - \bar{x})^2 = (kn-1) \sigma^2 \left\{ 1 - \frac{2}{(kn)(kn-1)} \sum_{u=1}^{kn-1} (kn-u) f_u \right\}$$

これに対応して、\$k\$ 個の連続した要素からなる一つの層内平方和の期待値をとるには、(3) の \$(kn)\$ の代りに \$(k)\$ を入れればよい。この結果は \$n\$ 個の層の何れにおいても同じであるから

$$(4) E(\text{層内の } S.S) = n(k-1) \sigma^2 \left\{ 1 - \frac{2}{k(k-1)} \sum_{u=1}^{k-1} (k-u) f_u \right\}$$

を得る。

公式(4)で \$(kn)\$ の代りに \$n\$, \$k\$ の代りに \$(kn)\$ を用いれば同様の計算から、特定の系統的標本内の平方和の期待値が得られる。

何と云へば、標本内には \$n\$ 個の要素があり、連続した要素間の相関は \$f_1, f_2, \dots\$ ではなく、\$f_k, f_{k+1}, \dots\$ である。この結果は \$k\$ 個の系統的標本の何れについても同じである。故に

$$(5) E(\text{系統的標本内の } S.S) = k(n-1) \sigma^2 \left\{ 1 - \frac{2}{k(n-1)} \sum_{u=1}^{n-1} (n-u) f_{ku} \right\}$$

5 無作為標本に対する平均的分散

有限母集団の平均に対する単純無作為、層化無作為、系統的標本の平均的分散を表わすのに σ_r^2 , σ_{st}^2 , σ_{sy}^2 なる記号を用いる。

この平均はこれまでの節で詳述した無限母集団からくり出したすべての有限母集団にわたってとられたものである。

無作為標本との比較は、こゝでの主目的ではないが興味ある問題であるので一語に論ずることとする。

1つだけの有限母集団の無作為標本に対する平均値の分散は多くの着者によって

$$(6) \frac{1}{n} \cdot \frac{(kn-n)}{(kn-1)} \cdot \frac{1}{kn} \sum_{i=1}^{kn} (X_i - \bar{X})^2$$

であることが示されている。こゝで \bar{X} は有限母集団の平均値である。(3)から

$$(7) \sigma_r^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{1 - \frac{2}{(kn)(kn-1)} \sum_{u=1}^{kn-1} (kn-u) f_u\right\}$$

6. 層化無作為標本の平均的分散

層化無作為標本の平均値を \bar{X}_{st} とすれば定義から \bar{X}_{st} の抽出分散 (sampling variance) は

$$(8) E(\bar{X}_{st} - \bar{X})^2$$

である。

まず1つの有限母集団のみに対するこの値の平均を考へる。 n 層の層の平均値を $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n$ とし、また各層から選ばれた要素を $X_{1j}, X_{2j}, \dots, X_{nj}$ とする。このとき

$$\sum_{j=1}^n X_{ij} = n\bar{X}_{st} \quad \text{で} \quad \sum_{i=1}^n \bar{X}_i = n\bar{X}$$

だから(8)は

$$(9) \frac{1}{n^2} E \left\{ (X_{1j} - \bar{X}_1) + (X_{2j} - \bar{X}_2) + \dots + (X_{nj} - \bar{X}_n) \right\}^2$$

とかける。

この有限母集団から得られる kn 個のすべての標本についての平均をとる。Cross product term はすべて消える。何故なら、例えば、 X_{ij} は $X_{21}, X_{22}, \dots, X_{2k}$ と同じ回数だけ現れるから、従って単一の有限母集団に対する分散は

$$(10) \frac{1}{kn^2} \sum_{i=1}^n \sum_{j=1}^k (X_{ij} - \bar{X}_i)^2$$

となる。この平方和は勿論単なる層内の平方和である。よって(4)から

$$(11) \sigma_{st}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{1 - \frac{2}{k(k-1)} \sum_{u=1}^{k-1} (k-u) f_u\right\}$$

7. 系統的標本の平均的分散

\bar{X}_{sy} を代表的な系統的標本の平均値とするとき、単一の有限母集団の分散は

$$(12) E(\bar{X}_{sy} - \bar{X})^2 = \frac{1}{kn} \left\{ n \sum (\bar{X}_{sy} - \bar{X})^2 \right\}$$

である。こゝで和は各層の系統的標本のすべてにわたる。標本間平方和は、(母集団の全平方和 - 標本内平方和) に等しいから、(12)は

$$(13) \frac{1}{kn} \sum_{i=1}^{kn} (X_i - \bar{X})^2 - \frac{1}{kn} (\text{系統的標本内SS})$$

に等しい。

すべての有限母集団について平均をとるための、オ1項及びオ2項にそれぞれ(13)と(5)を代入して

$$(14) \sigma_{sy}^2 = \frac{(kn-1)}{kn} \sigma^2 \left\{ 1 - \frac{2}{(kn)(kn-1)} \sum_{u=1}^{kn-1} (n-u) f_u \right\} - \frac{(n-1)}{n} \sigma^2 \left\{ 1 - \frac{2}{n(n-1)} \sum_{u=1}^{n-1} (n-u) f_u \right\}$$

を得る。

変形して

$$(15) \sigma_{sy}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{ 1 - \frac{2}{k(k+1)} \sum_{u=1}^{kn-1} (kn-u) f_u + \frac{2k}{n(k-1)} \sum_{u=1}^{n-1} (n-u) f_u \right\}$$

上述の式と記号は Madaw が用いたものと違っていることに注意。
Madaw は f と σ^2 を単一の有限母集団について考え、単一の有限母集団の抽出分散 (sample variance) を論じた。

8. 無作為標本と層化無作為標本の相対的な精度

まず残った一般的なことを述べる。

(7), (11), 及び (15) の三種類の標本の相対的な効率率は $\sigma_{st}^2, \sigma_{sy}^2$ に現れる f の一次関数にのみ関係するにわがわが違った。どの場合についても f の係数の和が 1 に等しいことは容易に証明できる。

無作為標本の場合この一次関数には、 \log が増加するにつれて直線的に減少し、標本の大きさとは無関係にただ有限集団の要素数 $N = (kn)$ にのみ依存する $\log(kn-1)$ までのすべての系列相関係数 (Serial correlation coefficient) が含まれる。

層化無作為標本では $\log(k-1)$ までの系列相関 (serial correlation) が現れる。ここで k は層内の要素数である。(15) に示した様に系統的標本に対する公式は二つの一次関数 (linear function) に分解される。

第1のものは、すべての係数が $(kn-1)/(k-1)$ 倍になっているを除外せば無作為標本に出てくるものと同じである。2番目は正の符号をもち、 \log が k の倍数であるような相関 (correlation) が含まれる。従ってこの公式には f についての制限は全然不要である。 f が正で単調減少の場合を考えれば次の Lemma が有用である。

Lemma $f_i (i=1, \dots, m)$ が正で、単調減少即ち

$$f_i > f_{i+1} > 0 \text{ で } (\alpha_1 + \alpha_2 + \dots + \alpha_m) = 0 \text{ ならば}$$

(16) $L = \alpha_1 f_1 + \alpha_2 f_2 + \dots + \alpha_m f_m > 0$ (可能な f のすべての組について) なる為の必要かつ十分な条件は

$$(17) \alpha_1 + \alpha_2 + \dots + \alpha_i > 0 \quad i = 1, 2, \dots, (m-1)$$

なることである。

(証明) 何故なら $f_i = f_{i+1} + \delta_i$ (ここで仮定から $\delta_i > 0$) とおくと $\delta_1, \delta_2, \dots, \delta_{m-1}$ の代りに次に f_1, f_2, \dots, f_{m-1} を代入してゆくと

$$(18) L = \alpha_1 \delta_1 + (\alpha_1 + \alpha_2) \delta_2 + (\alpha_1 + \alpha_2 + \alpha_3) \delta_3 + \dots + (\alpha_1 + \alpha_2 + \dots + \alpha_{m-1}) \delta_{m-1}$$

と成って、 $(\alpha_1 + \dots + \alpha_m) = 0$ だから最後の項 f_m が消える。

すべての δ_i について > 0 だから (17) が十分条件であることは明らか。またもし任意の i に対する δ_i の係数が負なら、 δ_i と正の値をとり δ_i を 0 とおくことによつて L を負にすることができる。これで必要性の条件が得られた。

系 (corollary). f_i が強い意味で単調即ち $f_i > f_{i+1}$ で α_i の少くとも一つが 0 でなければ (17) は L が 0 以上であるための十分条件となる。

(証明) なぜなら (18) の δ はすべて 0 より大であるから (17) によつて δ_i の係数はすべて正である、また δ の少くとも一つの係数は 0 より大でなければならぬが他のものは全部 0 と成ってもよい。
故に $L > 0$

さて f_u が単調減少のとき

$$(19) L(k) = \frac{2}{k(k-1)} \sum_{u=1}^{k-1} (k-u) f_u$$

が k の単調減少関数であることを示そう。

これは層化無作為標本の分散中に現れる一次関数である。

$$(20) L(k) - L(k+1) = \frac{2}{k(k-1)} \sum_{u=1}^{k-1} (k-u) f_u - \frac{2}{(k+1)k} \sum_{u=1}^k (k+1-u) f_u$$

$$(21) \quad \frac{2}{k(k^2-1)} \sum_{u=1}^k (k+1-2u) f_u$$

$L(k)$ と $L(k+1)$ の f_u の係数の和は 1 だから (21) の和は 0 である。従つて Lemma が適用できる。しかし係数は $u \leq (k+1)/2$ についてはすべて正で、 $u > (k+1)/2$ についてはすべて負だから、(21)

の係数の最初の \$i\$ 個の和が 0 よりも大きいことは明らかである。

故に

$$(22) \quad L(k) - L(k+1) > 0$$

さらに、系から、もし \$f_u\$ が強い意味で単調なら、\$L(k)\$ は強い意味で単調である。\$f_u\$ は負だからこの結果によって

$$(23) \quad 1 - \frac{2}{k(k-1)} \sum_{u=1}^{k-1} (k-u) f_u \leq 1 - \frac{2}{(nk)(nk-1)} \sum_{u=1}^{nk-1} (nk-u) f_u$$

を証明できる。

結局任意の大きさの標本について、層化無作為標本の平均的な分散は無作為標本のそれよりも大きくはならない。また、層化標本の無作為標本に対する効率 (efficiency) は層の大きさが小さくなるにつれて、即ち標本の大きさが大となるにつれて単調に増大する。勿論これらは予期されなかつたものではない。

(22) の方程式はまた、第 3 節で注意した単調減少な \$f\$ についての結果即ち層内の平均的な分散は、層の大きさが増すにつれて急速に増大するという性質をも説明するものである。何故なら、層内の平方和の自由度が \$n(k-1)\$ なら、上の式(4) は層内の平均的な分散が、

$$(24) \quad \sigma^2 \left\{ 1 - \frac{2}{k(k-1)} \sum_{u=1}^{k-1} (k-u) f_u \right\} = \sigma^2 \{ 1 - L(k) \}$$

であることを示す。

9. 系統的標本と無作為標本との比較

母集団の形について立入った限定を設けない限り、無作為標本と比較したときの系統的標本の効率について一般的な結果が得られないことは上述の研究から明らかである。Lemma を応用するため、公式(7), (11), (15) に現れる \$f\$ の一次指数の最初の \$i\$ 個の係数の和を求めてみる。簡単な方法でこれらの和が

$$\sum r = \frac{i(2nk-i-1)}{nk(nk-1)},$$

$$(25) \quad \sum_{st} = \frac{i(2nk-i-1)}{k(k-1)}, \quad 1 \leq i \leq (k-1)$$

$$1, \quad i \geq k$$

$$\sum_{sy} = \frac{i(2nr-i-1)}{nr(k-1)} - \frac{rk(2n-r-1)}{n(k-1)}$$

であることがわかる。ここで \$r\$ は \$(r+1)k \leq nk\$ なる整数。

Lemma から、\$\sigma_{sy}^2 \leq \sigma_{st}^2\$ ということには \$\sum_{sy} > \sum_{st}\$ が各 \$i\$ についてなり立つことを示せばよい。いまもし \$i < k\$ なら \$r\$ は 0 だから \$n=1\$ のとき以外は

$$(26) \quad \sum_{sy} > \sum_{st} > \sum r \quad i = 1, 2, \dots, (k-1)$$

\$n=1\$ のときはこの三つはすべて等しい。

しかし \$i\$ が \$k\$ の整数倍 \$rk\$ なる場合には

$$(27) \quad \sum r = \frac{r}{n} \left[1 + \frac{(n-r)k}{(nr-1)} \right], \quad \sum_{st} = 1, \quad \sum_{sy} = \frac{r}{n}$$

なることがわかるから

$$(28) \quad \sum_{st} > \sum r > \sum_{sy}$$

である。

結局 Lemma の条件は systematic sample に対しては満足されないから単調減少な \$f\$ をもつ母集団全体に対する一般的定理は存在しない。(26) と系から、\$f_u = 0, u > (k-1)\$ なる class に含まれるような任意の母集団に対しては、系統的標本は層化無作為標本よりも明らかに効率がよいことがわかる。一方(28)によれば、\$f\$ の最初の \$i\$ 個の係数が等しく残りが 0 であるような母集団については、系統的標本の分散は無作為標本のものより大となる。

これらの二つの結果を一語にすると、最初の \$j\$ 個の \$f\$ が等しく残りが 0 であるような母集団では大きさ \$j\$ の層をもつ系統的標本は比較し得る無作為標本より精度が悪く、また \$(j+1)\$ なる大きさの層をもつ標本は対応する層化無作為標本より精度がよい。その異なる母集団が実際に出てくる争いはないとしても、この結果は系統的標本の平均値の分散の標本の大きさに対するグラフ (graph) が、無作為標

本と同じ様な規則性は表わしそうにないことを示している。

10. コレログラムが上に凹る母集団

更に研究すれば系統的標本と無作為標本の相対的精度を定める決定的な因子は、一次階差 (first difference) ではなく、むしろ二次の階差 (second difference) であることが示される。

次の結果を証明してみよう。

定理

$$f_i > f_{i+1} > 0, \quad i = 1, 2, \dots, (kn-1)$$

で

$$\delta_i^2 = f_{i-1} + f_{i+1} - 2f_i > 0 \quad i = 2, 3, \dots, (kn-2)$$

なるようなすべての無限母集団については、任意の大きさの標本に対して

$$\sigma_{sy}^2 \leq \sigma_{st}^2 < \sigma_r^2$$

なり立つ。

また $\delta_i^2 = 0, \quad i = 2, 3, \dots, (kn-2)$ でない限り

$$\sigma_{sy}^2 < \sigma_{st}^2$$

但し $i = 2, 3, \dots, (kn-2)$

この結果は f_u の一次関数を二次の階差で表わし、二次階差に適用できる新しい Lemma を構成することによって証明できる。

もう一つの方法は、より簡単でまたおそろくよりわかりやすい。

f_u は単調減少だから section 8 の結果から $\sigma_{st}^2 \leq \sigma_r^2$ 。上の (12) において特定の有限母集団の系統的標本の平均値の分散は、

$$\frac{1}{kn} \sum_{i=1}^{kn} (X_i - \bar{X})^2 = \frac{1}{kn} (\text{系統的標本内の全 S.S.})$$

$$(29) \quad = \frac{1}{kn} \sum_{i=1}^{kn} (X_i - \bar{X})^2 - \frac{1}{n} (\text{系統的標本内の平均的 S.S.})$$

と表わすことができる。

層化無作為標本についてもこれに対応する方程式がなり立つ。というのはもし平均値 \bar{X}_{st} をもつ任意の層化無作為標本の要素を X_{ij} ,

X_{2j}, \dots, X_{nj} とすると

$$(20) \quad \sum (X_{ij} - \bar{X})^2 = \sum_{i=1}^n (X_{ij} - \bar{X}_{st})^2 + n(\bar{X}_{st} - \bar{X})^2$$

こゝで n 個のすべての標本について平均すれば、

$$(31) \quad \frac{1}{k} \sum_{k=1}^{kn} (X_k - \bar{X})^2 = (\text{平均的標本内 S.S.}) + nE(\bar{X}_{st} - \bar{X})^2$$

一番右の項は層化無作為標本の分散の n 倍だからこゝでも (29) と同様な結果が得られる。

結局系統的標本の平均的な「組内 (within)」平方和が層化無作為標本の「組内 (within)」平方和より大きいかあるいは等しい場合には

$$\sigma_{sy}^2 \leq \sigma_{st}^2$$

となる。いま (2) から、 (kn) の代りに n をおきかえればこれらのそれぞれの平均は

$$(32) \quad \frac{(n-1)}{2} E(X_{ij} - X_{lj})^2$$

に等しい。こゝで X_{ij} 及び X_{lj} はそれぞれ i 層及び l 層がうとった標本の要素で、平均は層の対のすべての可能な組についてとるものとする。

我々は層の対を固定して考え $l-i = u$ とおくことにする。系統的標本で i 層及び l 層の対にする要素は常に (kn) 要素だけ離れておから、

$$(33) \quad E_{sy} (X_{ij} - X_{lj})^2 = 2\sigma^2(1 - \rho_{kn})$$

層化無作為標本では二つの層の要素を一つ宛組合せる仕方のが数だけある。 $(kn-k+1)$ 要素離れている対は一つで $(kn-k+2)$ 要素だけ離れている対は二個等々……である。

対の個数は k まで直線的に増加しそれ以上では直線的に減少して $(kn+k-1)$ だけ離れた最後の対の間数は 1 個となっている。こゝから

$$(34) \quad Est (X_{ij} - X_{lj})^2 = 2\sigma^2 \left\{ 1 - \frac{1}{k^2} \sum_{i=-(k-1)}^{(k-1)} (k-|i|) \rho_{kn+i} \right\}$$

故に $\sigma_{sy}^2 < \sigma_{st}^2$ を完全に証明するには

$$(35) \sum_{i=(k-1)}^{(k-1)} (k-i) f_{kui} - k^2 f_{ku} > 0$$

を $u = 1, 2, \dots, (n-1)$ 即ち層の任意の粗に対して示せば十分である。これは

$$(36) \sum_{i=1}^{(k-1)} (k-i) (f_{kui+1} + f_{kui} - 2f_{ku}) > 0$$

とかけらる。しかし $\delta_{ku}^2 = f_{k,u-1} + f_{k,u+1} - 2f_{ku}$ を δ = 中央階差 (second central difference) とすると容易に

$$(37) f_{kui+1} + f_{kui} - 2f_{ku} = \sum_{j=-(i-1)}^{(i-1)} (i-|j|) \delta_{kuj}^2 > 0$$

を証明できる。何故なら仮定から $\delta_j^2 > 0, j = 2, 3, \dots, (kn-2)$ だから。このことは、任意の固定された層の対に対しては、系統的標本の要素間の分散が層化無作為標本の要素間分散よりも大きいから、等しいことを示している。これから全平均に対する結果も得られる。故に、 $\sigma_{sy}^2 < \sigma_{st}^2$ 、また更に標本が1つ丈の場合を除いて、 $\sigma_j^2 = 0$ でない限り明らかに $\sigma_{sy}^2 < \sigma_{st}^2$ である。

証明の本質的な部分は次の様である。i番目及びj番目の層の要素は、系統的及び層化無作為標本の何れにおいても、(ku)要素丈離れている。層化無作為標本の二つの要素が (ku+i) 要素だけ離れているとき、これらの相関は平均以下である。何故なら、 $f_{kui+1} < f_{ku}$ であるから、これはより多くの独立な情報を与えるからである。この要素間の分散は系統的標本の分散より $2\sigma^2(f_{ku} - f_{kui})$ だけ大きい。しかしこの様なことの起るのは要素が (ku-i) 丈離れていて分散が系統的標本のものより $2\sigma^2(f_{ku-i} - f_{ku})$ だけ小さいなる場合と同数だから釣合がとれる。f_{ku}が凹だから平均的な損失は釣合がとれるかあるいは利益の方が多い。

Section 9で論じた $f_u = f, u = 1, 2, \dots, j, f_u = 0, u > j$ なる母集団では、 $\delta_j^2 < 0, \delta_{j+1}^2 > 0$ 、他のuについては $\delta_u = 0$ が得られる。

第2の差の符号が反っていることは層の大きさがjとj+1の系統的標本の異なる性質を説明するものである。

上の定理は層化無作為標本に対する系統的標本の相対的精度が、nの単調函数で表わされることを証明するものでもなく、また、 σ_{sy}^2 がnの増加につれて減少するということすら証明するものではない。更だ何れの結果も成立たない様な母集団も存在する。これを次のsectionで述べることにする。実践的な応用を考えている場合には、 f_u が上に凹でなければならぬという様な条件が満たされないかもしれない。例えばこの条件はコレログラムが直線的 (linear) 即ち $f_u = (L-u)/L$ なる場合には満足される。この f_u の形は、Wald [6]が経済のデータに適用し得ることを考えた型のコレログラムである。凹になるという性質はOsborn [7]が森林及び土地利用調査に対して提案した $f_u = 0^{u/L}$ なる函数や、FisherとMarkensilが二つの観測所で測った通毎の雨量の相関を観測所間の距離を用いて表わした $f_u = \tan A (u^{1/2})$ なる関係式についても成立つ。実際もし f_u が正しくして表わされ、すべてのuについて連続なら、上に凹な函数はそれ自身自然的なものであることを示す。(naturally)

11. 直線的なコレログラム (linear correlogram)

コレログラムが(i) 直線的 (ii) 指数的 (exponential) な場合に得られる幾つかの結果を述べることは興味ある問題である。こののは、この両方の型は実際に出現する可能性の多い母集団だといふことがいわれているからである。

直線 (linear) の場合

$$(38) f_u = (L-u)/L, u \leq L; f_u = 0, u > L$$

もし $L > (nk-1)$ なるコレログラムは有限母集団の全範囲を通じて直線である。この場合には2次階差はすべて0だから、

$$\sigma_{sy}^2 = \sigma_{st}^2 < \sigma_y^2$$

なることが期待される。

$L < (nk-1)$ なるすべての第2階差は、正である σ_{st}^2 を除いては

すべて0となる。故に $\sigma_{yt}^2 < \sigma_{st}^2 < \sigma_r^2$ これらの各々の場合における結果は、(7), (11), (15) の基礎公式を初等的な方法で加え上げれば得られる。和をとることは詳しく述べない。

$L > (nk-1)$ なる場合は

$$(39) \sigma_{yt}^2 = \sigma_{st}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \frac{k+1}{jL} ; \sigma_r^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \frac{1}{jL} (nk+1)$$

となる。比 $\sigma_r^2 / \sigma_{yt}^2$ は $(nk+1)/(k+1)$ だから抽出比が大きい場合の外はこれは大体標本の大きさに等しい。このようにすべての無作為標本を上まわる非常に大きい効率 (efficiency) の利得 (gain) が得られる。

$L < (nk-1)$ なる公式はより複雑となる。まず k, L を考える。即ち抽出される割合が 10% 以下の場合である。 $N = nk$ なら

$$(40) \sigma_r^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{ \frac{jN(N-L) + (L^2 - 1)}{jN(N-1)} \right\}$$

$$(41) \sigma_{st}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{ \frac{jk(k-L) + (L^2 - 1)}{jk(k-1)} \right\} \quad k > L$$

$$(42) \sigma_{yt}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{ \frac{jN(k-L) + (L^2 - 1)}{jN(k-1)} \right\} \quad k > L$$

$\sigma_{yt}^2 < \sigma_{st}^2$ なることは明らか、また層化無作為標本に対する系統的標本の効率は標本の大きさが増すにつれて急速に増大することは容易に証明できる。

標本の大きさがもっと大きくなると $k \leq L$ となって (40) はそのまゝ成立する。一方、 σ_{st}^2 はこゝでは (39) と同じ公式で与えられる。 σ_{yt}^2 に対する公式はもっと複雑である。

L を k で割ったときの商の整数部分を g とし、残りを r とする。 $L = (gk + r)$ となるから、公式は

$$(42) \sigma_{yt}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{ \frac{gk(k^2-1) + jkr(n-g)(k-r) + r(r^2-1)}{jNk(k-1)} \right\}$$

中括弧内の分子の最後の二項は、 L が丁度 k で割り切れる時は0となることに注意せよ。また第二項の order は $nk = N$ で今の場合に

は、第1次より大きい影響を与えている。故に σ_{yt}^2 は L が k の倍数となるとき急激に低下する。実際 $L = gk$ なるときは (42') は

$$(42') \sigma_{yt}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \frac{(k+1)}{jN} \quad L = gk$$

となるから N が十分大きいときこの分散は0になる。 σ_{yt}^2 に対する公式 (39) の比較によって、 $L = gk$ なる場合には系統的標本と無作為標本の効率の比率は N/L であることがわかる。これは N が十分大きいときには無限に増加する。普通 r が0でない場合には、 N が大きいときの相対的効率の主要項は $(k^2-1)/kr(k-r)$ である。これは L と k の間の関係に従って幾分定期的に変動する。

説明のため $L = 10$ で、有限母集団が十分大きく、従って、 $1/n$ の項を無視し得る場合の数値を下に示す。

16. σ_{yt}^2 なる量は係数 σ^2/N を別にし、対応する分散を表わす。層化標本の分散は抽出比が大きくなるにつれて急速に減少する。一方系統的標本の分散は $k = 2, 5, 10$ で0に、相対的効率は ∞ となる。こうして中間の場合、即ち $k = 3, 4, 6, 7, 8, 9$ では分散及び相対的効率は抽出比に全然満足な関係を示さない。10% 以下の標本では、表に表してはいない様な場合も含めて、相対的効率が $k = 11$ の4から k が大きいときの1まで連続的に減少する。

12. 指数的コレログラム (Exponential correlogram)

指数函数 (exponential) $f(x) = e^{-\lambda x}$ の場合には結果はもっと規則的である。 f の一次函数 (linear function) がそれぞれ $(1-x)^{-2}$ の展開式の有限個の項からなっている。もし σ_r^2 の和が

$$(44) f(N, \lambda) = \frac{2}{N(N-1)} \left\{ \frac{(N-1)e^{-\lambda} - N + e^{-(N-1)\lambda}}{(e^{-\lambda} - 1)^2} \right\}$$

なら

$$(45) \sigma_r^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{ 1 - f(N, \lambda) \right\}$$

$$(46) \sigma_{st}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{1 - f(k, \rho)\right\}$$

$$(47) \sigma_{sy}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left\{1 - \frac{(N-1)}{(k-1)} f(N, \rho) + \frac{k(n-1)}{(k-1)} f(n, \rho)\right\}$$

なることがわかる。

標本の大きさが増すにつれて、系統的標本の分散は直線的に減少し同時に層化抽出 (stratified sampling) に対する相対的効率 (relative efficiency) も直線的に増大する。

効率の利得 (gain) の大きさを知るため反復が k が大きい場合をよめる。この場合相対的効率は実際には k, n, ρ の函数であるが、これは殆んど全部が (k, ρ) なる量のみによってきまる。即ち系統的標本の次々の層の調査対象 (item) 間の相関 (correlation), ρ のみ依存する。もし $\rho = (k, \rho)$ とすると、 $\sigma_y^2 = \sigma^2/n$,

$$(48) \sigma_{st}^2 = \frac{\sigma^2}{n} \left\{1 - \frac{2}{k} + \frac{2}{k^2} - \frac{2\rho}{k^2}\right\}$$

$$(49) \sigma_{sy}^2 = \frac{\sigma^2}{n} \left\{1 - \frac{2}{k} + \frac{2}{(k^2-1)}\right\}$$

を得る。

Table 1. コレログラムが直線の場合の因数 σ^2/n を除いた分散と系統的 (systematic) と層化無作為標本の相対的効率 (relative efficiency)

k	2	3	4	5	6	7	8	9	10	11	20
% sampled	50	33	25	20	17	14	12	11	10	9	5
V_{st}	10	.27	.50	.80	1.17	1.60	2.10	2.67	3.30	4.00	11.65
V_{sy}	0	.20	.40	0	.80	1.20	1.20	.80	0	1.00	10.00
V_{st}/V_{sy}	∞	1.33	1.25	∞	1.46	1.33	1.75	3.33	∞	4.00	1.16

Table 2 に幾つかの ρ に対する相対的効率を手立てである。 ρ は次々の層の調査対象 (item) 間の相関である。

この相対的効率には ρ が 1 に近づく (tend) するときの 2 なる極値がある。さうしてこれは ρ が 0 となるときの 1 にまでゆっくりと減少する。効率の利得 (gain) は ρ が $1/2$ をこえる場合には相当大である。

Table 2

ρ	.9	.8	.7	.6	.5	.4	.3	.2	.1
$\sigma_{st}^2 / \sigma_{sy}^2$	1.96	1.90	1.84	1.78	1.71	1.64	1.55	1.46	1.33

単一の標本からは系統的標本でも層化無作為標本でも誤差の予備推定値は得られないということを section 1 で注意しておいた。

これは誤差の推定が全然できないということではない。しかし任意の推定値は、抽出される母集団の形についてのはっきりした仮定にもとづいていなければならぬ。これらの仮定が誤っている場合には推定値が全然異なることになる。

例えば、コレログラムを指数型 (exponential) と仮定すると、 n, k が大きいとき公式 (47) あるいは (49) は 1 つの系統的標本から誤差を推定する場合の適切な基礎となり得る。系統的標本の次々の調査対象 (item) 間の相関 (correlation) は ρ の推定値となるから、この推定値が得られる。またもし $1/n$ の項が無視し得るなら、系統的標本の粗内 (within) の平均平方 (mean square) は σ^2 の不偏推定値であることがわかる。

指数型 (exponential) という仮定が正しいときは (49) にこれを代入して単一の系統的標本の分散の一致推定量 (consistent estimate) が得られる。さうして層化及び単純無作為抽出とこえる効率の利得も推定することが出来る。

By. W. G. Cochran A.M.S. Vol. 17 (1946)

引用文献

- [1] W. G. and L. K. Madow, On the theory of systematic sampling, *Annals of Math. Stat.*, Vol. 15 (1944), pp. 1-24.
- [2] H. Fairfield Smith, An empirical law governing soil heterogeneity, *Jour. Agr. Sci.*, Vol. 28 (1938), pp. 1-23.
- [3] R. J. Hessebm Statistical investigation of a sample survey for obtaining farm facts, *Sov. Agr. Exp. Sta., Res. Bull.*, No. 304, (1942).
- [4] R. C. Mahalanobis, On large scale sample surveys, *Roy. Soc. Phil. Trans.*, B. 231 (1944), pp. 329-451.
- [5] M. H. Hansen and W. N. Hurwitz, On the theory of sampling from finite populations, *Annals of Math. Stat.*, Vol. 14 (1943), pp. 333-362.
- [6] H. Wald, A study of the analysis of stationary time series, *Uppsala* (1938).
- [7] J. G. Osborne, Sampling errors of systematic and random surveys of covertype areas, *Amer. Stat. Assoc. Jour.*, Vol. 37 (1942), pp. 256-270.
- [8] R. A. Fisher and W. A. Mackenzie, The correlation of weekly rainfall, *Quart. Jour. Roy. Met. Soc.*, Vol. 48 (1922), pp. 234-245.

4. 系統的抽出の理論について

On the theory of systematic sampling. II

1. 要約および序

前論文 (I) では系統的抽出の問題を研究し、関連する分散をえた。母集団の型を幾つか想定し、これらについて系統的抽出計画の効率を単純および層化無依為標本の効率を比較した。系統的抽出計画では無依為抽出計画でのものと反対に、標本の大きさが増すと分散が増加することがわかった。この可能性は (II) において正しいことが証明された。

Chochran (3) によって与えられた系統的抽出計画の一つの研究手法は、母集団の要素を確率変数と仮定し、期待分散の一つの問題を転換することによってこの難点をある程度まで解決した。彼はこれらの確率変数のコレログラムが上に凸なら、系統的計画の期待分散は層化計画の分散より小となり、また多くの場合相当小さくなるということを示した。

この論文では前論文の結果を大きさが等しい場合、および等しくない場合の集落の系統的抽出に拡張することにする。二次元における系統的抽出についても若干の説明を行なった。

Sec. 2 ではサンプリングの多くの分野でかなり広く応用できる二つの定理を導いた。サンプリング理論の研究者の間では、これらの定理の成立は当然とされていたものであるが、しかしこれについて言及している文献は今まで存在しなかった。

Sec. 3 では、標本調査の設計の際、層間に負の相関を導くようにすべきだという注意 (I. P. 13) の内容を明らかにした。定理 1 では相関が負となるための十分条件を求めた。Sec. 4 の Lemma および定理 4 によつて、定理 3 が実際の場合に適用できることになる。これらの結果の応用として、我々は母集団のコレログラムが上に凸で、右層から 1 要素を無依為にとり出す抽出の形式が最適とな

るように層を定めれば、各層からの独立な無作為抽出よりも更に効率の良い系統的抽出計画を定めることができる。

Sec 5 とともににおいては、大部分これまでの節の一般定理の応用である集約の系統的抽出に対する色々な結果を導いた。大体において、この結果は系統的抽出の効率が無作為または層化無作為抽出のものより高くなることが予想される条件を、公式として示している(1)および(3)と性格を同じくするものである。しかし我々はこれらの公式を母集団のタイプを指定して応用することはしなかった。(1)および(3)から、このことは有用でもあり、またそのような研究は一般的論文の創りよりも重要なタイプの母集団またはデータについて進める方が一層貴重だということがすでにわかっている。

2. 無作為事象および条件付期待値

標本は殆んどの場合幾つかの段階で抽出される。たとえばあるシティから世帯の標本を抽出するには層々次のような二段抽出法が用いられる。

- a. 各ブロックの位置を示すシティの地図を手に入れ、その時点のものに修正する。
- b. この地図を用いてシティのブロックの標本を選ぶ(これがオノ一段)
- c. オノ一段で選んだブロックの世帯から世帯の副次標本を選ぶ(これがオニ一段)

この節では多段階抽出の平均値と分散の評価のための一般的方法を与えることにする。この方法は各段階で生ずる分散の寄与を同時にみせるという利点がある。その上ここに示す定理は興味が多段階抽出にない場合でも分散の計算においても有用なものである。この定理はサンプリングにおいて広い応用範囲をもつから一般的な用語で提出することにする。

ある操作を行った結果、確率 P_1, \dots, P_m で m 個の状態 A_1, \dots, A_m

が起りうるよ仮定できるとき、この操作の結果を確率事象 A^* という。ここで

$$P\{A^* = A_i\} = P_i \quad \sum_{i=1}^m P_i = 1$$

で $P\{A^* = A_i\}$ は「確率事象 A^* が A_i なる状態となる確率」という。

操作の一例はブロックの標本を抽出する手続きである。シティに N 個のブロックがあって、その中の n 個ずつのブロックの組がいずれも同じ機会に標本となるような方法で選出すると、 C 個の可能な標本がありうる。この場合 $m = C_n$ で、 C 個の可能な標本は「ブロックの標本を選んだ結果 A^* の m 個の状態である。さらにブロックの可能な標本がそれぞれ同じ機会に選ばれるときには

$$P\{A^* = A_i\} = \frac{1}{C_n} = \frac{1}{m}$$

である。

確率事象 A^* はまた確率変数の可能な値の一つで示すこともできる。確率変数 Z^* のとりうる値が z_1, \dots, z_m でその確率が P_1, \dots, P_m なら A^* が A_i となる状態を " $Z^* = z_i$ " なる場合と定義できる。

このように確率事象の概念には標本抽出の際にみられる二つのタイプの無作為性が含まれている。

Z^* を確率変数とする。このとき確率事象 A^* に関する Z^* の条件付期待値とはそのとりうる値が $E(Z^* | A_i); i=1, \dots, m$ で確率が P_i であるような確率変数 $E^*(Z^* | A)$ をいう。すなわち

$$P\{E^*(Z^* | A) = E(Z^* | A_i)\} = P_i = P\{A^* = A_i\}$$

ここで

$$(2.1) \quad E(Z^* | A_i) = \sum_{j=1}^m x_{ij} P_j(A_i)$$

で、 x_{ij} は A_i が起ったときに Z^* のとりうる N_i 通りの値のうち、 j 番目のものであり、

$$P_j(A_i) = P\{Z^* = x_{ij} | A_i\}$$

は A_i が起ったとき $Z^* = x_{ij}$ となる確率である。注意すべきことは

$$P_{ij} = P\{Z^* = x_{ij}\}$$

なら

$$P_{ij} = P\{X' = X_{ij}, A^* = A_i\}$$

となることである。というのは $X' = X_{ij}$ は A_i の生じたことを示すからである。そうすると

$$(2.2) \quad P_i \cdot P_j(A_i) = P_{ij}$$

である

ここで定理1と2を証明なしで述べておく。証明は容易である。

定理1. 確率変数 $E^*(X' | A)$ の期待値は $E X'$ である。すなわち

$$E\{E^*(X' | A)\} = E X'$$

$\sigma_{X'Y' | A}$ なる記号によって、とりうる値が $\sigma_{X'Y' | A}, i=1, \dots, k$ なる確率変数を表わす。ここで

$$\sigma_{X'Y' | A} = E\{(X' - E(X' | A_i))(Y' - E(Y' | A_i)) | A_i\}$$

であり

$$P\{\sigma_{X'Y' | A} = \sigma_{X'Y' | A_i}\} = P_i = P\{A^* = A_i\}$$

である。すなわち

$$\sigma_{X'Y' | A} = E^*\{(X' - E^*(X' | A))(Y' - E^*(Y' | A)) | A\}$$

さうに記号 $\sigma_{E^*(X' | A)E^*(Y' | A)}$ は "二つの確率変数 $E^*(X' | A)$ と $E^*(Y' | A)$ の共分散" を表わす。分散についての対応する定義は上のものを X' でおきかえることによってえられる。

定理2. X', Y' が確率変数なら

$$\sigma_{X'Y'} = E^* \sigma_{X'Y' | A} + \sigma_{E^*(X' | A)E^*(Y' | A)}$$

および

$$\sigma_{X'}^2 = E \sigma_{X' | A}^2 + \sigma_{E^*(X' | A)}^2$$

である

$P_{ij}, P_i, P_j(A_i)$ は指定されていないから定理1と2は任意の二層抽出計画で成立することに注意。定理1および2の多層抽出計画への拡張はすぐわかるが、実際にはこの定理を何回かにわけて適用する方が簡単なことが多い。定理1と2の適用を示すのはたやすいがしかしこれは本論文の目的には本質的なものでない。序に述べた通りこれらの二つの定理は長い間サンプリングの伝承とよびうる

ものの一つとなっていたものである。

3 層化抽出と負の相関ならびにその系統的抽出に対する応用

層化された母集団からのサンプリング計画を論ずる際には、仮定として、 X' が推定値で $X' = X_i + \dots + X_k$ (ここで X_j' は j 個の層の j 番目のものから生じた X' への寄与である) なら、確率変数 X_i' と X_j' が独立となるように抽出するものとする。

[1, p. 13] で、もし母集団が層化されており、異なる層からの寄与が負の相関をもつように要素を抽出すれば、この寄与が独立でもとも層内の共分散が同じ場合より推定値の精度は低くなることを注意した。もちろんこの結論は

$$\sigma_{X'}^2 = \sum_{i,j} \sigma_{X_i X_j}$$

からもし

$$(3.1) \quad C = \sum_{i \neq j} \sigma_{X_i X_j} < 0$$

なら $C = 0$ の場合より $\sigma_{X'}^2$ は小となることからすぐわかる $C < 0$ なら我々はこの標本計画は "負の相関" をもつとよぶ

容易にわかるように、任意の母集団はそれ自身一つの標本があると考えることが出来る。すなわち現実の母集団を定める力の造り出す可能な母集団からとられた標本である。抽出計画は支配力の知識である種の過去の経験をもととして選ぶことが多いから、特定の母集団の期待値や分散のみならず、同じ力で決定されるすべての可能な母集団上でのこれらの期待値をも考えるのが現実的である。標本計画の期待分散を考えることの有甲性は Cochran [3] が一つの例で示したところである。彼はそれ自身、確率変数であるところの母集団の要素 X_1, \dots, X_n を考え、 $E X_i = \mu$ および $E(X_i - \mu)^2 = \sigma^2$ と仮定した。彼の目的からいうと更に $u > 0$ のとき $E(X_i - \mu)(X_{i+u} - \mu) = \rho_u \sigma^2$ とおくと都合がよい。これによって、彼はコレロゾフムについての現実的な仮説、すなわち特定の母集団ではこれほど妥当とは考えられなかったところの " ρ_u は u の 数である" という

仮説を仮けることができた。このようにして彼は無作為および層化無作為計画と比較したときの系統的計画の期待効率についての一般的な結論をえた。

この論文では与えられた母集団に対する期待値や分散だけでなく、母集団の要素自身が確率変数であるという仮定のもとにこれらの期待値および分散の期待値をも考慮することにする。我々は母集団の要素を確率変数と考えた場合の期待値を表わすのに E なる記号を用い、また前と同じく特定の有限の母集団についての期待値を E と書くことにする。

$$E\sigma_{xi}^2 = \sum_{i,j=1}^k E\sigma_{xi}x_j$$

であるからもし $E C < 0$ ならこの計画は " 真の期待相関をもつ " という。

さて我々は標本計画中に真の相関や真の期待相関を導入したり、それを利用したりできる場合の標本計画を論ずることから始めよう。

簡単のためにまず層が二つの場合を考え、 X のとりうる値を x_1, \dots, x_n また Y のとりうる値を y_1, \dots, y_n とする。さらに

$$P\{X' = x_i\} = P\{Y' = y_i\} = P\{X = x_i, Y = y_i\} = p > 0$$

なるように抽出を行なうものとする。したがって $\sum_{i=1}^n p_i = 1$ であり、なら $P\{X', Y = y_j\} = 0$ である。上記の仮定から

$$(3.2) \quad \sigma_{X'Y'} = \sum_{i=1}^n p_i x_i y_i - \sum_{i,j=1}^n p_i p_j x_i x_j$$

が導かれる。

記号 Y_j の 0 はすべての i, j について $Y_j \geq 0$ であり、少くとも $i=j$ の i, j については $Y_j > 0$ なることを表わす。我々は $(x_i - x_j)(y_i - y_j) \geq 0$ なら組 (X) と (Y) は同様に順序づけられており、 $(x_i - x_j)(y_i - y_j) < 0$ ならこれらの組は順序が逆であるという。ただし (X) は x_1, \dots, x_n を表わし (Y) は y_1, \dots, y_n を表わす。このとき数値が逆の順序のときは $\sigma_{X'Y'} < 0$ 、また同じ順序なら $\sigma_{X'Y'} > 0$ というこ

と (2. P. 43) は直接容易に証明することができる。

これより幾分一般な結果は次の通りである

定理 3

$n \leq k$ とし

$$b = \sum_{i=1}^n \sum_{j=1}^k a_{ij} w_i z_j$$

を実数一次形式、また

$$c = \sum_{i=1}^n a_{ii} w_i$$

実一次形式とする。ただし $w_i > 0, z_i > 0$ かつ

$$\sum_{i=1}^n w_i = \sum_{i=1}^k z_i = 1$$

このとき $b > c$ なるための十分条件は

$$(3.3) \quad a_{ij} \geq a_{ii}$$

$k = n$ で $w_i = z_i$ なら $b > c$ となるのは

$$(3.4) \quad a_{ij} + a_{ji} \geq a_{ii} + a_{jj}$$

ときこめらる。

証明

$$b - c = \sum_{i=1}^n \sum_{j=1}^k a_{ij} (w_i z_i - w_i) + \sum_{i,j=1}^k a_{ij} w_i z_j$$

よって

$$1 - z_i = \sum_{j=1}^k z_j$$

だから

$$b - c = \sum_{i,j=1}^k (a_{ij} - a_{ii}) w_i z_j$$

がえられる。よって (3.3) が成りたてば $b > c$ 。同様にして $k = n$ かつ $w_i = z_i$ なら (3.4) が成りたてば $b > c$ となる。定理 3 の自明な拡張は、こゝでは必要がないので省略する。

組 (X) と (Y) が逆の順序のとき $\sigma_{X'Y'} < 0$ なることを証明するため、 $a_{ij} = x_i y_j$ および $z_i = w_i = p_i$ とおきえれば

$$(3.5) \quad a_{ii} + a_{jj} - a_{ij} - a_{ji} = (x_i - x_j)(y_i - y_j)$$

まうる。したがってこの値が逆に順序づけられていれば $\sigma_{x'y'} < 0$ であるから、二つの層は負の相関をもつ。

負の期待相関を考えるため

$$(3.6) \quad E \sigma_{x'y'} = \sum_{i=1}^k p_i \sigma_{ii} + \sum_{i \neq j} p_i p_j r_{ij}$$

なることに注意する。ここで $E x_i = p_i E y_i = V$ および $E(x_i - V)(y_i - V) = \sigma_{ij}$ と仮定する。したがって σ_{ii} は分散でなく共分散である。

$\sigma_{ij} = \sigma_{ji}$, $Z_i = W_i = p_i$ とおくと (3.4) が成立して、 $E \sigma_{x'y'}$ が負となるための十分条件、

$$(3.7) \quad \sigma_{ij} + \sigma_{ji} \geq \sigma_{ii} + \sigma_{jj}$$

がえられる。あるいはまた p_{ij} を

$$\sigma_x \sigma_y p_{ij} = \sigma_{ij}$$

によって定義すると、 $E \sigma_{x'y'} < 0$ なるための十分条件が

$$(3.8) \quad p_{ij} + p_{ji} \geq p_{ii} + p_{jj}$$

のようになれる。ただし $\sigma_x^2 = E(x_i - \mu)^2$, $\sigma_y^2 = E(y_i - \mu)^2$ とする。

単一要素の系統的抽出を考えよう。系統的抽出では何れ個の順序づけられた要素 $x_1, x_2, \dots, x_k, x_{1+k}, \dots, x_{2+k}, \dots, x_{i+(k-1)k}, \dots, x_{nk}$ からなる母集団を仮定し、その算術平均値を推定したいものとする。推定値としては

$$\bar{x}' = (x_1 + \dots + x_{nk}) / n$$

を用いる。ここで x'_i は x_1, \dots, x_k から無作為に選び、もし $x'_i = x_j$ ならば $x'_i = x_j + (i-1)k$, $i=2, \dots, n$ とする。よって \bar{x}' は n 層が

$$x_{i+(i-1)k}, \quad x_{i+(i-1)k}$$

からなっていて、

$$P\{x'_i = x_{j+(i-1)k}\} = P\{x'_i = x_{j+(i-1)k}, x'_j = x_{j+(j-1)k}\} = \alpha$$

かつ $\alpha + \beta$ のとき

$$P\{x'_i = x_{j+(i-1)k}, x'_j = x_{j+(j-1)k}\} = 0$$

なるような層化母集団よりの推定値と解釈できる。このとき、

$$\sigma_{x'_i x'_j} = \left(\frac{1}{n}\right) \sum_{i=1}^k x_{i+(i-1)k} \cdot x_{j+(j-1)k} - \bar{x}'_i \bar{x}'_j$$

ただし $\bar{x}'_i = \left(\frac{1}{n}\right) \sum_{i=1}^k x_{i+(i-1)k}$

ゆえに逆に順序づけられている二つの層は分散に対して負の寄与を与えることになる。しかしすべての層を負に順序づけることはできないからこれでは有用な結果はえられない。それで (1) のように C または σ_{ij} そのものに立かえて考えなければならぬ。しかしもし Cochran の仮定を認め、 $E \sigma_{x'y'}$ を考えると i 層 j 層について

$$f_{ij} = f_{(j-i)k} + \beta - \alpha$$

なる関係が導かれる (3.8) は

$$(3.9) \quad f_{(j-i)k} + (\beta - \alpha) + p_{(j-i)k} + (\alpha - \beta) \geq 2p_{(j-i)k}$$

となる。すなわち相関関数 f_{ij} は上方に凹でなければならぬ。これは Cochran が別の方法で証明したところのものである。 $E C$ を考えることによつて、ある種の平均的凹型関係をもつということだけが系統的抽出の分散を層化無作為抽出以下とするのにコレログラムに要求される条件だということが示しうる。

層の大きさが等しくない場合に相関が負となるための条件と系統的抽出に対するその応用

大きさに比例した確率による集落の系統的抽出 (SERC で論ずる) のように、定理3で取扱つて簡単な場合は直接あてはまらないことが多い。しかし定理3は次のようにして利用することができる。

x' のとりうる値を x_1, \dots, x_n y' のとりうる値を y_1, \dots, y_n とした

$$P\{y' = y_{i'} | x' = x_{i'}\} = p_{i' i'}$$

とする。ただし $k > n$ 。このときもし

$$(4.1) \quad f_{i' i'} = \sum_{j=1}^n y_j p_{i' j}$$

と定義すれば

$$f_{i' i'} = E(y_j | x' = x_{i'})$$

である。

yをその取りうる値がy₁, ..., y_mで、それからの確率がP_αである確率変数だとすると

$$y' = E^*(y'_0 | X')$$

で

$$\sigma_{X'y'_0} = \sigma_{X'y'}$$

ただし P_α = P{X' = X_α}

確率変数zとyについても定理3の成立つことは明らかである。結局我々には組X₁, ..., X_n およびy₁, ..., y_{m}が逆に順序づけられる、あるいは(3.7)が成立するために満足しなければならない条件付確率P_{β_{1α}}、およびy₀の値に対する制限が何であるかを定めることだけが必要である}

(3.5) のy₀ およびy₁ を代入すると、もし

$$(4.2) \quad (X_{j+1} - X_j) \sum_{\alpha=1}^n \frac{y'_0}{y'_1} (P_{\beta_{1\alpha}} - P_{\beta_{2\alpha}}) \leq 0$$

なら $\sigma_{X'y'_0} - \sigma_{X'y'_1} < 0$ である。

$$\sigma_{\alpha y} = E(X_{j+1} - \mu)(y'_j - \gamma)$$

とする。よって(3.7)に代入すると

$$(4.3) \quad \sum_{\alpha=1}^n (P_{\beta_{1\alpha}} - P_{\beta_{2\alpha}})(\sigma_{\alpha y} - \sigma_{\alpha y'}) \leq 0$$

あるいは

$$(4.4) \quad \sum_{\alpha=1}^n (P_{\beta_{1\alpha}} - P_{\beta_{2\alpha}})(P_{\beta_{1\alpha}} - P_{\beta_{2\alpha}}) \leq 0$$

なら $E \sigma_{X'y'_0} < 0$

である。(4.2) および(4.3)を用いるには次のよく知られた Lemma が要々有用である。

Lemma E₁ ≤ E₂ ≤ ... ≤ E_n ≤ 0 である E₁, ..., E_n について

$$\sum_{\alpha=1}^n E_{\alpha} E_{\beta} \leq 0 \text{ が成立すれば } \sum_{\alpha=1}^n E_{\alpha} E_{\beta} \leq 0, \alpha=1, \dots, n$$

層間の頁の相関または頁の期待相関を示すのに有用な他の定理を導くためにこの Lemma を使うことにする。

定理4. 一次形式

$$a = \sum_{i=1}^n \sum_{j=1}^m a_{ij} w_i z_j$$

について

$$\sum_{i=1}^n w_i \geq 0, \sum_{j=1}^m z_j \geq 0 \quad \alpha=1, \dots, n-1 \quad \alpha'=1, \dots, m-1$$

および

$$(4.5) \quad \sum_{i=1}^n w_i = \sum_{j=1}^m z_j = 0$$

とする。

$$\delta'_{ij} = a_{ij} - a_{i+1,j} - a_{i,j+1} + a_{i+1,j+1}$$

とおく。そうすると $\theta \leq 0$ なるための十分条件は $\delta'_{ij} \leq 0$ である。

証明 θ の w_n と z_m に(4.5)を代入すると

$$\theta = \sum_{i=1}^{n-1} \sum_{j=1}^{m-1} \delta'_{ij} w_i z_j$$

なることがわかる。ただし

$$\delta'_{ij} = a_{ij} - a_{i,m} - a_{n,j} + a_{n,m}$$

あるいは

$$E_j = \sum_{i=1}^{n-1} \delta'_{ij} w_i$$

と定義すれば

$$\theta = \sum_{j=1}^{m-1} E_j z_j$$

よって Lemma によって $\theta \leq 0$ なるための十分条件

$$E_1 \leq E_2 \leq \dots \leq E_{m-1} \leq 0$$

である。同様に

$$E_j - E_{j+1} \leq 0$$

なるための十分条件は

$$\delta'_{ij} - \delta'_{i,j+1} \leq \delta'_{i+1,j} - \delta'_{i+1,j+1}$$

である。よって証明で残されたものは

$$\delta'_{ij} = \delta'_{ij} - \delta'_{i,j+1} - \delta'_{i+1,j} + \delta'_{i+1,j+1}$$

を確かめることだけである。

我々は前の頁で層内の抽出が独立でなく、一つの層からの要素の

抽出が他層の抽出を決定してしまうような層化抽出は系統的抽出と同じものだという事を証明した。しかしこの場合層内の要素数は等しいことを仮定している。我々はここで層内の要素数が異なる場合についてこの標本抽出値を拡張してみよう。そうしてその過程の中で Lemma および定理4の用法を例示することにしよう。

さて N 個の要素 x_1, \dots, x_N からなる母集団が k 個の層に分類されており、その α 層には N_α 個の要素

$$x_{N_\alpha+1}, \dots, x_{N_\alpha+1+N_i}$$

が含まれるものとする。これらの要素を x_{i1}, \dots, x_{iN_i} とする。

これらの k 個の層からそれぞれ α 要素を一つづつ選ぶ。 α 層から選んだ要素を x_i と書く。母集団の算術平均値 \bar{x} の推定値として

$$\bar{x}' = \sum_{i=1}^k \frac{N_i}{N} x_i$$

をとる。このときよく知られているように抽出が各層から独立かつ無依存に行なわれるとすると

$$\sigma_{\bar{x}'}^2 = \sum_{i=1}^k \left(\frac{N_i}{N}\right)^2 \sigma_i^2$$

である。ただし σ_i^2 は x_i の分散すなわち α 層の分散である。

ここで普通に用いられるものと異なるもう一つの方法を考えてみよう。我々は一般性を失なうことなく $N_i > 1$ と仮定できる。(この方法は $N_i = 1$ なるいづれの層についても同じで、またすべての $N_i = 1$ であるかまたは唯一つの N_i を除く他のすべてが 1 に等しくなっているような母集団についても同じ結果がえられる。しかし少くとも二つの N_i が 1 と異なれば結果は違ってくる。)

我々はまず最初の層から無依存に一要素を選ぶ。 $x_1 = x_\alpha$ とおく。 $N_2 > 1$ と仮定して α 層から次のようにして一要素を選ぶ。 k_2 が整数のとき $N_2 t_2 / N_1 = k_2$ となるような正の整数 t_2 をとり N_2 に乗する。 α 層の各要素にサイズ t_2 の大きさ (measure) を与え、二組の累積和 $t_2, 2t_2, \dots, N_2 t_2$ および $k_2, 2k_2, \dots, N_1 k_2$ を作る。そうすると層 α の各要素に与えられたサイズ t_2 の大きさ

と層 1 の各要素に与えられたサイズ k_2 の大きさで層 1 と α は同じサイズをもつことになる。

以下に与える計算の一例として次のような簡単な場合を考える。 $N_1 = 3$ で $N_2 = 4$ と仮定する。このとき t_2 を 6 とすれば $k_2 = 8$ となる。我々は整数 1, 2, 3 の中の一つを等確率で選ぶ。もし整数 1 がえられたら α 層の α 要素が選ばれたのであるから 1 から 8 までの整数の中一数を等確率で抽出する。選ばれた整数が 1 と 6 の間の数であれば α 層の α 要素が選ばれる。同様にもし α 層の α 要素が選ばれたら 9 と 16 の間の整数から 1 数を等確率で選ぶ。この数が 9, ..., 12 なら α 層の α 要素をとり、13, ..., 16 なら α 層の α 要素をとる。

α 層に対する抽出手段の一般的方式は次の通りである。

β_0 を $(\alpha - 1)k_2 + 1 \leq \beta_0 t_2$ なる最小の整数とし、また β は $(\beta - 1)t_2 < k_2 \leq \beta t_2$ なる整数とする。1, ..., β の中から無依存に一つの整数を選んで β とする。

このときもし

$$(\alpha - 1)k_2 < (\alpha - 1)k_2 + \beta \leq \beta_0 t_2$$

なら α 層から β_0 番目の要素を選ぶ

$$\beta_0 t_2 < (\alpha - 1)k_2 + \beta \leq (\beta_0 + 1)t_2$$

なら $(\beta_0 + 1)$ 番目の要素を選ぶ。 そうしてもし

$$(\beta - 1)t_2 < (\alpha - 1)k_2 + \beta \leq \alpha k_2$$

なら α 層から β 番目の要素を選出する。

このようにすると容易に証明できるように層 α の各要素は等確率で抽出されることになる。ゆえにこの手段を各層に用いれば

$$\sigma_{\bar{x}'}^2 = \sum_{i=1}^k \left(\frac{N_i}{N}\right) \sigma_i^2 + \sum_{i \neq j} \frac{N_i N_j}{N^2} \sigma_{x_i x_j}$$

がえられる。

このタイプの抽出での $\sigma_{x_i x_j}$ を評価してみよう。いま

$$\sigma_{x_i x_j} = E(x_i - \bar{x}_i)(x_j - \bar{x}_j)$$

である。ここで \bar{x}_i は i 層の要素の算術平均値である。よって定理

である

$$\begin{aligned} \sigma_{x_i x_j} &= \frac{1}{N_i} \sum_{\alpha=1}^{N_i} E \{ (x_i - E(x_i | x_{i\alpha})) (x_j - E(x_j | x_{i\alpha})) \} \\ &\quad + \frac{1}{N_i} \sum_{\alpha=1}^{N_i} \{ (E(x_i | x_{i\alpha}) - \bar{x}_i) (E(x_j | x_{i\alpha}) - \bar{x}_j) \} \end{aligned}$$

$\sigma_{x_i x_j}$ の第1項は上記の抽出法では0になることが容易にわかる。さらに \bar{x}_i は条件付期待値の算術平均値だから問題は、この条件付期待値が、頁の相関または頁の期待相関をもつための条件を満たすかどうかを定めることに帰着する。

$E(x_i | x_{i\alpha})$ を $y_{i\alpha}$ で表わす。そうすると組 y_{i1}, \dots, y_{iN} と y_{j1}, \dots, y_{jN} が逆順かどうかを調べなければならぬ。いま

$$(y_{i\alpha} - y_{i\beta})(y_{j\alpha} - y_{j\beta}) = \sum_{g=1}^{N_i} \sum_{h=1}^{N_j} x_{ig} x_{jh} E_{i\alpha\beta} E_{j\alpha\beta}$$

である。ただし

$$E_{i\alpha\beta} = P\{x_i = x_{ig} | x_{i\alpha}\} - P\{x_i = x_{ig} | x_{i\beta}\}$$

である。 $\alpha < \beta$ ならこの抽出法にしたがって

$$\sum_{g=1}^A E_{i\alpha\beta} x_{ig} \geq 0 \quad \alpha = 1, \dots, N-1$$

で、また一方

$$\sum_{g=1}^{N_i} E_{i\alpha\beta} x_{ig} = 0$$

である。定理4で

$$\pi = N_i, \quad m = N_j, \quad w_g = E_{i\alpha\beta}, \quad z_h = E_{j\alpha\beta}, \quad u_{gh} = x_{ig} x_{jh}$$

とおく。そうすると

$$\delta_{gh} = (x_{ig} - x_{i,g+1})(x_{jh} - x_{j,h+1})$$

であるから、層間に頁の相関が存在するためには組 x_{i1}, \dots, x_{iN} と x_{j1}, \dots, x_{jN} が $\delta_{gh} \leq 0$ で表わされるような型の頁の順列のわれは十分である。同様に

$$\bar{\sigma}_{gh} = E(x_{ig} - \mu_i)(x_{jh} - \mu_j), \quad \mu_i = E x_{ig}$$

なり、頁の期待相関が生ずるための十分条件は

$$\bar{\sigma}_{gh} - \bar{\sigma}_{g,h+1} - \bar{\sigma}_{g+1,h} + \bar{\sigma}_{g+1,h+1} \leq 0$$

である。

もちろんこれらの関係はコレログラムが上に凹なら満足される。したがってもし母集団が上に凹なコレログラムをもつ N 個の確率変数 x_1, \dots, x_N からなっていると、これらの要素がどの層に分類されていても要素の出現の順序が交らない限り最適配分による標本要素の層化無作為抽出より分散の小さい推定値を与えるような標本要素の系統的抽出を計画することができる。もし最適配分のもとで一個以上の要素を層から選出するなら、目的の要素を系統的に抽出すれば十分である。もし最適配分だけでなく最適層化 (optimum stratification) をも用いて、各層から一要素だけを選出するものとするれば、この節で述べた方式にしたがう系統的抽出の分散は層化無作為抽出のものよりも大きくはならない。しかし層を考えないで系統的抽出に夢中になってもは小さくならないことに注意しなければならぬ。分散を減少させるためにはあらかじめおこなわなければならない手続きが存在するのである。

この例の方法によって母集団の大きさが標本サイズの倍数でない層目の要素の系統的抽出について答がえられるであろう。

5. 大きさに比例した確率を用いる集層の系統的抽出

集層抽出で確率をその大きさに比例させた場合には推定値の分散が層々大きく減少することが知られている(5) しかし確率を大きさに比例させて行なう幾つかの集層の抽出理論はこれまで研究されたいない。この節の目的はこの理論に対して何らかの寄与を与えることである。

大きさに比例した確率で集層を抽出するのに最も多く用いられる方法は次に述べるものと同様である。集層を C_1, \dots, C_M で表わし、これらの M 個の集層の h 番目のものに P_h なる確率を与える。次々

の累積和 $P_1, P_1 + P_2, P_1 + P_2 + P_3, \dots, P_1 + \dots + P_M$ を作る。これらの中から m 個の集落を選びたいものとする。 $\bar{P}_m = (P_1 + \dots + P_M) / m$ を計算し、次に $P_j \leq \bar{P}_m, j = 1, \dots, M$ と仮定して $1, \dots, \bar{P}_m$ の中から等確率で一つの整数を選ぶ。この整数を P' とし、 m 個の数 $P', P' + \bar{P}_m, P' + 2\bar{P}_m, \dots, P' + (m-1)\bar{P}_m$ を計算する。もし任意の整数 $i, i = 1, \dots, M$ に対し

$$(5.1) \quad P_1 + \dots + P_{i-1} + 1 \leq P' + (i-1)\bar{P}_m \leq P_1 + \dots + P_i$$

ならば集落 C_i を標本に選ぶ。 $P_i > \bar{P}_m$ なる集落はすべて自動的に標本内に含まれ、更にもしそのような集落がたとえば d 個あったとすると、これらの d 個の集落が標本に含まれた後に残る $M-d$ 個の集落に対して \bar{P}_{m-d} を計算し前と同様な手続を繰返す。

我々の用いる推定値の分散を小さくし、我々はこの推定値を層化抽出の推定値と考えることにする。推定値の分散はこのような解釈をしても容易に求められるが、おしで分散の推定が必要となるのでここで完全な式を少し繕って示しておく。

集落 C_1, \dots, C_{k-1} は

$$P_1 + \dots + P_{k-1} < \bar{P}_m \leq P_1 + \dots + P_k$$

を満たすものとする。このとき層1は集落 C_1, \dots, C_k からなるものとする。容易にわかる通り上の抽出方法をとれば

$$P\{C_k \text{ が層1から選ばれる } k < k_1\} = \frac{P_k}{\bar{P}_m}$$

$$P\{C_k \text{ が層1から選ばれる}\}$$

$$= \frac{\bar{P}_m - P_1 - \dots - P_{k-1}}{\bar{P}_m}$$

である。

さらに $C_{k_1}, C_{k_1+k_2}$ について

$$P_1 + \dots + P_{k_1+k_2-1} < 2\bar{P}_m \leq P_1 + \dots + P_{k_1+k_2}$$

が成立つとする。このとき層2は集落 $C_{k_1}, \dots, C_{k_1+k_2}$ からなるものとする。容易にわかるように上記の抽出方法を用いれば

$$P\{C_{k_1} \text{ が層2から選ばれる}\} = \frac{P_1 + \dots + P_{k_1} - \bar{P}_m}{\bar{P}_m}$$

$$P\{C_{k_1+k_2} \text{ が層2から選ばれる, } 1 \leq k_1 < k_2\}$$

$$= \frac{P_k}{\bar{P}_m}$$

$$P\{C_{k_1+k_2} \text{ が層2から選ばれる}\} = \frac{2\bar{P}_m - P_1 - \dots - P_{k_1+k_2-1}}{\bar{P}_m}$$

である。 $P_k \leq \bar{P}_m$ だから C_{k_1} が層1と層2の両方から選ばれることは不可能である。

一般に集落 $C_{k_1}, \dots, C_{k_1+k_2-1}, \dots, C_{k_1+k_2}, \dots, C_{k_1+k_2+k_3}$ について

$$(5.2) \quad P_1 + \dots + P_{k_1+k_2-1} < i\bar{P}_m \leq P_1 + \dots + P_{k_1+k_2+k_3}$$

なるものとする。もし層はこれらの $k_1+k_2+k_3$ 個の集落からなるから、我々は確率 $P_{k_1+k_2+k_3}, d = 0, \dots, k_2+k_3$ を次のように定義する。

$$P_{k_1+k_2+k_3} = P\{C_{k_1+k_2+k_3} \text{ が } i \text{ 層から選ばれる}\}$$

$$= \frac{P_1 + \dots + P_{k_1+k_2+k_3} - i\bar{P}_m}{\bar{P}_m}$$

$$P_{k_1+k_2+k_3} = P\{C_{k_1+k_2+k_3} \text{ が } i \text{ 層から選ばれる, } k_1, \dots, k_2+k_3 < k_1+k_2+\dots+k_3\}$$

$$(5.3) \quad = \frac{P_k}{\bar{P}_m} \quad d = k_1+k_2+\dots+k_3-1$$

$$P_{k_1+k_2+k_3} = P\{C_{k_1+k_2+k_3} \text{ が } i \text{ 層から選ばれる}\}$$

$$= \frac{i\bar{P}_m - P_1 - \dots - P_{k_1+k_2+k_3-1}}{\bar{P}_m}$$

すなわち

$$(5.4) \quad P_{i-1, k_{i-1}} + P_{i0} = \frac{P_{k_1 + \dots + k_{i-1}}}{\bar{P}_m}$$

に注意しておく。

いま母集団の要素を X_{ij} , $j=1, \dots, N_{ik}$ とし、
 九番目の集落の算術平均値を \bar{X}_{i0} とおく。 N_{ik} は通例未知だが大
 さきの尺度 P_{ik} はわかっているから、抽出確率を N_{ik} に比例させ
 ず P_{ik} に比例した確率で抽出を行なう。オ_i層の集落 C_{i0}, \dots, C_{ik}
 の大き

$$(5.5) \quad C_{i0} = C_{i, k_1 + \dots + k_{i-1}}$$

とおく

さらに集落の要素数を N_{i0}, \dots, N_{ik} また集落平均値を $\bar{X}_{i0}, \dots,$
 \bar{X}_{ik} とする。

ここで

$$N_{i0} = N_{i, k_1 + \dots + k_{i-1}}$$

$$(5.6) \quad \bar{X}_{i0} = \bar{X}_{i, k_1 + \dots + k_{i-1}}$$

だから

$$\bar{X}_{i0} = \bar{X}_{i-1, k_{i-1}}$$

また

$$N_{i0} = N_{i-1, k_{i-1}} \quad i=1, \dots, m$$

である。更

$$(5.7) \quad \bar{X}_{i0} = N_{i0} \bar{X}_{i0} / P_{i0} = \bar{X}_{i, k_1 + \dots + k_{i-1}}$$

と定義する。

我々はオ_i層の平均値を

$$(5.8) \quad \bar{X}_i = \sum_{j=0}^{k_i} P_{ij} \bar{X}_{ij} / \bar{P}_m$$

と定義しオ_i層の分散を

$$(5.9) \quad \sigma_i^2 = \sum_{j=0}^{k_i} \frac{P_{ij}}{\bar{P}_m} (\bar{X}_{ij} - \bar{X}_i)^2$$

とする。

このとき母集団の平均値と分散を

$$(5.10) \quad \bar{X} = \sum_{h=1}^M P_h \bar{X}_h / P$$

および

$$(5.11) \quad \sigma^2 = \sum_{h=1}^M \frac{P_h}{P} (\bar{X}_h - \bar{X})^2$$

と定義すると容易に

$$(5.12) \quad \bar{X} = \frac{1}{M} \sum_{i=1}^M \bar{X}_i$$

および

$$(5.13) \quad \sigma^2 = \frac{1}{M} \sum_{i=1}^M \sigma_i^2 + \frac{1}{M} \sum_{i=1}^M (\bar{X}_i - \bar{X})^2$$

が証明できる。

特性統計の不偏推定値

我々は X の推定値を求めようことを証明しよう。ただし

$$X = \sum_{i=1}^M \sum_{j=1}^{N_i} X_{ij}$$

すなわち X は母集団要素の総計である。 N は未知だから X の推定値
 として用いられるのは X と N の不偏推定値の比である。よく知られ
 ているように、この比は偏りをもつ。我々はこの比推定値の研究
 を行なうのでないから、この推定値の分散の近似式を導くこと
 はしない。しかしそれはここで与えた結果の単純な拡張でえられる
 ことは注意できる。

我々はこの推定値の一般式の次のようになることをみてみよう。
 もし母集団のオ_i層が選ばれば、それからの N_{ij} 個の要素の特
 性値統計を X_{ij} で表わす。さらにオ_i層からとられたものは同じ

ことだが、 i 層からとった集落より抽出された副次標本内の要素の統計を n_i とし、これらの要素の統計を X_i とかく。このようにして j 番目の集落が i 番目に選ばれたなら $n_i = n_j$ で $X_i = X_j$ である。我々は母集団統計 X の推定値 X^* を

$$(5.14) \quad X^* = K(X_1^* + \dots + X_m^*)$$

と定義する。

このときもし $K = P/mn$ で $n_i = n, N_i/P_i$ なら X^* が X の不偏推定値であることは容易にわかる。

推定値の分散

X の分散が計算できる。ここで

$$(5.15) \quad X^* = \bar{P}_m (\bar{X}_1^* + \dots + \bar{X}_m^*) \text{ および } \bar{X}_i^* = X_i^*/n$$

いま定理2によつて

$$(5.16) \quad \sigma_{X^*}^2 = E \sigma_{X^*|A}^2 + \sigma_{E^*}^2 (X^*|A)$$

である。ここで A^* は前に定義したものがある。我々は $E \sigma_{X^*|A}^2$ の詳細は行なわない。というのはこれには無作為または系統的抽出を用いる副次抽出法あるいは大きさ按比例した確率を用いる。副次抽出法に関する新しい問題は何も含まれていないからである。

(5.15) から

(5.17) $E^*(X^*|A) = \bar{P}_m (\bar{X}_1^* + \dots + \bar{X}_m^*)$
 がえられる。いいかえれば $E^*(X^*|A)$ は標本内の集落を全部調べた場合にえられる推定値である。(5.16) のオズ項を σ_B^2 と書けば

$$(5.18) \quad \sigma_B^2 = \bar{P}_m^2 \left\{ \sum_{i=1}^m \sigma_{\bar{X}_i^*}^2 + \sum_{i \neq j} \sigma_{\bar{X}_i^* \bar{X}_j^*} \right\}$$

である。

いま

$$(5.19) \quad \sigma_{\bar{X}_i^*}^2 = \sum_{\alpha=0}^{n_i} \frac{P_{i\alpha}}{\bar{P}_m} (\bar{X}_{i\alpha} - \bar{X}_i)^2 = \sigma_i^2$$

$i \neq j$ のときの $\sigma_{\bar{X}_i^* \bar{X}_j^*}$ を計算するため定理1を用いる

$$(5.20) \quad \sigma_{\bar{X}_i^* \bar{X}_j^*} = E(\bar{X}_i^* - \bar{X}_i)(\bar{X}_j^* - \bar{X}_j) = E\{(\bar{X}_i^* - \bar{X}_i) E^*(\bar{X}_j^* - \bar{X}_j | \bar{X}_i^*)\}$$

$E^*(\bar{X}_j^* - \bar{X}_j | \bar{X}_i^*)$ を計算するためまず

$$(5.21) \quad E^*(\bar{X}_j^* - \bar{X}_j | \bar{X}_i^*) \equiv E^*(\bar{X}_j^* - \bar{X}_j | C_i)$$

に注意する。ここで C_i^* は i 層の標本の集落としては C_{i0}, \dots, C_{in_i} の何れかが選ばれるという $n_i + 1$ 通りの可能な場合のある確率空間である。いまもし $C_{i\alpha}$ が i 層の集落なら

$$(5.22) \quad E(\bar{X}_j^* - \bar{X}_j | C_{i\alpha})$$

を計算する。我々はまず i 層から $C_{i\alpha}$ が選ばれることがわかっているとき j 層で標本の集落となり得る力ほどの集落があるかを定めることにする。 i 層 j 層の大きさはどちらも \bar{P}_m だから

$$P_{j0} + \dots + P_{j\beta_0-1} \leq P_{i0} + \dots + P_{i\alpha-1} < P_{j1} + \dots + P_{j\beta_1}$$

と

$$P_{j0} + \dots + P_{j\beta_0-1} < P_{i1} + \dots + P_{i\alpha} \leq P_{j1} + \dots + P_{j\beta_1}$$

なるような整数 β_0 と β_1 が存在する。ゆえにも i 層から $C_{i\alpha}$

が選ばれていることかわかっていれば、j層からは $C_{j\beta_0}, C_{j\beta_1}, \dots, C_{j\beta_r}$ のどれか1つを選ばなければならずまた

$$P\{C_{j\beta} \text{ が選ばれる} | C_{i\alpha} \text{ が選ばれる}\} = P_{j\beta}^i / P_{i\alpha}, \beta = \beta_0, \beta_0+1, \dots, \beta_r \text{ のとき } = 0 \text{ 他の場合}$$

なることがわかる。ただし、

$$P_{j\beta_0}^i = P_{j\beta_1}^i + \dots + P_{j\beta_r}^i - P_{i\beta_0}^i - \dots - P_{i\beta_{r-1}}^i$$

$$P_{j\beta}^i = P_{j\beta_0}^i, \beta = \beta_0+1, \dots, \beta_r-1$$

$$P_{j\beta_r}^i = P_{i\beta_r}^i + \dots + P_{i\beta_0}^i - P_{j\beta_0}^i - \dots - P_{j\beta_{r-1}}^i$$

で

$$\sum_{\beta=\beta_0}^{\beta_r} P_{j\beta}^i = P_{i\alpha}$$

である。

このとき

$$(5.23) \quad E\{(\tilde{x}_{ji} - \tilde{x}_{j\beta} | C_{i\alpha}) - \tilde{x}_{j\beta\alpha} - \tilde{x}_{j\beta}\}$$

は 0 である。

$$(5.24) \quad \tilde{x}_{j\beta\alpha} = \sum_{\beta=\beta_0}^{\beta_r} \frac{P_{j\beta}^i}{P_{i\alpha}} \tilde{x}_{j\beta}$$

よって (5.20) に代入すると

$$\sigma_{\tilde{x}_{ji} \tilde{x}_{j\beta}}^2 = E\{(\tilde{x}_{ji} - \tilde{x}_{j\beta})(\tilde{x}_{j\beta\alpha} - \tilde{x}_{j\beta})\}$$

ただし j層から $C_{i\alpha}$ が選ばれたときは、 $\tilde{x}_{j\beta\alpha} = \tilde{x}_{j\beta}$ である。それゆえ

$$(5.25) \quad \sigma_{\tilde{x}_{ji} \tilde{x}_{j\beta}}^2 = \sum_{\alpha=0}^{k_i} \frac{P_{i\alpha}}{P_m} (\tilde{x}_{i\alpha} - \tilde{x}_{j\beta})(\tilde{x}_{j\beta\alpha} - \tilde{x}_{j\beta})$$

がえられる。(5.23) を用いれば (5.25) の条件付期待値は求めらることは容易にわかるが、しかしそうしても簡略化とかが

化とかの点では何等うるところがない。

特別の場合には可能な標本を全部書き出して x_{ij} の分散共分散を求めることができよう。これを一般化するためには必要な記法を与えさえすればよい。

(5.18) に代入すると

$$\sigma_B^2 = \bar{P}_m^2 \left\{ \sum_{i=1}^m \sigma_i^2 + \sum_{i+j} \sigma_{\tilde{x}_{ij} \tilde{x}_{ij}} \right\}$$

がわかる。ただし $\sigma_{\tilde{x}_{ij} \tilde{x}_{ij}}$ は (5.25) で与えられたものである。

$$\sum_{i=1}^m (\tilde{x}_{ij} - \bar{x}) = 0$$

以下半角を用いければ

$$\sigma_B^2 = \bar{P}_m^2 \sum_{i=1}^m \sum_{\alpha=0}^{k_i} \frac{P_{i\alpha}}{P_m} (\tilde{x}_{i\alpha} - \bar{x})^2 + \bar{P}_m^2 \sum_{i+j} \sum_{\alpha=0}^{k_i} \frac{P_{i\alpha}}{P_m} (\tilde{x}_{i\alpha} - \bar{x})(\tilde{x}_{j\beta\alpha} - \bar{x})$$

あるいは一部を「簡化する前」の記法にもどして

$$(5.26) \quad \sigma_B^2 = \frac{P^2}{m} \sum_{h=1}^m \frac{P_h}{P} (\tilde{x}_h - \bar{x})^2 + \frac{P^2}{m} \sum_{i+j} \sum_{\alpha=0}^{k_i} \frac{P_{i\alpha}}{P} (\tilde{x}_{i\alpha} - \bar{x})(\tilde{x}_{j\beta\alpha} - \bar{x})$$

がえられる。

(4.26) の2番目の部分の項を組合わせると (1) でえた公式の拡張が容易にえられる。

σ_B^2 を少し違った形で表わせは

$$(5.27) \quad \sigma_B^2 = \frac{P^2}{m} \left\{ \sigma^2 - \sigma_{b.s}^2 + \sum_{i+j} \sum_{\alpha=0}^{k_i} \frac{P_{i\alpha}}{P} (\tilde{x}_{i\alpha} - \bar{x})(\tilde{x}_{j\beta\alpha} - \bar{x}) \right\}$$

となる。ここで

$$\sigma_{b.s}^2 = \frac{1}{m} \sum_{i=1}^m (\tilde{x}_i - \bar{x})^2$$

これは選ばれた標本を復元し (with replacing) . 大きさに比例した確率を用いて抽出する方法と比較したとき 効率と差がでてくる二つの原因を示すものである。(勿論 $P^2 \sigma^2 / m$ は

その前に選ばれたものを必ずもとに戻すようにして m 回の集

るを、大きさと比例した確率で抽出したときの $E^*(x'' | A)$ の分散であることは容易にわかる。

(5.26) と (5.27) をみると、単一要素の抽出の場合と殆んど同じ条件のもとで、大きさと比例した確率で行う系統的抽出は複元 P.P.S 抽出 (sampling with probability proportionate to size) よりも有効になることがわかる。詳細は省略する。これらは Lemma と定理 4 を適用すればえられる。この条件はまとめると次のようになる。P.P.S で系統的に標本を抽出して、二つの組

x_{1i}, \dots, x_{ki} および y_{ji}, \dots, y_{li} が単調で、一方が単調非増加、一方が単調非減少になっているとすると、 i 層と j 層間の共変動は負となるから、層から独立な抽出を行なうよりも有利である。

$$\sigma_{\alpha\beta}^0 = \sum (\bar{x}_{i\alpha} - E\bar{x}_{i\alpha})(\bar{x}_{j\beta} - E\bar{x}_{j\beta})$$

と定義すると、系統的 P.P.S 抽出の分散が、各層からの独立な P.P.S 抽出のものより小となるための concavity の条件は、 $\alpha < \beta$ とすると、

$$\sigma_{\alpha 1}^0 - \sigma_{\beta 1}^0 \leq \sigma_{\alpha 2}^0 - \sigma_{\beta 2}^0 \leq \dots \leq \sigma_{\alpha k}^0 - \sigma_{\beta k}^0 \leq 0$$

で与えられる。

6 等しい大きさの集団に対する系統的抽出

いま母集団が要素の集団からなっており、集団の大きさは等しい。すなわち同数の要素が含まれるものとしよう。はつきりさせるため母集団は M 個の集団からなり、各集団には N 個の要素が含まれるものとする。ただし $M = cm$, $N = kn$ である。このとき α 集団の α 要素を測定した特性の値を $x_{i\alpha}$ で表わし α 集団の全要素の合計を x_i で表わす。母集団の算術平均は \bar{x} から

$$M\bar{x} = \sum_{i=1}^M x_i$$

ここで

$$N_i \bar{x}_i = x_i$$

a 標本内の集団を完全に調査する場合

まず \bar{x} を \bar{x}' で推定したいものとする。ここで \bar{x}' は M 個の集団のうちから大きさと m なる系統的標本をとり、標本内の各集団の全要素を調べてえられる標本算術平均である。そうすると、

$$(6.1) \quad m\bar{x}' = \sum_{i=1}^m \bar{x}_i$$

と書ける。ここで \bar{x}_i は標本に選ばれた α 集団の平均値である。そうすると (1) から

$$\sigma_{\bar{x}'}^2 = \frac{\sigma_a^2}{m} \{1 + (m-1)\bar{f}_c\}$$

となる。ここで $M\sigma_a^2 = \sum_{i=1}^M (\bar{x}_i - \bar{x})^2$ 、また \bar{f}_c は (1.p.b) の \bar{f}_c で定義されるが x_i のところは \bar{x}_i で置きかえる。いま集団の無作為抽出の理論から

$$\sigma_a^2 = \frac{\sigma^2}{N} \{1 - (N-1)f\}$$

がえられる。この σ^2 は母集団分散すなわち

$$MN\sigma^2 = \sum_{i=1}^M \sum_{j=1}^N (x_{ij} - \bar{x})^2$$

で、 f は集団内要素の組内相関係数、すなわち

$$\sigma^2 f = \sigma_a^2 - \sigma_w^2 / N - 1$$

で与えられる。ただし

$$MN\sigma_w^2 = \sum_{i=1}^M \sum_{j=1}^N (x_{ij} - \bar{x}_i)^2$$

よって

$$(6.2) \sigma_{\bar{x}'}^2 = \frac{\sigma^2}{mN} \{1 + (N-1)f\} \{1 + (m-1)\bar{f}_c\}$$

となる。(6.2)の3つの因子の中 σ^2/mN は繰返しを許して大きさ mN の標本をとったときの分散で $1 + (N-1)f$ は無作為を用いたことによつて生ずる因子、また $1 + (m-1)\bar{f}_c$ は無作為を系統的に抽出したことによる因子である。

6. 層化と副次抽出

層化および副次抽出の可能性を考える場合には、可能な計画の数は極めて大きくなる。たとえば、母集団を層化して、抽出単位を大きさに比例した確率でぬき、次に系統的な副次抽出を行ない、再び系統的に副次抽出をして最後に愚作為な副次抽出を行なうというようにしてえた算術平均値の分散は簡単に計算できるであろう。しかし、そのような研究はそれを用いなければならぬ実際問題との関連をやるべきであろう。ここでは実際に生ずるであろう多くの可能性を考察するよりむしろ系統的標本の系統的副次抽出の結果のみを述べることにする。他の多くの計画の分散は定理1および2を用いれば容易に求められる。

いま系統的に選ばれた m 個の無落の標本の各々から系統的に n 個の要素を副次抽出するものとする。そうして \bar{x} の推定値を \bar{x}' とする。ただし σ_i 標本無落から選ばれた σ_i 要素を x'_{ij} とすると

$$\bar{x}' = \left(\frac{1}{m n}\right) \sum_{i=1}^m \sum_{j=1}^n x'_{ij} = \frac{1}{m} \sum_{i=1}^m \bar{x}'_i$$

で、

$$\bar{x}'_i = \left(\frac{1}{n}\right) \sum_{j=1}^n x'_{ij}$$

である。

定理2から再び

$$(6.3) \sigma_{\bar{x}'}^2 = \frac{\sigma^2}{mN} \{1 + (N-1)f\} \{1 + (m-1)\bar{f}_c\} + \frac{1}{M} \sum_{i=1}^M \frac{\sigma_i^2}{m n} \{1 + (m-1)\bar{f}_c\}$$

がえられる。ここで σ_i^2 は σ_i 無落内の分散で、 \bar{f}_c は (I.P.b) で定義された、 σ_i 無落内の平均的系列相関である。同様と系統的標本を抽出した標本内の m 個の無落を1つの母集団と考へて副次抽出を行なった場合にも、 \bar{x}' の分散は簡単に計算できる。これは *block* の標本を選び、標本 *block* のすべての世帯に追番号をつけ、このリストから系統的標本を選ぶような場合に生ずるものである。しかしここでの目的には(6.2)の分析が重要であるから、前にもどつて(6.2)を簡単に論じてみよう。

(6.2) からみよさせる最も重要な結論は、無落を系統的に抽出することは、それが望ましい場合でさえ無落を使うために分散が増加するという不利益を償いえないであろうということである。系統的な抽出によつて同様な相対的利得がえられるけれども、これらの利得は不等式

$$\{1 + (N-1)f\} \{1 + (m-1)\bar{f}_c\} < \frac{MN - mN}{MN - 1}$$

を成立させるほど大きくはない。

我々がこれまでに論じてきた問題は次の通りである。母集団の要素を確率変数と考えることによつて、異なる無落の要素間の平均的相関についてと同じく、単一の無落の要素間の平均的相関についても、同じ大きさの無落の系統的抽出の方が無落または個々の要素の愚作為または層化愚作為抽出より分散が小さくなるというるための条件がえられる。この解は直ちに求められるはずである。

c. 二次元の系統的抽出

二次元の系統的抽出の問題は *city* から *block* の標本を抽出するとか、*block* から *plot* の標本を抽出するとかの場合に

生じてくる。

city から block を選ぶ場合、問題を効果的に一次元に還元するため最も多く用いられる方法は、まず city の block またはその一部に、たとえば city の地図の右上隅から始めてオ1行の右から左に、次にオ2行の左から右というように番号をつけてゆくやり方である。そうするとこれらの block 番号の系統的標本、したがって block 自体が選ばれることになる。明らかにこの方法はもし隣接 block に高度の相関があるときは最も有効というわけにはゆかないはずである。というのは非現実的な可能性を例にとれば、可能な標本を city のブロックの列になってしまうかもしれないからである。オ2の二次元の系統的抽出の方法は、行の系統的標本および列の系統的標本を選ぶ、グリッドの標本をうるものである。この計画は行または列に沿って「肥沃度の傾斜」がある場合には非常に効率が悪くなる。

これら両者の方法の効率が悪いことの原因は、系統的標本の分散公式を調べれば判明する。もし上のように番号をつけてゆくとコログラフが上に凹とは考えられなくなるから、標本との食い違いが非常に大きくなるであろう。グリッド計画は、乗取の系統的抽出で系統的副次抽出を行なう場合の特別な形のものであるが、これについては(6.3)と \bar{P}_0 が負のときでさえ、域内相関係数 ρ が、 σ^2 の値を大にする程十分小さくなることを考えてみればよい。

明らかに(6.3)は可能な標本が ρ をできるだけ小さくするように定義されることを示している。正方形の面積には同じ処理をもつ Knut Vi 形の方形プロットを可能な標本であると定義すればどうなるであろうか、不規則な面積に対しても同様な可能な標本の定義が容易に与えられる。この問題はしかし

将来の研究にまたねばならない?)

By W. G. Madow (A. M. S. Vol. 1949, PP333-

この論文の引用文献の10(6)は著者の注意をひいた。その中のデータ、特にオ3表は上に述べた見解と一致している。

参 考 文 献

- (1) W. G. Madocw and L. H. Madocw "On the theory of Systematic sampling. I." A.M.S. Vol. 15 (1944) PP 1-27
- (2) L. H. Madocw "Systematic Sampling and its relation to other sampling design" J.A.S.A. Vol 41 (1946) PP 204-207
- (3) W. G. Cochran "Relative accuracy of systematic and stratified random samples for a certain class of populations" A.M.S. Vol 17. (1946) PP 164-177
- (4) G. H. Hardy, J. E. Littlewood and G. Polya Inequalities, Camb. Univ Press Lond & Newyork 1934. P43
- (5) M. H. Hansen an W. N. Hurwity "On the theory of sampling from finite populations A.M.S. Vol 14 (1943) PP 333-362
- (6) P. G. Homeyer & C. A. Black "Sampling replicated field experiments on oats for yield determinations" Soel. Sci. Soc. Proc. Vol 11 (1946) PP 341-344

5 二次元のサンプリングについて
Problems in Plane Sampling

1 要 約

1次元における系統的抽出 (systematic) と層化無作為抽出 (stratified random sampling) の相対的精度を考察した後、一次元抽出の抽出誤差 (linear sampling error) の推定の問題を論じた。

個體的な標本抽出 (sampling) の方法を述べて、これらの方法の精度を表わす式を導いた。これらの式を特に理論的にも実際的にも当て得ていると思われる相関函数 (Correlation function) に関連づけて、大標本に対して比較した。その結果多くの場合系統的抽出 (systematic sampling) は層化抽出 (stratified sampling) よりずっと精度のよいことがわかった。抽出誤差 (sampling error) を推定する方法を再び考察し、実例を与えた。この論文には抽出の行なわれる母集団の傾向 (trend) の問題に関しても幾つかの注意を与えてある。

2 一次元における系統的標本と層化無作為標本の精度

(accuracy systematic and stratified random samples in one dimension) W. G. Cochran (1) は X_1, X_2, \dots, X_n なる母集団から、単純無作為 (random) (ト)、層化無作為 (stratified random) (スト)、系統的 (systematic) (システム) 抽出を行なつて大きさ n の標本を抽出したときの平均値の分散を表す式を与えている。

母集団は要素 X_1, X_2, \dots, X_n が $E(X_i) = \mu, \quad E(X_i - \mu)^2 = \sigma^2, \quad E(X_i - \mu)(X_{i+u} - \mu) = \rho_u \sigma^2$ なる母集団から抽出されると仮定した。

ここで、 $u < v$ なるとき常に $\rho_u > \rho_v > 0$

とすると、

$$(1) \sigma_s^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{kn} \right) \left(1 - \frac{2}{kn} \frac{kn-1}{u} \right) (kn-u) \rho_u$$

$$(2) \hat{\sigma}_{st}^2 = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \left[1 - \frac{2}{kn(k-1)} \sum_{u=1}^{kn-1} (kn-u) \rho_u\right]$$

$$(3) \frac{\sigma_{st}^2}{S_y^2} = \frac{\sigma^2}{n} \left(1 - \frac{1}{k}\right) \cdot \left[1 - \frac{2}{kn(k-1)} \sum_{u=1}^{kn-1} (kn-u) \rho_u\right]$$

$$\therefore \rho_u + \frac{2-k}{n(k-1)} \sum_{u=1}^{kn-1} (kn-u) \rho_{kn-u}$$

なる式を導いた。

ρ_u の線形函数であるこれらの式を用いて、Cochran は幾つかのタイプのコレログラム (Correlogram) についてそれぞれの抽出方法の相対的な効率 (efficiency) を比較した。もし我々が

(a) 各 x_i が平均値 μ , 分散 σ^2 の母集団から抽出された標本で、

(b) μ_i は平均値 μ のまわりに、分散 σ^2 で分散する、

(c) $E\{(u_i - \mu)(u_j - \mu)\} = \rho_{ij} \sigma^2$ かつ

$$(d) \rho_u = \frac{1}{kn} \sum_{i=1}^{kn-u} \rho_i \cdot i + u$$

という仮定を置くならば、(1)、(2)、(3)、が成立するためには、右辺に、

$$\frac{1}{n} \left(1 - \frac{1}{k}\right) \cdot \frac{1}{kn} \sum_{i=1}^{kn} \sigma_i^2$$

なる定数を掛けなければならぬことは容易にわかる。従つて、Cochran の結果は色々な抽出方法の理論的な最大の相対効率を与えるものであることを忘れてはならない。ここで ρ_u は u だけ離れた標本間の平均的な相関である。この結果はおそらく、各点で平均的な変延度があり、標本がこの平均値のまわりで poisson 分散をしている様な昆虫変延度の標本調査で興味あるものである。

この場合相対的な変延度は、

$$\frac{1}{n} \left(1 - \frac{1}{k}\right) \cdot \frac{1}{kn} \sum_{i=1}^{kn} \mu_i^2 \sim \frac{1}{n} \left(1 - \frac{1}{k}\right) \mu^2$$

である。

我々が連続過程 (Continuous process) を標本抽出⁽¹⁾ するとき、 n が大なら (1)、(2)、(3) の積分による等式を得ることかできる。

$$(4) \frac{\sigma_{st}^2}{S_y^2} \sim \frac{\sigma^2}{n} \left[1 - \frac{2}{d^2} \int_0^d (d-u) \rho_u \delta u\right]$$

$$(5) \frac{\sigma_{st}^2}{S_y^2} \sim \frac{\sigma^2}{n} \left[1 - \frac{2}{d} \int_0^\infty \rho_u \delta u + 2 \sum_{u=1}^\infty \rho_{du}\right]$$

ここで ρ は標本の u だけ離れた次々の要素間の相関で d は標本間の平均距離である。従つて我々は次式を得る。

$$\frac{\sigma_{st}^2 - \sigma_{st}^2}{\sigma^2} \sim \frac{2}{d} \left[\int_0^d \frac{u}{d} \rho_u \delta u + \int_1^\infty \rho_u \delta u - d \sum_{u=1}^\infty \rho_{du} \right]$$

これは我々が層化無作為および系統的標本抽出の効率間の差のグラフを利用して研究に屢々用い得るものである。Fig 1 は4種類のコレログラムについてどの様に行なわれるかを示すものである。

連続的マルコフ過程 (Continuous Markov scheme) の場合には

$\rho_u = \rho^u$ だから

$$\frac{\sigma_{st}^2}{S_y^2} \sim \frac{\sigma^2}{n} \left[1 + \frac{2}{\log \rho d} + \frac{2}{(\log \rho d)^2} - \frac{2 \rho d}{(\log \rho d)^2}\right]$$

$$\frac{\sigma_{st}^2}{S_y^2} \sim \frac{\sigma^2}{n} \left[1 + \frac{2}{\log \rho d} + \frac{2 \rho d}{1 - \rho d}\right]$$

これは Cochran の結果と一致する。

3. 複製と誤差の推定 (Replication and the estimation of error) Yates [2] 系統的標本 (systematic sample) の誤差の推定にかゝる困難を指摘した。しかしこの点について上記の公式を用いて研究することはやり甲斐があると思われる。

単純無作為、層化無作為、系統的抽出で、 n が大きく k が一定と考えられるなら、平均値の推定量の分散は $\sigma^2 F(k)/n$ の形となる。ここで $F(k)$ は事実上 n と独立である。従つて、もし標本誤差を

(1) 実際において我々が連続過程から標本を抽出する場合には、これは殆ど n が大なる連続過程 (discontinuous process) とみなすときにのみ可能である。

推定する方法があれば、我々は標本抽出を行なう系列 (Series) を幾つかの部分 (即ちブロック) に等分して、各部分の平均値の誤差を求めるとともに、全平均値の誤差のより正確な推定量を得るため、これを結合 (Combine) することが可能となる。実際 k が非常に大きい場合には、これらの部分を無作為に選び、誤差の推定量を求めることによつて、観測値の個数を減少させることが考えられる。層化無作為抽出では、 k は N と完全に独立だから、我々は各層 (strata) からの誤差の推定量を結合することができる。このことから、各層あたり各個の無作為に選ばれた要素をとり、誤差の推定量を求めるのに、自由度 $k-1$ の分散の組を結合するという普通に行われている方法が導かれる。もし我々が、標本を排斥 (exclusive)、すなわちどの二要素も一致しないように作れば、標本平均の推定分散を求めるにはこの分散に $1 - k/N$ をかければよい。

同様にして、任意に選ばれた出発点から始まる十分な長さをもつ各個の系統的標本 (systematic sample) の組を用いて、系統的標本の平均値の分散を推定することができる。しかしこの標本抽出は、実際に用いるには一層難しいから他の方法を考える必要がある。我々の系統標本は、この系列を分割した部分、即ちブロックの各々の中では一定となるよう選ぶことができるから、この抽出方法には全部で各個の系統標本だけが含まれる。即ち、我々の各個の標本を等間隔で選ぶという Yates の提唱する方法に従うべきである。この場合それらはより大きい系統的標本の副次標本 (subsamples) である。

この後者の方法は簡単でもあり左範囲の計画にも取り入れることが出来るという長所を有するが、その使用には非常に慎重な考慮が必要である。離散的な場合を考へれば、我々は

$$(b) \sigma^2 \left(1 - \frac{2}{k-1} \sum_{u=1}^{\infty} \rho_{ku} + \frac{\sigma^2}{k-1} \sum_{u=1}^{\infty} \rho_{ku} \right)$$

を推定したい。しかし各個の等間隔系統的標本 (evenly-spaced

systematic sample) にもとづく分散の推定量には、 ρ_{ku}/k なる形の項だけしか含まれない。一方各個の無作為に選ばれた (randomly chosen) 系統的標本にもとづく分散の推定量には明らかに制限があり、多くの場合においてより代表的なものとなる。例えば $k=10, g=4$ とすると、我々は分散の推定量における観測された相関 (observing the correlation) ρ_1, \dots, ρ_{15} の相対的な出現度を比較することができる。

これの6つの例をオ/表にあげる。乱数は Fisher と Yates の表からとつた。 ρ_u と ρ_{16-u} との出現度数は等しいから、一語に示してある。この表は、最初の二つの無作為に選ばれた標本のように random 系統的標本に近い無作為標本でも、このコレログラムを系統的に抽出するのを避けられるかを示している。多くの場合、どちらの方法を用いても可成り良い結果の得られることは明らかである。しかし、後の方法を用いた方がより正確である。オ/表に示した標本を用いて色々な型のコレログラムに対する比較をオ2表で行なつた。初等階級の系統的標本の組を分解できるような多くの種類のコレログラムを理論的に仮定することも可能であるが、しかし我々は最終的には、経験から得られたコレログラムの型について決定を下さなければならぬ。この点については、二次元の標本抽出を考察し、あとで、もつと詳しく考えることにする。

オ/表 それぞれ $1/k$ 単位毎に抽出を行なつて4つの系統的標本をとるとき、分散推定値に現れる系列相関係数 $\rho_1, \rho_2, \dots, \rho_{15}$ の出現度

D	等間隔に抽出した4つの系統的標本	4つの系統的標本の無作為に選んだ出発点						頻度数
		4, 7 8, 12	3, 7 8, 12	3, 6 10, 13	4, 6 7, 12	2, 8 11, 15	2, 5 11, 16	
1	15	1	1		1			3
2	14				1		1	2
3	13	1		2	1	2		6
4	12	2	2	1		1	1	7
5	11	1	2				2	5
6	10			1	1	1	1	4
7	9		1	2	1	2	1	7
8	8	2			2			4

表2 系統的標本で推定した $\frac{2}{15} \sum_{u=1}^{15} P_u$ の値

P_u	等間隔抽出 の系統的 標本	系統的標本の無作為に選んだ出発点						平均値	期待値
		1	2	3	4	5	6		
$1-0.2 (u=1, 5)$	0.17	0.27	0.20	0.17	0.30	0.17	0.13	0.21	0.27
$1-0.1 (u=1, 10)$	0.53	0.62	0.58	0.53	0.60	0.53	0.53	0.57	0.60
2^{-u}	0.04	0.13	0.12	0.06	0.15	0.06	0.07	0.10	0.13
$2^{-u/4}$	0.58	0.66	0.64	0.60	0.66	0.60	0.60	0.63	0.65
Kendallの系列1	-0.14	0.03	0.00	-0.05	0.16	-0.05	-0.05	0.01	0.07

このような抽出誤差の計算は当該系列を分割した各部分即ちブロックの対応する要素間の相関は無視できると仮定している。例えばこの場合には P_{16} と上記は無視できる。この場合 $P_{16} = 1/16$ であるから結局 (6) において要求される項は

$$2 \left(\frac{1}{15} \sum_{u=1}^{\infty} P_u - \frac{16}{15} \sum_{u=1}^{\infty} P_{16u} \right) = 0.56$$

となつて我々が推定しようとする項 $\frac{2}{15} \sum_{u=1}^{15} P_u = 0.65$ と少し違つている。

4. 二次元での標本抽出法 (Method of sampling in 2 dimensions)

ある二次元空間から標本を抽出する方法の数は非常に多い。(3) 何故ならどの方向に対しても単純無作為、層化無作為、あるいは系統的抽出を用いることができるからである。この様に我々はこれらの方法の可能な組合せの何れをも考察することができる。例えば一方向に単純無作為で他の方向で系統的なとり方をする抽出は r, S といい記号で表わす。

更に我々は一方向が他に対して一直線になるような標本の組合せを (3) 一般に我々はこの二次元空間が矩形であると考え、任意の形の領域に対しても容易に同様な結論を導くことができる。

ることとできるし、また、独立にきめることとできる。添字 i は一直線にとつた aligned 標本を、また o は独立な標本を表わすのに用いる。例えばシステム r, S, o に従つて標本抽出ができるであろう。標本抽出の幾つかの方法の例を Fig 2 に与える。

Fig 2 幾々の相関函数について比較した系統的標本と層化無作為標本の効率の図示

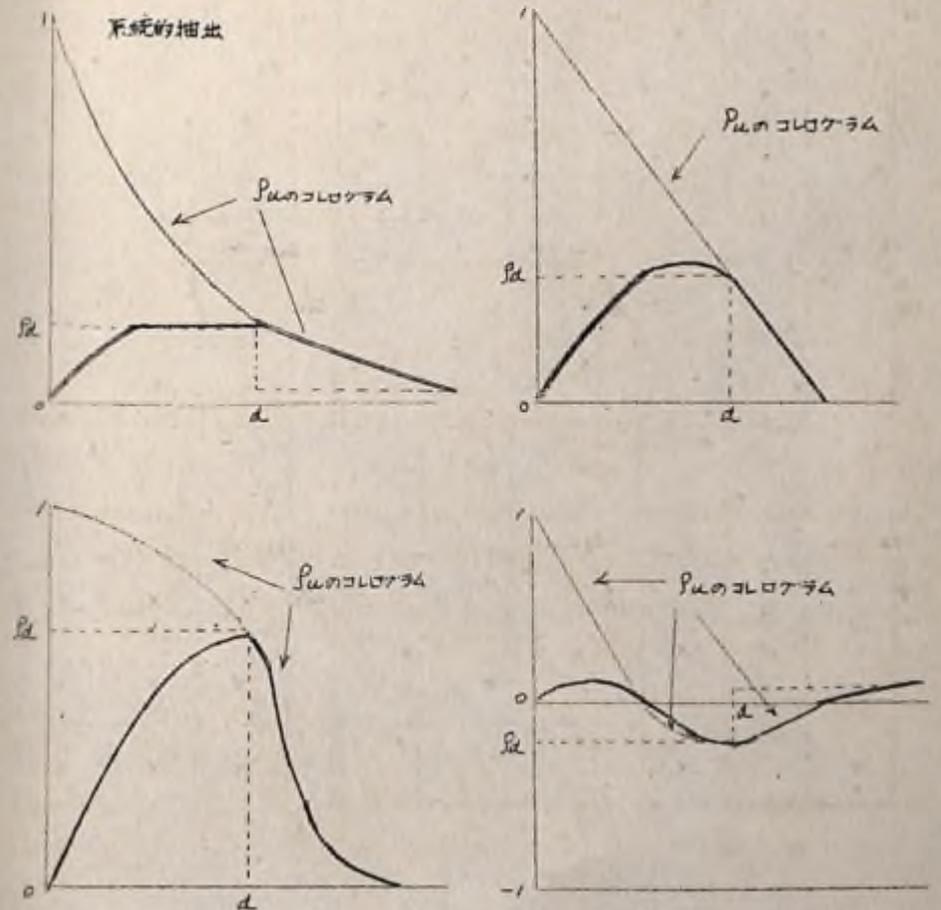
本線は函数
 $f(u) = u P_u / d$
 で点線は函数
 $S_2(u) = P_{id}$
 を与える。

$$0 \leq u \leq d$$

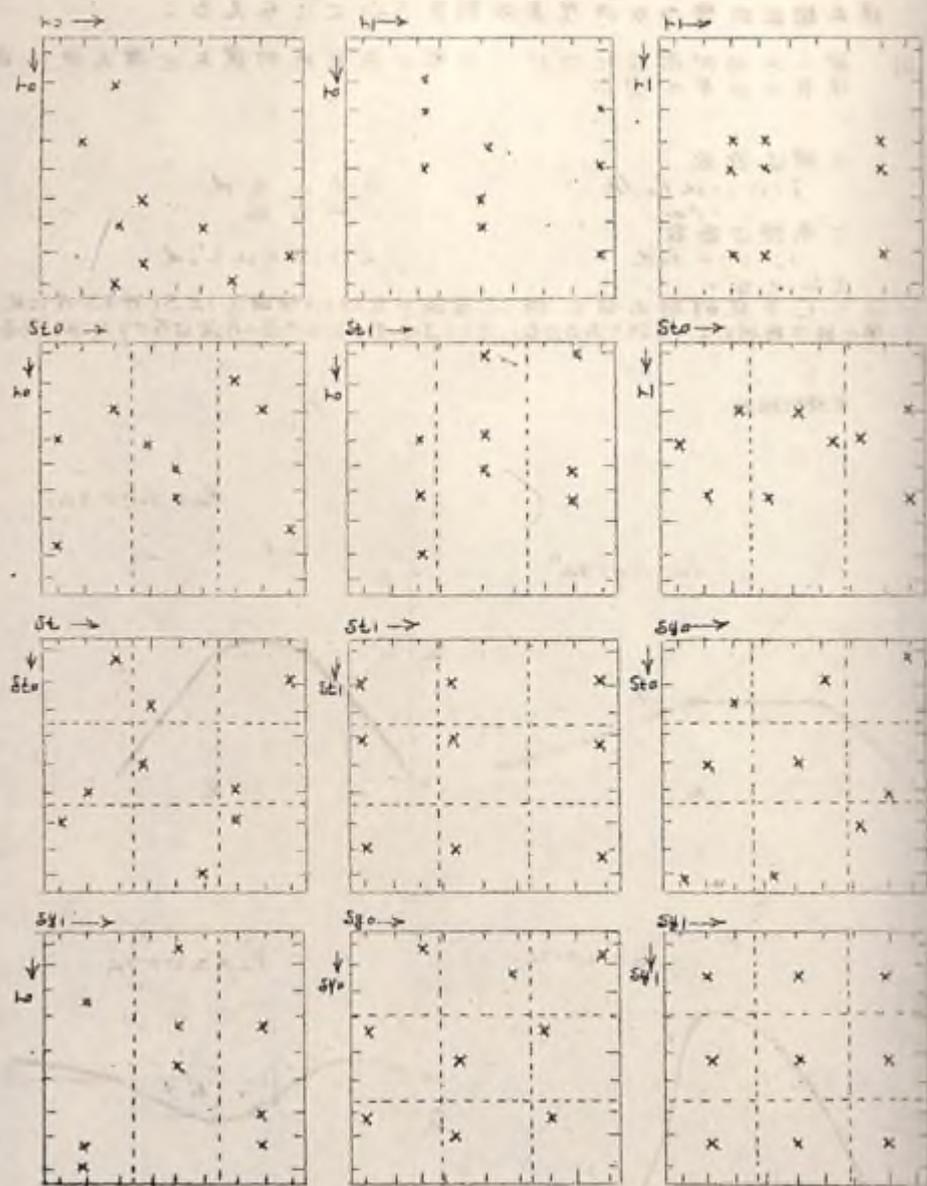
$$d \leq u$$

$$(i-1)d < u \leq id$$

このように系統的抽出は本線下の面積が点線下の面積より大きいから従つて、層化無作為抽出より有効であるかないか決まる。最も効率の高い方法は各グラフに示してある。



オ2区 四角の抽出法 この場合 $n_1 = n_2 = k_1 = k_2 = 3$



5 二次元の標本抽出の精度

(Accuracy of sampling in two dimensions)

我々は、二方向における平均間隔が k_1, k_2 となる様に、要素 X_{ij} ($i = 1, 2, \dots, n_1, k_1; j = 1, 2, \dots, n_2, k_2$) から抽出された、 n_1, n_2 個の要素からなる標本を考えよう。(これは仮想的無限母集団から抽出された1つの有限母集団を形成する。) これらのパラメータ (parameter) はもし必要なら抽出方法のあとに [] をつけて示すことにする。例えば r, S_{ij} ($n_1, k_1; n_2, k_2$)

X は考えている抽出方法で作られた標本平均を表わし、 X' でこの標本の構成因子を表わす。同様に X_{ij} は

$$E(X_{ij}) = \mu \quad E(X_{ij} - \mu)^2 = \sigma^2$$

$$E(X_{ij} - \mu)(X_{i+u, j+v} - \mu) = \rho_{ijuv} \sigma^2$$

なる母集団から抽出されると仮定する。

更に我々は

$$\rho_{uv} = \rho_{-u, -v} \text{ を}$$

$$\sum_i \sum_j \rho_{ijuv} = (k_1 n_1 - 1)(k_2 n_2 - 1) \rho_{uv}$$

なる関係で定義するため、 i と j のすべての可能な値にわたつて、 ρ_{ijuv} を平均することができる。これらの定義の目的は、有限母集団のパラメータを取扱う場合の困難性を、この母集団自身が無限母集団から得られた1つの標本であると考えることによって、除去することにある。Cochranも同様な手法を用いている。

5.2. 単純無作為抽出

$$\sigma^2(X) = \frac{1}{2} E(X_1 - X_2)^2 = E(X_1 - \mu)^2 - E(X_1 - \mu)(X_2 - \mu)$$

なることは容易に証明できる。ここで X_1 と X_2 は独立な標本である。同様にして

$$E(X_1 - \mu)(X_2 - \mu) = E(X'_1 - \mu)(X'_2 - \mu)$$

$$= \frac{\sigma^2}{k_1 k_2 n_1 n_2} \left[1 + \frac{1}{k_1 k_2 n_1 n_2} \sum_i \sum_j (k_1 n_1 - 1)(k_2 n_2 - 1) \rho_{ij} \right]$$

ここで二重和 (double summation) (4) は $u=v=0$ を除いて $|u| \leq k_1, |v| \leq k_2$ で与えられる範囲 S の上で存在する。ここで我々は種々な型の無作為標本について $E(X, -\mu)^2$ を評価しなければならぬ。

$$E(X, -\mu)^2 = \frac{\sigma^2}{n_1 n_2} \cdot \left[1 + \frac{n_1 n_2 - 1}{k_1 k_2 n_1 n_2 (k_1 k_2 n_1 n_2 - 1)} \sum \sum (k_1 n_1 - |u|) (k_2 n_2 - |v|) p_{uv} \right]$$

r_0, r_2 に対して

$$= \frac{\sigma^2}{n_1 n_2} \left[1 + \frac{n_1 - 1}{k_1 k_2 n_1 (k_1 k_2 n_1 n_2 - 1)} \sum \sum (k_1 n_1 - |u|) (k_2 n_2 - |v|) p_{uv} \right. \\ \left. + \frac{2(n_2 - 1)}{k_2 n_2 (k_2 n_2 - 1)} \sum_{v=1}^{k_2 n_2} (k_2 n_2 - v) p_{0v} \right]$$

r_1, r_0 に対して

$$= \frac{\sigma^2}{n_1 n_2} \left[1 + \frac{(n_1 - 1)(n_2 - 1)}{k_1 k_2 n_1 n_2 (k_1 k_2 n_1 n_2 - 1)} \sum \sum (k_1 n_1 - |u|) (k_2 n_2 - |v|) p_{uv} \right. \\ \left. + \frac{2(n_2 - 1)}{k_2 n_2 (k_2 n_2 - 1)} \sum_{v=1}^{k_2 n_2} (k_2 n_2 - v) p_{0v} + \frac{2(n_1 - 1)}{k_1 n_1 (k_1 n_1 - 1)} \sum_{u=1}^{k_1 n_1} (k_1 n_1 - u) p_{u0} \right]$$

r_1, r_1 に対して

故に

$$(7) \left\{ \begin{aligned} \sigma^2(r_0 r_0) &= \frac{1}{n_1 n_2} \left(1 - \frac{1}{k_1 k_2}\right) \sigma^2 \\ &\cdot \left[1 - \frac{1}{k_1 k_2 n_1 n_2 (k_1 k_2 n_1 n_2 - 1)} \sum \sum (k_1 n_1 - |u|) (k_2 n_2 - |v|) p_{uv} \right] \end{aligned} \right.$$

$$(8) \left\{ \begin{aligned} \sigma^2(r_1 r_0) &= \frac{1}{n_1 n_2} \left(1 - \frac{1}{k_1 k_2}\right) \sigma^2 \left[1 - \frac{k_1 k_2 n_2 - 1}{(k_1 k_2 - 1) k_1 k_2 n_1 n_2 (k_1 k_2 n_1 n_2 - 1)} \right. \\ &\cdot \sum \sum (k_1 n_1 - |u|) (k_2 n_2 + |v|) p_{uv} + \frac{2(n_2 - 1)}{k_2 n_2 (k_2 n_2 - 1)} \sum_{v=1}^{k_2 n_2} (k_2 n_2 - v) p_{0v} \end{aligned} \right.$$

(4) 特に断らない限り、一般に二重和 (double summation) は $u=v=0$ を除く係数が正となる狭い範囲上で存在する。

$$(9) \left\{ \begin{aligned} \sigma^2(r_1 r_1) &= \frac{1}{n_1 n_2} \left(1 - \frac{1}{k_1 k_2}\right) \sigma^2 \left[1 - \frac{k_1 k_2 (n_1 n_2 - 1) - 1}{(k_1 k_2 - 1) k_1 k_2 n_1 n_2 (k_1 k_2 n_1 n_2 - 1)} \right. \\ &\cdot \sum \sum (k_1 n_1 - |u|) (k_2 n_2 + |v|) p_{uv} + \frac{2k_1 (n_2 - 1)}{(k_1 k_2 - 1) n_2 (k_2 n_2 - 1)} \\ &\cdot \sum_{v=1}^{k_2 n_2} (k_2 n_2 - v) p_{0v} + \frac{2k_2 (n_1 - 1)}{(k_1 k_2 - 1) n_1 (k_1 n_1 - 1)} \sum_{u=1}^{k_1 n_1} (k_1 n_1 - u) p_{u0} \end{aligned} \right.]$$

なることは容易に証明できる。

5. 層化無作為抽出 (Stratified random sampling)
もし i 層で抽出される要素の平均値 \bar{x}_i が \bar{x}_j と独立なら、層別標本抽出の方法に対する分散を推定することかできる。
何故なら このとき

$$E(X - \bar{x})^2 = E(\bar{x}_i' - \bar{x}_i)^2 / n$$

であるから。ここで \bar{x} は標本抽出の行なわれる層別母集団の平均値である。故に

$$(10) \left\{ \begin{aligned} \sigma^2(sto r_0) &= \frac{1}{n} \sigma^2 \{ r_0 r_0 (1 - k_1 / n_1) \} \\ &= \frac{1}{n_1 n_2} \left(1 - \frac{1}{k_1 k_2}\right) \sigma^2 \left[1 - \frac{1}{k_1 k_2 n_1 n_2 (k_1 k_2 n_1 n_2 - 1)} \right. \\ &\cdot \sum \sum (k_1 |u|) (k_2 n_2 - |v|) p_{uv} \end{aligned} \right.$$

$$(11) \left\{ \begin{aligned} \sigma^2(sto r_1) &= \frac{1}{n_1 n_2} \left(1 - \frac{1}{k_1 k_2}\right) \sigma^2 \left[1 - \frac{1}{k_1 k_2 n_1 n_2 (k_1 k_2 - 1)} \right. \\ &\cdot \sum \sum (k_1 - |u|) (k_2 n_2 - |v|) p_{uv} \\ &+ \frac{2k_1 (n_2 - 1)}{(k_1 k_2 - 1) n_2 (k_2 n_2 - 1)} \sum_{v=1}^{k_2 n_2} (k_2 n_2 - v) p_{0v} \end{aligned} \right.$$

$$(12) \left\{ \begin{aligned} \sigma^2(sto r_2) &= \frac{1}{n_1 n_2} \left(1 - \frac{1}{k_1 k_2}\right) \sigma^2 \left[1 - \frac{1}{k_1 k_2 (k_1 k_2 - 1)} \right. \\ &\cdot \sum \sum (k_1 - |u|) (k_2 - |v|) p_{uv} \end{aligned} \right.$$

他の抽出方法の分散を推定するため、我々は(8) - (12) 式を弄びくのに用いた一般公式を使うことにする。

X'_i を標本 X の任意の要素とすると

$$\begin{aligned} (X - \bar{X})^2 &= \frac{1}{n_1 n_2} \left[\sum (X'_i - \bar{X})^2 - \sum (X'_i - \mu)^2 \right] \\ &= \frac{1}{n_1 n_2} \left\{ \sum (X'_i - \bar{X})^2 - \frac{n_1 n_2 - 1}{n_1 n_2} \sum (X'_i - \mu)^2 \right. \\ &\quad \left. + \frac{2}{n_1 n_2} \sum \sum (X'_i - \mu) (X'_j - \mu) \right\} \end{aligned}$$

よって

$$\begin{aligned} (13) \left\{ \begin{aligned} \sigma^2(X) &= E(X - \bar{X})^2 \\ &= \frac{k_1 k_2 n_1 n_2 - 1}{k_1 k_2 n_1 n_2} \sigma^2 \left(1 - \frac{1}{k_1 k_2 n_1 n_2} \sum \sum (k_1 n_1 - 1) \right. \\ &\quad \cdot (k_2 n_2 - 1) p_{uv} \left. - \frac{n_1 n_2 - 1}{n_1 n_2} \sigma^2 + \frac{n_1 n_2 - 1}{n_1 n_2} E(X'_i - \mu) (X'_j - \mu) \right. \\ &\quad = \frac{1}{n_1 n_2} \left(1 - \frac{1}{k_1 k_2} \right) \sigma^2 \left(1 - \frac{1}{k_1 k_2 n_1 n_2} \sum \sum (k_1 n_1 - 1) \right. \\ &\quad \cdot (k_2 n_2 - 1) p_{uv} + \frac{k_1 k_2 (n_1 n_2 - 1)}{k_1 k_2 - 1} \frac{E(X'_i - \mu) (X'_j - \mu)}{\sigma^2} \left. \right) \end{aligned} \right. \end{aligned}$$

このように $E(X'_i - \mu) (X'_j - \mu) / \sigma^2$ の推定が可能でありさえすれば(13)式はすべての抽出方法に対する誤差を与える。

1例として(12式)を算いてみよう。任意の層内要素 X'_i をとれば、オスの層内要素 X'_j は、 X'_i と同じ層の中に X'_j が占めることのできない $k_1 k_2 - 1$ 個の場所が存在するというを除いては、 X'_i に関して無作為に配置されるであろう。

かくして期待された相関 $E(X'_i - \mu) (X'_j - \mu) / \sigma^2$ は

$$(14) \left\{ \begin{aligned} &\frac{k_1^2 k_2^2 n_1 n_2 (n_1 n_2 - 1)}{k_1^2 k_2^2 (n_1 n_2 - 1)} \sum \sum (k_1 n_1 - 1) (k_2 n_2 - 1) p_{uv} \\ &- \frac{1}{k_1^2 k_2^2 (n_1 n_2 - 1)} \sum \sum (k_1 - 1) (k_2 - 1) p_{uv} \end{aligned} \right.$$

で与えられる。

(14) を(13)に代入すれば、 $St_0 St_0$ に対する分散の式(12)が得られる。

同様にして $St_1 St_1$ に対する式

$$\begin{aligned} (15) \left\{ \begin{aligned} &E(X'_i - \mu) (X'_j - \mu) \\ &= \frac{1}{k_1 k_2 (n_1 n_2 - 1)} \left(\frac{1}{k_1 k_2 n_1 n_2} \sum \sum (k_1 n_1 - 1) (k_2 n_2 - 1) p_{uv} \right. \\ &\quad - \frac{1}{k_1 k_2 n_1} \sum \sum (k_1 n_1 - 1) (k_2 - 1) p_{uv} \\ &\quad - \frac{1}{k_1 k_2 n_2} \sum \sum (k_1 - 1) (k_2 n_2 - 1) p_{uv} + \frac{1}{k_1 k_2} \sum \sum (k_1 n_1) \\ &\quad \cdot (k_2 n_2) p_{uv} + \frac{2(k_1 k_2 n_1 - 1)}{k_1 n_1 (k_1 n_1 - 1)} \sum_{u=1}^{k_1} (k_1 n_1 - 1) p_{0v} \\ &\quad - \frac{2(k_1 k_2 - 1)}{k_1 (k_1 - 1)} \sum_{u=1}^{k_1} (k_1 - 1) p_{uv} + \frac{2(k_1 k_2 n_2 - 1)}{k_2 n_2 (k_2 n_2 - 1)} \sum_{v=1}^{k_2} (k_2 n_2 - 1) p_{0v} \\ &\quad \left. - \frac{2(k_1 k_2 - 1)}{k_2 (k_2 - 1)} \sum_{v=1}^{k_2} (k_2 - 1) p_{uv} \right) \end{aligned} \right. \end{aligned}$$

このようにして我々はすべての型の層化無作為抽出について $\sigma^2(X)$ を計算することかできる。

5c. 系統的抽出 (Systematic sampling)

層化無作為抽出において用いたと同様互やり方で、系統的抽出の分散を計算するために(13)式を用いることかできる。

可能な3つの抽出方法に対する $E(X'_i - \mu) (X'_j - \mu)$ の値を下に掲げる。

5y1 5y1 については

$$(16) E(X'_i - \mu) (X'_j - \mu) = \frac{1}{n_1 n_2 (n_1 n_2 - 1)} \sum \sum (n_1 - 1) (n_2 - 1) p_{k_1 u, k_2 v}$$

Sy₁₀ については、

$$(17) \left\{ \begin{aligned} E(X'_i - u)(X'_j - u) &= \frac{1}{k_2^2 \pi_1 \pi_2 (\pi_1 \pi_2 - 1)} \Sigma \Sigma (\pi_1 - u_1) \cdot \\ &\quad (k_2 \pi_2 - v_1) P_{k_1, u_2} \\ &- \frac{2(k_2 - 1)}{k_2^2 \pi_2 (\pi_1 \pi_2 - 1)} \frac{k_2 \pi_2}{\Sigma_{u=1}^{k_2 \pi_2} (k_2 \pi_2 - u) P_{0, u}} \end{aligned} \right.$$

Sy₀ Sy₀ については

$$(18) \left\{ \begin{aligned} E(X'_i - u)(X'_j - u) &= \frac{1}{k_1 k_2 (\pi_1 \pi_2 - 1)} \left\{ \frac{1}{k_1 k_2 \pi_1 \pi_2} \Sigma \Sigma (k_1 \pi_1 - u_1) \cdot \right. \\ &\quad \cdot (k_2 \pi_2 - v_1) P_{u, v} - \frac{1}{k_1 k_2 \pi_1} \Sigma \Sigma (k_1 \pi_1 - u_1) (k_2 - v_1) P_{u, v} \\ &\quad - \frac{1}{k_1 k_2 \pi_2} \Sigma \Sigma (k_1 - u_1) (k_2 \pi_2 - v_1) P_{u, v} + \frac{1}{k_1 k_2} \Sigma \Sigma (k_1 - u_1) \\ &\quad \cdot (k_2 - v_1) P_{u, v} + \frac{k_1}{k_2 \pi_1} \Sigma \Sigma (\pi_1 - u_1) (k_2 - v_1) P_{k_1, u, v} \\ &\quad - \frac{2k_1}{k_2} \frac{k_2}{\Sigma_{v=1}^{k_2} (k_2 - v) P_{0, v}} + \frac{k_2}{k_1 \pi_2} \Sigma \Sigma (k_1 - u_1) (\pi_2 - v_1) P_{u, k_2, v} \\ &\quad \left. - \frac{2k_2}{k_1} \frac{k_1}{\Sigma_{u=1}^{k_1} (k_1 - u) P_{u, 0}} \right\} \end{aligned} \right.$$

(18) の導き方は(15)のそれと比較できるものである。

6 標本を直線的にとる場合の影響 (Effect of alignment)

直線的抽出 (alignment) の効果は、色々の標本の分散を調べるかあるいは直接(13)式を用いて研究することかできる。

無作為および層に無作為抽出での直線的抽出の効果は標本分散を

$$\Sigma \Sigma a_{uv} (P_{0, v} - P_{u, v}) + \Sigma \Sigma b_{uv} (P_{u, 0} - P_{u, v}) \quad \text{ここで } a_{uv} > 0, \quad b_{uv} > 0$$

なる量だけ増加させる。

これは単純減少の相関函数および、実際において出現する大多数の函数に対しては正である。この様に直線的抽出は普通無作為および層に無作為抽出の分散を増加させる。

系統的標本については問題はもつと複雑であるが、しかし大體

にまつて、分散は

$$\Sigma \Sigma a_{uv} (P_{k_1, u, k_2, v} - \bar{P}_{k_1, u, k_2, v})$$

なる量だけ増加する。ここで $a_{uv} \geq 0$ で、 $P_{k_1, u, k_2, v}$ は 0 でない u, v に対しては中心が $P_{k_1, u, k_2, v}$ 中心 $P_{0, k_2, v}$ なる長さ k_1 の直線上の平均値である。($v=0$ なるときも同様)

これが正となるか負となるかは、相関函数に關係してくるから、この問題は出現する相関函数の型について研究しなければならない。

7 極限における式 (Limiting forms)

連続的過程で、 π_1, π_2 が大きいとき、我々は一次元の抽出 (linear sampling) の場合と同様にして、 $\Sigma \Sigma P_{d_1, u, d_2, v}$ が収斂する限り、抽出分散の積を近似を得ることかできる。

即ち

$$(19) \sigma^2(r_0, r_0) = \sigma^2(st_0, r_0) \sim \sigma^2 / \pi_1 \pi_2$$

$$(20) \sigma^2(r, r_0) \sim \frac{\sigma^2}{\pi_1 \pi_2} \left[1 + \frac{2}{d_2} \int_0^\infty P_{0, u} du \right]$$

$$(21) \sigma^2(r, r_1) \sim \frac{\sigma^2}{\pi_1 \pi_2} \left[1 + \frac{2}{d_2} \int_0^\infty P_{0, v} dv + \frac{2}{d_1} \int_0^\infty P_{u, 0} du \right]$$

$$(22) \sigma^2(st, r_0) \sim \frac{\sigma^2}{\pi_1 \pi_2} \left[1 - \frac{1}{d_1^2 d_2} \int_{-\infty}^{d_1} \int_{-d_1}^{d_1} (d_1 - |u_1|) P_{u, v} du dv + \frac{2}{d_2} \int_0^\infty P_{0, v} dv \right]$$

$$(23) \sigma^2(st_0, st_0) \sim \frac{\sigma^2}{\pi_1 \pi_2} \left[1 - \frac{1}{d_1^2 d_2^2} \int_{-d_2}^{d_2} \int_{-d_2}^{d_2} (d_1 - |u_1|) (d_2 - |v_1|) P_{u, v} du dv \right]$$

$$(24) \left\{ \begin{aligned} \sigma^2(st, st) &\sim \frac{\sigma^2}{\pi_1 \pi_2} \left[1 - \frac{1}{d_1^2 d_2} \int_{-\infty}^{d_1} \int_{-d_1}^{d_1} (d_1 - |u_1|) P_{u, v} du dv \right. \\ &\quad + \frac{1}{d_1 d_2} \int_{-d_2}^{d_2} \int_{-\infty}^{\infty} (d_2 - |v_1|) P_{u, v} du dv + \frac{1}{d_1^2 d_2^2} \int_{-d_2}^{d_2} \int_{-d_2}^{d_2} (d_1 - |u_1|) \\ &\quad \cdot (d_2 - |v_1|) P_{u, v} du dv + \frac{2}{d_1} \int_0^\infty P_{u, 0} du - \frac{2}{d_1^2} \int_0^{d_1} (d_1 - u) P_{u, 0} du \\ &\quad \left. + \frac{2}{d_2} \int_0^\infty P_{0, v} dv - \frac{2}{d_2^2} \int_0^{d_2} (d_2 - v) P_{0, v} dv \right] \end{aligned} \right.$$

$$(25) \left\{ \begin{aligned} \sigma^2(Sy_1, t_0) &\sim \frac{\sigma^2}{n_1 n_2} \left[1 - \frac{1}{d_1 d_2} \sum_{-\infty}^{\infty} \sum_{-\infty}^{\infty} \rho_{uv} S_u S_v \right. \\ &\left. + \frac{1}{d_2} \sum_{u=-\infty}^{\infty} \int_{-\infty}^{\infty} \rho_{d_1, u, v} S_v - \frac{2}{d_2} \int_{-\infty}^{\infty} \rho_{0v} S_v \right] \end{aligned} \right.$$

$$(26) \sigma^2(Sy_0, Sy_1) \sim \frac{\sigma^2}{n_1 n_2} \left[\sum_{u=-\infty}^{\infty} \sum_{v=-\infty}^{\infty} \rho_{d_1, u, d_2 v} - \frac{1}{d_1 d_2} \sum_{-\infty}^{\infty} \sum_{-\infty}^{\infty} \rho_{uv} S_u S_v \right]$$

$$(27) \left\{ \begin{aligned} \sigma^2(Sy_0, Sy_0) &\sim \frac{\sigma^2}{n_1 n_2} \left[1 - \frac{1}{d_1^2 d_2} \sum_{-\infty}^{\infty} \sum_{-d_1}^{d_1} (d_1 - |u|) \rho_{uv} S_u S_v \right. \\ &- \frac{1}{d_1 d_2} \sum_{-d_2}^{d_2} \int_{-\infty}^{\infty} (d_2 - |v|) \rho_{uv} S_u S_v \\ &+ \frac{1}{d_1^2 d_2^2} \sum_{-d_2}^{d_2} \int_{-d_1}^{d_1} (d_1 - |u|) (d_2 - |v|) \rho_{uv} S_u S_v \\ &+ \frac{1}{d_2^2} \sum_{u=-\infty}^{\infty} \int_{-d_2}^{d_2} (d_2 - |v|) \rho_{d_1, u, v} S_v - \frac{1}{d_2^2} \int_{-d_2}^{d_2} (d_2 - |v|) \\ &\rho_{0v} S_v \\ &\left. + \frac{1}{d_1^2} \sum_{v=-\infty}^{\infty} \int_{-d_1}^{d_1} (d_1 - |u|) \rho_{u, d_2 v} S_u - \frac{1}{d_2} \int_{-d_2}^{d_2} (d_1 - |u|) \rho_{u0} S_u \right] \end{aligned} \right.$$

3. $\rho_{uv} = \rho_u \rho_v$ なる特別の場合 (particular case where $\rho_{uv} = \rho_u \rho_v$)

$\rho_{uv} = \rho_u \rho_v$ なる。これらの式の大部分は非常に簡単となることに注意せよ。

もし仮に Sy_u および St_v に対して同様に式で

$$Sy_u = 1 - \frac{2}{d_1} \int_0^{\infty} \rho_u S_u + 2 \sum_{u=1}^{\infty} \rho_{d_1, u}$$

$$St_v = 1 - \frac{2}{d_2} \int_0^{d_1} (d_1 - u) \rho_u S_u,$$

と置き、同様に

(5) これが正しい自己相関函数 (auto correlation function) であるための必要十分条件は ρ_u および ρ_v がともに自己相関函数となることである。

$$f_1 = \frac{2}{d_2} \int_0^{\infty} \rho_u \rho_v, \quad f_1' = \frac{2}{d_2^2} \int_0^{d_2} (d_2 - v) \rho_u \rho_v, \quad f_1'' = 2 \sum_{u=1}^{\infty} \rho_{d_2, u}$$

$$f_2 = \frac{2}{d_1} \int_0^{\infty} \rho_u \rho_v, \quad f_2' = \frac{2}{d_1^2} \int_0^{d_1} (d_1 - u) \rho_u \rho_v, \quad f_2'' = 2 \sum_{u=1}^{\infty} \rho_{d_1, u}$$

と置くことにすれば、例えは

$$(28) \sigma^2(Y, Y_0) \sim \frac{\sigma^2}{n_1 n_2} (1 + f_1)$$

$$(29) \sigma^2(t, t_0) \sim \frac{\sigma^2}{n_1 n_2} (1 + f_1 + f_2)$$

$$(30) \sigma^2(St_0, St_0) \sim \frac{\sigma^2}{n_1 n_2} (St_u St_v + St_u t + St_v)$$

$$(31) \sigma^2(St_1, St_1) \sim \frac{\sigma^2}{n_1 n_2} (St_u St_v + f_1 St_u + f_2 St_v)$$

$$(32) \sigma^2(Sy_1, Sy_1) \sim \frac{\sigma^2}{n_1 n_2} (Sy_u Sy_v + f_1 Sy_u + f_2 Sy_v)$$

$$(33) \sigma^2(Sy_0, Sy_0) \sim \frac{\sigma^2}{n_1 n_2} (St_u St_v + f_1' Sy_u + f_2' Sy_v)$$

これらの式から

$$(34) \left\{ \begin{aligned} \sigma^2(St_1, St_1) - \sigma^2(Sy_1, Sy_1) &\sim \frac{\sigma^2}{n_1 n_2} \\ &\cdot [(St_u St_v - Sy_u Sy_v) + f_1 (St_u - Sy_u) + f_2 (St_v - Sy_v)] \end{aligned} \right.$$

$$(35) \left\{ \begin{aligned} \sigma^2(Sy_1, Sy_2) - \sigma^2(St_0, St_0) &\sim \frac{\sigma^2}{n_1 n_2} \\ &\cdot \left[\{ (1 - S_{1, u}) (1 - S_{2, v}) - (1 - S_{1, u}) (1 - St_v) \} + f_1' Sy_u + f_2' Sy_v \right] \end{aligned} \right.$$

$$(36) \sigma^2(St_0, St_0) - \sigma^2(Sy_0, Sy_0) \sim \frac{\sigma^2}{n_1 n_2} [f_1' (St_u - Sy_u) + f_2' (St_v - Sy_v)]$$

が得られる。

(34), (35), (36) 式から、一次元の結果を用いて、二次元の

標本の分散の比較が可能となる。

実際の場合の大部分では、 σ は正であり $St_u \geq Sy_u, St_u \geq Sy_v$ だから

$$(37) \sigma^2(St_1, St_1) \geq \sigma^2(Sy_1, Sy_1) \geq \sigma^2(St_0, St_0) \geq \sigma^2(Sy_0, Sy_0)$$

となる。 $\rho_{d1u} = \rho_1^{1u}$ および $\rho_{d2v} = \rho_2^{1v}$ に対する $\sigma^2(St_0, St_0)/\sigma^2(t_0, t_0), \sigma^2(Sy_1, Sy_1)/\sigma^2(t_0, t_0), \sigma^2(Sy_0, Sy_0)/\sigma^2(t_0, t_0)$

の値を次の表に与える。与えられた数の標本に対して (d_1, d_2 一定) は、 $\rho_1 = \rho_2$ のとき $\sigma^2(St_0, St_0), \sigma^2(Sy_1, Sy_1)$ および $\sigma^2(Sy_0, Sy_0)$ が最小となることは容易にわかる。

表示された式は $\rho_1 = \rho_2 = 0$ に対しては値 1 をとり、 ρ_1 および ρ_2 が 1 に近づくに従って、それぞれ極限值 0, 2/3, 0 および 2 に近づく。

0.4 以上異なる ρ_1 と ρ_2 に対しては、 Sy_1, Sy_0 によつて組まれた格子 (grid) は単純無作為抽出法より効率が高いことは注目すべきことである。

しかし $\rho_{uv} = \rho_u \rho_v$ なる形の函数は *centrally-symmetric function* よりも現実において現れることか少ないようである。これは座標軸の並び方と独立である。この理由で、我々は幾何の型の函数を次に考察する。

9. Centrally-symmetric な相関函数

Leдебонт と Wehrte (3) は $\rho(u, v)$ が相関函数 (Correlation function) であるための必要十分条件を

$$(38) \rho(u, v) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cos(wu - \mu v) \delta F(w, \mu)$$

あるいは逆に

$$(39) \delta(w, \mu) = \frac{1}{(2\pi)^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cos(wu - \mu v) \rho(u, v) \delta u \delta v$$

のように与えている。

(6) 都市調査を行なう場合には、二点間の相関が street 内および street 間の相関からなるから、この函数が現実的となる。

表 3 異なる ρ_1 および ρ_2 の値に対する系統抽出と無作為抽出の効率の比較

$\rho_2 \backslash \rho_1$	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
0	1.000 1.000 1.000 1.000	1.000 1.222 1.000 1.000	1.000 1.500 1.000 1.000	1.000 1.857 1.000 1.000	1.000 2.222 1.000 1.000	1.000 3.000 1.000 1.000	1.000 4.000 1.000 1.000	1.000 5.667 1.000 1.000	1.000 9.000 1.000 1.000	1.000 19.000 1.000 1.000	1.000 ∞ 1.000 1.000
0.1		0.720 0.739 0.596 1.21	0.669 0.754 0.534 1.25	0.632 0.827 0.493 1.28	0.601 0.956 0.462 1.30	0.575 1.160 0.437 1.31	0.551 1.488 0.416 1.32	0.529 2.055 0.398 1.33	0.508 3.215 0.382 1.33	0.489 6.734 0.367 1.33	0.471 ∞ 0.354 1.33
0.2			0.609 0.706 0.462 1.32	0.565 0.721 0.416 1.35	0.529 0.798 0.380 1.39	0.497 0.914 0.354 1.41	0.469 1.134 0.328 1.43	0.443 1.532 0.307 1.44	0.419 2.362 0.289 1.45	0.396 4.911 0.272 1.46	0.375 ∞ 0.257 1.46
0.3				0.516 0.689 0.365 1.41	0.476 0.797 0.327 1.45	0.441 0.793 0.297 1.49	0.409 0.924 0.271 1.51	0.380 1.209 0.249 1.53	0.354 1.825 0.227 1.54	0.328 3.751 0.212 1.55	0.305 ∞ 0.196 1.55
0.4					0.432 0.680 0.298 1.50	0.394 0.702 0.256 1.54	0.360 0.787 0.229 1.57	0.329 0.983 0.206 1.60	0.300 1.437 0.185 1.62	0.272 2.900 0.167 1.63	0.247 ∞ 0.151 1.64
0.5						0.354 0.675 0.223 1.59	0.317 0.703 0.195 1.63	0.284 0.821 0.171 1.66	0.253 1.139 0.150 1.68	0.223 2.228 0.132 1.70	0.196 ∞ 0.115 1.71
0.6							0.279 0.671 0.167 1.67	0.243 0.712 0.142 1.71	0.210 0.908 0.121 1.74	0.180 1.699 0.102 1.76	0.151 ∞ 0.085 1.78
0.7								0.206 0.667 0.118 1.75	0.172 0.742 0.096 1.79	0.139 1.226 0.071 1.82	0.109 ∞ 0.059 1.84
0.8									0.134 0.667 0.074 1.84	0.102 0.763 0.055 1.87	0.070 ∞ 0.037 1.89
0.9										0.067 0.667 0.035 1.92	0.034 ∞ 0.018 1.95

Centrally - symmetric 相関函数について

$u = r \cos \theta, v = r \sin \theta$ とおけば $P(u, v) = P(r)$ で

$$f(u, v) = \frac{1}{(2\pi)^2} \int_0^{2\pi} \int_0^{\infty} \cos(r \sqrt{u^2 + v^2} \cos \theta) P(r) r dr d\theta$$

ここで $\theta_1 = \theta + \tan^{-1}(v/u)$

$$= \frac{1}{2\pi} \int_0^{\infty} J_0(rT) P(r) r dr \quad \text{ここで } T = \sqrt{u^2 + v^2}$$

従つて $P(u, v)$ が centrally - symmetric なら、更に $f(u, v)$ もそうなるから、

$$(40) f(T) = \frac{1}{2\pi} \int_0^{\infty} J_0(rT) P(r) r dr$$

および

$$(41) P(r) = 2\pi \int_0^{\infty} J_0(rT) f(T) T dT$$

が互に立つ。

故に $P(r)$ および $f(T)$ の適当な形を求めることができる。

このことから、式

$$\int_0^{\infty} J_0(yz) e^{-ay} dy = 1/(a^2 + z^2)^{1/2}, \quad a \gg 0$$

が有用である。何となれば、我々が式を導く場合、収斂性の問題とともに $P(r)$ の確率的性質によつて、制約を受けるのであるが、

$$\frac{1}{2\pi} \int_0^{\infty} J_0(yz) e^{-ay} dy \quad \text{と} \quad \frac{1}{2\pi} \int_0^{\infty} J_0(yz) dy$$

は $2\pi f(T)$ と $P(r)$ となる

る可能性をもつ函数である。故に例えは $a = \pi = 0$ とすればスペクトラル (spectral) および相関函数 (correlation function) として、 $1/2\pi T$ および $1/r$ が得られるがこれは収斂しない。

線型 (linear) な場合マルコフ過程 (markoff process) $P(u, v)$ のスペクトラル函数 (spectral function) は、一次元の Cauchy 分布、 $f(t) = 1/\pi (a^2 + t^2)$ である。

二次元の Cauchy 分布をスペクトラル函数としてとれば、

(7) 同様にして普通の Cauchy 分布は距離 a の点を中心としてすべての方向に作られた直線の一つ上の密度分布 (density distribution) と考えられるから、二次元の分布は、距離 a から始まる平面の一つ上の密度分布として考えることができる。

$$f(T) = a/2\pi (a^2 + T^2)^{3/2}$$

および

$$P(r) = -\frac{1}{\pi a} (e^{-ar}/r) = e^{-ar}$$

が得られる。即ち一般化された Cauchy 分布は一般化されたマルコフ過程のスペクトラル函数であることがわかる。勿論我々は

$$(42) P(u, v) = \exp\left[-\left\{\frac{u^2}{a^2} - \frac{2uv}{ae} + \frac{v^2}{e^2}\right\}^{1/2}\right]$$

で与えられる。(8) “楕円的”マルコフ過程 (“elliptical” Markoff process) を考えることもできる。

しかし $m=0$ とすれば計算を簡略化することができるから、これからは α_1, α_2 の測定単位を変えて過程 $P(r) = e^{-\alpha r}$ を研究することにする。

この過程は農業の圃場研究で経験される相関函数とそう大きくかけ離れている様には思えない。Osborne [4] は $P(u) = e^{-\lambda|u|}$ の可能な形式を指摘している。Mahalanobis [5] は 800 区画の水田に対する相関を計算した。その値をオメガ表に函数 e^{-r} の値とともに示しておく。

結果及び北東方向の値をもつ elliptical process が確かにあてはまりの良いことはわかるが Mahalanobis の値は、それぞれ標準誤差が約 0.035 であることに注意すれば、あてはまりは十分良いことわかる。

(10) 系統的抽出と層化無作為抽出の相対的効率

(The efficiencies of systematic and stratified random sampling)

最後の節で概念的に展開した相関函数は (19) - (27) において

(18) これらの点から、 $P(r) = e^{-\alpha r}$ は circular マルコフ過程とよばれるが、一方 $P(u, v) = P_1(u) P_2(v)$ および $P(u, v) = \exp\left\{-\left|\frac{u}{a} + \frac{v}{e}\right|\right\}$ は二次元 およびこの退化した (degenerate) マルコフ過程として知られている。

(19) このことは、広い範囲の値にわたつて Fairfield-Smith の法則と準一致一致する法則が数値的に得られることによつて支持される。

表4 観測された系列相関と centrally-symmetric 相関函数から求めた理論値の比較

距離 マイル	行		列		北-東		南-東	
	観測値	計算値	観測値	計算値	観測値	計算値	観測値	計算値
1	0.332	0.368	0.310	0.368	-	-	-	-
2	-	-	-	-	0.264	0.243	0.264	0.243
2	0.149	0.135	0.090	0.135	-	-	-	-
2√2	-	-	-	-	0.050	0.059	0.129	0.059
3	0.009	0.050	0.029	0.050	-	-	-	-
3√2	-	-	-	-	0.050	0.018	0.070	0.018
4	0.034	0.018	0.041	0.018	-	-	-	-
4√2	-	-	-	-	0.020	0.004	0.060	0.004

て用いるべきであるが、しかしこれらの函数は簡単には積分できない。

$$(43) \left\{ \frac{\alpha^2 (S_{to} S_{to}) - \alpha^2 (S_{yo} S_{yo})}{\alpha^2 (h_o h_o)} \sim \frac{1}{d_1} \int_{-d_1}^{d_1} \left(1 - \frac{|u|}{d_1}\right) F(u, d_2) \delta v \right. \\ \left. + \frac{1}{d_2} \left(1 - \frac{|v|}{d_2}\right) F(v, d_1) \delta u \right.$$

$$\text{ここで、} F(u, d_2) = \frac{2}{d_2} \left\{ \int_0^{d_2} \frac{v}{d_2} \rho_{uv} \delta v + \int_{d_2}^{\infty} \rho_{uv} \delta v - d_2 \sum_{v=1}^{\infty} \rho_{d_2, v} \right\}$$

$$F(v, d_1) = \frac{2}{d_1} \left\{ \int_0^{d_1} \frac{u}{d_1} \rho_{uv} \delta u + \int_{d_1}^{\infty} \rho_{uv} \delta u - d_1 \sum_{u=1}^{\infty} \rho_{d_1, u, v} \right\}$$

なることに注意すればもう一つの研究方法が可能である。

$F(u, d_2)$ および $F(v, d_1)$ はオズワルドの $(\alpha_{st}^2 - \alpha_{sy}^2) / \alpha_r^2$ について得られた式の拡張である。故にもし $F(u, d_2)$ および $F(v, d_1)$ がともに正の函数なら系統的抽出は層化無作為抽出よりも正確である。これの特別の場合 $\rho_{uv} = \rho_{1u} \rho_{2v}$ なるときに生ずる。しかし $\rho_{uv} = \exp\{-(u^2 + v^2)^{1/2}\}$ なるとき、 $F(u, d_2)$ は必ずしも正でない。何故なら u が増加するとき ρ_{uv} は v について凹 (Concave) となる。これによつて、 u が 0 から d_1 まで動くとき、 $F(u, d_2)$ は $+\infty$ からある未知の値 X まで動くことから、(43) の解釈が非常

に面倒になる。この値は $d_2 \gg d_1$ なるとき正で、 $d_1 \gg d_2$ なるときは抽出が二つの方向で比例しないから系統的抽出は層化無作為抽出より有効である。その上 $d_1 = d_2 = d_0$ で $d \rightarrow 0$ なら $F(u, d) \rightarrow \infty$ で、この場合も系統的抽出の効率の高いことは明らかである。この種の系統的抽出即ち $S_{yo} S_{yo}$ は非常に多くの場合無作為抽出より正確な結果を与える。

4. 抽出誤差の推定 (Estimation of sampling errors)

公式の (7) - (18) を見れば二次元抽出 (Plane sampling) においても一次元抽出の抽出誤差 (linear error) の推定に用いた原理が使えることがわかる。いま各標本が、独立な単位に分断され、その各々から個の層の何れか一つに位置させることかできるとすると、各個の反復に対しては誤差の自由度は $g-1$ となる。例えば、 $h_{oo}, h_{ot}, S_{to}, r_{oo}$ および S_{to}, r_{ot} の自由度はそれぞれ $g-1, g-1, g-1, g-1, g-1, g-1, g-1$ および $g-1, g-1$ となる。単一の標本は誤差の下限推定量を含む。しかし $S_{to}, S_{to}, S_{to}, S_{to}, S_{to}, S_{to}, S_{to}, S_{to}$ および $S_{yo}, S_{yo}, S_{yo}, S_{yo}, S_{yo}, S_{yo}, S_{yo}, S_{yo}$ の自由度は $g-1, g-1, g-1, g-1, g-1, g-1, g-1, g-1$ である。したがって、適切な誤差の推定量を求めると、反復が必要である。しかし我々は、標本を十分正確な誤差の推定量を与えるような数つかの部分に分割するという方法を用いることかできる。

更に系統的標本の組を明瞭にすることの可能性も考えることかできる。これは抽出誤差を推定するに際し層間隔で配置される。それに標本の各々の頂の間の相関が小さくない限り、各次の f を無視するとはつきりした偏りが生ずることがわかる。しかし Yates が指摘した通り、この方法によれば我々の抽出誤差の上限を求めることかできる。

これらの抽出法は例によつて下に示される。

例 (Example)

我々は系統的標本の抽出誤差を推定する三つの方法を考えることとする。

(1) 各々無作為に配置された系統的標本の組を用いる方法即ち、抽出対象を小地域 (sub-area) 即ちブロック (block) の系列に分割し、各ブロックにおいて幾つかの系統的標本を抽出する誤差の分散は各ブロックの系統的標本の分散から計算される。

(2) 無作為に配置された一組の系統的標本を用いる方法即ち、幾つかの系統的標本をとり、この地域を小地域即ちブロックに分割する。誤差分散は各ブロック内で系統的標本が占める割合の分散から計算される。

(3) 1つの系統的標本を用いる方法、即ち1つの系統的標本をとり、それをより広い間隔の幾つかの系統的標本に分割する。例えば、もとの間隔の4倍の標本を4つとると、この地域は幾つかのsub-areaに分割されるから、誤差の分散は各ブロック内の副次的系統的標本の割合の分散から計算される。これらの三つの方法は、平均値の推定は益々正確に、また抽出分散の推定は次第に偏りが大きくなり、また実際的な適用は順次容易となるから、幾々の抽出方法は母集団および結果の使用目的に従って変ってくる。例えば、抽出を次々に行えば、誤差の推定は改善されてくると思われるからはじめには大体の見当だけで良いかも知れない。

(2) もし幾々が系列を分割したときそれぞれの部分内の観測値の個数が非常に大きい様な連続な一次元母集団 (Continuous linear population) について抽出を行なうなら、(1)と(2)の方法はともに項 (term) 当り分散

$$\sigma^2 \left(1 - \frac{2}{d} \int_0^{\infty} p_u q_u + 2 \sum_{u=1}^{\infty} p_u du \right)$$

の正確な推定量を与えるであろう。

しかし(3)は正しい項当りの分散

$$\sigma^2 \left(1 - \frac{2g}{d} \int_0^{\infty} p_u q_u + 2 \sum_{u=1}^{\infty} p_u du / g \right)$$

の代りに σ^2 を推定する。

従つて(3)による抽出分散の推定量は、実際の分散が小さい場合でも、(1)、(2)の方法による推定量よりも一般に過大である。

(3) Kendall は480個からなる人為的な系列、

$u_{x+2} = 1/4 u_{x+1} - 0.5 u_x + E_{x+2}$ を作った。ここで E_x は-49から49までの矩形分布である。この系列では $\sigma^2 = 2379.81$ および $\rho^2 = 2535.11$ である。この系列を80項の6つの部分にわけた。これは何れも $n=5$, $k=16$, $g=4$ であるから誤差に対しては自由度18が利用できる。これのsamplingの結果を表5表に示す。

表5表 自己回帰模型に対する系統的標本

抽出誤差の三つの推定方法の比較

方法	自由度18にちとすく項当りの抽出分散の推定値 S^2	$E(S^2)$	項当りの真の抽出分散
(1)	3228	2170	2170
(2)	1872	2170	2167
(3)	3709	2577	423

この表の値は連続母集団の大標本に対する結論を保証するものである。

(3) 多数の一様性の実験 (uniformity trials) が行われ、システム (system) S_1, S_2, S_3, S_4 に対応して標本抽出が行われた。システム S_1, S_2 に対応する抽出について、各々二つの標本をとりて誤差を推定した。一方システム S_3, S_4 に対応する抽出では、(2) および (3) の方法による系列の各部分においてとられた4つの標本の組を比較して誤差を推定した。この抽出の結果を表6表に示す。一方試行の数が小さければ、結果に現れた筈の傾向は上に得られた結論と非常によく一致する。

14 母集団における trend (Trend in the population)

母集団から標本を抽出するとき或々は幾々 trend の問題に直面する。これは無作為抽出や層化無作為抽出による母集団平均値の推定の場合にはどう大きな影響を与えない。しかし系統的標本の初年よりこれによつて非常に大きな影響を受ける。いま一次元の抽出 (linear sampling) を考え、 X_i が最初の要素となっている際

本を S_i とかけば、標本 S_i の組は普通 n について単純であつて、 S_1 と S_n の差は大となる (大体 $X_i - X_n$ に等しい)
 Yates (1) はこの難点を克服する方法を提案した S_i を

$$\frac{1}{n-1} \left\{ \frac{i}{k} X_i + X_{i+k} + \dots + X_{i+(n-2)k} + \frac{k-i}{k} X_{i+(n-1)k} \right.$$

とせば、*trand* にもとづく系統的標本間の差は大部分除かれる。
 この方法を用いると情報を少し失ふことがすむ。特に短作命な連続的母集団 (Continuous random population) については、分散は σ^2/n でなく $(n - \frac{3}{2})\sigma^2 / (n-1)^2$ である。二次元の組 (Plane sampling) における対応する補正された標本は

$$S_{ij} = \frac{1}{(n_1-1)(n_2-1)} \left\{ \frac{i}{k_1 k_2} X_{i+k_1, j} + \dots + \frac{j(k_1-i)}{k_1 k_2} X_{i+(n_1-1)k_1, j} \right. \\
 + \frac{j}{k_1} X_{i, j+k_2} + X_{i+k_1, j+k_2} + \dots + \frac{(k_1-i)}{k_1} X_{i+(n_1-1)k_1, j+k_2} \\
 \left. + \frac{i(k_2-j)}{k_1 k_2} X_{i, j+(n_2-1)k_2} + \dots + \frac{(k_1-i)(k_2-j)}{k_1 k_2} X_{i+(n_1-1)k_1, j+(n_2-1)k_2} \right\}$$

これにも同様な損失がある。
 しかし *trand* がはつきりした形をとるのは大標本の場合である。
 この場合においては、端末補正 (end adjustment) による情報の損失は無視できるから上に得られた結論はそのまゝ成立する

表6 表三つの奇一性乗乗における種々の抽出法の初期標本の比較

出所 Cochran (1) 抽出法に依りて	Kalenkar (8)	Witte (9)	Wyman, Yates & Kazisima, Iyeg (10)
作功の種別	72	132	108
プロットの数	5	4	7
平均値	23.555	587.95	960
平均値の標準差	15.555	10.018	270.89
抽出形式	S_i, S_{ij}	S_i, S_{ij}	S_i, S_{ij}
抽出の推定法	$\frac{1}{16}$	$\frac{1}{9}$	$\frac{1}{18}$
抽出の数の例	1, 3, 2, 16, 2	1, 4, 3, 20, 6, 2	1, 4, 2, 15, 8, 2
平均値の分散推定値	23.140	586.54	275.29
分散推定値の自由度	9.763	5151.6	1320.15
	48	80	60
	23.435	598.65	266.72
	2.689	5772.7	799.29
	12	12	15
	23.323	7038.5	271.27
	4.889	12	1269.54
	12	12	15

★ もとの1500個のプロットから求めたもの

引 用 文 献

1. W. G. Cochran, "The relative accuracy of systematic and stratified random samples for a certain class of populations", *Annals of Math. Stat.*, Vol. 17 (1946) P. 164
2. F. Yates, "A review of recent statistical developments in sampling surveys", *Roy. Stat. Soc. Jour.*, Vol. 109 (1946) P. 12
3. G. Ledebant and P. Uchite "Mecanique alatoire", *Portugaliae physics*, Vol 1 (1945)
4. J. G. Osborne, "Sampling errors of systematic and random surveys of cover-type areas", *Am. Stat. Assoc. Jour.*, Vol. 37 (1942), P. 256.
5. P. C. Mahalanobis, "On large-scale sample surveys", *Roy. Soc. Phil. Trans.*, B. 231 (1944), P. 329
6. M. G. Kendall, "On the analysis of oscillatory time-series", *Roy. Stat. Soc. Jour.*, Vol. 108 (1945), P. 93
7. M. G. Kendall, *Contributions to the Study of Oscillatory Time-Series*. Nat. Inst. Econ. Soc. Res, 1946.
8. R. J. Kalamkar, "Experimental errors and the field-plot technique with potatoes", *Jour, Agr. Sci.*, 1932, P. 373.
9. G. A. Wiebe, "variation and correlation in grain-yield among 1500 wheat nursery plots", *Jour. Agr. Res.*, 1935.
10. Wynne Sayer and P. V. Krishna Iyer, "on some of the factors that influence the error of field experiment with special reference to sugarcane", *Ind. Jour. Agr. Sci.*, 1936. P 917.
11. W. G. Cochran, "Catalogue of uniformity trial data", *Roy. Stat. Soc. Suppl. Jour.*, Vol. 4 (1937), P. 233

12. F. Yates, "Systematic sampling" *Roy. Soc Phil. Trans.*, Vol. 241 (1948), P. 345.

By. M. H. Quenouille

A. M. S. Vol. 20 (1949) 51

ら、二次元の系統的抽出およびこれに関連する層化および無作為抽出について

Two dimensional systematic sampling and the associated stratified random sampling

緒言

今日系統的抽出法は理論および応用統計の両面において非常に興味をもたれている。この抽出方式については色々な修正方式とともによく知られている。Madour (1944), Cochran (1946) および Yates (1948) は、一次元の場合の幾つかの特別なタイプの系統的抽出を論じた。この論文では一次元の場合の Cochran の結果を二次元の場合に一般化することを考える。

抽出すべき母集団 (univers) がそれぞれ l 個の小区間 (cell) からなる $m \times l$ 個の行 (row) から構成されていて、全部で ml 個の層をなしていると仮定する*。この小区間は場合に応じて正方形でもよくまた矩形でもよい。座標を次のように定める。すなわち $y = Const$ なら行、また $x = Const$ なら列を数わす。このようにすると、 l 個の小区間の l 個の組である各行は x 方向で、 m 個の小区間の m 組からなる各列は y 方向にある。特定の層内において、 l 個の小区間中から ν 個を無作為に選ぶ。そうして選ばれたこれらの小区間が含まれる行から出発して、 l 番目ごとの小区間の行 (y 方向) を選ぶ。これらの選ばれた行の各々において無作為に抽出されたものと同じ列の小区間とこれから l 番目ごとの小区間 (x 方向に沿って) をとれば、大きさ $m \times \nu$ なる標本ができる。これは二次元の場合の系統的抽出の一種である。このような系統的抽出はまた次のようにしても定義できる。すなわち、特定の層内において l 個の中から ν 個の小区間を無作為に選ぶ。次に組内および層に対する配置の最初の層と同じになるようにして他の各層から同様な ν 個の小区間の組を選ぶ。我々はまた全部で大きさ $m \times \nu$ なる層化標本となるようにそれぞれの層内から大きさ ν の無作為標本を抽出することでも

とる。これとは別に全部で $m \times l$ 個の小区間の中から層化なしで $m \times \nu$ 個の区間を無作為に選ぶこともできる。これは単純無作為標本 (perfectly random sample) である。

我々はこのユニバースが $E(z_i) = \mu$ なる意味で齊一であり、また観測値は共通の分散と異なる空間的相関 (different space correlation) を有するものとする。ただし z は小区間の観測値である。 x および y 方向にそれぞれ u および v だけ離れた点の間の空間的相関係数を $\rho(u, v)$ とする。すなわち小区間 (x, y) を P とすればその収量と区間 $(x+u, y+v)$ の収量と組合せ、それらの可能なすべての組から $P(u, v)$ を計算するのである。

これはユニバースの有限性のためにある変動をもつかもしれない。もしこのユニバースが無限ならそのような影響は存在しない。しかしユニバースが有限でも、Cochran (1946) がしたようにそれを無限母集団から得られた標本と考えることができる。このとき無限母集団では、 $E(z) = \mu$, $V(z) = \sigma^2$ および

$E\{(z_{i+u, j+v} - \mu)(z_{ij} - \mu)\} = \rho(u, v)$ である。色々な u および v の値に対するこの $\rho(u, v)$ は三次元の図形となるから、これは *Convolopine* とよぶことができる。この図形は区間が無限小のときに連続となる。

この論文においては、無限母集団からとられたすべての有限母集団について平均したときの無作為、層化および系統的標本の総平均の分散の期待値を Cochran (1946) の用いた記号にしたかつてそれぞれ σ_n^2 , σ_{st}^2 , σ_{sy}^2 で表わす。

$\Delta_1, \Delta_2, \delta_1, \delta_2$ は

$\Delta_1 P(u, v) = P(u+1, v) - P(u, v)$

$\Delta_2 P(u, v) = P(u, v+1) - P(u, v)$

$S_1 P(u, v) = P(u+\frac{1}{2}, v) - P(u-\frac{1}{2}, v)$

$S_2^2 P(u, v) = P(u+1, v) + P(u-1, v) - 2P(u, v)$

$\Delta_1 \Delta_2 P(u, v) = \Delta_1 P(u, v+1) - \Delta_1 P(u, v)$

なる演算子 (operator) で S_2, S_2^2 および S_1^2, S_1^2 の定義も同様である。我々はまた $E, P(u, v) = P(u+1, v)$ および

$E_2 P(u, v) = P(u, v+1)$ なる E_1, E_2 を用いる。
 これらの E_i は "期待値" を表わす E , (添字のない) とは異なる。
 また $P(u, v) = P(u, -v)$ と $P(u, -v) = P(-u, v)$ は別のものであることに注意しなければならぬ。そうして $P(u, v) + P(-u, v)$ を $\psi(u, v)$ と書くことにする。これは u および v に関して対称である。 $\Delta_1 \psi$ および $\Delta_2 \psi$ は常に $\Delta_1 \psi(u, v)$ および $\Delta_2 \psi(u, v)$ を表わす。
 §1 においては $\sigma_x^2, \sigma_y^2, \sigma_{xy}^2$ を計算し、それらの相対的割合が γ と独立であり、したがって相対的効率 (relative efficiency) も γ と独立であることを示した。§2 では層化抽出が無作為抽出より効率が良くなるための幾つかの条件を多数の定理の形で述べ、また §3 では系統的抽出が層化抽出よりも有効となるための十分条件を定理4として与えた。したがって、直接この比較はしなかつたが、これらの結果を用いて系統的抽出が無作為抽出よりも有効となるための十分条件の組を見出すことができる。Cochran の結果の特別な場合として出てくる。しかしこれらは P が正または負の何れであつても、また γ より大きい γ に対しても成立するという点で Cochran (1946) の結果よりも多少改善されている。

§1

$$\sum_{j=1}^n (Z_j - \bar{Z})^2 = \frac{1}{n} \sum_{j=1}^n \sum_{k < j}^n (Z_j - Z_k)^2$$

が成立する。

Lemma I $\phi(i, j)$ 区割に対応する確率変数として Z_{ij} をもつような区割の N 個からなつて M 個の行を含むプロットにおいて

$$E \left\{ \sum_{i=1}^N \sum_{j=1}^M (Z_{ij} - \bar{Z})^2 \right\} = \frac{2\sigma^2}{MN} \left[\frac{MN(MN-1)}{2} - N \sum_{u=1}^M (M-u) P(0, u) - M \sum_{u=1}^{N-1} (N-u) P(u, 0) - \sum_{u=1}^{N-1} \sum_{v=1}^{M-1} (N-u)(M-v) \psi(u, v) \right]$$

ただし $i = 1, \dots, N; j = 1, \dots, M$.

証明

$$(1.1) \text{ によつて } \sum_{i=1}^N \sum_{j=1}^M (Z_{ij} - \bar{Z})^2 = \frac{1}{MN} \sum (Z_t - Z_t')^2$$

ここで和は、總言で述べたタイプの無限個のプロットからとられた標本であるところの有限個のプロット内の Z_t と Z_t' のすべての組合せの上にとらるものである。全部で MNC_2 個の組合せがあるが、これらは

- (i) 各行の中で NC_2 通りの組合せが全部で $M \cdot NC_2$ 通り
- (ii) 各列の中で MC_2 通りの組合せが全部で $N \cdot MC_2$ 通り
- (iii) (x, y) と $(x+u, y+v), (u \neq 0, v \neq 0)$ 間で $2^{MC_2} \cdot MC_2$ 通りのように分けられる。

いま

$$E \left\{ \sum_{t,t'} (Z_{t,t'} - \bar{Z})^2 \right\} = \frac{1}{MN} E \left\{ \sum (Z_t - u + u - Z_t')^2 \right\} = \frac{2\sigma^2}{MN} \sum \{ 1 - P(Z_t, Z_t') \} \dots (1.2)$$

ここで $P(Z_t, Z_t')$ は Z_t と Z_t' 間の相関係数で、可はもどきのすべての組合せ(すなわち全部で MNC_2 通り)にわたる。まず $\sum P(Z_t, Z_t')$ を求めてみよう。これは上に述べた3つの場合 (i) (ii) (iii) に対応する3つの式の和である。(i) の場合には、各行で $(u, 0)$ の距離のものは $(N-u)$ 区割ある。したがってすべての行についてのこのタイプの $\sum P(Z_t, Z_t')$ は $N \sum_{u=1}^{N-1} (N-u) P(u, 0)$ である。同様に (ii) の場合の $\sum P(Z_t, Z_t')$ は $N \sum_{v=1}^{M-1} (M-v) P(0, v)$ である。(iii) の場合は4点 $(x, y), (x+u, y), (x, y+v)$ および $(x+u, y+v)$ (ただし $u \neq 0, v \neq 0$) を考えれば、明らかにこれらの4点の中で (iii) の場合の異なる距離のタイプは2つのみ、すなわち (u, v) および $(-u, v)$ だけである。距離 $(-u, -v)$ と (u, v) は同じである。何となれば空間的相関はこの双方の距離に対して同一だからである。同様に距離 $(u, -v)$ は $(-u, v)$ と同等である。よつてこのような4点対は全部で $(N-u)(M-v)$ 組あるからこのタイプの全体について $\sum P(Z_t, Z_t')$ を加え上げれば、 $\sum_{u=1}^{N-1} \sum_{v=1}^{M-1} (N-u)(M-v) \psi(u, v)$ を得る。注意すべき点は

$$P(u, 0) = P(-u, 0) \text{ および } P(0, v) = P(0, -v)$$

である。

このときもすべての組合せ上での

$$\sum 1 \text{ は } MN C_2 \text{ である。}$$

よって(1.2)から Lemmaがえられる。

これは次のように書ける。

$$E \left\{ \sum_{i=1}^N \sum_{j=1}^M (z_{ij} - \bar{z})^2 \right\} = (MN-1)\sigma^2 \{1 - \phi(M, N)\} \dots \dots \dots (1.3)$$

ただし

$$\phi(M, N) = \frac{2}{M(MN-1)} \sum_{v=1}^{M-1} (M-v) P(0, v) + \frac{2}{N(MN-1)} \sum_{u=1}^{N-1} P(u, 0) + \frac{2}{MN(MN-1)} \sum_{u=1}^{N-1} \sum_{v=1}^{M-1} (N-u)(M-v) \psi(u, v) \dots \dots (1.4)$$

および

$$L(M, N) = 1 - \phi(M, N) \dots \dots \dots (1.5)$$

ゆえに(1.3)で $M = ml$ および $N = nk$ とおけば有限母集団の全体について

$$E \left\{ \sum (z - \bar{z})^2 \right\} = (mnlk-1)\sigma^2 \{1 - \phi(ml, nk)\} \dots \dots (1.6)$$

をうる。

同様に各層に対して(1.3) $M=l$ および $N=k$ とおけば

$$E \left\{ \sum (z - \bar{z})^2 \right\} = (lk-1)\sigma^2 \{1 - \phi(l, k)\} \dots \dots (1.7)$$

したがって mn 個の層に対する

$$E(\text{層内の}s.s.) = mn(lk-1)\sigma^2 \{1 - \phi(l, k)\} \dots \dots (1.8)$$

$Y=1$ なる各系統的標本については、(1.3) および (1.4) で

$M=m, N=n$ とおき、 $P(u, v)$ を $P(u_k, v_l)$ でまた ϕ を ϕ_{kl} とおきかえると

$$E \left\{ \sum (z - \bar{z})^2 \right\} = \sigma^2 \{1 - \phi_{kl}(m, n)\}$$

である。

ゆえに $Y=1$ なるとき

$$E(\text{系統的標本内}s.s.) = lk(mn-1)\sigma^2 \{1 - \phi_{kl}(m, n)\} \dots (1.9)$$

いま大きさ N の有限母集団からとられた大きさ n の無作為標本については

$$\text{標本平均値の分散} = \frac{1}{n} \cdot \frac{N-n}{N-1} \left\{ \frac{1}{N} \sum (z - \bar{z})^2 \right\}$$

である。ただし和は全有限母集団の上ごとの。

よってこの場合 $N = mnlk$ とおけば $n = mnv$ であるから、(1.6)を用いて

標本平均値の分散の期待値は

$$\sigma_{\bar{z}}^2 = \frac{1}{mnv} \frac{mn(lk-v)}{(mnlk-1)} \frac{1}{mnlk} (mnlk-1)\sigma^2 \{1 - \phi(ml, nk)\} = \frac{\sigma^2}{mnv} \left(1 - \frac{v}{lk}\right) \{1 - \phi(ml, nk)\} \dots (1.10)$$

となる。

更にそれぞれ lk 個の小区画からなる層が mn 個あるから $X_{11}, X_{12}, \dots, X_{1v}$ をその層からとった大きさ v なる標本とする。ここで X_{ij} は確率変数で、縦書きで述べた座標を表わすところの添字を付けない X と混同してはならない。

そうすると

$$\bar{z} = \frac{1}{mnv} \sum_{i=1}^{mn} \sum_{j=1}^v X_{ij} = \frac{1}{mn} \sum_{i=1}^{mn} \bar{X}_i, \text{ ただし } \bar{X}_i = \frac{1}{v} \sum_{j=1}^v X_{ij}$$
$$V(\bar{z}) = \frac{1}{m^2 n^2} \sum_{i=1}^{mn} V(\bar{X}_i)$$

したがって(1.7)から

$$E \left\{ V(\bar{X}_i) \right\} = \frac{1}{v} \frac{(lk-v)}{(lk-1)} \frac{(lk-1)}{lk} \sigma^2 \{1 - \phi(l, k)\}$$

よって

$$\sigma_{\bar{z}}^2 = E \left\{ V(\bar{z}) \right\} = \frac{\sigma^2}{mnv} \left(1 - \frac{v}{lk}\right) \{1 - \phi(l, k)\} \dots \dots (1.11)$$

いま系統的標本の分散の期待値 $\sigma_{\bar{z}}^2$ を求めるため $X_{11}, X_{12}, \dots, X_{1v}$ をその層からとった標本とする。標本総平均は $\bar{X}_{1j} = \frac{1}{mnv} \sum_{i=1}^{mn} \sum_{j=1}^v X_{ij}$

で、 $\sigma_{sy}^2 = E(\bar{x}_{sy} - \bar{x})^2$ である。ただし \bar{x} は有限母集団の総平均である。

ここで $\sum_{j=1}^{m\pi} x_{ij} / m\pi$ を ${}_j\bar{x}_{sy}$ と書くことにすると

$$(\bar{x}_{sy} - \bar{x}) = \frac{1}{\nu} \sum_{j=1}^{\nu} ({}_j\bar{x}_{sy} - \bar{x})$$

よって

$$(\bar{x}_{sy} - \bar{x})^2 = \frac{1}{\nu^2} \left\{ \sum_{j=1}^{\nu} ({}_j\bar{x}_{sy} - \bar{x})^2 + \sum_{\substack{j,j'=1 \\ (j \neq j')}}^{\nu} ({}_j\bar{x}_{sy} - \bar{x})({}_{j'}\bar{x}_{sy} - \bar{x}) \right\} \dots \dots \dots (1.12)$$

lk (ν 個)の系統的標本全体について (1.12) を平均する。

もし x_{ij} が第 i 層かうとうれた無作為標本 (観測可能な) なら、 ${}_j\bar{x}_{sy}$ は対応する系統的標本の平均値である。ゆえに lk 個の系統的標本全体での (1.12) の和において、 ${}_j\bar{x}_{sy}$ ($j=1, 2, \dots, m\pi$) は $lk-1$ ($\nu-1$ 回) 現れる。というのは大きさ ν の標本をとるためには残りの $(lk-1)$ 個の中から $(\nu-1)$ 個をけとり出せばよいからである。同様に、 ${}_j\bar{x}_{sy}$ および ${}_{j'}\bar{x}_{sy}$ (j, j') は同時に $lk-2$ ($\nu-2$ 回) 現れる。よって lk (ν 個) の標本全体について加え上げると $({}_j\bar{x}_{sy} - \bar{x})$ ($j=1, 2, \dots, lk$) は $lk-1$ ($\nu-1$ 回) 現れる。また $({}_j\bar{x}_{sy} - \bar{x})({}_{j'}\bar{x}_{sy} - \bar{x})$ (j, j') は $lk-2$ ($\nu-2$ 回) 現れてくる。しかるに $\sum_{t=1}^{lk} ({}_t\bar{x}_{sy} - \bar{x}) = 0$

$$\begin{aligned} \text{よって } \sum_{t=1}^{lk} ({}_t\bar{x}_{sy} - \bar{x})({}_t\bar{x}_{sy} - \bar{x}) &= 0 \quad \text{かつ} \\ \sum_{j=1}^{lk} \sum_{t=1}^{lk} ({}_j\bar{x}_{sy} - \bar{x})({}_t\bar{x}_{sy} - \bar{x}) &= - \sum_{j=1}^{lk} ({}_j\bar{x}_{sy} - \bar{x})^2 \end{aligned}$$

故に和においては $({}_j\bar{x}_{sy} - \bar{x})^2$ は $(lk-1)C_{\nu-1} - (lk-2)C_{\nu-2} = (lk-2)C_{\nu-1}$ 回でてくる。したがって求める (1.12) の平均は

$$\begin{aligned} \frac{(lk-2)}{(\nu-1)} \frac{1}{\nu^2} \sum_{j=1}^{lk} ({}_j\bar{x}_{sy} - \bar{x})^2 &= \frac{1}{m\pi\nu} \frac{(lk-\nu)}{lk(lk-1)} \sum_{j=1}^{lk} m\pi ({}_j\bar{x}_{sy} - \bar{x})^2 \\ &= \frac{1}{m\pi\nu} \frac{(lk-\nu)}{lk(lk-1)} \left\{ \sum_{i=1}^{m\pi} \sum_{j=1}^{lk} (x_{ij} - \bar{x})^2 \right\} \quad (\nu=1 \text{ 区3区2の系統的標本内1区}) \quad (1.13) \end{aligned}$$

ゆえに期待値をとると

$$\begin{aligned} \sigma_{sy}^2 &= \frac{\sigma^2}{m\pi\nu} \frac{(lk-\nu)}{lk(lk-1)} \left\{ (m\pi lk - 1) - (m\pi lk - 1) \phi(m\pi, lk) \right. \\ &\quad \left. - lk(m\pi - 1) + lk(m\pi - 1) \phi_{sy}(m, \pi) \right\} \end{aligned}$$

(1.6) および (1.9) から

$$= \frac{\sigma^2}{m\pi\nu} \frac{\nu}{(lk-\nu)} \left\{ 1 - \frac{(m\pi lk - 1)}{(lk-1)} \phi(m\pi, lk) + \frac{lk(m\pi - 1)}{(lk-1)} \phi_{sy}(m, \pi) \right\} \dots \dots \dots (1.14)$$

また (1.10) と (1.11) から

$$\sigma_r^2 = \frac{\sigma^2}{m\pi\nu} \left(1 - \frac{\nu}{lk}\right) \left\{ 1 - \phi(m\pi, lk) \right\} \dots \dots \dots (1.15)$$

$$\sigma_{st}^2 = \frac{\sigma^2}{m\pi\nu} \left(1 - \frac{\nu}{lk}\right) \left\{ 1 - \phi(lk) \right\} \dots \dots \dots (1.16)$$

ただし $\phi(M, N)$ は (1.4) で与えたものでありまた

$$\begin{aligned} \phi_{sy}(m, \pi) &= \frac{2}{m(m\pi - 1)} \sum_{j=1}^{m-1} (m-j) \rho(0, \nu l) + \frac{2}{\pi(m\pi - 1)} \sum_{u=1}^{\pi-1} \rho(u l, 0) \\ &\quad + \frac{2}{m\pi(m\pi - 1)} \sum_{u=1}^{\pi-1} \sum_{v=1}^{m-1} (m-u)(m-v) \psi(u l, \nu l) \end{aligned}$$

である

ゆえに相対的効率

$$\frac{1}{\sigma_r^2} : \frac{1}{\sigma_{st}^2} : \frac{1}{\sigma_{sy}^2} \text{ は } Y \text{ と独立で層化および correlated}$$

の性質によつてきまる。もし後層が既知なら前者はその目的によつて調整することができる。

また

ここで層化抽出が無作為抽出よりも層化となる、すなわち $\sigma_r^2 \gg \sigma_{st}^2$ なるための十分条件の組を求め、まず重要な 2.3 の Lemma を証明しておく。

Lemma II 条件

- (i) $\Delta_1 \psi(u, v) \leq 0$
- (ii) $\Delta_2 \psi(u, v) \leq 0$
- (iii) $2\psi(u+1, v) \gg \psi(u, v+1) + \psi(u+1, v+1)$

が満足される時、 l および k のすべての値に対して

$$L(l+1, k) \geq L(l, k)$$

である。

証明

(1.4)および(1.5)の M を l で、 N を k でおきかえれば $L(l, k)$ がえられる。

そこで $l' = l+1, A = lk-1, A' = l'k-1$ および

$$\beta_{0v} = 2 \left(\frac{l-v}{2A} - \frac{l'v}{l'A'} \right), \quad d_{u0} = \frac{k-u}{AA'}, \quad d_{uv} = \frac{(k-u)}{k} 2 \left(\frac{l-v}{2A} - \frac{l'v}{l'A'} \right) \quad (2.1)$$

ただし $u \neq 0$ および $v \neq 0$, とおくと

$$L(l+1, k) - L(l, k) = \sum_{v=1}^k \beta_{0v} P(0, v) + \sum_{u=1}^k \sum_{v=0}^{l-1} d_{uv} \psi(u, v) \dots (2.2)$$

かえられる。

いま $\psi(u, v) = -\Delta_2 \psi(u, v) + \psi(u, v+1)$ であるから

$$\sum_{v=0}^l d_{uv} \psi(u, v) = \sum_{v=0}^{l-1} S_{uv} \Delta_2 \psi(u, v) + S_{ul} \psi(u, l) \dots (2.3)$$

$$\text{ただし } S_{uv} = \sum_{j=0}^v d_{uj}$$

$$\text{また } \sum_{v=1}^k \beta_{0v} P(0, v) = -\sum_{v=1}^{l-1} S'_{0v} P(0, v) + S'_{0l} P(0, l) \dots (2.4)$$

ただし

$$S'_{0v} = \sum_{j=1}^v \beta_{0j}$$

よって

$$S'_{0v} = v \left\{ \gamma \left(\frac{1}{2A} - \frac{1}{l'A'} \right) + \frac{(k+1)}{l'AA'} \right\}$$

ただし $\gamma = l-1-v$

ゆえに $\gamma \geq 0$, すなわち $v \leq (l-1)$ のとき

$$S'_{0v} \geq 0 \dots (2.5)$$

および

$$S'_{0l} = -\frac{(k-1)}{AA'} \dots (2.6)$$

いま(2.5)から $v \leq (l-1)$ になるとき

$$S_{uv} = d_{u0} + \sum_{j=1}^v d_{uj} = \frac{k-u}{AA'} + \frac{k-u}{k} S'_{0v}$$

は正であるから(2.6)を用いて

$$S_{ul} = \frac{(k-u)}{kAA'} \dots (2.7)$$

$$L(l+1, k) - L(l, k) = -\sum_{u=1}^k \sum_{v=0}^{l-1} \frac{(k-u)}{AA'} \Delta_2 \psi(u, v)$$

$$- \sum_{u=1}^k \sum_{v=0}^{l-1} \frac{(k-u)}{k} S'_{0v} \Delta_2 \psi(u, v) - \sum_{v=1}^{l-1} S'_{0v} \Delta_2 P(0, v) + \sum_{u=1}^k \frac{(k-u)}{kAA'} \psi(u, l) - \frac{(k-1)}{2AA'} \psi(0, l) \dots (2.8)$$

いま

$$\sum_{u=1}^k \frac{(k-u)}{kAA'} \psi(u, l) - \frac{(k-1)}{2AA'} \psi(0, l) = \sum_{u=1}^k \frac{(k-u)}{kAA'} \{ \psi(u, l) - \psi(0, l) \} \\ = \sum_{u=1}^k \frac{(k-u)}{kAA'} \sum_{i=0}^{u-1} \Delta_1 \psi(i, l) \dots (2.9)$$

$$= \sum_{u=0}^k \frac{(l-u)(k-u)}{2kAA'} \Delta_1 \psi(u-1, l) + \sum_{u=1}^k \frac{(k-u)}{2AA'} \Delta_1 \psi(u-1, l) \dots (2.10)$$

$\Delta_1 \psi \leq 0$ になるとき(2.10)の最初のところからは0に等しいかまたは0より大である。しかし他の部分は正にならないから $\Delta_2 \psi \leq 0$ および $\Delta_1 \psi \leq 0$ のとき(2.8)から

$$L(l+1, k) - L(l, k) = Q + \sum_{u=1}^k \frac{(k-u)}{2AA'} \Delta_1 \psi(u-1, l) \\ - \sum_{u=1}^k \sum_{v=0}^{l-1} \frac{(k-u)}{AA'} \Delta_2 \psi(u, v) \dots (2.11)$$

よって $Q \geq 0$ である。

次に

$$\sum_{v=0}^{l-1} \Delta_2 \psi(u, v) = \psi(u, l) - \psi(u, 0)$$

を代入する。

よって(2.11)から

$$L(l+1, k) - L(l, k) = Q + \sum_{u=1}^k \frac{(k-1)}{2AA'} \{ 2\psi(u, 0) - \psi(u, l) - \psi(u-1, l) \} \dots (2.12)$$

いま $\Delta_2 \psi \leq 0$ になるとき

$$2\psi(u,0) - \psi(u,l) - \psi(u-1,l) \geq 2\psi(u,l-1) - \psi(u,l) - \psi(u-1,l)$$

ゆえにすべての u の値に対して

$$2\psi(u+1,v) \geq \psi(u,v+1) + \psi(u+1,v+1)$$

$2\psi(u,l-1) - \psi(u-1,l) - \psi(u,l) \geq 0$. すなわち (2.12) およびこれから (2.11) は 0 または正である。

よって Lemma は証明された。

同様にして次の Lemma を証明することかできる。

Lemma (III) 条件

(i) $\Delta_1 \psi(u,v) \leq 0$

(ii) $\Delta_2 \psi(u,v) \leq 0$

(iii) $2\psi(u+1,v) \geq \psi(u+1,v)$

が満たされるとき、 l および l のすべての値について

$$L(l, l+1) \geq L(l, l)$$

である。

定理 1

(i) $\Delta_1 \psi \leq 0$

(ii) $\Delta_2 \psi \leq 0$

(iii) $2\psi(u,v+1) \geq \psi(u+1,v) + \psi(u+1,v+1)$

(iv) $2\psi(u+1,v) \geq \psi(u,v+1)$

が満たされるすべての無限母集団においては、

任意の大きさの標本について

$$\sigma_{st}^2 \leq \sigma_r^2$$

が成立する。また上の4つの場合のそれぞれについて等号が成立

しなければ

$$\sigma_{st}^2 < \sigma_r^2$$

である。

証明

これらの3つの条件は Lemma I のすべての条件を満足している。

よって

$$L(l, l) \leq L(l+1, l) \leq \dots \leq L(m, l)$$

すなわち

$$1 - \phi(l, l) \leq 1 - \phi(m, l)$$

この場合の l は勿論層内の全区間数であるから

$$\sigma_{st}^2 \leq \sigma_r^2$$

で定理がえられる。

いまもし $\Delta_1 \psi \geq \Delta_2 \psi$ なら (iii) が満足される。ゆえに更に強い条件の値を次のように与えることかできる。

系 I u -方向 (別えは) に対つた平行なストリップで層化がおこなわれているすべての無限母集団で、もし $\Delta_2 \psi \leq \Delta_1 \psi \leq 0$ なら

$$\sigma_r^2 \leq \sigma_{st}^2$$

で、またこれらの場合の各々において等号が成立しない限り

$$\sigma_r^2 > \sigma_{st}^2$$

である。

系 II 上記において $l = 1$ すなわち v -方向に沿つた一次元の面場でも $\Delta_2 \psi(u,v) \leq 0$ または $\Delta_2 \psi(0,v) \leq 0$ なら、

$u=0$ であるから

$$\sigma_r^2 > \sigma_{st}^2$$

である。この場合は Cochran (1946) が $\sqrt{v} = 1$ について論じている。これらの定理はすべて l, m, n, l の値と無関係な条件

を与えており、したがつてどのような層化方式が有効であるかということには考えていない。しかし著者 (1949) はパラメータ

の関係に依存する他の条件の組を与えた。ここではそれについて論ずることにする。この場合にもある種の層化の任意性

(some degree of freedom in stratification) を得た。

これは次の定理のように述べられる。

定理 3 $\sqrt{l} \leq l \leq l-1$ および $(n-1)l = (m-1)l$ なる型のすべての層化方法に対して、もし

(i) $\Delta_1 \psi \leq 0$

(ii) $\Delta_2 \psi \leq 0$

(iii) $\Delta_1 \Delta_2 \psi \geq 0$

なら

$$\sigma_{st}^2 \leq \sigma_r^2$$

で、これらの3の場合の各々において等号が成立しない限り

$$\sigma_{st}^2 < \sigma_r^2$$

である。

このパラメトリックな関係は勿論 $k=l, m=n$ なら成立する。

この定理は次の Lemma N を用いて説明する。

Lemma N 条件

(i) $u \leq k, v \leq l$ のとき

$$\sum_{i=0}^u \sum_{j=0}^v d_{ij} \geq 0 \text{ ただし } \sum_{j=0}^k \sum_{i=0}^l d_{ij} = 0$$

(ii) $\Delta_2 \psi(k, v) \leq 0$

(iii) $\Delta_1 \psi(u, l) \leq 0$

(iv) $\Delta_1 \Delta_2 \psi(u, v) \geq 0$

が満たされるとき

$$\sum_{i=0}^k \sum_{j=0}^l d_{ij} \psi(i, j) \geq 0$$

である。

証明

$$\begin{aligned} \sum_{i=0}^k \sum_{j=0}^l d_{ij} \psi(i, j) &= - \sum_{j=0}^l \sum_{u=0}^{k-1} S_{u,j} \Delta_1 \psi(u, j) + \sum_{j=0}^l S_{k,j} \psi(k, j), \text{ かつ } S_{k,j} = \sum_{i=0}^k d_{ij} \\ &= \sum_{u=0}^{k-1} \sum_{v=0}^{l-1} d_{uv} \psi(u, v) - \sum_{u=0}^{k-1} T_{u,l} \Delta_1 \psi(u, l) - \sum_{v=0}^{l-1} T_{k,v} \Delta_2 \psi(k, v) + T_{k,l} \psi(k, l) \end{aligned}$$

$$\text{ここで } T_{uv} = \sum_{j=0}^v S_{uj} = \sum_{i=0}^u \sum_{j=0}^v d_{ij}$$

これで Lemma は証明された。

ここで定理3を証明しよう。

$$\begin{aligned} L(l+1, k+1) - L(l, k) &= \sum_{v=1}^k d_{0v} \psi(0, v) + \sum_{u=1}^k d_{u0} \psi(u, 0) + \sum_{u=1}^k \sum_{v=1}^l d_{uv} \psi(u, v) \\ &= \sum_{i=0}^k \sum_{j=0}^l d_{ij} \psi(i, j) \end{aligned}$$

$$k' = k+1, l' = l+1, A = lk-1 \text{ および } A' = l'k'-1$$

$$d_{00} = 0, d_{0j} = \frac{d_j}{lA} - \frac{l'j}{l'A'}, \quad d_{i0} = \frac{k-i}{kA} - \frac{k'-i}{k'A'}$$

また $i \neq 0, j \neq 0$ のとき

$$d_{ij} = \frac{2(k-i)(l-j)}{lkA} - \frac{2(k'-i)(l'-j)}{l'k'A'}$$

とおく。

そうすると

$$T_{uv} = \sum_{i=1}^u d_{i0} + \sum_{j=1}^v d_{0j} + \sum_{i=1}^u \sum_{j=1}^v d_{ij}$$

よつて

$$\begin{aligned} 2T_{uv} &= \frac{u(k+x-1)}{kA} + \frac{v(l+y-1)}{lA} + \frac{uv(k+x-1)(l+y-1)}{lkA} \\ &\quad - \frac{u(k'+x)}{k'A'} - \frac{v(l'+y)}{l'A'} - \frac{uv(k'+x)(l'+y)}{l'k'A'} \end{aligned}$$

ただし $x = (k-u), y = (l-v)$ である。

これから $\sqrt{l} \leq k \leq l^2$ なら $T_{uv} \geq 0$ および併に T_{kl} が 0 になるから、 $\sqrt{l} \leq k \leq l^2$ なら

$$\sum_{i=0}^k \sum_{j=0}^l d_{ij} = 0 \text{ および } \sum_{i=0}^u \sum_{j=0}^v d_{ij} \geq 0$$

が証明できる。

よつて Lemma N を適用すれば

(i) $\Delta_1 \psi(u, l) \leq 0$

(ii) $\Delta_2 \psi(k, v) \leq 0$

(iii) $\Delta_1 \Delta_2 \psi(u, v) \geq 0$

なら

$$L(l+1, k+1) - L(l, k) \geq 0$$

ゆえに

(i) $\Delta_1 \psi \leq 0$

(ii) $\Delta_2 \psi \leq 0$

(iii) $\Delta_1 \Delta_2 \psi(u, v) \geq 0$

(iv) $(m-1)l = (n-1)k$

(v) $\sqrt{l} \leq k \leq l^2$

なら

$$L(l, k) \leq L(l+1, k+1) \leq \dots \leq L(ml, nk)$$

である。証明終り。

§ 3

我々はここで系統的抽出が層化抽出より有利となるための十分条件の超を求めることにする。

定理 4

(i) $S_1^2 P(u, v) \geq 0$ (ii) $S_2^2 P(u, v) \geq 0$

なるすべての無反母集団において、任意の大きさの標本に対して

$$\sigma_{sy}^2 \leq \sigma_{st}^2$$

が成立する。また上の二つの場合の各々において等号が成立しない限り

$$\sigma_{sy}^2 < \sigma_{st}^2$$

である。

証明

相対的効率 $\sigma_{st}^2 / \sigma_{sy}^2$ は $V > 1$ であるから、これを 1 にとることかできる。すなわち $m\pi$ 個の層の各々から 1 つの標本をとる。

(1/3) から $V = 1$ のとき $(\bar{X}_{sy} - \bar{X})^2$ の平均値すなわち σ_{sy}^2 の推定値は

$$\frac{1}{m\pi l k} \left\{ \sum_{i=1}^{m\pi} \sum_{j=1}^{lk} (X_{ij} - \bar{X})^2 = \frac{1}{m\pi} \times (\text{系統的標本内の平均的 SS.}) \right\} \dots (3.1)$$

$m\pi$ 個の層からとられた系統的標本を $X_{1j}, X_{2j}, \dots, X_{m\pi j}$ とし、またその平均値を \bar{X}_{st} とおく。 X_{ij} は i 層からとられた標本である。

そうすると

$$\bar{X}_{st} = \sum_{i=1}^{m\pi} X_{ij} / m\pi$$

i 層の平均値を \bar{X}_i 、総平均を \bar{X} とおくと

$$\bar{X}_i = \frac{1}{lk} \sum_{j=1}^{lk} X_{ij}, \quad \bar{X} = \frac{1}{m\pi} \sum_{i=1}^{m\pi} \bar{X}_i, \quad \sigma_{st}^2 = E(\bar{X}_{st} - \bar{X})^2$$

このとき

$$\sum_{i=1}^{m\pi} (X_{ij} - \bar{X})^2 = \sum_{i=1}^{m\pi} (X_{ij} - \bar{X}_{st})^2 + m\pi (\bar{X}_{st} - \bar{X})^2$$

$(lk)^{m\pi}$ 個のすべての標本について平均すれば

$$\frac{1}{lk} \sum_{i=1}^{m\pi} \sum_{j=1}^{lk} (X_{ij} - \bar{X})^2 = (\text{標本内の平均的 SS.}) + m\pi (\bar{X}_{st} - \bar{X})^2$$

$$\therefore \sigma_{st}^2 = E(\bar{X}_{st} - \bar{X})^2 = \frac{1}{m\pi lk} E \left\{ \sum_{i=1}^{m\pi} \sum_{j=1}^{lk} (X_{ij} - \bar{X})^2 - \frac{1}{lk} (\text{標本内平均的 SS.}) \right\} \dots (3.2)$$

標本を $X_1, X_2, \dots, X_{m\pi}$ 、その平均値を \bar{X} とおくと標本内の S.S は

$$\sum_{i=1}^{m\pi} (X_i - \bar{X})^2 = \frac{1}{m\pi} \sum_{\substack{i,j=1 \\ i>j}}^{m\pi} (X_i - X_j)^2 \dots (3.3)$$

これは系統的または層化標本内の平均的 S.S である。

系統的小よび層化標本内の平均的 S.S の期待値をそれぞれ

$$\Sigma_{sy}, \Sigma_{st} \text{ とすると (3.1) および (3.2) から } \Sigma_{sy} \geq \Sigma_{st} \text{ なら } \sigma_{st}^2 \geq \sigma_{sy}^2$$

である。

(3.3) から、 Σ_{sy} および Σ_{st} は層の対 (i, j は層を表わす) のすべての組合せについて $E(X_i - X_j)^2$ なる項を加え上げたものに $\frac{1}{m\pi}$ を乗じたものである。

ゆえに Σ_{sy}, Σ_{st} は規則に比較することかできる。

E_{sy} および E_{st} で系統的小よび層化標本に対する期待値を表わす。また X_i, X_j はそれぞれ $(i_1 - j_1) = u$ および $(i_2 - j_2) = v$ なるような (i_1, i_2) および (j_1, j_2) 層に属するものとする。そうすると

$$E_{sy} (X_i - X_j)^2 = E_{sy} (X_{i_1 - \mu + 1} - X_{j_1})^2 = 2\sigma^2 \{ 1 - P(u, v) \} \dots (3.4)$$

ここで X_j は X_i が含まれるはきまり、逆も成立する。

$E_{st} (X_i - X_j)^2$ については、この場合 l^2 対 l^2 通りの対の組合せが可能である。したがって我々はこれらのすべての組合せについて平均値をとる。

($l-1$) 個の対は距離 ($lu + l$) だけ離れており、($l - |j_1|$) は ($lv + j$) だけ離れていることかわかる。

$l^2 k^2$ 通りの組合せについて平均すると

$$E s_t (X_i - X_j)^2 = 2 \sigma^2 \left\{ 1 - \frac{1}{k^2 l^2} \sum_{i=1}^{(k-1)} \sum_{j=1}^{(l-1)} (k-i)(l-j) P(u_{k+i}, v_{l+j}) \right\} \dots (3.5)$$

をうる。

(3.4) と (3.5) によつて、 $(\sum s_y - \sum s_x)$ は

$$\frac{1}{k^2 l^2} \sum_{i=1}^{(k-1)} \sum_{j=1}^{(l-1)} (k-i)(l-j) P(u_{k+i}, v_{l+j}) - P(u_k, v_l) = \frac{1}{k^2 l^2} (T_1 + T_2 + T_3 + T_4)$$

のような項の和に比例することかわかる。

そこで

$$T_1 = \sum_{i=1}^{(k-1)} \sum_{j=1}^{(l-1)} (k-i)(l-j) \{ P(u_{k+i}, v_{l+j}) + P(u_{k+i}, v_{l-j}) + P(u_{k-i}, v_{l+j}) + P(u_{k-i}, v_{l-j}) \}$$

$$T_2 = k \sum_{j=1}^{(l-1)} (l-j) \{ P(u_k, v_{l+j}) + P(u_k, v_{l-j}) \}$$

$$T_3 = l \sum_{i=1}^{(k-1)} (k-i) \{ P(u_{k+i}, v_l) + P(u_{k-i}, v_l) \}$$

$$T_4 = -lk(lk-1) P(u_k, v_l)$$

である。

いま、 T_1 から $4P(u_k, v_l)$ を、 T_2, T_3 からそれぞれ $2P(u_k, v_l)$

を引いて、 T_4 にその分を加えると

$$T'_1 = \sum_{i=1}^{(k-1)} \sum_{j=1}^{(l-1)} (k-i)(l-j) \{ P(u_{k+i}, v_{l+j}) + P(u_{k+i}, v_{l-j}) + P(u_{k-i}, v_{l+j}) + P(u_{k-i}, v_{l-j}) - 4P(u_k, v_l) \}$$

$$T'_2 = k \sum_{j=1}^{(l-1)} (l-j) \{ P(u_k, v_{l+j}) + P(u_k, v_{l-j}) - 2P(u_k, v_l) \}$$

$$T'_3 = l \sum_{i=1}^{(k-1)} (k-i) \{ P(u_{k+i}, v_l) + P(u_{k-i}, v_l) - 2P(u_k, v_l) \}$$

となり、この差引いた分の合計 $lk(lk-1)P(u_k, v_l)$ を T_4 に加えると 0 になる。

よつて式は $T'_1 + T'_2 + T'_3 \dots (3.6)$

となる。

いつものように

$$S^2 = \Delta^2 E^{-1} \equiv (E^{\frac{1}{2}} - E^{-\frac{1}{2}})^2 \text{ とすると}$$
$$\sum_{j=-(i-1)}^{(i-1)} (i-1-j) S^2 E^j \equiv E^i + E^{-i} - 2$$

が証明できる。

よつて

$$\sum_{j=-(i-1)}^{(i-1)} (i-1-j) S^2 \phi(u_{k+j}) = \phi(u_{k+i}) + \phi(u_{k-i}) - 2\phi(u_k) \dots (3.7)$$

ゆえに T'_1 においては

$$P(u_{k+i}, v_{l+j}) + \dots + P(u_{k-i}, v_{l-j}) - 4P(u_k, v_l) = (E^i + E^{-i} - 2) \{ P(u_k, v_{l+j}) + P(u_k, v_{l-j}) + 2(E^{\frac{j}{2}} + E^{-\frac{j}{2}} - 2) P(u_k, v_l) \}$$

(3.7) を用いて

$$= \sum_{d=-(i-1)}^{(i-1)} (i-1-d) \{ S_1^2 P(u_{k+d}, v_{l+j}) + S_2^2 P(u_{k+d}, v_{l-j}) \} + 2 \sum_{t=-(j-1)}^{(j-1)} (j-1-t) S_2^2 P(u_k, v_{l+t})$$

同様に T'_2 については

$$P(u_k, v_{l+j}) + P(u_k, v_{l-j}) - 2P(u_k, v_l) = \sum_{t=-(j-1)}^{(j-1)} (j-1-t) S_2^2 P(u_k, v_{l+t})$$

また T'_3 については

$$P(u_{k+i}, v_l) + P(u_{k-i}, v_l) - 2P(u_k, v_l) = \sum_{d=-(i-1)}^{(i-1)} (i-1-d) S_1^2 P(u_{k+d}, v_l)$$

よつて

(i) $S_1^2 P(u, v)$ (ii) $S_2^2 P(u, v) \geq 0$ なら T'_1, T'_2, T'_3 はいずれもすべての u, v について正である。

ゆえにこれらの条件のもとでは $\sum s_y \geq \sum s_x$ であるから $\sigma_{sy}^2 \leq \sigma_{sx}^2$ この差は明らかに上の3つの場合の各々について解きが成立するときのみ 0 となる。よつて定理は証明された。

この節と前節の定理の証明から $\sigma^2 \geq \sigma_{sy}^2$ などの条件をためることができる。ある場合には、層化した系統的抽出法を採用しなければならぬことがあるから、効率の悪さの程度

(degree of inefficiency) を求めなければならない。
これは容易ではないが、§1の(1.14), (1.15), (1.16) から求めることができる。

ここでは1つの例を考えてみよう。

$$P(u, v) = \sum_{i=1}^p A_i e^{-\lambda_i |u| - \mu_i |v|}$$

ただし $\sum_{i=1}^p A_i = 1$ で λ_i および μ_i は正である。
とする。明らかに $P(u, v) = P(u, v)$ であるから u, v の正の値についてのみ考えればよい。
これは自然な条件であることが示されている (Ghosh, 1949)

いま $\psi(u, v) = 2P(u, v)$
よつて $\Delta_1 \psi(u, v) = 2 \sum_{i=1}^p A_i e^{-\lambda_i u - \mu_i v} (e^{-\lambda_i} - 1)$
しかるに $\lambda_i \leq 0$ のとき

$$e^{-\lambda_i} \leq 1$$

ゆえに正の恒尙の場合

$$\Delta_1 \psi(u, v) \leq 0$$

$\Delta_2 \psi(u, v)$ についても同様である。

しかし λ_i および μ_i の値がわからなければ定理1の条件が満足されるかどうかはわからない。ただしこのすべての値について $\lambda_i = \mu_i$ になるときにのみ定理1の系が適用できる。

更に

$$\Delta_1 \Delta_2 \psi(u, v) = 2A_i e^{-\lambda_i u - \mu_i v} (e^{-\lambda_i} - 1) (e^{-\mu_i} - 1)$$

でこれは正である。

よつて定理3が適用できる。すなわち $\sqrt{r} \leq k \leq l^2$ および $(n-1)k = (m+1)l$ なる層化抽出は、ユニバースが期待平均値 (expected mean) に関して一様であるとしても無作為抽出より
善悪は $\Delta_1 \psi(u, v) \leq 0, \Delta_2 \psi(u, v) \leq 0, S_x^2 P(u, v) \geq 0$ および $S_y^2 P(u, v) \geq 0$ のとき $\sigma_r^2 > \sigma_{sc}^2 \geq \sigma_{sy}^2$ が成立し、また上の4つの場合の各々において等号がなり立つときは $\sigma_r^2 = \sigma_{sc}^2$ で、またそのときには $\sigma_{sc}^2 = \sigma_{sy}^2$ なることも証明した。これらの結果は *Science and Culture* の次号に発表する予定である。

同様に

$$\Delta_2^2 \psi(u, v) = 2 \{ P(u+1, v) + P(u-1, v) - 2P(u, v) \} \\ = 2 \sum_{i=1}^p A_i e^{-\lambda_i u - \mu_i v} (e^{\lambda_i/2} - e^{-\lambda_i/2})^2$$

でこれは P が正のとき正である。

同様に

$$S_2^2 \psi(u, v) \geq 0$$

ゆえに定理4を適用すれば、系統的抽出は層化抽出より有効なることがわかる。

この論文の印刷中、この中で論じた結果の幾つかが Quenouille (1949) によつてもえられているという注意があつた。看看は1949年にこれらの結果を *Science and Culture* に発表した。

幾つかの有益な助言を与えられた Prof. B. N. Ghosh にお礼を言上げる。

By A. C. Das (*Sankhya* vol. 10, 1950 より)

参考文献

1. Cochran, W. G. (1946) Relative accuracy of systematic and stratified random for a certain class of populations. *A. M. S.* 17, 164-177.
2. Das, A. C. (1949): Two dimensional systematic, stratified and random sampling Proc. Ind. Sci. Congress, 36th Session Part III, 6
3. Ghosh, B. N. (1949): On a particular type of natural field Proc. Ind. Sci. Congress, 36th Session, part III, 7.

- 4 Madow, W. G. and S. H. (1944): On the theory of systematic sampling I. *A. M. S.* 15, 1-24.
- 5 Queneuille, M. H. (1949): problem in plane sampling *A. M. S.* 20, 355-375
- 6 Yates, F. (1948): Systematic sampling, *Phil. Trans. Roy. Soc.* 241 (A), 345-377

7 種々の形式の二重抽出法による
推定値の誤差について

On errors of estimates in various
type of double sampling procedure

1 緒言

二重抽出法なる用語は二つの抽出調査を含むようなサンプリング方法に適用されるようになってきた。他の多くの種類のサンプリングと同じく、費用の減少と精度の増加はこの種のサンプリングの主な利益でもある。

Keyman (1938)は、補助変量のオ1の大標本を用いてオ2の(主たる)特性の変動を小さくするように母集団をグループわけし、特性に相関がある場合には比較的小数のオ2標本から主特性の良好な推定値を求めうるような抽出方法を与える。

他のよく知られた種類の二重抽出法は、二つの特性についてとられた大きさNのオ1標本を用いて主特性yの他のXに対する回帰を定め、補助特性Xのみについて観察された大きさNのオ2標本を用いて主特性yの推定値を求めめるものである。この方法は特に、主特性の調査に非常に費用がかかるのに対して、それと相関を有する補助変量は容易に測定できるというような場合に適用できる。Cochran (1943)は Snedecor and King (1942)およびC. Bose (1943)によって与えられた線型回帰の仮定のもとに、この種の推定値の分散に対する抽出公式 (sampling formulas) の例を与えた。補助変量か1つの場合の特別な形の非線型回帰に対する推定値の分散の式は Bose と Goryen (1946) によって算びかれた。

しかしながら、この種の二重抽出法の推定値の精度は、一つでなく多くの相関ある補助変量を含めることによつて向上させることが可能であろう。B. Ghosh (1947)は、回帰が線型で無作為抽出を行なう場合は推定値は不偏であることを示しまた多くの補助変量にせよ推定値の分散の近似公式をえている。

前に述べたように二重抽出法は補助変数が1つでもあつても多数であつても、それ自身が標本の精度を増加させる1つの方法なのである。したがつて更に他の精度の増すことのでわかつていゝサンプリング方法と組合せて用いることが可能であらう。すなわちオノおよびオ2標本の単位の選択には色々な方法を用いることができる。例えばオノ標本を無作為に、またオ2標本を系統的に選ぶというように、その上費用と精度の考慮から、ある時は予め定められたものをオノ標本としあるいは特に選ばれた値の組にとることも是認されるであらう。

この論文では色々な形式の二重抽出法に対する推定値および推定値の分散の式を求めた。補助変数が1つの場合と2つ以上の場合を別々に取扱つた。殆んどの場合線型回帰を仮定したか、2例だけは非線型回帰を仮定して考案した。オノ標本の補助変数の期待値が他の変数の期待値の一定倍になつてゐるように修正された場合を考案した。

オノ、オ2標本に対する標本単位の最適配分の問題は、Schumaker and Chapman (1942) の採用した方針にしたがつて、色々なタイプの中の一つに対する最適個数を導くことによつて論じておいた。

附録においてまだ知られていないと思われるある種の興味ある結果、たとえば多変量正規母集団での回帰係数の同時分布、および標本変動行列の代表要素の期待値などを導びいた。これとともに、すでに知られてゐる偏回帰係数の分布も直交座標を用いるもう一つの方法で導びいた。

2. 多変量補助変数の組 — サンプリングにおける色々な場合

2.1 X_n — ランダム, y_n — ランダム, X_N — ランダム:

補助変数が X_1, X_2, \dots, X_k なる組からなつてゐるとき、線型回帰を仮定すると、無作為抽出 n の y の母集団平均値の推定値は

$$Y = \bar{y}_n + \sum_{i=1}^k b_{ni} (\bar{X}_{n_i} - \bar{X}_{N_i}) \quad (2.11)$$

で与えられる。ここで $\bar{y}_n, b_{ni}, \bar{X}_{n_i}$ はオノ標本から、ま

\bar{X}_{N_i} はオ2標本から導びかれる。

多変量正規母集団を仮定し

$$E(y) = \eta \quad E(X_i) = \xi_i (i=1, 2, \dots, k)$$

$$V(y) = \sigma_y^2 \quad V(X_i) = \sigma_{ii} = \sigma_i^2$$

$$\text{Cov}(X_i, X_j) = \sigma_{ij} = \sum_{l=1}^k \sigma_{li} \sigma_{lj}$$

$$\text{Cov}(y, X_i) = \sigma_{yi} = \sum_{j=1}^k y_j \sigma_{ij} \sigma_j$$

$$E(b_{ni}) = b_{ni} = -\frac{R_{yi}}{R_{YX}} \frac{\sigma_y}{\sigma_i}$$

ただし R_{ij} は行列式 $|R_{ij}|$ ($i, j=1, 2, \dots, k$) の $\sum_{i=1}^k$ の余因子、 R_{YX} は行列 A の結果を用いて直ちに次のような関係がえられる。

$$V(b_{ni}) = \sigma_y^2 \cdot X \cdot E(C_{ii}) = \frac{\sigma_y^2 \cdot X \cdot \sigma_{ii}}{n-k-2} = \frac{\sigma_y^2 (1-R_{Y,12\dots k}) \sigma_{ii}}{n-k-2}$$

$$\text{Cov}(b_{ni}, b_{nj}) = \sigma_y^2 \cdot X \cdot E(C_{ij}) = \frac{\sigma_y^2 \cdot X \cdot \sigma_{ij}}{n-k-2} = \frac{\sigma_y^2 (1-R_{Y,12\dots k}) \sigma_{ij}}{n-k-2}$$

ここで $\sigma_y^2 \cdot X$ は X を固定したときの y の変異分散で、 $R_{Y,12\dots k}$ は X_1, X_2, \dots, X_k 同の重相関係数

また行列 (σ_{ij}) は (σ_i^2) の逆行列である。

したがつてオノ標本内の y と X 、オ2標本内の X をすべてランダムに選ぶ場合の二重抽出においては

$$E(Y) = \eta + \sum_{i=1}^k \beta_{ni} (\bar{X}_{n_i} - \bar{X}_{N_i}) = \eta$$

$$V(Y) = E\left\{ (\bar{y}_n - \eta) + \sum_{i=1}^k b_{ni} (\bar{X}_{n_i} - \bar{X}_{N_i}) - \sum_{i=1}^k b_{ni} (\bar{X}_{n_i} - \bar{X}_{N_i}) \right\}^2$$

$$= \frac{\sigma_y^2}{n} + \left(\sum_{i=1}^k \sum_{j=1}^k \sigma_{ij} \beta_{ni} \beta_{nj} + \frac{\sigma_y^2 \cdot X}{n-k-2} \right) \left(\frac{1}{n} + \frac{1}{N} \right) - 2 \sum_{i=1}^k \beta_{ni} \frac{\sigma_{yi}}{n}$$

$$= \frac{\sigma_y^2}{n} + \left(\frac{1}{n} + \frac{1}{N} \right) \sum_{i=1}^k \sum_{j=1}^k \beta_{ni} \beta_{nj} \sigma_{ij}$$

$$+ \frac{\sigma_y^2 (1-R_{Y,12\dots k})}{n-k-2} \cdot \left(\frac{1}{n} + \frac{1}{N} \right) \sum_{i=1}^k \sum_{j=1}^k \sigma_{ij} \sigma_{ij} - 2 \sum_{i=1}^k \beta_{ni} \frac{\sigma_{yi}}{n}$$

が成立する。

いま $\sum_{i=1}^k \sigma_{ij} \sigma_{ij} = 1$ $\therefore \sum_{i=1}^k \sum_{j=1}^k \sigma_{ij} \sigma_{ij} = k$

$$\sum_{i=1}^k \beta_{ni} \sigma_{ij} = \sigma_{ij}$$

$$\sum_{i=1}^k \sum_{j=1}^k \beta_{ni} \beta_{nj} = \sum_{i=1}^k \beta_{ni} \sigma_{ij} = \sum_{i=1}^k \beta_{ni} - \frac{R_{ij} \sigma_{ij}}{R_{ij} \sigma_{ij}} \sum_{i=1}^k \beta_{ni} \sigma_{ij} \sigma_{ij}$$

$$= \frac{-\left(\sum_{i=1}^k \sum_{j=1}^k Z_{ij} R_{ij} \sigma_{ij}^2\right)}{R_{ij}} = \left(\frac{R - R_{ij}}{R_{ij}}\right) \sigma_{ij}^2$$

$$= \sigma_{ij}^2 R^2_{y,12 \dots k}$$

よって

$$V(Y) = \frac{\sigma_y^2}{n} + \frac{k}{n-k-2} \sigma_y^2 (1 - R^2_{y,12 \dots k}) \left(\frac{1}{n} + \frac{1}{N}\right)$$

$$+ \sigma_y^2 R^2_{y,12 \dots k} \left(\frac{1}{N} - \frac{1}{n}\right)$$

$$= \sigma_y^2 (1 - R^2_{y,12 \dots k}) \left(\frac{1}{n} + \frac{k}{n-k-2} \left(\frac{1}{n} + \frac{1}{N}\right)\right) \frac{R^2_{y,12 \dots k}}{N} \quad (2.13)$$

(2.13)式は各々の分散 σ_{ij} と独立であることに注意せよ

(2.2) X_n -ランダム, Y_n -ランダム, X_N -システムティック:
 ここで更に一般の場合, すなわちオノ標本の X とオノ標本の Y はランダムに選ぶが, オノ標本はシステムティックな方法で選ぶことにしても, 推定値 Y については(2.11)と同じ関係が成立する.

$$E(Y) = \eta$$

$$E\left\{(\bar{X}_{Np} - \bar{z}_i)(\bar{X}_{Nq} - \bar{z}_j)\right\} = \frac{1}{N^2} E\left\{\sum_{p=1}^N \sum_{q=1}^N (X_{Nip} - \bar{z}_i)(X_{Njq} - \bar{z}_j)\right\}$$

$$= \frac{\sum_{p=1}^N \sum_{q=1}^N \sum_{i,p,j,q}^{(XX)} \sigma_{ij}}{N^2}$$

ここで $Z_{ij}^{(XX)}$ は($p=1, \dots, N; q=1, \dots, N$)は母集団におけ

る X_{Nip} と X_{Njp} の相関を表わす.

$$V(Y) = E\left\{\left(\bar{y}_n - \eta\right) + \sum_{i=1}^k \beta_{ni} (\bar{X}_{Ni} - \bar{z}_i) - \sum_{i=1}^k \beta_{ni} (\bar{z}_i - \eta)\right\}^2$$

$$= \frac{\sigma_y^2}{n} + \sum_{i=1}^k \sum_{j=1}^k \left[\beta_{ni} \beta_{nj} + \frac{\sigma_{ij}^2}{n-k-2}\right] \left[\frac{\sum_{p=1}^N \sum_{q=1}^N \sum_{i,p,j,q}^{(XX)} \sigma_{ij} \sigma_{ij}}{N^2} + \frac{\sigma_{ij}}{n}\right]$$

$$- 2 \sum_{i=1}^k \beta_{ni} \frac{\sigma_{ij}}{n}$$

$$= \frac{\sigma_y^2}{n} + \sum_{i=1}^k \sum_{j=1}^k \left[\beta_{ni} \beta_{nj} + \frac{\sigma_y^2 (1 - R^2_{y,12 \dots k}) \sigma_{ij}}{n-k-2}\right] \times$$

$$\left\{ \frac{\sum_{p=1}^N \sum_{q=1}^N \sum_{i,p,j,q}^{(XX)} \sigma_{ij} \sigma_{ij}}{N^2} + \frac{k \sigma_y^2 (1 - R^2_{y,12 \dots k})}{n(n-k-2)} - \frac{\sigma_y^2 R^2_{y,12 \dots k}}{n} \right\}$$

$$= \frac{\sigma_y^2 (1 - R^2_{y,12 \dots k})}{n} \frac{n-2}{n-k-2} + \sum_{i=1}^k \sum_{j=1}^k \frac{\sum_{p=1}^N \sum_{q=1}^N \sum_{i,p,j,q}^{(XX)} \sigma_{ij} \sigma_{ij}}{N^2}$$

$$\times \left[\beta_{ni} \beta_{nj} + \frac{\sigma_y^2 (1 - R^2_{y,12 \dots k}) \sigma_{ij}}{n-k-2}\right] \dots \dots \dots (2.22)$$

2.1で考察したのは, これの

$$Z_{ipjp}^{(XX)} = Z_{ij} \quad p=q \text{ のとき}$$

$$= 0 \quad p \neq q \text{ のとき}$$

であるような特別な場合になっている.

(2.3) X_n -固定, Y_n -ランダム, X_N -ランダム

次に我々は, オノ標本の X とは固定するが, オノ標本の X とオノ標本の Y はランダムに選ぶというオノの場合を考える. Y の母集団平均値の推定値はこの場合も(2.11)と同じである:

$$E(Y_{np}) = \alpha_n + \sum_{i=1}^k \beta_{ni} X_{nipo} \quad (p=1, 2, \dots, n)$$

と推定すると

$$E(\bar{y}_n) = \alpha_n + \sum_{i=1}^k \beta_{ni} (\equiv \eta' \text{ とおく})$$

よって

$$E(Y) = \eta' + \sum_{i=1}^k \beta_{ni} (\bar{z}_{ni} - \bar{X}_{ni}) = \eta \quad \dots \dots \dots (2.31)$$

$$V(Y) = E \left\{ (\bar{y}_n - \eta) - \sum_{i=1}^k (\bar{x}_{ni} - \bar{z}_i) (b_{ni} - \beta_{ni}) + \sum_{i=1}^k b_{ni} (\bar{x}_{ni} - \bar{z}_i) \right\}^2$$

$$= A + \frac{B}{2} + \frac{C}{N} \dots \dots \dots (2.32)$$

ただし

$$A = C_{ix}^2 \left\{ \sum_{i=1}^k \sum_{j=1}^k (\bar{x}_{ni} - \bar{z}_i) (\bar{x}_{nj} - \bar{z}_j) C_{ij} \right\}$$

$$B = C_{ix}^2$$

$$C = \sum_{i=1}^k \sum_{j=1}^k \sigma_{ij} \left\{ C_{ij} \sigma_{yx}^2 + \beta_{ni} \beta_{nj} \right\}$$

系総費用 T が

$$T = d + \beta n + \gamma N \dots \dots \dots (2.33)$$

ただし、 β, γ はテータから推定されるパラメータである。
 なる形で表わされるとすれば、与えられた費用 T に対して分散を最小とするような n および N の最適値は

$$\frac{\partial V}{\partial n} + \lambda \frac{\partial T}{\partial n} = 0$$

$$\frac{\partial V}{\partial N} + \lambda \frac{\partial T}{\partial N} = 0$$

すなわち

$$\frac{Bn}{\sqrt{B\beta}} = \frac{\gamma N}{\sqrt{C\gamma}} = \frac{T-d}{\sqrt{B\beta} + \sqrt{C\gamma}}$$

なる方程式を満足しなければならぬから

$$\left. \begin{aligned} n &= \frac{T-d}{\beta} \frac{\sqrt{B\beta}}{\sqrt{B\beta} + \sqrt{C\gamma}} \\ N &= \frac{T-d}{\gamma} \frac{\sqrt{C\gamma}}{\sqrt{B\beta} + \sqrt{C\gamma}} \end{aligned} \right\} \dots \dots \dots (2.34)$$

である。

(2.4) X_n - 固定, Y_n - ランダム, X_N - システムティック

ここではオノ標本の X は固定するか、 Y はランダムで、オノ標本の X には相関があるというオノの場合を考える。ここでは (2.11)

は成立して

$$E(Y) = \eta + \sum_{i=1}^k \beta_{ni} (\bar{z}_i - \bar{x}_{ni}) = \eta \dots \dots \dots (2.41)$$

$$V(Y) = E \left\{ (\bar{y}_n - \eta) - \sum_{i=1}^k (\bar{x}_{ni} - \bar{z}_i) (b_{ni} - \beta_{ni}) - \sum_{i=1}^k b_{ni} (\bar{x}_{ni} - \bar{z}_i) \right\}^2$$

$$= \frac{\sigma_{yx}^2}{n} + \sum_{i=1}^k \sum_{j=1}^k E \left\{ \frac{\sum_{p=1}^N \sum_{q=1}^N (X_{ni,p} - \bar{z}_i) (X_{nj,q} - \bar{z}_j)}{N^2} \right\} E(b_{ni} b_{nj})$$

$$+ \sum_{i=1}^k \sum_{j=1}^k (\bar{x}_{ni} - \bar{z}_i) (\bar{x}_{nj} - \bar{z}_j) C_{ij} \sigma_{yx}^2$$

$$= \frac{\sigma_{yx}^2}{n} + \frac{1}{N^2} \sum_{i=1}^k \sum_{j=1}^k \left\{ \beta_{ni} \beta_{nj} + C_{ij} \sigma_{yx}^2 \right\} \left\{ \sum_{p=1}^N \sum_{q=1}^N L_{ip, jq}^{(xx)} \sigma_i \sigma_j \right\}$$

$$+ \sum_{i=1}^k \sum_{j=1}^k (\bar{x}_{ni} - \bar{z}_i) (\bar{x}_{nj} - \bar{z}_j) C_{ij} \sigma_{yx}^2 \dots \dots \dots (2.42)$$

(2.42) から $V(Y)$ はオノ標本の X の範囲にわたり、しかも平均値 \bar{x}_{ni} が母集団平均 \bar{z}_i ($i=1, 2, \dots, k$) に近い程小となることかわかる。

明らかに (2.32) の結果は (2.42) から

$$L_{ip, jq}^{(xx)} = z_i \quad \text{オノのとき}$$

$$= 0 \quad \text{オノでないとき}$$

3. 補助変数が1つのとき - サムプリングにおける色々な場合
 (3.1) X_n - 固定, Y_n - システムティック, X_N - ランダム:
 我々はここで補助変数が1つしかないという特別な場合の二重相関性を考察する。 X はオノ標本では固定するか Y には相関があり、オノオノ標本の X はランダムに送るものとする。

$$E(X_{ni}) = \bar{z}_i \quad (i=1, 2, \dots, N)$$

$$E(Y_{ni}) = d_n + \beta_n X_{ni} \quad (i=1, 2, \dots, n)$$

$$V(Y_{ni} | X_{ni}) = \sigma_{yx}^2$$

$$\text{Cor}(Y_{ni}, Y_{nj} | X_{ni}, X_{nj}) = C_{ij} \sigma_{yx}^2 \quad (3.1)$$

このとき

$$Y = \bar{y}_n + b_n (\bar{X}_N - \bar{X}_n) \dots \dots \dots (3.11)$$

$$E(b_n) = E \left\{ \frac{\sum_{i=1}^n (X_{ni} - \bar{X}_n) (y_{ni} - \bar{y}_n)}{\sum_{i=1}^n (X_{ni} - \bar{X}_n)^2} \right\}$$

$$= \beta_n$$

$$E(Y) = (d_n + \beta_n \bar{X}_n) + \beta_n (\bar{Z} - \bar{X}_n)$$

$$= d_1 + \beta_n \bar{Z} \dots \dots \dots (3.12)$$

が成立する。

さらに

$$V(b_n) = \frac{\sigma_{yx}^2 \sum_{i=1}^n \sum_{j=1}^n (X_{ni} - \bar{X}_n) (X_{nj} - \bar{X}_n) z_{ij}^{(YY)}}{\left\{ \sum_{i=1}^n (X_{ni} - \bar{X}_n)^2 \right\}}$$

$$\text{cov}(b_n, \bar{y}_n) = \frac{1}{n} E \left\{ \frac{\sum_{i=1}^n y_{ni} (X_{ni} - \bar{X}_n) \sum_{j=1}^n (y_{nj} - d_n - \beta_n X_{nj})}{\sum_{i=1}^n (X_{ni} - \bar{X}_n)^2} \right\}$$

$$= \frac{\sigma_{yx}^2 \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)} (X_{ni} - \bar{X}_n)}{n \sum_{i=1}^n (X_{ni} - \bar{X}_n)^2}$$

$$V(\bar{y}_n) = \frac{\sigma_{yx}^2}{n} \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)}$$

である。

よって

$$V(Y) = E \left\{ (\bar{y}_n - d_n - \beta_n \bar{X}_n) + b_n (\bar{X}_n - \bar{Z}) + (\beta_n - \beta_n) (\bar{Z} - \bar{X}_n) \right\}^2$$

$$= \sigma_{yx}^2 \left[\frac{\sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)}}{n^2} + \frac{\sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)} (X_{ni} - \bar{X}_n) (X_{nj} - \bar{X}_n)}{\left\{ \sum_{i=1}^n (X_{ni} - \bar{X}_n)^2 \right\}^2} \right]$$

$$\times \left\{ \frac{\sigma_x^2}{N} + (\bar{Z} - \bar{X}_n)^2 \right\} + \frac{2(\bar{Z} - \bar{X}_n) \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)} (X_{ni} - \bar{X}_n)}{n \sum_{i=1}^n (X_{ni} - \bar{X}_n)^2} + \frac{\sigma_{yx}^2 \beta_n^2}{N} \dots \dots \dots (3.13)$$

(3.13)の特別な場合、すなわち

$$z_{ij}^{(YY)} = 1 \quad (i=j, \dots, n)$$

$$z_{ij}^{(YY)} = 0 \quad (i \neq j)$$

なる場合が $d=1$ のときの (2.32) である。 $V(Y)$ の推定値を求めるには相関係数 $z_{ij}^{(YY)}$ について $|i-j|=1$ なるとき $z_{ij} = z_{ji}$ というようなある仮定を設けなければならぬ。

(3.2) X_n -固定、 y_n -システマティック、 X_n -システマティック

ここではオノ標本の X は固定するか、オノ標本の X とオノ標本の y には相関があるものとする。

このときも (3.11) と (3.12) は成立するから

$$\text{cov}(X_{ni}, X_{nj}) = z_{ij}^{(XX)} \sigma_x^2 \quad (i, j = 1, 2, \dots, N)$$

と仮定すれば

$$V(Y) = \sigma_{yx}^2 \left[\frac{\sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)} + \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)} (X_{ni} - \bar{X}_n) (X_{nj} - \bar{X}_n)}{n^2} \right]$$

$$\times \left\{ \frac{\sigma_x^2}{N^2} \sum_{i=1}^N \sum_{j=1}^N z_{ij}^{(XX)} + (\bar{Z} - \bar{X}_n)^2 \right\} + \frac{2(\bar{Z} - \bar{X}_n)}{\sum_{i=1}^n (X_{ni} - \bar{X}_n)^2}$$

$$\times \left[\sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(YY)} (X_{ni} - \bar{X}_n) \right] + \frac{\sigma_{yx}^2 \beta_n^2}{N^2} \sum_{i=1}^N \sum_{j=1}^N z_{ij}^{(XX)} \dots \dots \dots (3.21)$$

が成立する。

$$z_{ij}^{(XX)} = 1 \quad (i=j, \dots, N)$$

$$z_{ij}^{(XX)} = 0 \quad (i \neq j)$$

ならこの式は (3.13) に帰着する。

(3.3) X_n -システマティック、 y_n -システマティック、 X_n -システマティック

我々はここで更に一般な場合、すなわちオノ標本内の X と y およびオノ標本内の X かすべてそれら自身の間で相関を有するという場

合を取扱う，この場合のYの期待値および分散は近似公式でさえ非常に複雑である。ここで(3.1)は成立する。オノ標本のXとYについて多変量正規分布を仮定し、

$$E(X_{ni}) = E(X_{N_i}) = \bar{x}$$

$$E(Y_{ni}) = \bar{y}$$

$$\left. \begin{aligned} \text{Cov}(X_{ni}, X_{nj}) &= z_{ij}^{(xx)} \sigma_x^2 \\ \text{Cov}(Y_{ni}, Y_{nj}) &= z_{ij}^{(yy)} \sigma_y^2 \\ \text{Cov}(X_{ni}, Y_{nj}) &= z_{ij}^{(xy)} \sigma_x \sigma_y \end{aligned} \right\} (i, j = 1, 2, \dots, n)$$

$$E(X_{ni} - \bar{x})(X_{nj} - \bar{x})(Y_{nk} - \bar{y}) = 0 \quad (i, j, k = 1, 2, \dots, n)$$

とおけば

$$E(X_{ni}, X_{nj}, Y_{nk}) = \left\{ \bar{x}^2 + z_{ij}^{(xx)} \sigma_x^2 \right\} \bar{y} + \bar{y} \sigma_x \sigma_y \left\{ z_{ik}^{(xy)} + z_{ij}^{(xy)} \right\}$$

であるから

$$\begin{aligned} E\left\{ \frac{\sum_{i=1}^n X_{ni} Y_{ni}}{n} - \frac{\sum_{i=1}^n X_{ni}}{n} \frac{\sum_{i=1}^n Y_{ni}}{n} \right\} \\ E(b_n) &= \frac{E\left\{ \frac{\sum_{i=1}^n X_{ni}^2}{n} - \frac{(\sum_{i=1}^n X_{ni})^2}{n} \right\}}{\left\{ n(n-1) - \sum_{i \neq j=1}^n z_{ij}^{(xx)} \right\}} \\ &= \frac{\sigma_y}{\sigma_x} \cdot \frac{\left\{ n(n-1) \bar{z} - \sum_{i \neq j=1}^n z_{ij}^{(xy)} \right\}}{\left\{ n(n-1) - \sum_{i \neq j=1}^n z_{ij}^{(xx)} \right\}} = b_n \text{ (とおく)} \end{aligned}$$

ここで $z_{ij} = z_{ij}^{(xy)}$ ($i, j = 1, 2, \dots, n$) が成立する。

さらに、

$$\begin{aligned} E\left\{ \frac{\sum_{i=1}^n X_{ni}}{n} \frac{\sum_{i=1}^n X_{ni} Y_{ni}}{n} \right\} &= \sum_{i=1}^n \sum_{j=1}^n E(X_{ni} Y_{ni} X_{nj}) \\ &= n^2 \bar{x}^2 \bar{y} + \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(xx)} \sigma_x^2 \bar{y} + n^2 \bar{x} \bar{y} \sigma_x \sigma_y + \sum_{i=1}^n \sum_{j=1}^n z_{ij} \\ &+ \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(xy)} \bar{y} \sigma_x \sigma_y \end{aligned}$$

$$\begin{aligned} E\left\{ \left(\frac{\sum_{i=1}^n X_{ni}}{n} \right)^2 \frac{\sum_{i=1}^n Y_{ni}}{n} \right\} &= \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n E(X_{ni} X_{nj} Y_{nk}) \\ &= n^2 \bar{x}^2 \bar{y} + n \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(xx)} \sigma_x^2 \bar{y} + 2n \sum_{i=1}^n \sum_{j=1}^n z_{ij}^{(xy)} \bar{y} \sigma_x \sigma_y \end{aligned}$$

よって

$$\begin{aligned} E\left\{ \frac{\sum_{i=1}^n X_{ni} Y_{ni}}{n} - \frac{(\sum_{i=1}^n X_{ni})^2}{n} \frac{\sum_{i=1}^n Y_{ni}}{n} \right\} \\ E(b_n \bar{X}_n) &= \frac{n E\left\{ \frac{\sum_{i=1}^n X_{ni}^2}{n} - \frac{(\sum_{i=1}^n X_{ni})^2}{n} \right\}}{n E\left\{ \frac{\sum_{i=1}^n X_{ni}^2}{n} - \frac{(\sum_{i=1}^n X_{ni})^2}{n} \right\}} \\ &= \bar{y} \frac{\sigma_y}{\sigma_x} \frac{\left\{ n(n-1) \bar{z} - \sum_{i \neq j=1}^n z_{ij}^{(xy)} \right\}}{\left\{ n(n-1) - \sum_{i \neq j=1}^n z_{ij}^{(xx)} \right\}} \\ &= E(\bar{X}_n) E(b_n) \\ \therefore E(Y) &\approx \bar{y} + b_n (\bar{x} - \bar{x}) = \bar{y} \quad \text{----- (3.31)} \end{aligned}$$

また

$$\begin{aligned} E(b_n^2) &= \frac{E\left\{ \frac{\sum_{i=1}^n (X_{ni} - \bar{x})(Y_{ni} - \bar{y})}{n} \right\}^2}{E\left\{ \frac{\sum_{i=1}^n (X_{ni} - \bar{x})^2}{n} \right\}^2} \\ &= \frac{E\left\{ \frac{\sum_{i=1}^n \sum_{j=1}^n X_{ni} X_{nj} Y_{ni} Y_{nj}}{n} \right\}}{E\left\{ \frac{\sum_{i=1}^n \sum_{j=1}^n X_{ni}^2 X_{nj}^2}{n} \right\}} \end{aligned}$$

ここで X'_{ni}, Y'_{ni} は対応する変量のその平均値からの偏差を表わす。

(3.31)のときの $E(X'_{ni}, X'_{nj}, Y'_{ni}, Y'_{nj})$ は4変数の積乗母函数(この分析は多変量正規分布であると仮定されている)から置ちに求められて

すれは、

$$y = a + \beta(x - \bar{x}_n) + \gamma(x^2 - \bar{x}_n^2)$$

と書ける。ここで \bar{x}_n と \bar{x}_n^2 はオノ標本の x と x^2 の平均値。すなわち

$$\bar{x}_n = \frac{\sum_{i=1}^n x_{ni}}{n} \quad \text{および} \quad \bar{x}_n^2 = \frac{\sum_{i=1}^n x_{ni}^2}{n}$$

である。このとき a, β, γ を推定するための正規方程式は

$$\bar{y}_n = a_n$$

$$\sum (x - \bar{x}_n) y = b_n \sum (x - \bar{x}_n)^2 + c_n \sum (x - \bar{x}_n) (x^2 - \bar{x}_n^2)$$

$$\sum (x^2 - \bar{x}_n^2) y = b_n \sum (x - \bar{x}_n) (x^2 - \bar{x}_n^2) + c_n \sum (x^2 - \bar{x}_n^2)^2$$

ここで a_n, b_n および c_n はそれぞれ a, β, γ の推定値である。上式から

$$b_n = C_{11} \sum (x - \bar{x}_n) y + C_{12} \sum (x^2 - \bar{x}_n^2) y$$

$$c_n = C_{21} \sum (x - \bar{x}_n) y + C_{22} \sum (x^2 - \bar{x}_n^2) y$$

と書ける。ここで C_{11}, C_{12}, C_{22} はすべて求められる。

オノ標本の x を固定し、オノ標本の y とオノ標本の x をランダムと仮定すると

$$V(a_n) = \frac{\sigma_y^2}{n}, \quad \text{Cov}(a_n, b_n) = \text{Cov}(a_n, c_n) = 0$$

$$V(b_n) = C_{11} \sigma_y^2, \quad \text{Cov}(b_n, c_n) = C_{12} \sigma_y^2$$

$$V(c_n) = C_{22} \sigma_y^2$$

この場合 y の母集団平均値は

$$Y = \bar{y}_n + b_n (\bar{x}_n - \bar{x}_n) + c_n (\bar{x}_n^2 - \bar{x}_n^2) \dots \dots (3.51)$$

である。ここで \bar{x}_n と \bar{x}_n^2 は \bar{x}_n, \bar{x}_n^2 と同様な意味をもつ。

$$E(\bar{x}_n) = \mu_1', \quad E(\bar{x}_n^2) = \mu_2'$$

とおけば

$$E(Y) = a + \beta(\mu_1' - \bar{x}_n) + \gamma(\mu_2' - \bar{x}_n^2) = \eta \dots \dots (3.52)$$

および

$$\begin{aligned} V(Y) &= E \left\{ (\bar{y}_n - a) + (b_n - \beta)(\bar{x}_n - \mu_1') - (c_n - \gamma)(\bar{x}_n^2 - \mu_2') \right. \\ &\quad \left. + \beta(\bar{x}_n - \mu_1') + (\gamma - \beta)(\bar{x}_n^2 - \mu_2') \right\}^2 \\ &= \frac{\sigma_y^2}{n} + C_{11} \sigma_y^2 \left\{ \frac{\mu_2' - \mu_1'^2}{N} + (\bar{x}_n - \mu_1')^2 \right\} + \beta \frac{\mu_2' - \mu_1'^2}{N} \\ &\quad + C_{22} \sigma_y^2 \left\{ \frac{\mu_2' - \mu_1' \mu_2'}{N} + (\bar{x}_n^2 - \mu_2')^2 \right\} + \gamma^2 \frac{\mu_2' - \mu_1'^2}{N} \\ &\quad + 2C_{12} \sigma_y^2 \left\{ \frac{\mu_2' - \mu_1' \mu_2'}{N} + (\bar{x}_n - \mu_1')(\bar{x}_n^2 - \mu_2') \right\} \\ &\quad + 2\beta\gamma \frac{\mu_2' - \mu_1' \mu_2'}{N} \dots \dots (3.53) \end{aligned}$$

ここで

$$\mu_r' = E \left\{ \frac{\sum_{i=1}^n x_{ni}^r}{n} \right\}$$

(3.5a) (3.5) の上に述べた結果は y と x 間の非線形関係に P 次の放物線すなわち

$$y = a + \beta^1(x - \bar{x}_n) + \beta^2(x^2 - \bar{x}_n^2) + \dots + \beta^p(x^p - \bar{x}_n^p)$$

ここで $\bar{x}_n^k = \frac{\sum_{i=1}^n x_{ni}^k}{n}$ (β の添字は巾ではない。)

を仮定した場合にも拡張できる。

係数 $a, \beta^1, \beta^2, \dots, \beta^p$ は $P+1$ 個の正規方程式を解いて得た $a_n, b_n^1, b_n^2, \dots, b_n^p$ で推定でき、

$$b_n^k = C_{k1} \sum (x - \bar{x}_n) y + C_{k2} \sum (x^2 - \bar{x}_n^2) y + \dots + C_{kP} \sum (x^P - \bar{x}_n^P) y$$

($k = 1, 2, \dots, P$)

と書ける。

ここで $V(a_n) = \frac{\sigma_y^2}{n}, \text{Cov}(b_n^i, b_n^j) = C_{ij} \sigma_y^2 (i, j = 1, 2, \dots, P)$

大きさ N の x の他の独立な標本をとつたときの y の母集団平均の不偏推定値は $Y = \bar{y}_n + b_n^1(\bar{x}_n - \bar{x}_n) + b_n^2(\bar{x}_n^2 - \bar{x}_n^2) + \dots + b_n^p(\bar{x}_n^p - \bar{x}_n^p) \dots (3.51a)$

$$\begin{aligned} E(Y) &= \eta \\ V(Y) &= \frac{\sigma_y^2}{n} + \sigma_y^2 \sum_{i=1}^P \sum_{j=1}^P C_{ij} \left\{ \frac{\mu_{i+j}' - \mu_i' \mu_j'}{N} + (\bar{x}_n^i - \mu_i')(\bar{x}_n^j - \mu_j') \right\} \\ &\quad + \sum_{i=1}^P \sum_{j=1}^P b_n^i b_n^j \frac{\mu_{i+j}' - \mu_i' \mu_j'}{N} \dots (3.53a) \end{aligned}$$

となる。

(3.6) 非変型回帰の場合: X_n -ランダム, Y_n -ランダム, Z_n -ランダム

ここで Y_n の誤差の X_n も固定されることなくランダムにとられるものとするは前と同じく

$$Y = \bar{y}_n + b_n (\bar{X}_N - \bar{X}_n) + c_n (\bar{X}_N^2 - \bar{X}_n^2) \dots \dots (3.61)$$

である。

$$E(\bar{X}_n) = E(\bar{X}_N) = \mu'_1$$

$$E(\bar{X}_n^2) = E(\bar{X}_N^2) = \mu'^2_2$$

$$E(b_n) = \alpha$$

とおき、大標本近似 (b_n が \bar{X}_n と独立、かつ c_n が \bar{X}_n^2 と独立) を仮定すると、 Y の期待値は近似的に

$$E(Y) \approx \alpha + E(b_n)E(\bar{X}_N - \bar{X}_n) + E(c_n)E(\bar{X}_N^2 - \bar{X}_n^2) \approx \alpha \dots (3.62)$$

で与えられる。

$$\text{また } V(b_n) = \sigma^2_{YX} E(C_{11}) = \sigma^2_{4b} \quad (\text{とおく})$$

$$V(c_n) = \sigma^2_{YX^2} E(C_{22}) = \sigma^2_{cc} \quad (\text{とおく})$$

$$\text{Cov}(b_n, c_n) = \sigma_{YX} \sigma_{YX^2} E(C_{12}) = \sigma_{bc} \quad (\text{とおく})$$

ここで

$$C_{11} = \frac{\sum (X - \bar{X}_n)^2}{\sum (X^2 - \bar{X}_n^2)}$$

で、 C_{22} および C_{12} も同様である。

ここで

$$E\left\{\sum (X - \bar{X}_n)^2\right\} = E\left\{\sum X^2 - \frac{(\sum X)^2}{n}\right\} = (n-1)(\mu'^2_2 - \mu'^2_1)$$

$$E\left\{\sum (X^2 - \bar{X}_n^2)\right\} = E\left\{\sum X^2 - \frac{(\sum X^2)^2}{n}\right\} = (n-1)(\mu'^4_4 - \mu'^2_2)$$

$$E\left\{\sum (X - \bar{X}_n)(X^2 - \bar{X}_n^2)\right\} = E\left\{\sum X^3 - \frac{(\sum X)(\sum X^2)}{n}\right\} = (n-1)(\mu'_3 - \mu'_1 \mu'_2)$$

$$E(C_{11}) \approx \frac{\text{分子の期待値}}{\text{分母の期待値}}$$

$$\approx \frac{1}{(n-1)} \frac{(\mu'^4_4 - \mu'^2_2)}{\mu_2(\mu'^4_4 - \mu'^2_2) - (\mu'_3 - \mu'_1 \mu'_2)^2}$$

$$E(C_{22}) \approx \frac{1}{(n-1)} \frac{\mu_2}{\mu_2(\mu'^4_4 - \mu'^2_2) - (\mu'_3 - \mu'_1 \mu'_2)^2}$$

$$E(C_{12}) \approx \frac{1}{(n-1)} \frac{(\mu'_3 - \mu'_1 \mu'_2)}{\mu_2(\mu'^4_4 - \mu'^2_2) - (\mu'_3 - \mu'_1 \mu'_2)^2}$$

ここで

$$\mu_2 = \mu'^2_2 - \mu'^2_1$$

$$V(Y) \approx E\left\{(\bar{y}_n - \alpha) + b_n(\bar{X}_N - \mu'_1) - c_n(\bar{X}_n - \mu'_1) + c_n(\bar{X}_N^2 - \mu'^2_2) - c_n(\bar{X}_n^2 - \mu'^2_2)\right\}^2 \approx \frac{\sigma^2_Y}{n} + \left(\frac{1}{n} + \frac{1}{n}\right) \left\{\mu_2(\sigma_{22} + \beta^2) + (\mu'^4_4 - \mu'^2_2)(\sigma_{cc} + \delta^2) + 2(\mu'_3 - \mu'_1 \mu'_2)(\sigma_{bc} + \beta\delta)\right\} \dots (3.62)$$

研究中、親切な援助と助言を与えられた C. R. Rao 博士にお礼を述べ、さらに原稿の作成と有益な示唆を与えられた Matkhai 氏に謝意を表す次第である。

附 録 A

《 S_{ij} 》を各変量正規母集団の標本分散行列とし、《 C_{ij} 》を対応するその逆行列とする。このとき次の結果を得る。

$$\left. \begin{aligned} E(C_{ij}) &= \frac{\sigma_{ij}}{n-k-2} \quad (i, j) \\ E(C_{ii}) &= \frac{\sigma_{ii}}{n-k-2} \end{aligned} \right\} \quad (A.1)$$

(証明) 行列《 C_{ij} 》の逆元 (右から出発して

$$E(C_{kR}) = E(S_{kR}) = E \left\{ \begin{vmatrix} S_{11} & S_{12} & \dots & S_{1, k-1} \\ S_{21} & S_{22} & \dots & S_{2, k-1} \\ \dots & \dots & \dots & \dots \\ S_{k-1,1} & S_{k-1,2} & \dots & S_{k-1, k-1} \end{vmatrix} \div \Delta \right\}$$

ここで

$$\Delta = \begin{vmatrix} S_{11} & S_{12} & \dots & S_{1k} \\ S_{21} & S_{22} & \dots & S_{2k} \\ \dots & \dots & \dots & \dots \\ S_{k1} & S_{k2} & \dots & S_{kk} \end{vmatrix}$$

よって

$$E(C_{kR}) = \frac{1}{\pi} E \left\{ \frac{(t_{11} t_{12} \dots t_{k-1, k-1})^2}{(t_{11} t_{22} \dots t_{kR})^2} \right\}$$

$$= \frac{1}{\pi} E \left(\frac{1}{t_{kR}} \right)$$

ここでもは直円座標 (rectangular coordinate) を表わす (Mahalanobis, Bose and Roy (1937))
他の t を積分してつたときの t_{kR} の分布は

$$\text{Const } e^{-\frac{\pi}{2} \left\{ T_{kR}^2 t_{kR}^2 \right\}} (t_{kR})^{\frac{\pi-k-1}{2}} dt_{kR}$$

で、この積分は

$$\text{Const} \int_{-\infty}^{\infty} e^{-\frac{\pi}{2} \left\{ T_{kR}^2 t_{kR}^2 \right\}} (t_{kR})^{\frac{\pi-k-1}{2}} d(t_{kR}^2) = 1$$

と与えられる。

$$\therefore E \left(\frac{1}{t_{kR}} \right) = \text{Const} \int_{-\infty}^{\infty} e^{-\frac{\pi}{2} \left\{ T_{kR}^2 t_{kR}^2 \right\}} (t_{kR})^{\frac{\pi-k-4}{2}} d(t_{kR}^2)$$

$$= \left(\frac{\pi}{2} T_{kR} \right)^{\frac{\pi-k}{2}} \pi \left(\frac{\pi-k-2}{2} \right)$$

$$= \frac{\pi \left(\frac{\pi-k}{2} \right)}{\left(\frac{\pi}{2} T_{kR} \right) \left(\frac{\pi}{2} T_{kR} \right)^{\frac{\pi-k-2}{2}}}$$

$$= \frac{\pi}{\pi-k-2} T_{kR}$$

$$= \frac{\pi}{\pi-k-2} \sigma_{kR}^2$$

すなわち

$$E(C_{kR}) = \frac{\sigma_{kR}^2}{\pi-k-2}$$

で一般に

$$E(C_{ij}) = \frac{\sigma_{ij}^2}{\pi-k-2}$$

次に

$$E(C_{k, k-1}) = E \left\{ \begin{vmatrix} S_{11} & \dots & S_{1, k-2} & S_{1k} \\ S_{21} & \dots & S_{2, k-2} & S_{2k} \\ \dots & \dots & \dots & \dots \\ S_{k-1,1} & \dots & S_{k-1, k-2} & S_{k-1, k} \end{vmatrix} \div \Delta \right\}$$

$$= -\frac{1}{\pi} \left\{ \frac{t_{k-1, k} (t_{11} t_{22} \dots t_{k-2, k-2}) (t_{11} t_{22} \dots t_{k-1, k-1})}{(t_{11} t_{22} \dots t_{kR})^2} \right\}$$

$$= \frac{1}{\pi} E \left\{ \frac{t_{k-1, k}}{t_{k-1, k-1} t_{kR}^2} \right\}$$

$$A_{11} = \frac{\pi}{2} T_{kR} \quad A_{12} = A_{21} = \frac{\pi}{2} T_{k-1, k}$$

$$A_{22} = \frac{\pi}{2} T_{k-1, k-1} \quad t_{kR} = X, \quad t_{k-1, k} = X_1, \quad t_{k-1, k-1} = X_2$$

と置く。

このとき t_{kR} , $t_{k-1, k}$ および $t_{k-1, k-1}$ の同時分布は

$$C \left\{ \exp -\frac{\pi}{2} \left[T_{kR}^2 t_{kR}^2 + (T_{k-1, k-1}^2 t_{k-1, k-1}^2 + 2 T_{k-1, k}^2 t_{k-1, k} t_{k-1, k-1} + T_{k-1, k}^2 t_{k-1, k}^2) \right] \right\} \times$$

$$t_{kR}^{n-k-1} t_{k-1, k}^{n-k-2} dt_{kR} dt_{k-1, k-1} dt_{k-1, k}$$

と与えられる。ここで積分 C は

$$C \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp -\frac{\pi}{2} \left[\dots \right] t_{kR}^{n-k-1} t_{k-1, k}^{n-k-2} dt_{kR} dt_{k-1, k-1} dt_{k-1, k} = 1$$

あるいは t_{kR} で積分して

$$C_1 \frac{\Gamma(\frac{n-k}{2})}{A_{11}^{\frac{n-k}{2}}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\{-A_{22}x_2^2 + 2A_{21}x_2x_3 + A_{11}x_3^2\} x_2^{k-1} dx_2 dx_3 = 1$$

すなわち

$$C_1 \frac{\Gamma(\frac{n-k}{2})}{A_{11}^{\frac{n-k}{2}}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\{-A_{11}(x_3 + \frac{A_{21}}{A_{11}}x_2)^2 + A_{33}x_2^2\} x_2^{k-1} dx_2 dx_3 = 1 \quad \dots (A.11)$$

で与えられる。

$$\text{そこで } A_{33} = \frac{A_{11}A_{22} - A_{21}^2}{A_{11}}$$

$$y_2 = x_3 + \frac{A_{21}}{A_{11}}x_2$$

$$y_3 = x_2$$

よおけは

$$\frac{\partial(y_2 - y_3)}{\partial(x_2, x_3)} = 1$$

よって (A.11) は

$$C_2 \frac{\Gamma(\frac{n-k}{2})}{A_{11}^{\frac{n-k}{2}}} \frac{\Gamma(\frac{1}{2})}{A_{11}^{\frac{1}{2}}} \frac{\Gamma(\frac{n-k+1}{2})}{A_{33}^{\frac{n-k+1}{2}}} = 1$$

すなわち

$$C_2 = \frac{\begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} \frac{n-k+1}{2}}{\Gamma(\frac{n-k}{2}) \Gamma(\frac{1}{2}) \Gamma(\frac{n-k+1}{2})}$$

よなる。

$$\text{いま } E\left\{-\frac{t_{k-1, k}}{t_{k-1, k-1}, t_{k, k}}\right\} = E\left\{-\frac{x_3^2}{x_2 x_1^2}\right\}$$

$$= C_1 \frac{\Gamma(\frac{n-k-2}{2})}{A_{11}^{\frac{n-k-2}{2}}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-A_{11}y_2^2 - A_{33}y_3^2} y_3^{n-k-1} \left(\frac{A_{21}}{A_{11}}y_2 - y_3\right) dy_2 dy_3$$

$$= C_2 \frac{\Gamma(\frac{n-k-2}{2})}{A_{11}^{\frac{n-k-2}{2}}} \left\{ \frac{A_{21}}{A_{11}} \frac{\Gamma(\frac{1}{2}) \Gamma(\frac{n-k+1}{2})}{A_{11}^{\frac{1}{2}} A_{33}^{\frac{n-k+1}{2}}} + 0 \right\}$$

第2項は奇函数で範囲は $-\infty$ から ∞ までであるから

$$= C_2 \frac{\Gamma(\frac{n-k-2}{2})}{(A_{11}A_{33})^{\frac{n-k+1}{2}}} \Gamma(\frac{1}{2}) \Gamma(\frac{n-k+1}{2}) A_{21}$$

$$= \frac{\Gamma(\frac{n-k-2}{2})}{\Gamma(\frac{n-k}{2})} A_{21}$$

$$= \frac{\pi}{n-k-2} T_{k-1, k}$$

$$= \frac{\pi}{n-k-2} \sigma_{k, k-1}$$

$$\therefore E(C_{k, k-1}) = \frac{\pi k k - 1}{n-k-2}$$

よって一般に

$$E(C_{ij}) = \frac{\sigma_{ij}}{n-k-2} \quad (i=1, 2, \dots, k; j=1, 2, \dots, k)$$

附 録 B

多変量正規分布に従うとき、 $x_i, i=1, 2, \dots, k$ の母に
対する回帰係数 (b_{ij}) の同時分布は

$$\frac{\Gamma(\frac{n}{2}) |\sigma_{ij}|^{\frac{n-1}{2}} (\sigma_{ij})^{\frac{k}{2}} \prod_{i=1}^k d b_{ij}}{\Gamma(\frac{n-k}{2}) \pi^{\frac{k}{2}} |\sigma_{ij} + (b_{ji} + (b_{ji} - \beta_{ji})(b_{ji} - \beta_{ji}))|^{\frac{n}{2}}}$$

で与えられる。

証明、各 $x_i (i=1, 2, \dots, k)$ が一定の場合 b_{ij} の多変分布は

$$\frac{|\Sigma_{ij}|^{\frac{1}{2}}}{(\sigma_{ij} \sqrt{2\pi})^k} e^{-\frac{1}{2\sigma_{ij}^2} \sum_{i=1}^k \sum_{j=1}^k S_{ij} (b_{ji} - \beta_{ji})(b_{ij} - \beta_{ij})} \frac{d b_{ij}}{d b_{ij}}$$

$$\left[\text{ただし } \sigma_{YX}^2 = \frac{1}{\sigma_{YY}} \right]$$

で与えられる。

一方Xの分布は Wishart

$$\frac{1}{2^{(n-1)k/2} \pi^{k(k-1)/4}} \frac{1}{\Gamma(\frac{n-k}{2})} e^{-\frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \sigma_{ij}^{-1} S_{ij}} \frac{\pi^{-k/2} \prod_{i=1}^k \Gamma(\frac{n-k+1}{2})}{|S_{ij}|^{\frac{n-k}{2}}} dS_{ij}$$

に従う。よってXが変動するときの b_{Yi} の分布は

$$\text{Const} \prod_{i=1}^k db_{Yi} \left\{ \dots \int e^{-\frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \left\{ \sigma_{ij}^{-1} + \frac{(b_{Yi} - \beta_{Yi})(b_{Yi} - \beta_{Yj})}{\sigma_{YX}^2} \right\} S_{ij}} |S_{ij}|^{\frac{n-k}{2}} \prod_{i=1}^k \prod_{j=1}^k dS_{ij} \right. \\ \left. = C \frac{\prod_{i=1}^k db_{Yi}}{|S_{ij}|^{\frac{n-k}{2}}} \right.$$

で与えられる。

$$\begin{aligned} \sigma_{ij}^{-1} &= \sigma_{ij}^{-1} + \frac{(b_{Yi} - \beta_{Yi})(b_{Yj} - \beta_{Yj})}{\sigma_{YX}^2} \\ &= \sigma_{ij}^{-1} + (b_{Yi} - \beta_{Yi})(b_{Yj} - \beta_{Yj}) \sigma^{YY} \end{aligned}$$

また

$$C = \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-k}{2})} \frac{|\sigma_{ij}^{-1}|^{\frac{n-1}{2}} (\sigma^{YY})^{\frac{k}{2}}}{\pi^{\frac{k}{2}}}$$

附 録 C

X_1, X_2, \dots, X_k が多変量正規分布に従うときの $b_{k, k-1, 1, 2, \dots, k-2}$ (すなわち X_{k-1} の X_k に対する回帰係数) の分布は

$$\frac{1}{a \beta(\frac{n-k+1}{2}, \frac{1}{2})} \left\{ 1 + \frac{(b_{k, k-1, 1, 2, \dots, k-2} - \beta_{k, k-1, 1, 2, \dots, k-2})^2}{a^2} \right\}^{-\frac{n-k+1}{2}} \times db_{k, k-1, 1, 2, \dots, k-2} \quad \dots (C.1)$$

$$\text{ここで } a = \frac{\begin{vmatrix} \sigma_{kk} & \sigma_{k, k-1} \\ \sigma_{k-1, k} & \sigma_{k-1, k-1} \end{vmatrix}}{\sigma_{kk}}$$

証明 $b_{k, k-1, 1, 2, \dots, k-2} = \frac{t_{k-1, k}}{t_{k-1, k-1}}$

$t_{k-1, k-1}$ および $t_{k-1, k}$ の同時分布は

$$\text{Const. exp} \left\{ -\frac{n}{2} \left[T^{-1} t_{k-1, k-1}^2 + 2T^{-1} t_{k-1, k-1} t_{k-1, k} + T^{-1} t_{k-1, k}^2 \right] \right\} t_{k-1, k-1}^{-n} dt_{k-1, k-1} \times dt_{k-1, k} \dots (C.11)$$

$$u = \frac{t_{k-1, k}}{t_{k-1, k-1}}, \quad v = t_{k-1, k-1} \quad t_{k-1, k-1}$$

とかくと変換の Jacobian は

$$\frac{\partial(u, v)}{\partial(t_{k-1, k-1}, t_{k-1, k})} = 2 \frac{t_{k-1, k}}{t_{k-1, k-1}} = 2u$$

となる。

(C.11) から u と v の同時分布は

$$\text{Const. exp} \left\{ -\frac{nv}{2} \left[T^{-1} u^2 + 2T^{-1} u + \frac{T^{-1}}{u} \right] \right\} v^{\frac{n-k+1}{2}} u^{-\frac{n-k+1}{2}} du dv \dots (C.12)$$

なることがわかる。

u について 0 から ∞ まで積分し、 $T^{-1} = \sigma_{ij}$ および

$$E(b_{k, k-1, 1, 2, \dots, k-2}) = \beta_{k, k-1, 1, 2, \dots, k-2} = \frac{\sigma_{k, k-1}}{\sigma_{kk}}$$

なることに注意すると、 u の分布は

$$\text{Const.} \cdot \left\{ 1 + \frac{(b_{k, k-1, 1, 2, \dots, k-2} - \beta_{k, k-1, 1, 2, \dots, k-2})^2}{a^2} \right\}^{-\frac{n-k+1}{2}} db_{k, k-1, 1, 2, \dots, k-2}$$

となる。ここで

$$a = \frac{\begin{vmatrix} \sigma_{kk} & \sigma_{k, k-1} \\ \sigma_{k-1, k} & \sigma_{k-1, k-1} \end{vmatrix}}{\sigma_{kk}} \cdot \frac{1}{\sigma_{kk}}$$

で、積分は $-\infty$ と ∞ の範囲で積分すると

$$a \beta\left(\frac{n-k+1}{2}, \frac{1}{2}\right)$$

なることがわかる。

by K.C. Seal (Sankhya vol. 11, 1951)

文 献

Bose, C. (1933) Note on sampling error in the method of double sampling. *Sankhyā* 6. 329-30

Bose C. and Gayen A. K. (1946) Note on the expected discrepancy in the estimation (by double sampling of a variate in terms of a concomitant variate when there exists a nonlinear regression between the two variates. *Sankhyā* 8, 73-74

Cochran, W. G. (1939) The use of analysis of variance in enumeration by sampling. *J. A. S. A.* 34, 490-96

Ghosh, B. (1947). Double sampling with many auxiliary variables. *Cal. stat. Assoc. Bull.* 1. 91-93

Mahalanobis, P. C. Bose, R. C. and Roy, S. N. (1937): Normalization of statistical variates and the use of rectangular coordinates in the theory of sampling distributions *Sankhyā* 3. 8-40

Neyman, J. (1938): Contribution to the theory of sampling human populations *J. A. S. A.* 33, 101-116

Schumacher, F. X. and Chapman, R. A. (1942): Sampling methods in forestry and range management (186-189)

Snedecor, G. W. and King, A. J. (1942) Recent development in sampling for agricultural statisticians. *J. A. S. A.* 37. 99-100

8 プロットの大きさにとむなう

収量の分散の変動について

On the variation of yield variance with plot size

研究した問題は、プロットの大きさおよび形状にとむなう収量分散の変動状態の把握から、収量密度の空間的共分散函数 (Spatial covariance function) を求めることである。収量のプロットについて得られた結果を3に示してある。またプロットの幾何学的形状に対する収量分散の依存関係についても、非常に小さいプロットおよび非常に大きいプロットの結果を4で得られた。共分散が距離の大きいところでべき則 (Power law) に従う場合に特別な注意を払った。

1 緒 言

よく知られているように収量分散の観測された変動をプロットの大きさおよび形状について説明するためには、プロット内の任意の二点の収量密度の間に相関が存在するという可能性を考慮することが必要である。(我々は均等と定常的 (stationary) な場合だけに限定する。この場合には収量密度の期待値はその地域内で一定である。) 更にこの空間的相関は幾々二点間の距離が増すにつれて、比較的ゆるやかに即ち指数的というよりは距離の巾函数として減少しなければならないように思われる。

同様な性質は約き糸の直径、洪水位 (flood height) (Feller, 1951) および母集団標本からの偏差などの観測で示される。この論文の計算はこれらの場合にもあてはまるものであるが、具体性をわたせるため我々は引続きプロットおよび収量の問題で話をすることにする。(たとえ時として二次元の場合だけに限定されないことを示すために、"プロット"の代りに"領域 (region)"なる用語を用いるとしても)。

最も普通に行なわれる計算は、一定の大きさおよび形状をもつプロット

トおよび与えられた空間的共分散函数 (Spatial covariance function) $\rho(S)$ に対して収量の分散を評価するものである。しかし一般にはおそらくプロットの幾学的形態の函数としての収量分散の知識から Covariance function を決定する筈の計算の方が立ちそうに思われる。この従な方法によれば実験結果を用いてかくとし $\rho(S)$ の部分的推定を得ることができよう。我々はこの問題を3で考えることにする。

同例の一つは Mellin 変換を用いて容易に解くことができる。従つて距離の中函数に従つて減少する共分散

$$\rho(S) \sim C|S|^{-\gamma} \quad (S \text{ が大のとき}) \quad \dots\dots (1)$$

は昔芳なしに得られる。

観測される多くの共分散函数が (1) の形であることを示す証拠が数多く存在する。(Fairfield Smith, 1938)。収量変動に対するより単純な線型模型は市法則を予想するものでなく、むしろ一層急速に減少する

$$\rho(S) \sim C|S|^{-\gamma} e^{-\lambda|S|} \quad (S \text{ が大のとき}) \quad \dots\dots (2)$$

なる型のものを要求するという意味においてこれらは一層興味をいくものである。我々はこの問題を34において更に追求するであろう。

一次元の大きい区間での“収量”の分散はその共分散函数が (1) の形なら漸近的に区間の長さの中にも比例することが容易に証明できる。しかし二次元あるいはそれ以上の次元で類似する結果を得ることはどう容易ではない。この問題も同様に33および34で取扱うことにする。

3 収量の分散に対する公式

我々は平面上の点を表わすのにベクトル座標を用いる。せうして面積 A 、収量 $y(\Omega)$ のプロットもしくは領域 Ω を考える(平面とか面積などの用語を用いるが、この論文の大部分の公式は n 次元空間 ($n=1, 2, 3, \dots$) にあてはまる)。

最終的には距離 S だけ離れた二点になるように収量する二つの領域 Ω_1 と Ω_2 を考えよう。共分散函数を

$$\rho(S) = \lim_{A_1, A_2 \rightarrow 0} \frac{\text{cov}[y(\Omega_1), y(\Omega_2)]}{A_1 A_2} \quad \dots\dots (3)$$

で定義し、この極限が存在する場合のみを考えることにする。

$y(\Omega)$ の母集団分散を $V(\Omega)$ とすると、(3) の定義から

$$V(\Omega) = \int_{\Omega} \int_{\Omega} \rho(x_1 - x_2) dx_1 dx_2 \quad \dots\dots (4)$$

が得られる。 $\rho(S)$ が原点で連続なら (4) から一つの直接的な結論を導くことができる。

$$V(\Omega) \sim \rho(0)A^2 \quad (\Omega \text{ が小なるとき}) \quad \dots\dots (5)$$

即ち小さい領域(即ちそれらの次元のすべてにおいて小さい領域)では V は面積の二乗に比例する。これはすべての状況について V が面積に比例するとき、空間的相関 (spatial correlation) が 0 の場合と対比さるべきものである。

一次元の場合で Ω が長さの線分からなっているときには (4) は

$$V(a) = 2 \int_0^a (a-x) \rho(x) dx \quad \dots\dots (6)$$

のように簡単になる。

Ω が x 座標軸に平行な a, b なる区間を有する平面上の矩である場合に同様にして

$$V(a, b) = 4 \int_0^a \int_0^b (a-x)(b-y) \rho(x, y) dx dy \quad \dots\dots (7)$$

と取る。(6) の積分は一般的 $\rho(S)$ について計算できるが (7) およびその他の二次元の公式は、普通あまり一般的でない形の函数 $\rho(S)$ に対してしか計算できない。

ここでスペクトル密度函数 (spectral density function)

$$F(u) = \int_{-\infty}^{\infty} e^{iu \cdot x} \rho(x) dx \quad \dots\dots (8)$$

および面積的特性函数 (areal characteristic function)

$$G(u) = \int_{\Omega} e^{iu \cdot x} dx \quad \dots\dots (9)$$

を挿入すれば

(4) 式は

$$V(\Omega) = \frac{1}{(2\pi)^n} \int_{-\infty}^{\infty} |G(u)|^2 F(u) du \quad \dots\dots (10)$$

と書ける。この公式は往々 (4) 式よりも便利である。というのは

積分の速度が簡単で種々の波数 (wave number) w に対して与えられる重み (weights) が全くはつきりしており、steepest descents法のような近似法を直接用いることができるからである。簡単な形状のプロットに対する函数 $G(w)$ は直ちに計算できる。即ち a, b をもつ矩形に対しては

$$G(w) = \frac{4 \sin(a w_1) \sin(b w_2)}{w_1 w_2} \dots \dots \dots (11)$$

で、半径 a 、中心が原点にある円形に対しては

$$G(w) = \frac{2 \pi a}{w} J_1(a w) \dots \dots \dots (12)$$

である。但し J_1 は order 1 の Bessel 函数である。また w_1, w_2 は w の成分で、 w はその絶対値である。

3 逆の問題 (The inverse problem)

実用上からは、(4) の関係を逆にして covariance function を与えられた Ω について観測できると思われる収束分散 $V(\Omega)$ で表わせれば有用である。この逆関係は、線型および矩形プロットに対応する (5) および (6) の特別な場合には非常に簡単に求められる。これらはそれぞれ、

$$p(x) = \frac{1}{2} \frac{\partial^2 v(x)}{\partial x^2} \dots \dots \dots (13)$$

$$p(x, y) = \frac{1}{4} \frac{\partial^4 v(x, y)}{\partial x^2 \partial y^2} \dots \dots \dots (14)$$

である。しかしこれらと同じように取扱いやすい場合が少しある。ここでは主として一般的な型の一つに問題を限定することにする。即ち共分散は等方性 (isotropic) であると仮定する。従って $p(s)$ は s のみの函数であるから

$$p(s) = \bar{p}(s) \dots \dots \dots (15)$$

と書ける。相似的にのみ (即ちそのすべての大きさを同じ比で変化させることによつて変化することのできる一定の形をしたプロット) についてのみを取扱うものと仮定する。与えられたプロットの形状

についてはプロットの大きさをその最も大きい寸法即ち x で規定すると便利である。従つて我々は

$$V(\Omega) = V(x) \dots \dots \dots (16)$$

と書くことにする (このように円形プロットでは x は直径に等しくまた矩形プロットでは対角線に等しい)。この場合プロットのそれそれの形状に対応して、プロット内で無作為に送んだ二点間の距離の分布を表わす次の逆函数 $K(s)$ が存在する。

即ち
$$V(x) = \int_0^1 K(s) p(s) ds \dots \dots \dots (17)$$

である。こゝでもしすべての一次元の寸法を $x:1$ に増加させれば、この積分の無限小要素は、 $x^2:1$ の割合で増加するであろうから全く一般的に

$$V(x) = x^{2n} \int_0^1 K(s) p(xs) ds \dots \dots \dots (18)$$

V および p の間の関係 (18) は Mellin 変換を適用できる形であるから

$$\bar{V}(\gamma) = \int_0^\infty V(s) s^{\gamma-1} ds \dots \dots \dots (19)$$

$$\bar{p}(\gamma) = \int_0^\infty p(s) s^{\gamma-1} ds \dots \dots \dots (20)$$

$$\bar{K}(\gamma) = \int_0^\infty K(s) s^{\gamma-1} ds \dots \dots \dots (21)$$

これを用いると (18) は
$$\bar{V}(\gamma - 2n) = \bar{K}(\gamma - \gamma) \bar{p}(\gamma) \dots \dots \dots (22)$$
 となる。

我々は §4 において変換 (19) および (20) が存在し、かつ (22) が成立つような γ の共通的な範囲を考慮することにする。暫らくの間適当な範囲が存在すると仮定して、(22) で与えられる \bar{p} の式の逆変換をとることによつて $p(s)$ の解を直接求めることができる。実際 (18) に対する希望する反転関係 (inverse relation) が

$$p(x) = x^{-2n} \int_0^1 L(s) V(xs) ds \dots \dots \dots (23)$$

なら
$$\bar{p}(\gamma) = \bar{L}(\gamma - 2n) \bar{V}(\gamma - 2n) \dots \dots \dots (24)$$

である。ここで $\bar{L}(Y)$ は $\bar{K}(Y)$ と類似をもつものである。(22)と(24)を比較すれば

$$\bar{L}(Y) = \frac{1}{\bar{K}(Y - n\pi)} \dots\dots\dots (25)$$

であることがわかる。

計算を正確に実行し得るためには(21)の積分を求めることができかつ(25)で与えられる $\bar{L}(Y)$ の反転 (invert) できねばならない。これが可能となる小数の場合の一つはどちらかというといふ非現実的な円形プロットの場である ($n = \text{直径}, \pi = 2$)。円については直接法から

$$\bar{K}(s) = \frac{1}{2} [s \cos^{-1}(s) - s^2 \sqrt{1-s^2}] \dots\dots\dots (26)$$

を得る。但し

$$\bar{K}(Y) = \frac{\Gamma(\frac{1}{2}) \Gamma\{\frac{1}{2}(Y+2)\}}{4(1+Y) \Gamma\{\frac{1}{2}(Y+5)\}} \dots\dots\dots (27)$$

$$\bar{L}(Y) = \frac{4(-3) \Gamma\{\frac{1}{2}(Y+1)\}}{\Gamma(\frac{1}{2}) \Gamma\{\frac{1}{2}(Y-2)\}} \dots\dots\dots (28)$$

しかし $\bar{L}(Y)$ を反転しようとするとき及々は発散積分を得る。これは $P(X)$ は $V(X)$ だけでなくその数係数 $V^{(j)}(X)$ でも表さるべきであるという事実を示すものとして受け取らねばならない。(13)および(14)の結果をみればこのことは驚くには当らない。もし(23)に対して更に一般な関係

$$P(X) = X^{-2\lambda} \int_0^1 \sum_{j=0}^{\infty} L_j(s) [(Xs)^j V^{(j)}(Xs)] ds \dots\dots\dots (29)$$

を代入すれば

$$[s^{Y+j-1} V^{(j)}(s)]_0^1 = O(j-0, 1, 2, \dots) \dots\dots\dots (30)$$

なる限り

$$\bar{L}(Y) = \sum (Y-1)(Y-2)\dots(Y-j) L_j(Y) \dots\dots\dots (31)$$

なることがわかる。

いま(28)の $\bar{L}(Y)$ は

$$\pi^{-1} \Gamma(\frac{1}{2}, \frac{Y-3}{2}) [(Y-1) - (Y-1)(Y-2) + (Y-1)(Y-2)(Y-3)] \dots\dots\dots (32)$$

と替ける。これは(29)および Beta 函数の積分表示を考れば

$$P(X) = \frac{2X^{-\lambda}}{\pi} \int_0^1 [(Xs) V'(s) - (Xs)^2 V''(Xs) + (Xs)^3 V'''(Xs)] \frac{s^{-\lambda} ds}{\sqrt{1-s^2}} \dots\dots\dots (33)$$

なる関係に対応するものである。

実用的に最も興味のあるのは矩形の場合であろう。しかし $K(s)$ はこの形に対して容易に計算できるのに反して $\bar{K}(Y)$ はそうでない。

4. 巾法則に従う共分散 (Power-law Covariances)

まず(22)の関係を打ちかえることにする。もし函数 $f(s)$ が正数 s の小さい値および大きい値に対してそれぞれ $O(s^{-\lambda})$ および $O(s^{-\rho})$ なるその Mellin 変換 $\bar{f}(Y)$ が

$$0 < \text{Re}(Y) < \rho \dots\dots\dots (34)$$

に対して確かに存在する。そうして変換 $\bar{f}(Y)$ は λ および ρ において単純な極 (pole) を有する。

いま原点で $f(s)$ が $O(1)$ (この点で連続と仮定すれば) で、かつ Y の大きい値に対して $O(s^{-\lambda})$ と仮定すれば $\bar{f}(Y)$ は $0 < \text{Re}(Y) < \lambda$ なる範囲で存在する。

小さい s に対しては $K(s)$ は半径 s の n 次元の球の表面積に比例するから、それは $O(s^{n-1})$ である。

一方 $s > 1$ に対しては 0 となる。結局 $\bar{K}(Y)$ は $(1-n) < \text{Re}(Y) < \infty$ に対して存在する。

故に(22)式の右辺の値は二つとも

$$0 < \text{Re}(Y) < \min(n, \lambda) \dots\dots\dots (35)$$

なる範囲で存在する。(22)式は \bar{K} および \bar{P} で \bar{V} を定義しているから、これらの \bar{V} の値についてもこの関係が成立つ。よって $\bar{V}(Y)$ は

$$-2\pi < \text{Re}(Y) < \min(-\pi, \lambda - 2\pi) \dots\dots\dots (36)$$

で存在し、またこれらの \bar{V} の極端な値では単純な極 (Simple pole) をもつ。

故にもし $f(s)$ が $s^{-\lambda}$ あるいは

$$V(X) = O(X^{2\lambda-\lambda}) = O(A^{2-\lambda/\pi}) \dots\dots\dots (37)$$

より急速に減少しまた $f(s)$ が $s^{-\lambda}$ ($\lambda < n$) と同程度に減少する

なら、小プロットの割合

$$V(X) = O(X^{2\pi}) = O(A^2) \dots\dots\dots (38)$$

大プロットでは

$$V(X) = O(X^\pi) = O(A) \dots\dots\dots (39)$$

を得る。言い換えれば、Vが小プロットのプロット面積に比例するかあるいはプロット面積よりも急速に増加するかの二つの場合のけじめをつけるのが巾法測 (power law) $[P(S) \sim (S^{-\pi})]$ である。

(22)の関係は反転に用いることができる。もし $V(X)$ が大面積に対して A^μ と同様に増加することからすれば、大きい距離に対して $P(S)$ が $S^{\pi(\mu-2)}$ と同程度に減少することが結論できる。

注意しておかねばならないことはこれらの条件はプロットの形が一定の場合にのみ成立するという点である。何となれば一般に V は面積の函数であるとともに形の函数でもあるのだから。しかし、P. Fairfield Smith (1938) によって与えられた幾つかの結果は、この形に対する依存関係は小さくて、 V は殆んどプロット面積だけから決定し得ることを示しているように思われる。このことは多分プロットの形が特に極端な場合、例えばあまり細長くない狭な場合にのみ正しいであろう。Fairfield Smith の結果はまた (37) 式に関連する性質について非常に納得のゆく証明を与えているという意味で興味をひくものである。彼は面積 A のプロット内の単位面積当り収量の分散が曲線 $A^{-0.74P}$ で非常によく近似できることを見出した。指数を $-\frac{3}{4}$ でおきかえると共分散函数

$$f(S) \sim \text{const} S^{-\frac{3}{2}} \dots\dots\dots (40)$$

に対応する

$$V(A) \sim \text{const} A^{2-\frac{3}{2}} = \text{const} A^{\frac{5}{2}} \dots\dots\dots (41)$$

を得る。(勿論これらの等式は無限に小さい A および S については成立しない)。

観測においてよく見られるこの結果は、三つの興味ある可能性の根拠を与えるものである。即ち

(a) 共分散は $S^{-\frac{3}{2}}$ に従って減少するという性質をもっている。

(b) 収量の分散をプロット面積以上の速さで減少させる理論率が小である。

(c) 観測された指数は簡単な有理数である。

これまで上記の型の共分散を導びくようなある地域上の収量の変動を表わす簡単な線型の模型が提出されたことは一度もない。

(Whittle, 1954; Hewson, 1955 を見よ)。例えば、 (X, Y) のまわりのすべての点の収量に対して均等に収量密度 (yield density) $\Sigma(X, Y)$ を関係づける模型

$$\left\{ \left(\frac{\partial}{\partial X} \right)^2 + \left(\frac{\partial}{\partial Y} \right)^2 - X^2 \right\} \Sigma(X, Y) = E(X, Y) \dots\dots (42)$$

は共分散函数

$$f(S) = \text{const} \cdot SK_1(HS) \sim \text{const} \cdot S^{\frac{1}{2}} e^{-KS} \quad (X, S \text{ が大きいとき}) \dots\dots (43)$$

を導びく。(ここで K_1 は修正された Bessel 函数を表わす)

折糸の直径の変動に対する一次元の問題で、D. R. Cox (未発表と思われる研究についての私信で) は大きい S に対して事実上巾法測 (power law) を予想する模型を提案している。

多次元の模型で任意の範囲の S に対して巾法測の共分散を予想する着目の知っているのは騒音 (turbulence) のだけである (Batchelor, 1953, p. 122 をみよ)。これらの模型では次元の試験から指数は簡単な有理数値でなければならぬことが示される。

騒音の見解では農業研究で観測される巾法測に従う減少を十分説明できる模型は、Cox の折糸の模型および騒音の模型に共通する二つの性質を備えていなければならないように思われる。

即ち

(a) 線型でないこと。

(b) 収量 (収量密度 (yield density) は空間的座標 (spatial co-ordinates) の函数であると同時に、時間的函数としても考えなければならないこと。

である。(Boxの場合では“時間”は紡ぎ糸の受けた smoothing の数で不連続である。)

このような模型としては、土壌の傾向 (soil trend) における肥沃度の勾配 (fertility gradients) は時の経過とともに拡散過程 (diffusion process) のために平滑化 (Smoothing) されてしまうことをあげ得るであろう。この過程は一定の値以上の勾配のみが減少するという範囲内では非線型である。

指数入の値については当然検討しなければならない。Fairfield Smith の推定した指数 $\lambda = 0.747$ が有理数 $\frac{3}{2}$ ($\lambda = \frac{3}{2}$ に対応する) に非常に近いということは偶然であろうか? 確かにこの一例だけから結論を出すのは早計にすぎる。しかし Fairfield Smith は、色々な種類の作物の斉一性試験 (uniformity trial) のテータから推定した λ の値をあげている。推定された 39 個の λ の値のヒストグラムでは主な傾点か $\lambda = 0.41 - 0.50$ の間があり、また λ の傾点および上の打切り点 (cut-off point) は $\lambda = 0.71 - 0.80$ にある。このことは少くとも $\lambda = \frac{1}{2}$ および $\frac{3}{4}$ なる値 ($\lambda = 1, \frac{3}{2}$ に対応する) はある方法で区別されることを示している。

更に多数のテータあるいは理論からの推測が与えられない限りおそらくこれ以上のことは言いえないであろう。

by P. Whittle

参 考 文 献

Batchelor, G. K (1953) *The theory of homogeneous turbulence* Cambridge Univ. press.

Fairfield, Smith H (1938) *J. Agric. Sci.* 28, 1-23

Feller (1951) *The asymptotic distribution of the range of independent random variables* Ann, Math. Stat. 22, 427-32

Hoine, V. (1955) *Models for two-dimensional stationary stochastic processes* Biometrika, 42, 170-8.

Whittle, P. (1954), *on stationary process in the plane.* Biometrika, 41, 434-49

9 系統的抽出研究の展望

A Review of the Literature of Systematic Sampling

要 約

系統的抽出に関する文献は、実際面の必要性とそれらの充足 (satisfaction) に関する統計的手法の供給との相互作用を示している。系統的抽出の問題は一般に認められているよりはるかに屢々生じ、また統計的手法の多くも、満足という点から程よいということから、この争論は一層の進歩に対する強い刺激となっている。

結 言

この展望の目的は系統的抽出に関する主要論文の見直しを与えることである。この分野におけるサンプリングの大部分の進歩における現代的な性格に照らして、この展望は大体年代順にまとめている。これによって発表当時の実際的問題の解決を助けるために与えられた手法の、要求に対する相互関係を調べるのが可能になる。この点よりすれば、西子分析の方法が心理学者の周辺から発展してきたのと全く同様に、系統的抽出法が林業および土地利用に関する問題から発展してきたことは注目し得る。Stephen の論文 (1948) はそれ自身が展望の形になっているから、最初に取り上げるのに都合がよい。

この論文のオズ部における二つの文章は重要だから次に引用する。
1. "現代のサンプリングの方法は、個体を無作為かあるいはある種の系統的な手続によって送り出すプロセスの上に成立している"

2 系統的抽出は単純な規則で場合をある便宜的な順序で数え、標本としては又番目ごとの '場合'、または測定された間隔でのなにか同様なパターンを用いて '場合' を選出する。
ある種の物理的および工学的科学においては、系統的抽出の形は

長い間用いられてきた。他の応用は林業、農業、畜産、および気象学においてなされてきた。しかし経済および社会学の分野の研究の場合 - 例えば Bowley, Caradog, Jones, Hilton および Kaiser の - には、標本抽出に用いるリストが前もつてある順序になっていて、その大部分はランダムかあるいは少なくとも全体のある division 内にあるということが根本的な特徴となっている。一方他の分野における対象は、確かにランダムでない順序に配列されていることがわかつている。したがって同様な標本抽出の手続によって現れる二種類の系統的抽出が存在することになるが、この二つの解決は全く異なるものである。経済の分野において Bowley や他の人々の用いた方法は屢々単一無作為といわれているもので、むしろ層化無作為抽出といつた方がよい。この用語の一般的なことは、この展望の最後において簡単にふりかえることにするものの一つである。

Hasel (1938) は林業における数種のサンプリングの文献をまとめているが、その中で標本内に抽出誤差を推定するのに必要な情報が含まれるためのよくいわれている条件を述べている。しかし生計する林木の単例 '毎木調査 (Censusing)' によって調査され、サンプリングは同時に行なわれる。この標本は普通ある形式の系統的送らえらるから、抽出誤差の正しい推定のための条件を欠いている。この困難を解決するため Hasel は補助標本を無作為に選ぶことを提案した。しかし彼はこの方法を実行するには費用がかかることになっている。その困難にもかかわらず、彼は管理面および回収の際に非常に都合がよいということから、系統的標本を支持しているが、系統的標本の誤差に対する方法を発展させることが實際上極めて重要であると指摘している。

次の二つの論文は土地利用の分野に関する研究である。Pearlfoot (1942) の論文は、土地利用型の分布の標本データが満足すべきものであるかどうかを推定する問題に 'Traverse Survey' の方法を適用した。彼は全 traverse intercept が大きいと traverse line の間隔が狭くなるために、平均平方誤差の平方根が小さくなること

を確かめた。これは形式的には抽出比の増加による標本サイズの増大と同じことである。外業調査テーカーでチェックされたのと同じ色々な間隔の *traverse line* による推定値間の相関から、色々な層数が計算された。これらは "traverse サンプルング" の正確さを予感するための *traverse* 精度図表の最初のものとなった。

Osborne (1942) の論文はその用語の層数密な解釈の上で立つて、系統的抽出の問題と取り組んだ最初の試みと思われる。色々な種類の林地によって被覆される面積を推定するための系統的ラインサンプルングの背景に対して、彼の論文は無作為および系統的に選ばれた標本からえられた標本推定値の正確度を報告し、後者の抽出誤差を推定する方法を与えている。彼は系統的にとられたテーカーを無作為標本の公式で評価すれば、抽出誤差の推定値には必ずず偏りが生ずると結論している。これとともに地域をブロックに細分するという解決方法もとれなかつた。というのは、このようにしても齊一な層の組を作ることができなかつたからである。実際、変数は位置の連続函数と考えられる変化を受けているのである。しかつてサンプルングの問題は言葉のあてはめに耐するものとなり、観測された値の間に存在するかもしれない内部的な相関を考慮しなければならぬことになる。

よつてこの初期の段階では、系列相関が問題の中心的部分を占めていて、空間または時間において順序づけられた系列の一次的分析につながっているということが明らかにされたのである。系統的抽出法の近代的な解釈を過去の散漫なる解釈と区別する問題はこの見方に存在するのである

理論の発展

1944年に系統的抽出法を取扱った最初の論文 (Madow and Madow, 1944) が発表された。それ以来望君 William and Silian Madow の名前は引続く理論の発展の各段階を通じて一まわり目立つていた。この論文において引用されたそれ以前の文献が

Hansen と Hurwitz (1943) による有限母集団からのサムプリングの研究だけであつたことは興味がある。

この論文は単一要素の層化および無作為母集団からの系統的抽出を論じている。後の論文 (Madow, 1949) では大きさが等しい場合、等しくない場合の *cluster* (集団) に対する理論の比較を取扱つていて、この最初の論文では "サンプルング理論の大きい欠陥は、完全な無作為標本をとるべきかをきめる同等の統計的方法のないことである" と考えられている。系統的抽出の手続は、多くの可能な選び方を特に除外した無作為抽出の手続として定義されている。もう一つの面からは、系統的抽出計画は集団抽出の一形式であるが、母集団における要素の順序の知識が、級内相関を求めるのに用いられるところが異なつていて、このために Hansen and Hurwitz の研究が引用されている。このようにして、母集団分散およびある種の系列相関の知識から、系統的標本の抽出誤差を評価できる。3つの基本的な結果が述べられている。

- 1. 系列相関が正の和なら系統的抽出は抽出分散が大きいという意味で無作為抽出より悪い。
- 2. 系列相関が負の和なら系統的抽出は無作為抽出よりも良い。
- 3. 系列相関の和が近似的に0なら、二つの抽出手続の間には大きい違いがない。

有限母集団からの系統的標本の数は小さいから、標本平均値が正偏分散すると仮定することは疑わしい。結局我々は、母集団の要素が、確率変数の単一の観測値であるという考えに従うことになる。この概念と理論の展開に完全に組み入れた Cochran (1946) の論文まで、このような考え方は採択しなかつたことに注意せよ。

一つだけの単位の状態化の系統的抽出の場合、系統的標本の平均値 (\bar{x}) は母集団平均値の無偏推定値であり、その分散は

$$\text{Var}(\bar{x}) = \frac{\sigma^2}{n} (1 + (n-1) \bar{r}_k) = \frac{\sigma^2}{n} \left\{ 1 - \sum_{k=1}^m p_{km} \right\}$$

であることが示された。ここで r_k はそれぞれ k 要素からなるクラス

の数で、 k_m は系列相関のラグ(すれ)である。また ρ_m^2 は母集団分散で、 ρ_m は級内相関係数である。筆者は単一の標本から得られた偏りのある分散および系列相関係数の推定値は、2つ以上の標本を用いることのみによつて修正できることを注意している。このような方法は前述のように *Hasel* (1938) によつても提案されているが、費用がかさむと思われる。(ρ_m の) 推定値 γ_{km} の偏りは無視できる程度であることを示して、 σ^2 の推定を解決する方法が幾つか与えられているが、これらはかなり複雑である。

層化された場合の系統的抽出については、二つの基本的な形式が考えられている。

- (i) 抽出比(または間隔)がすべての層で一定のとき。
- (ii) 可変抽出比(または間隔)を用いたとき 標本平均値は層平均値の重みづけ平均であるが、分散についての考察から次のような実際的問題が生ずる。
 - 1 異なる層の対応する *item* が正の相関をもつときは、抽出の効率以外の考慮が重要でない限り、抽出間隔を一定とすることは得策でない (*Madou, W. G. 1949* を参照)。
 - 2 対応する *item* が負の相関をもつときは、抽出間隔を一定とすると抽出分散が小さくなる。
 - 3 可変抽出比を用いた場合には、層間で独立な抽出を行なうと層の対の間の共分散が消えて、平均値の分散の式はよく知られた層化抽出のものとなる。

系統的抽出の分散の無作為抽出のそれに対する比の式が与えられ、 k が k' に比して大きく、かつ、

$$\sum \rho_{km} < - (n-1) / 2 (kn-1) \approx - \frac{1}{2} k$$

なるときには、前者の分散が後者より小さくなることが示された。

次にこの論文はテータの同期性の効果を考察している。右をテータの同期とすると系統的標本の層平均値間の相関は $+1$ であるから、無作為標本のほうがよい。しかし、同期が ± 1 にならざる層また

は位置の差が奇数となるような層については、層間の相関が -1 となり、また位置の差が偶数となる層の前では $+1$ となるから、系統的標本の分散はおそらく小となるであろう。興味ある結果は、母集団の型 (*form*) が線型るとき、系統的標本は系統の傾向 (*Trend*) の影響を取除くという点で無作為標本よりはるかに有効であるが、層化無作為標本よりは劣っているということである。

この論文の簡潔な説明が筆者の一人 (*Madou, S. H. 1946*) によつて与えられている。この研究方法は集落抽出においてはさうにはつきりしている (*Hansen and Hurwitz* をみよ) が、その最も単純なのが、集落内で全然別次抽出を行なわないで、唯一つの集落を抽出するものである。系統的抽出計画を採用するために必要な修正は明瞭に述べられている。この論文に与えられたオスの例は、色々な標本サイズに対する系統的標本の分散が標本サイズの増加にもなつてて、並進に減少しないという意味において *consistent* でないことを示している。したがつて系統的抽出計画の効率は不定なということになる。

理論の次の発展を論ずる前に、3つの論文を検討しておくのが有用であろう。これらはいずれも重要な実際的問題を論じている。

最初の論文 (*Wadley 1945*) は昆虫の母集団の推定—単位面積あたりの数—あけまたは負換体の問題に関するものである。もし全所に空間的な差のあることかわかつていれば、系統的標本は無作為標本より代表性がある。系統的抽出理論の発展の実際面からの必要性と理論的発展の相互作用は次の引用文中に示されている。

「……系統的標本の抽出効率を標準的な方法では正しく推定できない。……そのような標本の分散について光を与える方法は同様なものでも有用である。Madou and Madou はこの問題を論じたが、それらの結果はまだ現場での利用に適当な形式では与えられていない。」

次のものは (*Lleming and Simmons, 1946*) 層内で系統的抽出を行なう層化母集団からの興味あるサムプリングの例である。

統計的な考察がなされていないのは、その時(1946)は理論におけるそのような発展がまだ十分利用できる程になつていなかったことによるものである。この種類に属するオ3の論文(Bayley and Hammersley, 1946)は、1946年の1月に Royal Statistical Society, Research Section でもたれた時系列における自己相関に関するシンポジウムの討論で筆者が与えた寄与から生れたものである。ここで展望を行なっている幾つかの論文の間につながりをつけるには、Bayley and Hammersley の論文のオ11番目の式は、Madow and Madow のオ4の式と形式的に同じだということを示せる。これはまた非常に長い形の式にも適用するのである—これについては Yule and Kendall の Introduction to the Theory of Statistics, 14th Edition, Page 404 の公式 17.11 および Yule (1945) の J. R. S. S. 108, 208 における時系列についての論文を参照されたい。Bayley and Hammersley が研究した問題は、対空兵器のテストの連続的記録をとる器械の系統的な抽出に関するものであった。

系統的抽出におけるコレログラム

Cochran の論文(1946)は、コレログラムの型およびその系統的抽出におよぼす影響についての最も重要な進歩である—系統的標本の分散がある程度まで系列相関のみに依存することが起こされるであろう。この論文はまた母集団の要素を確率変数と考えることによつて算びかれる期待分散について研究するという考えを展開していることでも重要である。

筆者は相関する要素のグループ内分散が、グループの大きさとともに増加する、すなわち系列相関が存在するような母集団模型を考察した。この系列相関が正なら、コレログラムの一般的な形は単調減少函数(上に凹)である。平均的にいうと、層化無作為標本の効率(単純無作為標本のそれ以下となる)はなく、またその相対的効率は標本サイズの単調増加函数であるが、系統的標本の場合にはこのよ

うな一般的結果はえられない。Madow の論文の一部をなすこの結果の確認は、コレログラムが上に凹なら、平均的にいつて系統的標本は任意の標本サイズの層化標本より正確だという点で拡張されている。ついでにいえばコレログラムに対するこのような制限は、実際にはありえない程厳格なものではないといえる Wold (1938) は経済データに線型のコレログラムを用いたし Osborne も林業および土地利用調査で同じ型の函数を用いている。また Fisher (1922) は二つの気象観測所における逐雨量間の相関の研究に幾分似通つた見のものを用いている。

Cochran は系統的小よび無作為抽出計画の相対的精度をきめる決定的因子は、系列相関のオ1階差でなくむしろオ2階差の方であると述べている。このことは Madow and Madow (1944) の与えた幾つかの不等式の解釈をわかりさせるものであるが、それでも系統的標本の層化無作為標本に対する相対的効率、標本サイズの単純函数であるとは推論できない。期待分散についての彼の研究方針にしたがつて、Cochran は無作為、層化無作為および系統的標本の平均値の分散の式を与え、またそれらの相対的効率か式中の系列相関の線型函数とどのように密接な関係にあるかを示した。

この論文は

- (i) 指数型
- (ii) 線型

のコレログラムをもつ母集団についての各種の標本抽出の幾つかの結果を述べている。

こゝにまた理論の発展の次の主要段階について述べる前に、系統的抽出法の使用を報告している一群の主要論文に眼を通しておくのが便利である。

最初のもの (Vordskoy and Crump, 1948) は畜産の分野の論文である。筆者は

- (1) 無作為な日 (random day)
- (2) 連続した日 (consecutive day)

(3) 区間ごとの日 (interval day)

による産卵 — サンプルング方式からえた結果を比較し、この三つの方法間の正確さの差は僅かであると結論した。異なる日について産卵の数の産卵数をチエックしなければならないとき(このような事情は大規模な養蚕場ではよくある。)は系統的標本(区間ごとの日)の方が僅かに精度がよかつた。

このグループに属するオスの論文(Finney, 1948)は一般の刊行物および *risk of event* において屢々書述されている多くの点を強調しているのが重要である。立木の材積測定の場合、全数測定に要する時間と労力が大きいから実際の価値の大きい抽出方法を考えることが必要となる。この論文は Hasel (1938) と大体同じように、林業問題に適用した場合のサンプルングの一般的概観と、あとの論文(Finney, 1949)の序論として書かれている。これらの論文のオスのもの (Shaul and Myburgh, 1948) は、アメリカの人口の標本センサスを行なうための計画を論じたものである。与えられた抽出誤差に対して必要とされる標本サイズの計算公式を引用しているが、単純無作為抽出を仮定したときのものであることを注意している。そうして層化標本の方が精度がよいと断言している。しかし実際には、それはリストから系統的標本をとる結果を無作為標本として取扱うことにより精度を増加させるように計画したのであつた。これは明らかに Bowley その他の新しい研究法の現代版であつて、我々がこの展望の中で考えてきたような厳格な意味での系統的標本でないことは明らかである。この点については Gray and Corlett (1950) の論文が、これらの方法を幾つかの点で改善していることを注意しておく。

英国学派の寄与

ここにおいてこの国(英)の統計家による主な寄与を概観することが出来る。初期の文献はほとんど北アメリカで発表されている。吾輩のためこれ以前は研究上の連絡が阻げられていたが、し

しこの中の一つの結果は今日実用上最も有用な(特に二次元の系統的抽出の分野における)発展の幾つかを含むと考えられる一群の論文である。1948年9月前後の時点で書かれた論文を考える場合は、どれが先かということ順序をつけることは困難である。というのはそれか色々な定期刊行物での出版争奪によることが多いからである。この困難を避ける一つの方法は、関連する論文を一まとめにして考えることであらう。そこでまず Yates の研究(1946, 1948 および 1949)に向するものをとりあげることにする。

主要論文は“系統的抽出法”と題して *Phil. Roy. Soc. (A)* (1948) に発表された。その結果の幾つかおよびもつと一般には論文(Yates, 1946)の一節で簡単に説明されておりまた

Sampling Methods for Census and Survey (1949) なる書物に要約と展望が与えられている。最後の書物は僅かの間にこの分野の標準的な教科書となつたものである。1948年の Yates の論文は多くの理由のために系統的抽出論の進歩において劃期的重要性をもつものであつた。これはその殆んどが直接的意味で読者に実際の局面を示し、また計量的な値の系列に向する系統的標本の有用性を与えたということによるものである。しかしこの論文の取扱つているのは計量的データだけではない。というのはその一節で屢々一二つの値をとる函数の標本として表わされている。この系統的抽出を取扱っているからである。筆者は気象学を含めて多くの分野における実践面の応用を引用しているが、ついでにいうと、最近の論文は系統的抽出の問題の所在は認められているが、その解決は殆んどなされていないという感じを受けるであらう。その連判は Carruthers (1949) の論文である。そこで筆者は明らかに気温の読み取りに系統的抽出法を用いることを推奨しているがしかし毎日一定の時間単一の観測値をとることから生ずる困難を解決する手段は何もとつていないのである。

Yates は一次元の系統的抽出を取扱い、抽出誤差を推定する一つの方法を展覧した。無作為に配置された系統的標本に固有な傾向

(trend)の誤差を除去するために、彼は前の論文(Yates, 1946)の中で見通して系統的抽出の結果のみからは抽出誤差の完全に信頼できる推定値は作れないことを示し、また補足の証明を与えるために、 $1/2$ および $1/2$ 間隔で追加標本をとる方法を述べた。系統的抽出の実行は次のような母集団類型について研究された。

- (1) 二つの値をとる函数、すなわち属性
- (2) 正規分布する資料
- (3) 一項の自己回帰函数

Cochran (1946) にならって Yates も期待分散の概念を用いた。彼はまた何もわかっていないかまたは周期性のある資料について系統的標本を用いることの危険性について彼の考え方を強調している。

区別された区間に対して二つの函数をサンプリングする場合、属性を計する区間の割合が小さいかまたは1に近い場合、あるいは区間の大部分が抽出間隔に比べて小さいようなときには、系統的抽出と無作為抽出形式のどれを送るかということはそう重要でなくなる。しかしもし区間の大部分が抽出間隔より大なら、系統的抽出は大体においてこの型の無作為抽出よりも正確と思われる。この中間の場合には更に研究しなければならないが、しかし同様な結果が成立し得ると思われる。抽出対象が正規分布にしたがうときは、垂直的空中写真からのサンプリング — この分野はまだ十分発展していない — において二次元の類推の成立すること示される。

研究されたオ3の場合は単一項の自己回帰函数(単純マルコフチェーン)

$$y_{t+1} = by_t + a_{t+1}$$

のサンプリングである。この種の函数の性質は、連続変量の型とは本質的に異なることが指摘されている。そのコレログラムは指数型で、この場合も Cochran (1946) の論文のときに引用した文献とつながりがある。もし対象が実際に自己回帰であることがわかっていれば、系統的標本の分散を推定することは容易である。

しかしこのようなことは殆んどなく、もし真の自己回帰からのへだたりが多項の変動(long term variation)によるものであれば、その影響は僅かである。しかしもしこのへだたりが無作為変動が加えられたためのものであれば非常に極端な過少推定が生ずる。この危険を避けるには、部分的(Partial)系統的標本によって補助情報を求めることが提案される。非常に興味ある抽出研究について幾つかの結果が与えられ、またこの分野の地理論の実験例が提出される。

英国学派のオ2の主な寄手は故 A. E. Jones の論文である。この研究は、1948年5月の不幸かつ悲愴な Dr. Jones の事故により M. G. Kendall (1948) によつてうけつがれた。この論文で考えられたのは、連続パラメーターに依存する1次元の斉一な母集団からとられた確率変数の平均値を系統的標本で推定することであつた。この問題の核心は標本の構成要素を次々の値の間の相関度によつて連続な基準で配分する最良の方法を求めることである。この自己回帰模型は Yates (1948) の研究したものと同様であつた。Jones はもしこの相関が非常に小なら、 n 個の標本点を

$$T/(n+1), \quad 2T/(n+1), \dots$$

の距離で配分することを提案した。ここで T は抽出すべきストリップの長さである。一方相関が0.25より大なら最良の配分は距離が $T/2n, \quad 3T/2n, \quad 5T/2n$ のときのものである。どちらの場合でも、最適な分布は点か等距離すなわち系統的配置のときにえられる。

先の論文で Kendall (1948) はもつと集約なサンプリングと抽出された strip の拡張との間の本質的な違いを指摘し、後者で端の区間の長さにあまり神経質になることは不必要な手数をかけることだと述べている。この注意は Yates (1949) の書物のノクテ頁にもとり入れられている。

一方これまでにみてきたそれ以前の論文の多くは、抽出率の増加に対して形式的に同一であるところの集約な系統的抽出の効果については、解説されているが、これと相補う抽出された "strip" の拡張は全然ふれられていない。コレログラムが指数型でないようなも

つと一般な場合について、Kendallは同じ距離でとった観測値を同一の重みで考えてえられた結論は誤りであると述べている。しかしこのために必要とされる一般的方法から直接的な解がえられるとは思えない。さらにKendallは対象に著しい周期性がある場合、系統的抽出を使用してはならないという強い基準をあげて述べている。

Finney (1949)の論文は前年に発表されたつと一般な論文の続編をなすものである。これは森林調査において精度が高いとされている系統的抽出法を、色々な出所の多数のデータをもとにして吟味検討した研究である。この論文もまた統計的手法の相互作用が関連する問題に導入される仕方のしるしを与えるものである。

最近の進歩

W. G. Madocによるこの分野の研究のオニの紹介が1949年の終りに発表された。この論文において、彼は前の論文の結果を集約の大きさか等しい場合と等しくない場合の系統的抽出に拡張した。これはまた二次元の系統的抽出についての幾つかの説明が含まれているが、しかしこの問題に関する最も重要な研究はQuenouille (1949)の論文である。集約の大きさか等しい場合の系統的抽出には提案されている方法が二つすなわち

- (i) 抽出された集約内の全要素を調べる
- (ii) 抽出された集約内で層化および副次抽出を行なう。

あるか、この特別な設計は、人口母集団よりもむしろ自然的母集団のサンプリングに適しているように思われる。最も重要な結論はたとえそれが非統計的な理由によるものであっても、集約を系統的に抽出する方法は集約を用いることによって生ずる分散の増加を減らすほど有効なものでないということである。

二次元の系統的抽出の問題については、二つの点か引用されている。(1) City内のブロックに“渦巻型”に番号をつけてゆくやり方は、ブロック間に相関がある場合には有効な方法でない。

ある系統的標本はCityのブロックの列となることもある。コレログラムは上に凹となり

(2) グリッド上の行および列からなる系統的標本は行または列に沿って肥沃度の傾斜があるときには有効にならないであろう。しかし上に注意したように、この展望のせいで結びともなる論文は、二次元の系統的抽出の問題を取扱う段階を定められるものである。

Quenouille (1949)はその“平面上のサムプリングについて”と題する論文において系統的抽出の問題を実際と非常に近いところまで近づけた。これはそれまで殆んど直観的に取扱われてきたものである。Cochranの論文(1946)の中のある式について、重ね合わせた反動をつけ加える必要のあることを注意してから、彼は無作為抽出の分散に対する比で表わした系統的小よび層化無作為計画の抽出分散の差をコレログラムとともに図上でどのように用いられはこれらの抽出方式の相対的効率を調べることかできるかを示した。幾つかの有用な図表を用いて、彼は産産距離の一方または両方に平行または独立な標本——の場合には非常に多くの組合せかある——についてだけでなく無作為、層化無作為、系統的抽出方式の二次元サムプリングの問題を説明した。二つの無作為抽出方式の場合には平行に単位をとると分散分散が大きくなる。予感されるように系統的抽出の場合の分散の増加、減少はコレログラムの形によつてきまる。産産距離について産産の式か与えられ、また各産産距離に沿つた内部相関の異なる値に対して、二つの異なる抽出方式の相対的効率の表か与えられている。色々な場合を論じての結論は、二次元の独立な系統的抽出は無作為抽出より正確な結果を与えるということである。標本方式の抽出誤差の推定に提案された方法は平面方式の場合にも採用できる

- (1) 互いに独立に配置された系統的標本の組、誤差は各ブロック内の系統的標本の分散から計算される。
- (ii) 無作為に配置した系統的標本の一組、誤差分散は系統的標本

の一部から計算される。

(iii) 一つの系統的標本を用いて、これを分割して更に間隔の異なる系統的標本を幾つかとらえる。抽出分散は各副次標本の分散から計算する。

これらの3つの方法は、平均値の推定の場合はこの順序に精度がよくなるが、その反面抽出分散の推定にはこの順に偏りが大きくなるから、次第に実用性が低くなる。したがって方法の選択はある程度まで問題の性質に支配される。

結 語

過去10年の間、系統的抽出を実用に供しようとする直観的試みに対して、多くの理論的裏づけがなされてきた。引続き理論の一面の進歩を促すためには、実際的問題が非常に重要だということが多くの面において強く示されている。しかし、我々は実際には互いに僅かしか類似していないにもかかわらず同一の名前でよばれている2種類の抽出方式があつて、用語の上で根本的な困難が生じていることをみてきた。ある人々は一方を進めて、ある段階で標本抽出をとり入れない抽出方式はすべて系統的抽出とよぶべきだといふかもしれない。これは疑いもなく連続パラメータをもつ母集団からの集落抽出の場合をも含めるものであろう (Kendall J. R. S. S. 5, 1948, 230. para. 11)。しかしこれは我々がこの論文でみてきた抽出方式とは明らかに異なるものである。抽出方式の標準的な分類ができるまでにはなお時日を要するものと思われる。

by William R. Buckland

(J. R. S. S. B. Vol. 13, 1951より)

系統的抽出に関する重要文献

Bayley, G. V., and Hammersley, J. M. (1946)
 "The effective number of 'independent' observations in an autocorrelated time series," J. R. S. S. Suppl. 8, 184

Carruthers, N. (1949), "Accuracy of mean of n temperature observations" Met. Mag., 78, 65*

Cochran, W. G. (1946), "Relative accuracy of systematic and stratified random samples for a certain class of population", Ann. Math. Stat., 17, 164

Deming, W. E., and Simmons, W. R. (1946), "On the design of a sample for dealer inventories?" J. A. S. A. 41, 16*

Fonney, D. J. (1948), "Volume estimation of standing timber by survey" Forestry, 21, 179*

— (1949), "Random and systematic sampling in timber surveys" Forestry, 22, 64*

Fisher, R. A. and Mackenzie, W. A. (1922), "The correlation of weekly rainfall" Q. J. Met. Soc., 48, 234.

Gray, P., and Corlett, F. (1950), "Sampling for the Social Survey" J. R. S. S. (A), 113, 150.

Hanson, M., and Hurwitz, W. N. (1943), "On the theory of sampling from finite populations," Ann. Math. Stat., 14, 333.

Hasel, A. A. (1938), "Sampling error in timber surveys," J. Agric. Res., 57, 713.*

Jones, A. E. (1948), "Systematic sampling of continuous parameter populations," Biomet., 35, 285

Kendall, M. G. (1948), "A continuation of Dr. Jones' paper" Biomet., 35, 291

Madow, L. H. (1946), "Systematic sampling and its relation to other sampling designs," J. A. S. A. 41, 204.

Madow, W. G. (1949), "Systematic sampling", *I. A. M. S.*, 20, 339
 — and Madow, L. H. (1944) "On the theory of systematic sampling", *I. A. M. S.*, 15, 1.
 Nardskog, A. W. and Crump, S. L. (1948), "Systematic and random sampling for estimating egg production in poultry", *Biometrics*, 4, 223*
 Osborne, J. G. (1942), "Sampling errors systematic and random surveys of cover-type areas", *J. A. S. A.* 37, 256*
 Prudfoot, M. J. (1942), "Sampling with transverse traverso lines", *J. A. S. A.*, 37, 265*
 Queneville, M. H. (1949), "Problem in plane sampling", *A. M. S.*, 20, 355
 Shaul, J. R. H. and Myburgh, C. A. L. (1948) "Sampling survey of African population Southern Rhodesia", *Population studies*, 2, 339*
 Stephan, F. F. (1948), "History of the uses of modern sampling Procedures", *J. A. S. A.*, 43, 12
 Wadley, F. M. (1945), "An application of the poisson series to some problem of enumeration", *J. A. S. A.* 40, 93*
 Wald, H. (1938), *A study in Analysis of Stationary Time Series*, Uppsala: Almqvist & Wiksells.
 Yates, F. (1946), "A review of recent statistical developments in sampling and sampling surveys", *J. R. S. S.*, 109, 12
 — (1948), "Systematic sampling, *phil. Trans.* (A), 241, 347
 — (1949), *Sampling Methods for Census and Surveys*. London: Griffin & Co.

*印は主として応用を目的とした論文である。

系統的抽出に関する重要文献

Bayley, G. V. and Hammersley, J. M. (1946) "the effective number of 'independent' observations in an autocorrelated time series" *J. R. S. S. Suppl.* 8, 184
 Carruthers, N. (1949) "accuracy of mean of 'n' temperature observations" *Met. Mag.* 78, 65*
 Cochran, W. G. (1946) "Relative accuracy of systematic and stratified random samples for a certain class of population" *Ann. Math. Stat.* 17, 164
 Deeming, W. E. and Simmons, W. R. (1946) "On the design of a sample for dealer inventories" *J. A. S. A.* 41, 16*
 Finney, D. J. (1948) "Volume estimation of standing timber by survey" *Forestry*, 21, 179*
 — (1949), "Random and systematic sampling in timber surveys" *Forestry*, 22, 64*
 Fisher, R. A. and Mackenzie, W. A. (1922) "The correlation of weekly rainfall" *J. J. Met. Soc.* 48, 234
 Gray, P. and Corlett, F. (1950) "Sampling for The Social Survey" *J. R. S. S.* 113, 150
 Hansen, M. and Hurwitz, W. Z. (1943) "On the theory of sampling from finite populations" *Ann. Math. Stat.* 14, 333
 Hasel, A. A. (1938) "Sampling error in timber surveys" *J. Agric. Res.* 57, 713*
 Jones, A. E. (1948) "Systematic sampling of continuous

parameter populations" *Biom.* 35. 283

Kendall. M. G. (1948) "A continuation of Sir Jones' paper" *Biom.* 35. 291

Madow W. G. (1949) "Systematic sampling and its relation to other sampling designs" *J. A. S. A.* 41. 204

Madow W. G. (1949) "Systematic sampling: II." *A. M. S.* 20. 333

— and Madow. J. H. (1944) "On the theory of systematic sampling: I." *A. M. S.* 15. 1

Nordskog. A. W. and Crump. S. L. (1948) "Systematic and random sampling for estimating egg production in poultry" *Biometrics.* 4. 223*

Osborne. J. G. (1942) "Sampling errors of systematic and random surveys of cover-type areas" *J. A. S. A.* 37. 256*

Proudfoot. M. G. (1942) "Sampling with transverse travers lines" *J. A. S. A.* 37. 265*

Zuenouille. M. H. (1949) "Problem in plane sampling" *A. M. S.* 20. 355

Schul. T. R. H. and Myburgh. C. A. L. (1948) "Sample survey of African population of Southern Rhodesia", *Population Studies.* 2. 339*

Stephan F. F. (1948). "History of the uses of modern sampling procedures" *J. A. S. A.* 43. 12

Wadley T. M. (1945) "An application of the Poisson series to some problem of enumeration" *J. A. S. A.* 40. 93*

Wold. H. (1938) "A study in analysis of

Stationary Time Series. Uppsala: Almqvist & Wiksell.

Yates F. (1946) "A review of recent statistical developments in sampling and sampling surveys" *J. R. S. S.* 109. 12

— (1948) "Systematic Sampling" *Phil Trans (A).* 241. 347

— (1949) *Sampling Methods for Census and Survey* London: Griffin & Co

*印は主として応用を目的とした論文である。

10 系統的抽出の理論 III

中心化された系統的抽出とランダムスタートの系統的抽出との比較

On the theory of systematic sampling III Comparison of centered and random start systematic sampling⁽¹⁾

1 要約

えられた主な結果は次の通りである。

母集団のコレログラムが単調減少なら、中心化された系統的抽出 (Centered systematic sampling) はランダムスタートの系統的抽出より効率が高い、これとともに、母集団が中心化された系統的抽出はランダムスタートの系統的抽出より効率が高いことがわかったが、しかしこの場合層化無作為抽出の効率はこのいづれより高い場合のあることが容易に調整できる。このように、Cochran (1) の証明した、ランダムスタート系統的抽出の効率が層化無作為抽出より高くなるような場合 (例えばコレログラムが上に凹で減少する) には中心化された系統的抽出はランダムスタートの系統的抽出より効率が高くなる。

1 緒言 考察した抽出のタイプ

この論文では中心化された系統的抽出方法の理論を考える。よく知られている通り、このような標本抽出の方法は長い間実証的望望性を認められてきたものである。中心化された系統的抽出の理論は端末補正を施したランダムスタート系統的抽出に対しても成立する (Yates (5))、というのは端末補正の方法によつてランダムスタート系統的抽出は實際上中心化された系統的抽出に帰着するからである。

論文に用いる方法は Cochran (1) および筆者 (3) (4) の以前の論文のものにしたがったので、証明と記法を簡潔な形で与えてお

く。

母集団の要素は X_1, X_2, \dots, X_N で $N = kn$ である。目的は大きさ n の標本から母集団の算術平均値 \bar{X} を推定することである。

ランダムスタート系統的抽出の推定値 \bar{X}_{sy} は、 X_1, \dots, X_n の中から単純抽出法で 1 要素を選びどのあとの k 番目ごとの要素を標本に含めることにしてえた n 個の要素の算術平均値である。これらの異なる n 個の標本の算術平均値を $\bar{X}_{s1}, \dots, \bar{X}_{sn}$ と書く。ここで X_i はどの k の要素が X_i となっている標本の平均値である。 \bar{X}_{sy} の分散を母集団の要素で表わして σ_{sy}^2 と書く。 k が奇数なら中心化された系統的抽出の推定値 \bar{X}_c は $\bar{X}_{(k+1)/2}$ で、 k が偶数なら我々は任意に $\bar{X}_c = \bar{X}_{k/2}$ と定義する。(実際には k が偶数なら、人によつて $k/2$ または $(k+2)/2$ どちらかを適時に送んだり、あるいは要素 $X_{k/2}, X_{(k+2)/2}, \dots, X_{(k+2)/2}$ を送り出す上記の方法の代りに別な形の標本の抽出法をとるかも知れない。たとえば単調な母集団では $X_{k/2}, X_{(2k+2)/2}, X_{3k/2}, X_{(4k+2)/2}$ をとる方が望ましい。我々のいまの目的については、最良の抽出の型を決定することはあまり重要でない) \bar{X}_c のまわりの \bar{X}_c の平均平方を σ_c^2 と書く。

層化無作為抽出の場合には我々は層を構成するため $X_{1+(j-1)k}, X_{2+(j-1)k}, \dots, X_{jk}$ を考える。ただし $j=1, \dots, n$ 。よつてそれぞれを個々の要素からなる n 個の層があることになる。我々は n 個の層の n 個の要素が単純抽出法で送られるものとする。この標本平均を \bar{X}_{st} と書き、また母集団要素で表わしたときの \bar{X}_{st} の分散を σ_{st}^2 と書く。

我々は、母集団の要素が常数と考えられる場合の期待値の記号を E で表わし、また母集団要素が確率変数と考えられる場合のものを E と書くことにする。

以下においてはコレログラムという用語を多少広義に用いるので、ここで用語について述べておく。 X_1, \dots, X_N が順序のついた確率変数列で、派字が k だけ異なる二つの確率変数の相関係数を ρ_k とするとき (たとえば $\rho_2 = \sigma_{X_1, X_3} / \sigma_{X_1} \sigma_{X_3}$)、この相関が k の

みでまざるものとするは、函数 $f(s) = \rho_s$, $s = 1, \dots, N-1$ は
種々この列のコレログラムとよばれる。通例確率変数は等平均、等
分散と仮定される。しかし我々が以下でこの用語を使うときは、

$s = 1, \dots, k-1$ のみに依存する仮定した観測値 X_{1s}, \dots, X_{ks} の期待値のみを指
すことにする。よつて確率変数が等しい平均値をもつなら、我々の
論述は普通のコレログラムを指すことになるが、それ以外では、我
々の述べる条件は確率変数の平均値の等しいことを仮定しない。

3 単調な母集団

Hotelling と Solomons (2) は Z_1, Z_2, \dots, Z_q に対し
て次の不等式の成立つことを証明した。

$$(3.1) \quad \frac{q \text{ (メテイアン - 算術平均)}^2}{\sum_{i=1}^q (Z_i - \bar{Z})^2} \leq 1$$

ただしこの項もすべて有限で、分母は 0 にならないものとする。
この場合 q が奇数ならメテイアンは普通に定義されるものであるが、
 q が偶数で、 Z と Z が二つの中央値であれば、メテイアンは
 $Z \leq \text{メテイアン} \leq Z'$ なる任意の値にとることかできる。(細かい
証明は q が奇数の場合しか与えられていないが q が偶数のときもす
ぐ証明できる)

母集団が単調ならそれは $\bar{X}_1 \leq \bar{X}_2 \leq \dots \leq \bar{X}_k$ か $\bar{X}_1 \geq \bar{X}_2 \geq \dots$
 \bar{X}_k のどちらかである。よつて k が奇数なら \bar{X}_c は $\bar{X}_1, \dots, \bar{X}_k$ のメ
テイアンで、 k が偶数なら \bar{X}_c は $\bar{X}_{k/2} \leq \bar{X}_c \leq \bar{X}_{(k+2)/2}$ または
 $\bar{X}_{k/2} \geq \bar{X}_c \geq \bar{X}_{(k+2)/2}$ である。それで (3.1) は

$$(3.2) \quad \frac{(\bar{X}_c - \bar{X})^2}{\sigma_{SY}^2} \leq 1$$

となる。($\bar{X}_c - \bar{X}$) は母集団の要素が確率変数でないときの \bar{X}
のまわりの \bar{X}_c の平均平方誤差だから次の定理が証明されたことに
なる。

定理 1 母集団が単調なら、中心化された系統的抽出はランダム
スタートの系統的抽出より良好である。

勿論母集団が単調で標本が十分大きければ、層化無作為抽出の方
が中心化された系統的抽出より効率が高くなる。なせなら後者の推
定値のもつ偏りは標本の大きさが増しても十分な速さで 0 に収斂し
ないからである。

しかし実際には母集団が単調なときでさえ中心化された系統的抽
出の方が層化無作為抽出より効率が高くなることが多い。なせこの
ようなことが起るかを調べるため、 σ_c^2 の平均的分散と平均的共分
散の項を定義してみよう。

$$\bar{X}_j = \frac{1}{k} \sum_{i=1}^k X_{i+(j-1)k} \quad j = 1, 2, \dots, n$$

とおく。そうすると

$$\sigma_c^2 = S + C$$

である。ここで

$$S = \frac{1}{n^2} \sum_{j=1}^n (X_{a+(j-1)k} - \bar{X}_j)^2$$

$$C = \frac{1}{n^2} \sum_{\substack{j,m=1 \\ j \neq m}}^n (X_{a+(j-1)k} - \bar{X}_j)(X_{a+(m-1)k} - \bar{X}_m)$$

で、 k が奇数なら $a = (k+1)/2$, k が偶数なら $a = k/2$ である。
我々は S, C をそれぞれ σ_c^2 の平均的分散項、平均的共分散項とよ
ぶことにする。

Hotelling と Solomons の結果から

$$(X_{a+(j-1)k} - \bar{X}_j)^2 \leq \frac{1}{k} \sum_{i=1}^k (X_{i+(j-1)k} - \bar{X}_j)^2 \quad j = 1, \dots, n$$

が成立つ。よつて $S \leq \sigma_{St}^2$ 。このように $C < \sigma_{St}^2 - S$ なら
 $\sigma_c^2 < \sigma_{St}^2$ である。実際の場合平均的共分散の項 C は C について
の上の条件を満足する位小さい値であることが多い。

4 コレログラムが単調減少であるような母集団 (実際には $k/2$ 以
下で単調減少と仮定しているようなものをいうのである。)

この節では次の記法が必要である。特に断らない限り、 i と j は 1 から k までのすべての整数値をとるものとし、 m は 1 から k までのすべての整数値をとるものとする。また δ は 1 から $k-1$ までの 1 から $k-1$ まで、それに E は 1 から $\frac{k-1}{2}$ ($k-1$) までのすべての整数値をとるものとする。

(証明の中で k は奇数と仮定する。 k が偶数のときは基本的な結果は変わらないが一層の複雑さと記法が導入される。) 我々はここで母集団の要素が確率変数と仮定し、また

$$E X_{A+(j-1)k} X_{B+(m-1)k} = \mu_{(m-j)k+\delta}$$

とおく。ここで $\delta = B - A$ 、 $j \leq m$ である。よって $-(k-1) \leq \delta \leq (k-1)$

定理2 上記の条件のもとで

$$(4.2) \quad E \sigma_{Sy}^2 - E \sigma_C^2 = \frac{4}{\pi^2 k^2} \sum_E E \sum_{j=0}^{n-1} (\mu_{jk+E} - \mu_{(j+1)k-E})$$

が成立する。

もし $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{nk-1}$ で不等号が1つでも成立てば中心化された系統的抽出はランダムスタートの系統的抽出より効率が良い。しかし $\mu_1 \leq \mu_2 \leq \dots \leq \mu_{nk-1}$ で不等号が1つでも成立てば逆である。

定理2の証明の前にその意味を少し考えてみよう。実際、母集団の要素がこれらの距離と無関係に同一の積の期待値(4.1)をもつ。すなわち、 $\mu_1 = \mu_2 = \dots = \mu_{nk-1}$ なら(4.2)から

$$E \sigma_{Sy}^2 = E \sigma_C^2$$

かえられる。もし母集団の全要素の期待値を同じと仮定すれば、上記の命題は期待積 (expected product) よりはむしろ系列共分散 (serial covariance) についていえるものである。たとえば母集団の要素の期待値が同じで無相関なら

$$E \sigma_{Sy}^2 = E \sigma_C^2$$

上述の、中心化された系統的抽出の効率からランダムスタートの系統的抽出より低くなるような条件は、現実には殆んど満たされない筈である。しかし実際には、コレログラムが不規則なためランダム

スタートの系統的抽出の方が、中心化された系統的抽出よりも効率の高くなることによくある。

証明(4.2)の証明は面倒だが、困難ではない。まず次の二つの

Lemma

Lemma 1 $f(i-k)$ を整数 i および k の差の函数とすると

$$(4.3) \quad \sum_{i=k} f(i-k) = k f(0) + \sum_{\delta} (k-\delta) [f(\delta) + f(-\delta)]$$

である。同じ様にして

$$(4.4) \quad \sum_{i=k} f(i-k) = k f(0) + 2 \sum_{\delta} (k-\delta) f(\delta)$$

証明は省略する。

Lemma 2 $\delta = |i-k|$ とおく。そうすると(4.3)なら

$$(4.5) \quad E \bar{X}_i \bar{X}_k = \frac{\mu_{\delta}}{n} + \frac{1}{n^2} \sum_{\delta} (n-\delta) (\mu_{jk+\delta} + \mu_{jk-\delta})$$

で、また

$$(4.6) \quad E \bar{X}_C^2 = \frac{\mu_0}{n} + \frac{2}{n^2} \sum_{\delta} (n-\delta) \mu_{jk}$$

が成立する

証明 ここで $X_{i+(j-1)k}$ を X_{ji} と書く。

$$\bar{X}_i = (1/n) \sum_{\delta} X_{ji} \text{ だから}$$

$$\begin{aligned} E \bar{X}_i \bar{X}_k &= \frac{1}{n^2} \sum E X_{ji} X_{jk} + \frac{1}{n^2} \sum_{j < m} E X_{ji} X_{mk} + \frac{1}{n^2} \sum_{j > m} E X_{ji} X_{mk} \\ &= \frac{1}{n} \left\{ \mu_{\delta} + \sum_{\delta} \frac{n-\delta}{n} (\mu_{jk+\delta} + \mu_{jk-\delta}) \right\} \end{aligned}$$

よって(4.5)は証明された。(4.6)はその特別な場合である。

定理2の証明にもとつて

$$\Delta = E \sigma_{Sy}^2 - E \sigma_C^2$$

$$\Delta = \frac{1}{k} \sum_{\delta} E \bar{X}_C^2 - E \bar{X}_C^2 + 2 E \bar{X}_C \bar{X} - 2 E \bar{X}^2$$

かえられる。Lemma 2によつて $E\bar{X}_i^2$ は i と独立だから

$$\Delta = 2E\bar{X}_i\bar{X} - 2E\bar{X}^2$$

をうる。ここで (4.4) と (4.5) から $i=C$ にとり、 k について平均すると

$$E\bar{X}_i\bar{X} = \frac{1}{nR} \left\{ \mu_0 + 2 \sum_{\delta} \frac{n-\delta}{n} \mu_{\delta R} \right\} + \frac{2}{nR} \left\{ \sum_E \mu_E + \sum_{\delta} \frac{n-\delta}{n} (\mu_{\delta R+E} + \mu_{\delta R-E}) \right\}$$

さらに

$$E\bar{X}^2 = \frac{1}{n^2} \sum_{i,k} E\bar{X}_i\bar{X}_k$$

で、Lemma 2により $E\bar{X}_i\bar{X}_k$ は $|i-k|$ にのみ関係するから Lemma 1によつて

$$E\bar{X}^2 = \frac{1}{nR} \left\{ \mu_0 + 2 \sum_{\delta} \frac{n-\delta}{n} \mu_{\delta R} \right\} + \frac{2}{nR} \left\{ \sum_{\delta} \frac{k-\delta}{k} \left(\mu_{\delta} + \sum_{\gamma} \frac{n-\gamma}{n} (\mu_{\delta R+\gamma} + \mu_{\delta R-\gamma}) \right) \right\}$$

である。そうすると

$$\Delta = \frac{4}{nR} \sum_E \frac{E}{R} (\mu_E - \mu_{k-E}) + \frac{4}{nR} \sum_{\delta} \frac{n-\delta}{n} \sum_E \frac{E}{R} (\mu_{\delta R+E} + \mu_{\delta R-E} - \mu_{(\delta+1)R-E} - \mu_{(\delta+1)R+E})$$

である。いま

$$\sum_{\delta} (n-\delta) (\mu_{\delta R+E} - \mu_{(\delta+1)R+E}) = -n\mu_{-E} + \sum_{\delta=0}^{n-1} \mu_{\delta R+E}$$

でまた

$$\sum_{\delta} (n-\delta) (\mu_{\delta R-E} - \mu_{(\delta+1)R-E}) = n\mu_{R-E} - \sum_{\delta=0}^{n-1} \mu_{(\delta+1)R-E}$$

である。

よつて

$$\sum_{\delta} \frac{n-\delta}{n} \sum_E \frac{E}{R} (\mu_{\delta R+E} + \mu_{\delta R-E} - \mu_{(\delta+1)R-E} - \mu_{(\delta+1)R+E}) = - \sum_E \frac{E}{R} (\mu_E - \mu_{k-E}) + \frac{1}{n} \sum_E \frac{E}{R} \sum_{\delta=0}^{n-1} (\mu_{\delta R+E} - \mu_{(\delta+1)R-E})$$

となる。これで (4.2) は証明された。 $1 \leq E \leq (k-1)/2$ だから、 $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{nA-1}$ なら $\mu_{\delta R+E} \geq \mu_{(\delta+1)R-E}$ である。よつてコレログラムが単調減少なら $\Delta \geq 0$ でありまたコレログラムが単調減少かつ一定でなければ $\Delta > 0$ である。

5 注

この論文の結果は二次元のサムプリングおよび乗算のサムプリングにも容易に拡張できるこれらの問題は次の論文で論ずることにする。

興味あることは、もし $E\bar{X}_i^2 + (i-1)k$ が i および i に独立でないとして仮定すれば、 $E\bar{X}_i^2 + (i-1)k$ について更に仮定を設けないと上記の結果が成立しないことである。

by W. G. Madow A. M. S. 1953 (pp. 105-)

参 考 文 献

- 1 W. G. Cochran "Relative accuracy of systematic and stratified random sampling for a certain population" A. M. S. vol 17 (1946) pp. 164-177
- 2 H. Hotelling and L. M. Solomon "Limits of measure of skewness" A. M. S. vol 3 (1932) pp. 141-142
- 3 W. G. Madow and L. H. Madow "On the theory of systematic sampling I" A. M. S. vol 15 (1944) pp. 1-24
- 4 W. G. Madow "On the theory of systematic sampling, II" A. M. S. vol 20 (1949) pp. 333-354
- 5 F. Yates "Systematic sampling" Philos. Trans. Roy. Soc. London Ser. A. vol 241 (1948) pp. 345-377

// 平面における確率過程について On stationary process in the plane

空間における定常過程 (stationary process) の抽出理論は定常時系列 (stationary time series) と全く同じようには論じられない。というのは、時系列 (time series) の変量は過去の値のみによって影響されるのに及し、空間的過程はすべての方向の値によって定まるのだから。この点については §2 スタームにおいて詳論した。§7 で展開した推定と検定の理論は §8 で小麦とオレンジのデータの二様性について適用した。最後の節は幾つかの特殊な二次元の過程の研究にあてた。

緒 言

画物実験、森林および収量調査あるいは居住地域 (populated area) の標本調査等における地形的相関 (topographic correlation) の騒乱効果については長く知られているところであって、二次元の確率過程 (two dimensional stochastic process) に関するこのデータについても認められるものである。物理学者も高次元の過程 (例えば乱流 (turbulence) や分子系) を取扱うが、実際には物理学者の方が主要な研究者であった。

我々の提起した過程が定常 (stationary) と考えられるのはかくとも一次元以上の意味においてである。しかしこの近似は屢々十分満足であるから定常的な過程は研究する価値がある。

二次元の理論の大部分は時系列の研究で用いられるもの、形式的な拡張にすぎないから、求めるだけ簡単にこれらの見直しをつけておこう。しかし上の要約で注意した様に興味のある一つの新しい側面が現れるので、これを詳細に考察することにする。

一般に二次元で特定の模型 (model) を研究する場合には技術的な面の数学的取扱いが難かしいということだけのために、一次元の

場合より困難が大きいものである。ある場合には初等函数の範囲を全く逸脱してしまう様にも感える。

我々はこの障害を指摘するであろうが、しかし幸いなことに多くの目的についてはこれらの困難を避けることができる。多くの応用では全く非決定論的な過程 (purely non-deterministic process) を考えるだけで十分であるから、我々はこの種の過程、特に線型自己回帰過程 (linear autoregressive process) のみに問題を限定する。

2 ライントランセクト

まずライントランセクト (line transect) (即ち地域上に設定した直線で、観測値はこの直線に沿って等距離にとる)。トランセクトの観測値は時系列の場合と全く同じく、一次元の過程から生成 (generate) されていると考えることができる。しかしこの二つの場合には重要な違いが存在する。すなわち時系列ではどの瞬間についても過去と将来の間に自然的な区別があり、またある瞬間における観測値は過去の値のみによって定まる。すなわちこの従属性は唯一方向 (すなわち右と向き) のみに進む。しかしトランセクトの場合にはこのような二つの方向の区別はなく、両方の側に対する従属性が存在する。我々は距離に対するより一般的に二次元の側を考えることができる。距離の任意の点で施した小量の肥料は、終局的にはすべての方向の土壌の肥沃度に影響を与える。(勿論、例外もありうる：たとえば風場の傾斜がひどくて、施肥した点より下側の地域のみが影響される場合のように。)

いま観測値と“誤差”の変量をそれぞれ z_t と ϵ_t ($t = \dots, -2, -1, 0, 1, 2, \dots$) と書けば、最も簡単な現実的な時系列の模型はおそらく一次の自己回帰 (autoregression)

$$z_t = a z_{t-1} + \epsilon_t \dots \dots \dots (1)$$

であろう。

しかしトランセクトでは模型 (1) は退化する場合のもので、一

つの方向についてのみ従属性が見られる。非退化 (non-degenerate) なトランセクト模型の最も簡単なものは、

$$z_t = a z_{t-1} + b z_{t+1} + \epsilon_t \dots \dots \dots (2)$$

であろう。ここで a, b があまり大きくてはならないということは直観的に明らかである。我々は (1) の一方向の自己相関

(unilateral autoregression) と區別して (2) を二方向の自己相関 (bilateral autoregression) と呼ぶ。一方向型 scheme

に対するこのような仮定を極めるほどの様な結果が得られるかはパラメーター (parameter) の推定を考える場合に明瞭になつてくる。

(1) のパラメーター a は残差平方和 (residual sum of square) $\sum_t (z_t - a z_{t-1})^2$ を最小にすれば (consistent) な推定値が得られるにもかかわらず、(2) の a, b を求めるにはたゞ、

$$U = \sum_t (z_t - a z_{t-1} - b z_{t+1})^2 \dots \dots \dots (3)$$

を最小にしただけでは無意味な結果がえられる。このことは z_{t+1} の値と z_t が相関している場合には z_t の条件付 (conditional) の平均値 $E(z_{t+1} | z_t)$ であることは向道しているということから説明できる。形式的には、 E_t から z_t への変換 (transformation) の Jacobian が、(1) のような一方向の相関 (unilateral relation) と違って (2) の相関では上でないということができる。しかし最小二乗推定の正しい方程式は U を最小にすれば得られるということが示される。(これは必ずしも一般的な二次元の場合について証明する) ここで我々は

$$\log k = -\frac{1}{2\pi} \int_0^{2\pi} \log (ae^{i\omega} + b\bar{e}^{i\omega})(a\bar{e}^{-i\omega} + b e^{i\omega}) d\omega \dots \dots (4)$$

で与えられるパラメーターのある種の函数である。

3 二方向 (bilateral) scheme の形式的な性質

二次元の過程を論ずる準備としてトランセクト模型の性質を簡単に調べてみよう。

一般的な二方向の線型自己相関 (bilateral linear autoregression)

$$L(T) z_t = \epsilon_t \dots \dots \dots (5)$$

を考えよう。ここで $L(T) = \sum a_j T^j$ で T は変換の作用素 (translation operator)

$$T z_t = z_{t+1}$$

である。

(5) 式の解は

$$z_t = \frac{\epsilon_t}{L(T)} = \sum b_j \epsilon_{t+j} \dots \dots \dots (7)$$

で与えられる。ここで b_j は $\{L(e^{i\omega})\}^{-1}$ の Fourier 展開の $e^{i\omega}$ の係数である (Bartlett, 1946, P 60 を見よ)。Fourier 変換を用いれば、Scheme (5) のスペクトル函数 (spectral function) は同様にして

$$F(\omega) = \frac{\sigma^2(\epsilon)}{L(e^{i\omega}) L(e^{-i\omega})} \dots \dots \dots (8)$$

であることがわかる (Lob, 1944, Llanjell, 1946) 自己共変動 (autocovariance)

$$\phi(j) = \text{cov}(z_t, z_{t+j}) \quad (j=0, \pm 1, \pm 2, \dots) \dots \dots (9)$$

は $F(\omega)$ の Fourier 展開の係数として得られる。したがって、簡単のために、 $e^{i\omega} = z, \sigma^2(\epsilon) = v$ と置くことにすれば

$$F(\omega) = \frac{v}{L(z) L(z^{-1})} \dots \dots \dots (10)$$

多くの本質的な点が含まれているから、展開式 (7) (8) が成立するための条件と考えることも自然である。(5) が "過去のみに対する従属" を表わす一方向の (unilateral) scheme を表わすとすると $|z|=1$ での $\{L(z)\}^{-1}$ の Laurent 展開にはこの指数に正数は含まれてはならないから、時系列自己相関 (time series) の安定性 (stability) に対する普通の条件が導かれる。すなわち $L(z) = 0$ のすべての根は単位円の内部にある (Wold, 1938) しかし、もし二方向 (bilateral) scheme を承認する場合に必要なことは、係数が何も制限されていなければ $\{L(z)\}^{-1}$ の Laurent 展開が収束することである。この場合の安定性の条件は非常に弱くなつていて、必要なのは、 $L(z) = 0$ の根がどれも単位円上にあ

つてはならないということである。

しかしこれは特定のL(z)が定常であることを保証するための条件にすぎない。

L(z)は(5)なる関係式では従属変数であるε_tであることが望ましい。

(例えば関係式(1)をε_{t+1} = ε_{t+1} - ε_tと変形したときのε_{t+1}ではなく)これにはL(z)にzのある整数乗をかければよい(あるいは同じことであるがε-系列(sequence)を転移(translation)によって再定義すればε_tはε_tと対応がつく)というのは、zが単位円周上を正の方向に一周すればL(z)は複素平面上で原点のまわりを一周する(ここでzは整数でなければならぬ)したかつてε-系列を一周階進して転移(translation)すれば正規化(normalization)できる。ゆえにL(z)を正規化作用素(normalized operator) T⁻¹L(z)即ちL*(z)でおきかえることができる。原点はlog L*(z)のbranch pointでないから(zの0でないときのL(z)の場合と同じく)log L(z)はzのLaurent級数に展開でき、また作用素L*(z)は

$$L^*(z) = e^{\sum_{j=1}^{\infty} T_j z^j} \dots \dots \dots (11)$$

と表わされる。

これらは、この正規化が行なわれているものと仮定する。したかつて(11)は常に可能である。

4 推定の非確定性

二方向の(bilateral) schemeを承認すると、Wold (1938)が移動平均(moving average)について示したのと同様な困難が自己相関を求める際に生ずる。これから、オーサーPの自己相関から生成される一連の自己共変動が与えられたものとするれば自己相関を定めることができる。すなわち(5)の多項式L(z)を定めることができる。ここで(10)を満足するL(z)を求めることが必要である。ここでF(w)は与えられた自己共変動で決定され

る。このようなL(z)のとり方は、2^P通りある。何故なら特定のL(z)の根をα₁, α₂, ..., α_PとすればF(w)の2^P通りの可能な表現

$$F(w) = \frac{Const}{|(z^1 - \alpha_1)(z^1 - \alpha_2) \dots (z^1 - \alpha_P)|^2} \dots \dots (12)$$

に対応して2^P通りの可能な有限自己相関(finite autoregression)が存在する。

これらの二つのschemeは一方向であるが、ただ違うのは(時間(time)軸が正の方向のみに向かっているという点である。

同様な論法はあるデータの組にP次の自己回帰を最小自乗法であてはめる場合にも適用できる。というのはこの場合にもまた自己共変動から計算を行うからである。これは2^P個のあてはめの自己相関の一つから選べるのと同じようにして他からも容易に求めることができる。

例として、もう一度二方向のscheme(2)を考える。このschemeから3の最後の断のようにして正規化されれば、方程式

$$a - z + b z^2 = 0 \dots \dots \dots (13)$$

の根は単位円の内部にあり他は外部にある。この根をα, β⁻¹と置くことにする。

(|α| < 1, |β| < 1)。原数AとBを

$$(z - \alpha)(z - \beta) = z^2 + Az + B \dots \dots \dots (14)$$

で定義すれば自己相関

$$\epsilon_t + A \epsilon_{t-1} + B \epsilon_{t-2} = \epsilon_t \dots \dots \dots (15)$$

を作ることもできる。これは一方向であるけれども(2)と同じ自己相関を生ずる。

積分(4)を計算すればa, bは

$$\left[\frac{1 + \sqrt{1 - 4ab}}{2} \right]^{-2} \sum (\epsilon_t [a \epsilon_{t-1} - b \epsilon_{t-2}]^2) \approx \left[\frac{1 + \sqrt{1 - 4ab}}{2} \right]^{-2} (1 + a^2 + b^2) (c_0 - 2(a+b)c_1 + 2ab c_2) \dots \dots (16)$$

を最小にすれば推定されることかわかる。

ここで C_s は $\log s$ で観測された自己共変動であるパラメータ A および B を変換すれば (16) 式は

$$(1/A^2+B^2)C_0+2(A+AB)C_1+2BC_2 \approx \text{Const.} \sum (\xi_{t+1}+A\xi_{t-1}-B\xi_{t-2})^2 \dots (17)$$

に比測することかわかる。これは *Scheme* (2) と (15) が同等であることを他の面から表わすものである。この同一性の実際的な意義は取扱いにくい (16) 式を最小にする代りにより簡単な (17) 式を最小にすればよいということにある。したがって A および B の推定値は (14) の関係を用いて、 $A : B$ のそれから計算することができる。

二方向型の *scheme* を効果的に一方向型に帰着させるような場合には、そのような *scheme* を導入する必要はないように思えるかもしれない。実際このやり方は、我々が二方向の観測にどれ位多くの面割かつ不確実なパラメータがあるかということを考える場合には程かに望ましくないことかわかる。しかし二次元の舞台には、いま考えた一次元の舞台と違って、一方向の *scheme* に帰着させるには非常に多くの面倒な問題が生ずることかわかるであろう。このようにすべての方向への従属性を正確に導入することは不可理である。またそれ以上の気が立っている。すなわち実際になく対応する多方向の *scheme* は、たとえそれが推定値等の形式的な作業のような場合でも同等な一方向の観測を用いて行えば最も簡単である。

5 二次元の過程の一般性

平面上のいかなる点においても、ある値となる様な変量の連続的な過程を考えることができる。それは普通、変量か矩型 of 平面上の格子でのみ観測される様なものである。我々は主として数学的により簡単でまたより実際の興味の多い離散的な場合を取扱うことにする。

観測値および誤差変量をそれぞれ ξ_{st} および ϵ_{st} ($s, t = \dots, -2, -1, 0, 1, 2, \dots$) で表わす。連続過程を論ずる場合、この記号は (x, y) および $\epsilon(x, y)$ (x, y はすべて実数とする) に変える。

我々かこれから主として取扱う特別な模型は二次元の観測自己回帰

である。これを我々は

$$L(T_s, T_t) \xi_{st} = \epsilon_{st} \dots \dots \dots (18)$$

と書くことにする。ここで T_s, T_t は

$$T_s \xi_{st} = \xi_{sT, t} \quad T_t \xi_{st} = \xi_{s, t+1} \dots \dots \dots (19)$$

および

$$L(T_s, T_t) = \sum_j \sum_k a_{jk} T_s^j T_t^k \dots \dots \dots (20)$$

で定義される。転移作用素 (*translation operator*) である。方程式 (17) - (19) に対応して

$$\xi_{st} = \frac{\epsilon_{st}}{L(T_s, T_t)} = \sum_j \sum_k b_{jk} \epsilon_{sT^j, t+T^k} \dots \dots (21)$$

$$F(W, W_2) = \frac{1}{L(Z_1, Z_2) L(Z_1^{-1}, Z_2^{-1})} \dots \dots \dots (22)$$

$$= \sum_j \sum_k \phi(j, k) Z_1^j Z_2^k \dots \dots \dots (23)$$

をうる。

ここで b_{jk} は $[L(Z_1, Z_2)]^{-1}$, ($Z_1 = e^{iW}, Z_2 = e^{iW_2}$) の Fourier 展開の Z_1^j, Z_2^k の係数。 v は ϵ の分散、 $\phi(j, k)$ は ξ_{st} と $\xi_{sT^j, t+T^k}$ の共変動 (*covariance*)。また $F(W, W_2)$ は (23) によつて定義される数値に対応するスペクトル函数である。

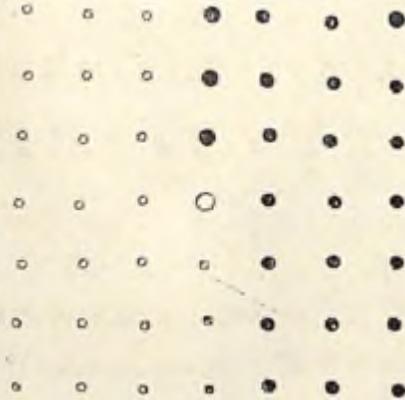
前と同じく、(21) と (22) が成立つための必要かつ十分な条件は、任意の Z_1, Z_2 に対し $L(Z_1, Z_2)$ が 0 でなく、同時に $|Z_1| = 1, |Z_2| = 1$ を満足することである (少くとも自己回帰 (18) が有限であるときに)。

よって、あるいは Z_2 の何れかを単位円周に沿つて動かすとき、複素平面の原点が $L(Z_1, Z_2)$ の円内に入らない様に作用素 (*operator*) L が正規化されていると仮定する。(もし Z_1 が単位円の外側をまわり、 Z_2 が特定の値をとるとき、 L が原点とその円周に含まなければ任意の他の Z_2 の値に対してもその誤差率は起らない。何となればもしその様なことが起れば仮定に反して L が原点を滑つて動くような Z_2 のある中間的な値がある筈である。

6 二次元過程の一方角模型による表現

今日よく知られているタイプの論法 (Wienerによる。例えは Wiener, 1949, 978) によつて、与えられた自己相関の組 (その場合 ξ_{st} は ξ_{su} ($u>t$) および ξ_{vw} ($w>s$, w は無条件) との自己相関、として表わされる) を生成する過程は唯一通りに求めうることを証明しよう。すなわち Fig. 1 の格子において、白い点の値は黒点の値によつて表わしうるのである。与えられたコレログラム (autocorrelogram) は完全にどんなものでもよいというわけにはゆかないが、必要条件は非常にゆるやかである。

Fig. 1



その条件は自己相関が純非一決定的過程 (non-deterministic process) によつて生成されるということである。これはとりわけ $\log F(W_1, W_2)$ が Fourier 展開でき、したかつて $F(W_1, W_2)$ が

$$F(W_1, W_2) = \exp \left\{ \sum_i \sum_k d_{ik} Z_i^+ Z_k^* \right\} \dots (24)$$

のように表わされることを意味する。

ここで函数

$$P(Z, Z_2) = \exp \left\{ \frac{d_{00}}{2} + \sum_{k=1}^{\infty} d_{0k} Z_2^k + \sum_{j=1}^{\infty} \sum_{k=-\infty}^{\infty} d_{jk} Z_j^+ Z_k^* \right\} \dots (25)$$

を定義する。

$$F(W_1, W_2) = \frac{1}{P(Z, Z_2) P(Z_1^+, Z_2^*)} \dots (26)$$

だから自己相関

$$P(T_s, T_t) \xi_{st} = \xi'_{st} \quad (\sigma^2(\xi') = 1) \dots (27)$$

は要求されたスペクトル函数 $F(W_1, W_2)$ をもつ。その上 $P(T_s, T_t)$ は (25) 式での Z_1 と Z_2 との同じ相指数の T_s, T_t を含む。ゆえにこの自己相関は要求された形のものである。(27) が自己相関として正しく書かれうるためには、 σ^2 の条件を緩める。即ち $P(e^{i\omega}, e^{i\omega'})$ は Fourier 展開が可能ということである。

(27) 式はトランセクトの場合の一方角表示に対応する。しかしそれは便利なものではないから、推定のような問題については一般にほとんどの模型について研究する方がよい。1つの点で二次元の場合には一次元の場合と同じことが成立しない。すなわち有限自己相関 (finite autoregression) の一方角模型による表現はまた有限自己相関であるということである。

このことは次の簡単な例で確かめられる。有限自己相関

$$(1+B^2) z_{st} = B(z_{s+1,t} + z_{s,t+1} + z_{s,t-1}) + \epsilon_{st} \dots (28)$$

は無限一方角模型による表現

$$\begin{aligned} z_{st} = & B z_{s+1,t+1} - B^2 z_{s+1,t} + B(1-B^2) \sum_{j=0}^{\infty} B^j \epsilon_{s+1,t-j} \\ & + \epsilon_{st} \dots (29) \end{aligned}$$

更に一方角模型による表現の実験的有用さは、それがパラメーターの交換 (change) を単純化するということである。(同様なパラメーターの交換は普通積分 (2) の計算から、あるいはその二次元での同等性から示されるものではあるが) 二次元模型の大部分ではそれが明らかでないものであつても、非常に複雑で、それをやる価値はあまりない。例えは模型

$$\epsilon_{st} = \alpha (\epsilon_{s+1,t} + \epsilon_{s-1,t} + \epsilon_{s,t+1} + \epsilon_{s,t-1}) + \epsilon_{st} \dots (30)$$

の一方角模型での表現を計算すれば、 $\log \{1/2(z_1 + z_1^+ + z_2 + z_2^+)\}$

の Fourier 展開の係数を計算しなければならぬ。これは通常積分で表わすより他に方法がないのである。

7 抽出理論

この節においては、離散的過程の一般な推定式、適合度の検定およびパラメーターの推定値の近似的な分散および共変動を導く抽出理論を述べる。理論の展開の大要は、多方向性に特別な注意を払わなければならぬという事実を別として、前に出した時系列解析 (Whittle, 1951-3) と似通っている。考察した特別の過程は定非一決定的でその $[F(\omega_1, \omega_2)]$ は任意の実数 ω_1, ω_2 に対して 0 とならないようなものである。したがってこの過程は自己回帰として表現可能である。このように最尤基準で得られた推定式は、変量が正規分布をしない場合には最小二乗推定方程式と考えることができる。しかしこのような場合に残りの結果が成立する限界はより広い研究の課題である。

m 個の観測値 ε_{st} ($s = 1, 2, \dots, m; t = 1, 2, \dots, n$) の組があるものとする。そうするとラグ (lag) j, k の経験的共変動は

$$C_{jk} = \frac{1}{m\pi} \sum_{s=1}^{m-j} \sum_{t=1}^{n-k} \varepsilon_{st} \varepsilon_{s+j, t+k} \dots (31)$$

である。我々は端末効果 (end-effect) は一概に無視することにするから、(31)では $(m-j)(n-k)$ で割る代りに $m\pi$ で割った。更に次の量

$$\begin{aligned} f(\omega_1, \omega_2) &= \frac{1}{m\pi} \left\{ \left[\sum_{s=1}^m \sum_{t=1}^n \varepsilon_{st} \cos(s\omega_1 + t\omega_2) \right]^2 + \left[\sum_{s=1}^m \sum_{t=1}^n \varepsilon_{st} \sin(s\omega_1 + t\omega_2) \right]^2 \right\} \\ &= \sum_{j=-m}^m \sum_{k=-n}^n C_{jk} \cos(j\omega_1 + k\omega_2) \\ &= \sum_{j=-m}^m \sum_{k=-n}^n (c_{jk} \cos(j\omega_1 + k\omega_2) + d_{jk} \sin(j\omega_1 + k\omega_2)) \dots (32) \end{aligned}$$

を定義する。これは (23) と比較すると経験的スペクトル函数であることがわかる。またこれは実際 Schuster のペリオドグラム

この節の終りて我々は次の基本的な結果を証明することにする。

すなわち、もし変量 ε_{st} が、平均値 0 で正規に分布し、スペクトル函数 $F(\omega_1, \omega_2)$ をしつと考えられる型の定常過程から生成されたものであれば、これらの同時的尤度 (joint likelihood) は端末効果を無視すれば

$$P(\varepsilon) = \frac{1}{(2\pi V)^{\frac{1}{2} m n}} \exp \left\{ -\frac{m\pi}{8\pi^2} \int_0^{2\pi} \int_0^{2\pi} \frac{f}{F} d\omega_1 d\omega_2 \right\} \dots (33)$$

で与えられる。ここで

$$V = \exp \left\{ \frac{1}{4\pi^2} \int_0^{2\pi} \int_0^{2\pi} \log F d\omega_1 d\omega_2 \right\} \dots (34)$$

(33) の対数をとれば F のパラメーターの最尤推定値は

$$\hat{F} = \frac{1}{4\pi^2} \iint \left(\log F + \frac{f}{F} \right) d\omega_1 d\omega_2 \dots (35)$$

を最小にしてえられることがわかる。変量が正規分布と考えられない場合でも、(35)式を最小とすれば最小二乗推定値がえられる。(Whittle, 1953, p.132)

実際は (35)式は推定方程式の最も便利な形を導くものではない。これはペリオドグラムよりむしろ自己共分散による \mathcal{L} の式

$$\mathcal{L} = \frac{1}{4\pi^2} \iint \log F d\omega_1 d\omega_2 + \sum_j \sum_k C_{jk} C_{jk} \dots (36)$$

によつてうるからである。

ここで係数 C は

$$\sum_j \sum_k C_{jk} Z_1^j Z_2^k = \frac{1}{F(\omega_1, \omega_2)} \dots (37)$$

で与えられる。

いま (18) の形の確率差方程式 (stochastic difference equation) によつて生成される過程の特別な場合を考えることにすれば

$$F = \frac{v}{L(Z_1, Z_2) L(Z_1', Z_2')} \dots (38)$$

(36) にこの関係式を代入すれば

$$\mathcal{L} = \log v + \log k + v / m n v \dots (39)$$

ここで

$$U = \sum_{T=1}^m \sum_{t=1}^n \epsilon_{St}^2 = \sum_{T=1}^m \sum_{t=1}^n \{L(T_s, T_b) \epsilon_{St}\}^2 = \sum_j \sum_k C_{jk} C_{jk} \dots (40)$$

および
$$\log k = \frac{1}{4\pi^2} \iint \log \{L(Z_1, Z_2) L(Z_1', Z_2')\} dw_1 dw_2 \dots (41)$$

無関係なパラメーター (irrelevant parameter) v に用いる (39) の最小値は $\log (kU/mn)$ であるから L のパラメーターに関して最小とすべき量は $\log U$ である。

これは今の様なタイプの研究においては最も有用な結果である。最小二乗推定値は普通の通り、'残差平方和倍' したパラメーター L の函数を最小とすることによって得られる。

(41) から $-\frac{1}{2} \log k$ は $(L, L'; L, L')$ の二重 fourier 展開の絶対項 (absolute term) と考えられることかわかる。屢々問題となるのはその評価である。

F は未知函数 $\theta_1, \theta_2, \dots, \theta_n$ (θ はこれらの中に含まれている。普通はその値は未知である) このとき $\theta_1, \theta_2, \dots, \theta_n$ の最小二乗推定値の近似的共変動行列が

$$\frac{2}{mn} \left[\frac{1}{4\pi^2} \iint \frac{\partial \log F}{\partial \theta_j} \frac{\partial \log F}{\partial \theta_k} dw_1 dw_2 \right]^{-1} \dots (42)$$

であるというよく知られた論法がえられる。

(Whittle, 1953, p. 35; ここでの場合と同様な応用について) ここでまた p 個のパラメーターにあてはめを行なったときの L の最小値を $(kU)_n$ と仮定し、 q 個のパラメーターを追加してあてはめるときこの量が $(kU)_{n+q}$ になるものとする。そうすると、最初の n 個のパラメーターの仮説が正しいとすると

$$\psi^2 = (mn - n - q) \log \frac{(kU)_n}{(kU)_{n+q}} \dots (43)$$

は近似的に自由度 q の χ^2 分布をすることを用いて、 q 個のパラメーターを補足したことによってえられた改善を決定できる。

この論法も前に用いたものである (Whittle 1952, 1953)

(33) の証明

考えている過程は一方模型による表現 (27) をもつものとする。 mn 個の残差 ϵ'_{St} ($s=1, 2, \dots, m; t=1, 2, \dots, n$) の同時的頻数函数 (joint frequency function) は

$$h(\epsilon') = \frac{1}{(2\pi)^{\frac{1}{2}(mn)}} \exp \left\{ -\frac{1}{2} \sum_s \sum_t (\epsilon'_{St})^2 \right\} \dots (44)$$

である。(27) が同準であるような一次変換を行なつて、次のような ξ_{St} の函数の式をうる。

$$\begin{aligned} P(\xi) &\approx \frac{A^{mn}}{(2\pi)^{\frac{1}{2}(mn)}} \exp \left\{ -\frac{1}{2} \sum_s \sum_t \{P(T_s, T_b) \xi_{St}\}^2 \right\} \\ &\approx \frac{A^{mn}}{(2\pi)^{\frac{1}{2}(mn)}} \exp \left\{ -\frac{mn}{2} \sum_j \sum_k C_{jk} C_{jk} \right\} \\ &\approx \frac{A^{mn}}{(2\pi)^{\frac{1}{2}mn}} \exp \left\{ -\frac{mn}{8\pi^2} \iint \frac{1}{F} dw_1 dw_2 \right\} \dots (45) \end{aligned}$$

ここで A は (27) の ϵ_{St} の係数である。(24) (25) から

$$A = e^{-\frac{1}{2} d_{00}} = \exp \left\{ -\frac{1}{8\pi^2} \iint \log F dw_1 dw_2 \right\} \dots (46)$$

これを (45) に代入して (33) が得られる。

(45) の近似的な等式の性質はしつと正確に説明することかできる。(45) の右辺を P' と書けば、得られた関係は $\log n$ と $\log n'$ が近似的に等しいということである。(端末効果の無視によって exponent の項を無視することになる)

この種の近似は全く十分なものである。その理由は尤度の対数に対する相対的なオーダー n^{-1}, n^{-1} の項を加えても、 mn の大きな値では最小推定値あるいは尤度比の有意点には大きい影響がない。

8 数値例

研究さるべき二組のデータはいずれも我々が論じてきた特別な模型のそのままの例 (picture book example) になつていないが、しかしこれらの理由によつて更に価値あるものである。

オ1のデータは Mercer と Hall (1911) が実行した小麦 (Wheat) の収量の一様性試験に用いるものである。

これは 20 x 25 の矩形に配列された 500 個の 11 x 10.82 ft のプロットを含むもので、プロット合計が観測値となっている。相関面 (Correlation field) のオ1象限とオ4象限が Table 1 に与えられている。

オ2とオ3象限は

$$\phi(-j, -k) = \phi(j, k)$$

なる関係で埋めることができる。ここで象のつくことは、南北軸 (S) に沿った相関は東北軸 (E) に沿った相関よりもかなり大きいこと、少くともプロットが正方形でないということが多分その説明となるものである。

相関は一般に北西方面より北東方向で高く、はつきりした方向の影響があることを示している。もう一つの注目すべき事実は原点から離れるにしたがって相関が著しく減少しないで最低にまで下つてから再び増加することである。距離が小さい場合にはこのような効果を生み出している近くの植物間の競合を考慮することもできるが、しかしこの場合もつと真実らしく思える説明は耕転された圃場でよく見られる "肥沃度の波" が存在するということである。(例えは Neyman, 1952, p. 75)。

Table 2 に考察したき々な scheme を、あてはめの係数および対応する λ , λ_0 の値とともに要約した。

λ_0 はすべて 0.7 に近いから、我々は大きさ λ は次にいうことができる。すなわち、新しいパラメーターの導入によつて λ_0 が Δ だけ減少したなら、(43) 式によつて

$$3.841 < 500 \log \left(\frac{0.7 + \Delta}{0.7} \right), \text{ すなわち } \Delta > 0.0054$$

なら Δ は 5% 水準で有意である。

二つあるいは三つのパラメーターを導入したとき、対応する限界は 0.0084 と 0.0110 である。このようにして、仮説 1 と 3 は適合度

に有意な差がない。しかし 1 と 4 とは非常に明瞭な差がある。

1 と 5 のような仮説はこの方法では正理に比較することかできない。なぜなら何れの仮説と他のもの、特別な場合ではないからである。しかし読者は 1 のあてはまりは 5 のそれよりはるかに優れていると言いたいように思うであろう。したがって、この結論はそれぞれを順次しつと一般な仮説 Γ と比較すれば証明される。

この結果で意外に思える点は簡単な一方向の scheme 1 が対称な二次 (second order) scheme 5, 6 よりずつとよくデータに適合していることである。この節の論題の大きな部分は、空間的変遷は一般に多方向的であるということであつたが、我々の最初の例では一方向型が有力である。

このことの理由は、Table 1 の相関を再検討すれば明らかとなる。これらの相関はラグが 1 つ導入されると急激に小さくなる (S 軸上で 0.52 も軸上で 0.29 に) しかしこれはラグが大きくなるにつれて非常にゆるやかにしか減少しない (二単位のラグに対応する数値は 0.41 と 0.15 である)。しかし λ_0 で示すように、5 または 6 のような scheme のトレログラムは、原点から右の方へなだらかに下降していて原点上原点ではその導関数は 0 となっている (cf. Fig. 2)。したがって 5 と 6 も観測値に適合しないことは当然である。

Table 1 小麦のデータに対する自己相関

τ	$S=0$	$S=1$	$S=2$	$S=3$	$S=4$
-3	0.1880	0.1602	0.1509	0.1296	0.1352
-2	0.1510	0.0234	0.0020	-0.0137	-0.1039
-1	0.2923	0.1853	0.1349	0.0788	0.0378
0	1.0000	0.5252	0.4055	0.3639	0.3561
1	0.2923	0.2352	0.1799	0.1205	0.1399
2	0.1510	0.1235	0.0999	0.0749	0.0859
3	0.1880	0.1935	0.2483	0.2415	0.2284

観測されたコレログラムのこの性質については少なくとも二つの説明が可能である。

Table 2 あてはめた模型

模型No	$L(T_s, T_t)$	k	U	kU
1	$1 - 0.488 T_s - 0.202 T_t$	1	0.6848	0.6848
2	$1 - 0.483 T_s - 0.179 T_t$	1	0.6940	0.6940
3	$1 - 0.492 T_s - 0.211 T_t + 0.019 T_s T_t$	1	0.6845	0.6845
4	$1 - 0.402 T_s - 0.168 T_t - 0.172 T_s^2 - 0.092 T_t^2$	1	0.6564	0.6564
5	$1 - 0.159 (T_s + T_s^{-1} + T_t + T_t^{-1})$	1.1240	0.6508	0.7304
6	$1 - 0.213 (T_s + T_s^{-1}) - 0.102 (T_t + T_t^{-1})$	1.1332	0.6217	0.7045
7	$1 - 0.488 T_s + 0.030 T_s^{-1} - 0.202 T_t - 0.034 T_t^{-1}$	0.9843	0.6816	0.6709

基となる scheme は例えば土地の傾斜、主風の方角等のために一方向的であるかもしれない。逆に我々は観測値か成育の点観測値 (point observation) ではなく、ある地域上で加え上げ (integrate) られ成育の観測値だということをおぼろげに忘れない。このような加え上げはオラケの自己共変動を他のものよりも大きくさせるであろうから、6のような scheme のコレログラムでこのような歪曲が行なわれれば、観測されたコレログラムとは違ったものとなる。しかしそのような効果を考えることは当面の問題と離れすぎるから、我々は観測値の表面上の値でとり扱うこととしよう。scheme 1 は距離の増加に伴う相関の低下、増大を説明しないから、全く不適當である。一次 (first order) の項と同方向の二次 (second order) の項を入れればかなりあてはまりがよくなる。scheme 4 を見よ。

1~4のような scheme では $k=1$ だから最小二乗回帰法で直接にあてはめが可能である。

しかし5~7のような多方向的 scheme ではあてはめはもつと面倒である。

U は (40) によれば、直接計算できる (41) から k を求める場合には新しい問題が生ずる。ここでの計算には他にもつとよい方法があるかもしれないが、たまたまありふれた級数展開の方法と数値積分を用いた。このようにして scheme

$$\epsilon_{st} = \alpha \epsilon_{s+1,t} + \beta \epsilon_{s-1,t} + \gamma \epsilon_{s,t+1} + \delta \epsilon_{s,t-1} + \epsilon_{st} \dots (47)$$

については

$$\log(k) = \log \{ 1 - \alpha Z_1 - \beta Z_1^{-1} + \gamma Z_2 + \delta Z_2^{-1} \} \text{の absolute term の } -2 \text{ 倍}$$

$$= \sum_{j=1}^{\infty} \sum_{k=0}^j \frac{(2j)!}{j! k! (j-k)!} (\alpha\beta)^k (\gamma\delta)^{j-k} \dots (48)$$

$\alpha\beta = \gamma\delta$ なるときは

$$\log(k) = \sum_{j=1}^{\infty} \frac{1}{j} \binom{2j}{j} (\alpha\beta)^j \dots (49)$$

となる。

展開式 (48) (49) は $\alpha\beta$ $\gamma\delta$ が非常に小さい場合には有用である。しかしこれらの値がその最大値 $1/16$ に近づくとき、収斂が緩慢となる。

Table 3 には $\theta = \sqrt{\alpha\beta}$ の幾つかの値に対する (49) の $\log(k)$ の値を示した。これらの値は数値積分で求めたか、変化は非常にゆるやかで十分公正な開示が可能である。

Table 3 $\log(k)$ の修正係数の例

θ	0.00	0.05	0.10	0.15	0.20	0.22	0.25
$\log(k)$	0.0000	0.0076	0.0120	0.0170	0.0228	0.2656	0.4406

Table 2 の模型は単にパラメーターの trial value を代入し、対応する kU を計算してあてはめた。そうして改善された値は、核物体 (nucleoid) から近似によってその最小値の位置をきめるといいうようにしてこの操作を繰返した。この作業は逐次であるが、

しかし、データそのものか選出された努力を考ればそう不当なものではない。

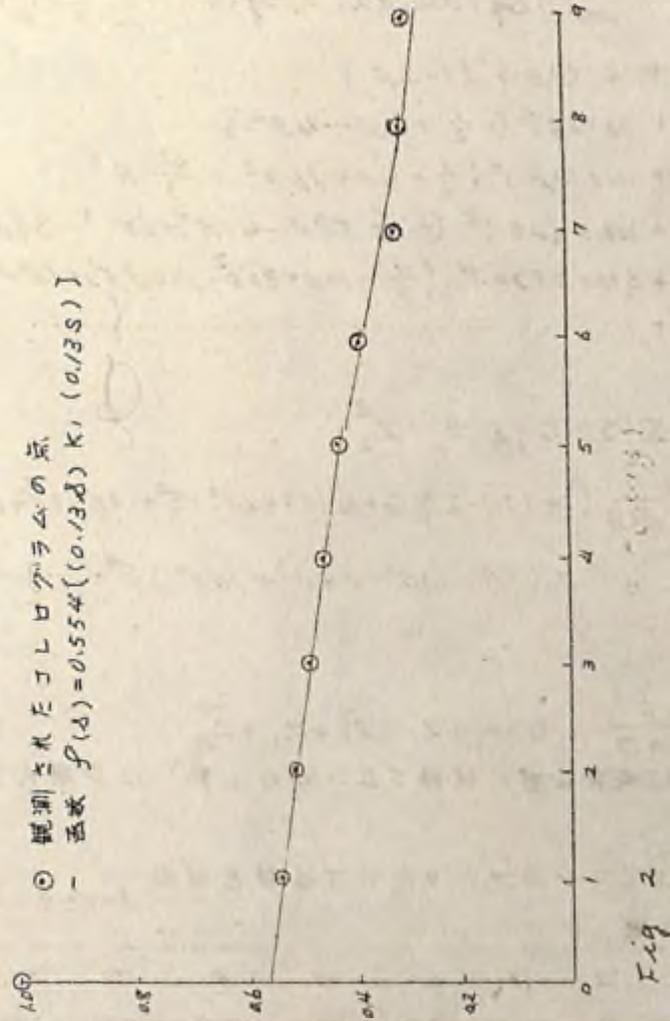
我々のオスの例は *Batchelor* と *Reed* (1924) が 1000 本のオレンジについて行った一様性試験 (*uniformity trial*) から得られたもので樹は 20 x 50 の矩形格子に配置されている。圃場の相関 $\rho(S, t)$ は *Table 4* に、また S 軸に沿った相関は *Fig 2* にプロットされている。二つの最良の相関はもとの圃場の高度の対称性を反映して、非常に似通った形態を示している。興味のある点は隣接した S, t の値を通つて動くとき大部分の相関 $\rho(S, t)$ の値はなめらかに変化しているのに、原点では非特にはつきりした不連続性が存在することである。

このように相関面 (*Correlation surface*) は假点で大きい突起をした広いドーム (*dome*) のまわりを滑らかに取りまいたような形をしている。これは個々の樹の収量 Y_{st} が

$$Y_{st} = \epsilon_{st} + \eta_{st} \dots \dots \dots (50)$$

のように表わされることを示している。ここで ϵ_{st} は滑らかなコレログラムをしつ過程にしたかう。一方異なる η_{st} は互いに (あるいは ϵ_{st} と) 無相関であるから η のコレログラムはただ原点の "突起" のみからなっている。(50) の自然的な解釈が心に浮かんでくる: すなわち各樹のまわりの土地の肥沃度を表わし、多分 (30) のような単純な対称的な過程 (*symmetric process*) にしたがい、また η は個々の樹に圃場の結実性を表わすということである。コレログラムを見たとき、読者は土壌と木の交異性は全葉動に対してそれぞれ 56 および 44% の寄与を与えているということもおれぬ。しかし前の例と同じく樹はまわりの土壌の肥料を "integrate" するかも知れぬと いうことを忘れてはならない。

(30) のように単純な自己相関をあてはめるのと置つて (50) + (30) のような複合した *scheme* をあてはめるのは非特に面倒である。



$$\text{Var}(\eta) = A \quad \text{Var}(\epsilon) = B \quad \text{とする}$$

$$F(W_1, W_2) = A + \frac{B}{(1 + \delta(Z_1 + Z_1' + Z_2 + Z_2'))^2} \dots \dots (51)$$

をうる。

推定の目的に必要な量は $\log V$ と F^{-1} を展開したときの係数である。(36)を見よ)

我々は

$$\begin{aligned} \log V = & \frac{1}{4\pi^2} \iint \log F dw_1, dw_2 = \log(A+B) + \sum_{j=1}^{\infty} \frac{1}{j} \binom{2j}{j} \alpha^{2j} t \\ & + 4(\alpha\beta)^2(1-2\beta) \\ & + 36(\alpha\beta)^4(-\frac{1}{2} + 4\beta - 4\beta^2) \\ & + 400(\alpha\beta)^6[\frac{1}{3} - 6\beta + 16\beta^2 - \frac{32}{3}\beta^3] \\ & + 4900(\alpha\beta)^8[-\frac{1}{4} + 8\beta - 40\beta^2 + 60\beta^3 - 32\beta^4] \\ & + 63504(\alpha\beta)^{10}[\frac{1}{5} - 10\beta + 80\beta^2 - 224\beta^3 + 256\beta^4 - \frac{128}{5}\beta^5] \\ & + \dots \dots \dots (52) \end{aligned}$$

および

$$\begin{aligned} F^{-1} = & \sum \sum C_{j,k} z_1^j z_2^k \\ = & \frac{1}{A+B} \{ (1 + (2\beta-2)S + (1-5\beta+4\beta^2)S^2 + (4\beta+2\beta^2+8\beta^3)S^3 \\ & + (-\beta+13\beta^2-28\beta^3+16\beta^4)S^4 + \dots) \dots \dots (53) \end{aligned}$$

もうる。

ここで

$$\beta = \frac{A}{A+B}, S = \alpha(z_1 + z_1^{-1} + z_2 + z_2^{-1})$$

この scheme は単純な自己回帰でないから、 F^{-1} は多項式でなく、

Table 4 オレンジのデータに対する自己相関

$\lambda=0.9, \tau=4.4$

δ の値

τ	0	1	2	3	4	5	6	7	8	9
-4	0.3912	0.3902	0.3771	0.3581	0.3370	0.2923	0.2549	0.2546	0.2440	0.2131
-3	0.4609	0.3956	0.3270	0.3696	0.3484	0.3353	0.3137	0.2834	0.2894	0.2870
-2	0.6667	0.4059	0.3982	0.3470	0.3243	0.3225	0.3250	0.2801	0.2721	0.2752
-1	0.5462	0.4669	0.4336	0.3834	0.3880	0.3761	0.3495	0.2914	0.2567	0.2606
0	1.0000	0.5403	0.5052	0.4800	0.4585	0.4233	0.3960	0.3246	0.3210	0.3150
1	0.5462	0.4658	0.4280	0.3751	0.3914	0.3622	0.3716	0.2948	0.2958	0.2925
2	0.4667	0.3858	0.3799	0.3459	0.3603	0.3475	0.3401	0.2905	0.2766	0.2230
3	0.4609	0.4190	0.3716	0.3764	0.3509	0.3336	0.3369	0.2765	0.2433	0.2741
4	0.3912	0.3543	0.3882	0.3392	0.3537	0.3446	0.3305	0.2893	0.2685	0.2321

また係数 $C_{j,k}$ も有限でない。
 したがって今の場合、我々は(36)の和 $\sum \sum C_{j,k} C_{j,k}$ の十分な近似をうるには大体 τ/λ までの相関を含まなければならぬことがわかる。しかし相関は僅かも $\tau=4$ までしかえられない。これ以上続けて求めることは非常に困難である。

一般に scheme が単純な自己回帰なものでは観測されたコレログラムが非常に急速に小さくなる場合に応じて最小二乗法を採用することができる。今の問題のような場合にはむしろせよくなめではめの方法、例えば観測された方程式と理論的自己相関係数を用いなければならないようになる。

一方減少の率がゆるやかたということは、この離散的な模型が連続な模型で近似しうるという利点を目指す。これを単純化するのに最も多く用いられる方法である。このように scheme (30) において λ が 0.25 に接近した値なら、これは自己共変関数 (autocovariance function)

$$f(\tau) = \text{const} \cdot K_1(K\tau) \dots \dots \dots (54)$$

を τ 連続な scheme で近似しうることを 9 において示されている。

ここで τ は考えられた二点間の距離であつて K_1 は修正された Bessel 函数で

$$K = \sqrt{\left(\frac{\tau}{2} - 4\right)} \dots \dots \dots (55)$$

である。コレログラム $(F; g, s)$ の薄い線は帯数 $K\tau_1$ と τ_2 で一致するように調整したときの函数 (55) の値を示す。 $\tau_2 = 0.2489$ のとき $K = 0.13$ であることがわかる。この一致は注目し得る。しかしこれは相当割引して考えなければならない。というのは帯数 K の単純減少函数も端点 (end-point) さえ一致するようにすれば観測値に相当よくあてはまるものだからである。

しかし例えば同じような型の指数曲線をあてはめるなら、指数曲線は中間であまり下りすぎて全然あてはまりが悪い。

9. 特別な過程

データの組にあてはめを行うとき模型の性質についてそう厳密に知っている必要はない(前節の例を参照)。しかしこれらの性質を知ることは優れた模型をうるための最も近道である。この場合「性質」の内容は、普通ゴレログラムと同義である。なぜなら与えられたデータの性質を改む場合はゴレログラムを用いることが最も多いからである。

ここで我々は二、三の特殊な模型とそのゴレログラムをどちらかというど無系統な形で論じてみよう。というのは厳密な性質の結果も幾つか得ることかできるからと思う。

興味ある問題で最も簡単なものは解

$$E_{st} = \sum_{j=0}^s \sum_{k=0}^t (j+k) \alpha^j \beta^k E_{s+j, t+k} \dots (57)$$

である。

E_{st} を独立変数とするこの過程が定常であるためには

$$|\alpha| + |\beta| < 1$$

でなければならぬ(55をみよ)

自己共変動は

$$(1 - \alpha Z_1 - \beta Z_2)^{-1} (1 - \alpha Z_1^{-1} - \beta Z_2^{-1})^{-1}$$

から生成され、相関は

$$p(s, 0) = A^s \quad p(0, t) = \beta^t$$
$$p(s, t) = \left(\sum_{i=0}^s \binom{s+t-i}{t} \alpha^{s-i} \beta^i A^i + \sum_{k=0}^t \binom{s+t-k}{s} \alpha^{s-k} \beta^k \beta^k \right) \dots (58)$$

となる。

$$\left. \begin{aligned} \text{ここで } A &= \frac{1 + \alpha^2 - \beta^2 - \alpha}{2\alpha}, & B &= \frac{1 + \beta^2 - \alpha^2 - \beta}{2\beta} \\ \Delta &= \sqrt{\{(1 + \alpha + \beta)(1 + \alpha + \beta)(1 - \alpha + \beta)(1 - \alpha - \beta)\}} \end{aligned} \right\} \dots (59)$$

である

(58)式は $\alpha, \beta > 0$ について成立つ。

Scheme (56) の連続の場合への類推は一次確率差方程式 (first-order stochastic differential equation) となり形式的には

$$\left(\alpha \frac{\partial}{\partial X} + \beta \frac{\partial}{\partial Y} + \gamma \right) E(X, Y) = E(X, Y) \dots (60)$$

のように書かれる。

しこの軸に対して再 t_{un} (β/α) だけ回転させた新しい座標軸 $X=0, Y=0$ に対して平面を回転させると、この関係は

$$\left(\sqrt{(\alpha^2 + \beta^2)} \frac{\partial}{\partial X} + \gamma \right) E'(X, Y) = E'(X, Y) \dots (61)$$

のように書けることかわかる。

ここで $E'(X, Y) = E(X, Y)$ $E' = (X, Y) = E(X, Y)$ である。

すなわち scheme (60) は t_{un} (β/α) の方向に沿って走る Markov 過程の系列を考慮することかできる。これは互いに独立である。このことは一次の模型 (first order scheme) の退化 (degenerate) する性質を系している。

最も簡単な二次の (second order) scheme (すなわち最も簡単な非退化 scheme) は系統的自己回帰

$$E_{st} = \alpha (E_{s+1, t} + E_{s-1, t} + E_{s, t+1} + E_{s, t-1}) + E_{st} \dots (62)$$

である。

$$\left[\Delta_s^2 + \Delta_t^2 + \left(\Delta - \frac{1}{\alpha} \right) \right] E_{st} = E_{st} \dots (63)$$

のように書かれる連続的関係に対する類似の stochastic laplace equation は (ここで Δ は central difference operator)

$$\left[\left(\frac{\partial}{\partial X} \right)^2 + \left(\frac{\partial}{\partial Y} \right)^2 - K \right] E(X, Y) = E(X, Y) \dots (64)$$

となることは明らかである。

この連続的な関係は (関係式 (56), (60) などの場合と同じく) 本質的に技巧的な性質がないため古くから取扱われている。

Scheme (62) の正確な結果は解られている。(Van der Pol

& Bremner, 1950; Stöhr, 1950) がこれらは簡単ではない。公式 (Titchmarsh, 1948, p 201) から

$$\frac{1}{4\pi^2} \iint \frac{e^{i(xw_1 + yw_2)} dw_1 dw_2}{(w_1^2 + w_2^2 + k^2)^{\mu+1}} = \left(\frac{r}{2k}\right) \frac{K_\mu(kr)}{\Gamma(\mu+1)} \dots (65)$$

$$(r = \sqrt{x^2 + y^2}),$$

であるから scheme (64) は

$$\xi(x, y) = \int \int_{-\infty}^{\infty} \epsilon(x+x', y+y') K_\mu(kr) dx' dy' \dots (66)$$

$$\phi(x, y) = \frac{Y}{2k} K_\mu(kr) \dots (67)$$

となる。

ここで K_ν は ν 級の修正された $\text{order } \nu$ の Bessel 函数である。

$\lim_{r \rightarrow 0} K_\mu(kr) = 1$ だから

(67) に対応する相関係数は

$$\rho(r) = K_\mu(kr) K_\mu(kr) \dots (68)$$

である。

scheme (64) は一般論 (general) second-order 確率差分方程式 (stochastic difference equation) の特別の場合である。通常の楕円的 (elliptic) 放物線的 (parabolic) 双曲線的形式 (hyperbolic form) への分類にしたがえば、(64) は原点を中心とする円形式 (circular form) と考えることができる。残りの second order scheme を研究することは理論的にも実際的にも甚だ興味がある。

相関函数 (68) は一次元の場合の指数式 $e^{-|x|}$ と同じように、二次元での "elementary" 相関と考えられるから興味がある。どちらの相関曲線も単調減少であるか、(68) は原点で平でありまたその減少の率が指数式よりも遅いという点で異なっている。

指数式はそれ自身一次元での当然の差次であることが証明されており、また観測された曲線も同様な単調減少を示すという大きな理

由から、二次元のコレログラムを指数式の和で表わすという多くの試みかなされてきた。

しかし前節で明らかのように、二次元の場合 指数式は必ずしもよくないし また最後の節の例によつて、 K_μ 函数は指数型よりも観測値によくあてはまることか示されている。

指数的相関をもつような二次元の過程も作れるか、それは非常に人工的なものである。例えば Matern (1947) は相関函数 $\exp - 2\sqrt{x^2 + y^2}$ に対応するスペクトル函数は $(w_1^2 + w_2^2 - 2^2)^{-3/2}$ であることを示した。

このようなスペクトル函数をもつ最も簡単な過程は形式的に

$$\left[\left(\frac{\partial}{\partial x}\right)^2 + \left(\frac{\partial}{\partial y}\right)^2 - 2^2 \right] \xi(x, y) = \epsilon(x, y) \dots (69)$$

と書くことができるが このような関係を等しく自然的なメカニズムを具体的に示すことは難しい。

by P Whittle

(Biometrika, 1954, vol 41, p 434-449 21)

参 考 文 献

Bartlett, M. S (1946). Stochastic processes. Monograph N. Carolina lecture notes.

Batchelor, S. L and Reed, H. S (1924) Relation of the variability of yields of fruit trees to the accuracy of field trials. J Agric Res 12, 295-83

Samuel P J (1946) contribution to discussion on stochastic processes J. R Statist Soc B. 8 88-90

Dool J S (1944) The elementary Gaussian processes
 Am Math Statist. 15, 229-82.

Matern. B. (1947). Metoder att upphatta noggrannheten
 vid linje-och punktavsering medd.
 Skogsforskalnst. 36 no 1.

Mercer. W. B and Hall. A. H (1911) the experimental
 error of field trials. J Agric.
 Sci. 4, 107-32

Neyman J (1952) Lectures and Conferences on
 Mathematical Statistics and probability
 Graduate School U S Dep of Agriculture.

Quenouille. M. H. (1949) problems in plane sampling.
 Ann. Math Statist 20 355-75

Stohe. A. (1950) Über einige lineare partielle differenzgleichungen
 mit konstanten Koeffizienten Math. Nachr 3, 295-315

Fitchmarsh. E. C (1948) introduction to the Theory of Fourier
 integrals oxford university press.

Van der pol B. and Bremmer. H (1950) Operational Calculus.
 Cambridge University press

Whittle. P (1952) Tests of fit in time series Biometrika, 39, 309-18.

Whittle. P (1953) The analysis of multiple stationary time series
 J. R. Statist. Soc B 15, 125-39.

Wiener, N. (1949) The Extrapolation interpolation and
 Smoothing of stationary time-series. New York Wiley.

Williams. R. M. (1954) The choice of sampling interval for systematic
 samples from population with stationary correlation.
 (Unpublished paper)

Wald, H. (1938). A Study in the Analysis of Station-
 ary Time Series. Upsala.

論者注

Whittle の上記論文 (p 59-60) および Patankar の論文 (1954
 Bion. pp 450-462) ではともに Mercer & Hall (1911)
 の小麦のデータを理論の説明に用いている。一見したところでは
 Whittle の table の (i) $s=0$ の列と (ii) $t=0$ の行の
 系列相関は Patankar の相関の (i) "行に沿って、もとの全デー
 タについての結果" (table 5.2, p 462 の最後の列) を (ii)
 の列) と対応すべきように思われる。しかし着者の調べたところ、
 計算された値は異なる定義のものである。すなわち

$i = 1, 2, \dots, 25$ 東西方向の列を示す。
 $j = 1, 2, \dots, 20$ 南北方向の行を示す。

これは Whittle の K_{30} に対する共変動は

$$\sum_{i=1}^{25} \sum_{j=1}^{20} X_{ij} X_{i,j+5} - \sum_{i=1}^{25} \sum_{j=1}^{20-5} X_{ij} X_{i,j+5} \sum_{i=1}^{25} \sum_{j=1}^{20-5} X_{ij} / \{25(20-5)\}$$

で、分母についても同じである。しかし Patankar の用いた対応す
 る共変動は

$$\sum_{i=1}^{25} \left\{ \sum_{j=1}^{20-5} X_{ij} X_{i,j+5} - \sum_{j=1}^{20-5} X_{ij} X_{i,j+5} \times \sum_{i=1}^{25} X_{i,j+5} / (20-5) \right\}$$

で、分母でも同じである。
 このように Whittle の系列相関は全変動 (total variation)
 と共変動から計算されているが、Patankar の相関は列内 (または
 行内) 変動にもとずいている。予知される様に Whittle の相関は
 行および列内の変動だから Patankar のものより大きい。

12 系統的標本の平均値の分散

The Variance of the Mean of Systematic Samples

緒言

系統的標本抽出の方法はこれまで色々な分野、特に林業および生態学者によつて用いられてきた。これは、現場での便利さと調査資料を地図作成に用いる場合に無作為抽出法より優れていることによるものである。これらの用法は Hasel (1938) や Osborn

(1942) などの研究によつて正当化されたものである。彼等は無作為、層化無作為および系統的調査の抽出誤差を、これらの三方法で分析するためにとられた詳細な森林調査のデータについて研究した。これらの分析が、系統的抽出の効率が一様に最も高いということが示された。Osborne の論文では実験結果と妥当な一致を示す平均平方系列相関 (mean square serial correlation) を用いる一つの方法が与えられてはいるが、大体において彼等は系統的標本の誤差を推定する問題を未解決のまま残している。

ごく最近 Finney (1948) も森林調査のデータを用いて、系統的標本 (妥当な精度を得るために群の重ね合せ法 (overlapping method) によつて計算した) と無作為標本、および層当り一つまたは二つの要素をもつ層化無作為標本の分散を比較した。これらの場合において Finney は系統的標本の分散は層当り1個の観測値をもつ同じ大きさの層化無作為標本の分散と少ししか変わらないが、層の数を半分にして層当り2個の観測値をとつた層化標本で得られる分散よりは非常に小さくなることを示した。これは補足的な情報なしで分散の不偏推定値を与える最も有効な方法であると思われる。色々な模型 (model) に対する系統的標本の抽出分散の更に進んだ研究が有用であるように思われる。

Yates (1948) は多数の特別な場合について抽出分散を研究し、補助情報がある場合に系統的標本の平均値の分散を与える式あるいは

はある種の仮定のもとで誤差の限界を与える式を示した。Jowett (1952) は、観測値が定常過程 (stationary process) から導かれたという仮定のもとに系統的標本の分散を決定する方法を与えた。これは同じ分野における Cochran (1946) および Quenouille (1949) の研究から導かれたものである。

この論文においては、Jowett の方法に類似したやり方を展開するか、それよりは幾分強い仮定によるものである。我々はまた観測データのそれより小さい間隔および等しい間隔でとつた標本の分散式を導くであろう。

(これは Jowett の方法には含まれていない)。これは分散が単一の標本から導かれる場合も含まれる。これは要々言われている原則すなわち系統的標本それ自身からは誤差の推定値は得られないことと矛盾するものではない。というのは母集団について行なつた仮定は補助情報と同等であるからである。他にも数多くの仮定を続けることができたが、特定のものを送ぶについてはそれが多くの種類のデータにあてはまりそうだという証明がなければならぬ。またあつて考察する場合において実験と一致する結果を導くことが可能である。

この研究は、直線上にとつた点で、地上の種々の地被物 (various types of cover) の割合 (%) を定め観測を行なうとき、この点の間隔をどのように定めるかということについて行なつた North Canterbury Catchment Board (New Zealand) の担当者との討論から生れたものである。

この理論の主要部は、観測値の組が後述の離散的および連続的変数の何れについても成立する模型にあてはまりさえすれば、時間又は空間上で一様にとられた任意の観測値の系列に対して適用することからである。

便宜のためこれからはいつも観測値をとり直線の全体を意味するトランセクトと、その上の点についてを考へることにする。実験的研究については §4 に詳述する。

2. 分散の公式

トランセクトは長 \$kn\$ の点からなっているとす。標本は各番毎の点をとることによって抽出されるから、可能な標本の数は \$k\$ となる。おし点での変量を \$X_i\$ (\$i = 1, 2, \dots, kn\$) とする。母集団から抽出された系統的標本の平均値の抽出分散 \$\sigma_s^2(n, k)\$ (標本平均値の母集団平均値 \$\bar{x}\$ からの偏差の平均平方として定義される) の公式は、連続する三つの論文 (Madow, 1949, 1953) の最初のものにおいて Madow & Madow (1944) が与えた。

これは平方和の普通の分割と。

$$\sum_{i=1}^{kn} (X_i - \bar{x})^2 = \sum_{i=1}^{kn} \sum_{j>i} (X_i - X_j)^2 / kn$$

なる関係から都合のよい形で並びくことかできる。

すなわち、これから

$$\sigma_s^2(n, k) = \frac{1}{(kn)^2} \sum_{i=1}^{kn} \sum_{j>i} (X_i - X_j)^2 - \frac{1}{kn} \sum_{s=1}^{n-1} \sum_{i=1}^{k(n-s)} (X_i - X_{i+ks})^2$$

$$= \left(\frac{f(0)}{n} + 2\bar{f}_n \right) - \left(\frac{f(0)}{kn} + 2\bar{f}_{kn} \right) \dots \dots (1)$$

$$\left. \begin{aligned} \text{ここで } f(0) &= \sum_{i=1}^{kn} (X_i - a)^2 / kn \\ \bar{f}_n &= \frac{1}{n} \sum_{s=1}^{n-1} \frac{1}{kn} \sum_{i=1}^{k(n-s)} (X_i - a)(X_{i+ks} - a) \\ \bar{f}_{kn} &= \frac{1}{kn} \sum_{s=1}^{kn-1} \frac{1}{kn} \sum_{i=1}^{kn-s} (X_i - a)(X_{i+s} - a) \end{aligned} \right\} \dots (2)$$

ここで \$a\$ は任意の定数である。

\$k\$ が無限大に \$t \to 0\$ するものとし、また \$s = i/kn\$、これに対応する \$X_i\$ の値を \$X(s)\$、および \$\phi(t)\$ を

$$\phi(t) = \int_0^{1-t} (X(s) - a)(X(s+t) - a) ds \dots \dots (3)$$

と定義する。

極限において

$$\bar{f}_n \rightarrow \frac{1}{n} \sum_{s=1}^{n-1} \phi\left(\frac{s}{n}\right) = \bar{\phi}$$

$$\bar{f}_{kn} \rightarrow \int_0^1 \phi(t) dt = \bar{\phi}$$

$$f(0) \rightarrow \phi(0)$$

であるから

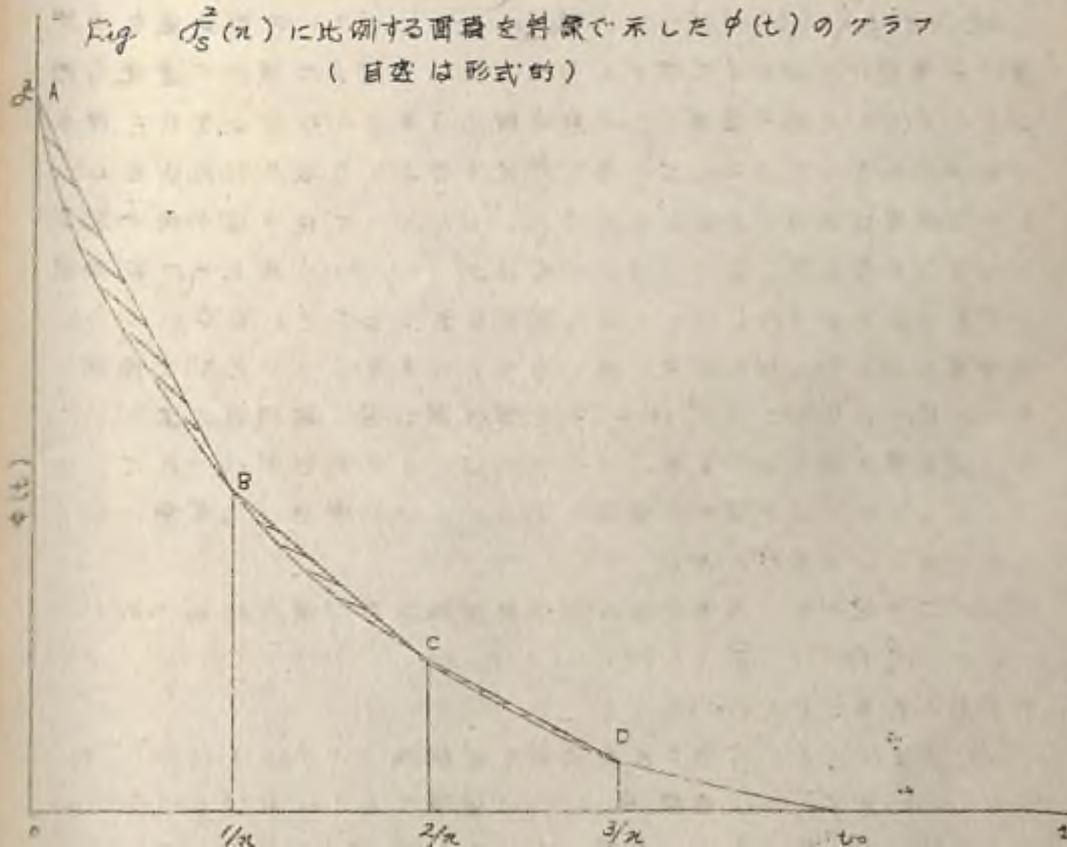
$$\sigma_s^2(n, k) \rightarrow \frac{\phi(0)}{n} + 2\bar{\phi} - 2\bar{\phi} = \sigma_s^2(n) \dots \dots (4)$$

\$\sigma_s^2(n)\$ のもう一つの形が Euler-Maclaurin の定理から得られる。この定理によればもし \$\phi(t)\$ が \$(2q+1)\$ 次までの連続な導関数をもてば

$$\sum_{s=0}^n \frac{1}{n} \phi\left(\frac{s}{n}\right) = \int_0^1 \phi(t) dt + \frac{1}{2n} (\phi(0) + \phi(1)) - \sum_{j=2}^q \frac{B_{2j}}{(2j)!} \frac{1}{n^{2j}} \phi^{(2j)}(0) - \phi^{(2q+1)}(1) + (-1)^{q+1} \frac{1}{n^{2q+2}} \int_0^1 P_{2q+1}(x) \phi^{(2q+1)}\left(\frac{x}{n}\right) dx$$

が得られる。

Fig. \$\sigma_s^2(n)\$ に比例する面積を斜線で示した \$\phi(t)\$ のグラフ (自変は形式的)



* Jonett (1952) はこの公式を少し変った形で与えている。Zuenille (1949) の与えたのは \$k\$ が \$n\$ と同じ位大きいときの式である。

ここで $\phi(x)$ は x の偶函数、 B_{2v} は Bernoulli の数で表わした $P_{2v+1}(x) = \sum_{e=1}^{\infty} \frac{\sin(2e\pi x)}{2^{2e} (2\pi)^{2e+1}}$ である。

この級数が収斂すると仮定すると

$$\sigma_S^2(x) = \frac{\phi'(1) - \phi'(0)}{6\pi^2} - \frac{\phi^{(3)}(1) - \phi^{(3)}(0)}{360\pi^4} + \frac{\phi^{(5)}(1) - \phi^{(5)}(0)}{15120\pi^6} + \dots \quad (5)$$

が得られる。

$\phi(x)$ が必要な次数まで微分可能であるという仮定は、一般には特別なトランセクト内の観測値によつて満足されない。しかし Cochran (1946) に従つて、特定のトランセクトの観測値を多変量の母集団から抽出した標本と考えるなら、個々の値が不連続な場合でも $\phi(x)$ の期待値等(この期待値は母集団から抽出された標本の全体にわたつてとられているが希望するような微分可能性をもつ)という仮定は正当であることが多い。したかつて我々は今後分散等の期待値を考えることにする。公式は $\sigma_S^2(x)$, $\phi(x)$ 等をその期待値へ置きかえても支障ないし、また混乱を生ずるおそれもないから、期待値に対しても同じ記号を用いることにする。この区別を強調したいところでは $\sigma_S^2(x)$, $\phi(x)$ 等の標本値、期待値のようになつた二つの値を考えることにする。この区別は (5) 式が用いられているようなところでは本質的に重要であるが他の場合には実際にあまり重大でないと思われる。

可換な方法から、大きさ n の無作為標本の平均値の抽出分散が

$$\sigma_r^2(n) = \frac{1}{n} (\phi(0) - 2\bar{\phi}) \dots \dots \dots (6)$$

で与えられることがわかる。

(4) および (6) に含まれる函数を単調減少の $\phi(x)$ について Fig 1 に示す。 $\bar{\phi}$ は曲線 $\phi(x)$ 下の面積で与えられ $\phi(0)/2n + \bar{\phi}$ は多角形 $OABCD \dots$ の面積であるから $\sigma_S^2(n)$ は斜線で示した面積の二倍に等しい。

* このグラフによる方法は Jowett (1952) が与えたものと同じである。しかしこれは時系列からとれているという Jowett の仮定は含んでいないことに注意せよ。

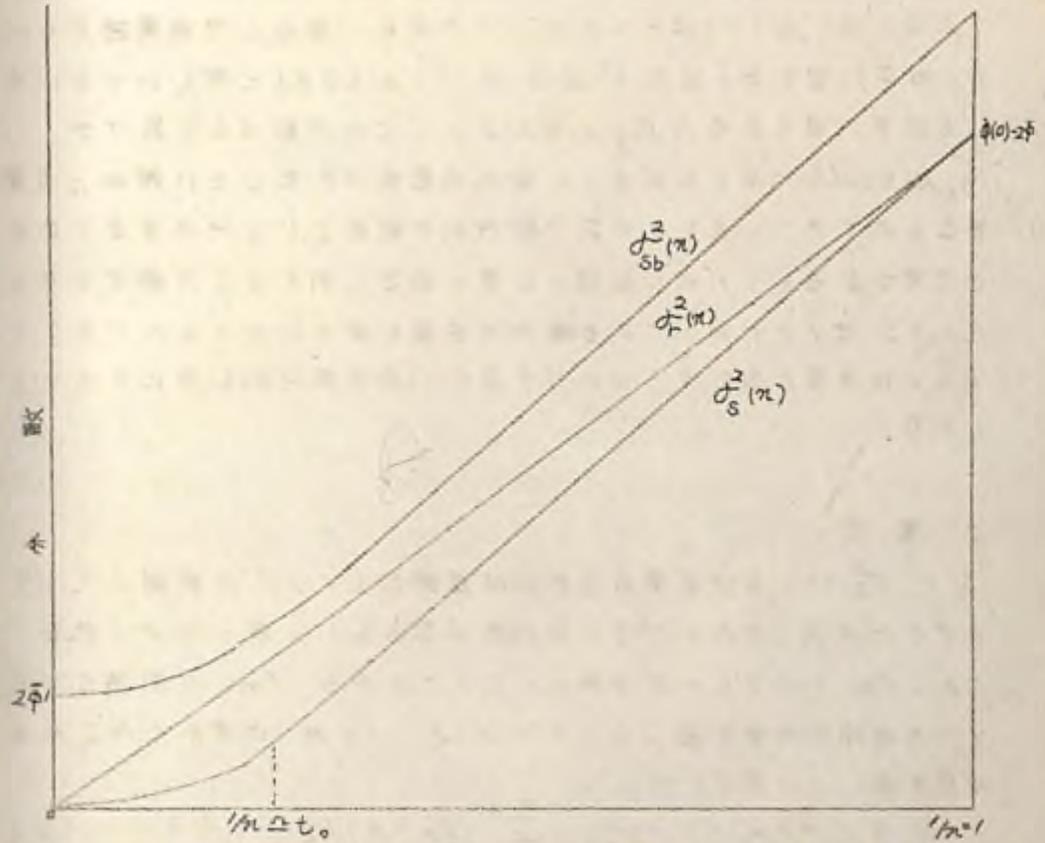


Fig 2 単調減少の $\phi(x)$ に対する相対的な分散の値

(5) から、もし展開が正しければ、まず $\sigma_S^2(n)$ は観測値間で距離 $(1/n)$ の二乗で変化する。一方 $\sigma_r^2(n)$ は $1/n$ とともに一次的に (linearly) 変化する従つて十分大きい n に対しては $\sigma_r^2(n)$ は $\sigma_S^2(n)$ よりも大きいから、系統的標本は無作為標本より有利である。もし $\phi(x)$ が 0 かまたはある点 x_0 を越えると一次的に変化するなら、抽出点 $(1/n)$ の間隔 h も。以上になるとき、 $\sigma_S^2(n)$ は $1/n$ とともに一次的に増加する (Fig 2 を見よ)

分散 $\sigma_S^2(n)$ は、同一トランセクトから抽出された全標本の平均値の変動性を表わす尺度である。 $\int_0^1 \mu(x) dx$ に対する標本平均値の偏差の平均平方を $\sigma_{Sb}^2(n)$ と表わすと

$$\sigma_{Sb}^2(n) = \sigma_S^2(n) + 2\bar{\phi}' \dots \dots \dots (7)$$

である。但し $\mu(t)$ はトランセクトの点を抽出した母集団点における平均値でありまた ϕ' は a が $\int_0^1 \mu(t) dt$ に等しいとおいたときの ϕ の値である (Fig 2 をみよ)。この分散はもし我々が $\int_0^1 \mu(t) dt$ に生じた疑わしい変化の意味づけをしたい場合には重要なものであり、また Σ の共分散行列が依然としてそのままであると仮定できる。これには証明が必要である。例えばこの論文で考えたトランセクトにおいてある植物の品種を徐々に他のものでおきかえると被覆度が変化するのはかりでなく、共分散行列も変化するかもしれない。

3 推定

3.1 $\sigma_0^2(x)$ を決定するためには理論的 $(0, 1)$ の範囲のすべての点での $\phi(t)$ がわかっている必要はない。我々は u/m 点を $(0, 1/m)$ の中から無作為にとり、これから $1/m$ の間隔ごとにとった m 点での観測値 X_r ($r = 1, 2, \dots, m$) の標本からこれを推定する。この標本に対して

$$\psi(u/m) = \frac{1}{m-u} \sum_{r=1}^{m-u} (X_r - a)(X_{r+u} - a) \dots \quad (8)$$

を定義する。 $\psi(u/m)$ は $\phi(u/m)$ の不偏推定値である。

もし $\phi(t)$ に対して特別な形 (例えば指数関数) を仮定できればむしろ我々は、 $\psi(u/m)$ が非偏に多数の点にもとづくことが保証される十分小さい u の値に対しては、その式をこれらの観測された $\psi(u/m)$ にあてはめ、因又は級数展開の式から $\sigma_0^2(x)$ を求めるのがよい。

これが不可能な場合にも、小さい u の値に対する $\phi(t)$ を定めることはむづかしくないから、結局級数式に必要な $t=0$ での微係数も容易に決定できる。しかし u が 1 に近くなると、含まれる項の数が少なくなりすぎて十分な推定値を得られなくなる。グラフによる方法にも同様な欠点があるが、それに加えてもし可能ならすべての点で $\phi(t)$ を計算するという手段を省略することが望ましい。

もし $\phi(t)$ がある値 t_0 で直線なら、グラフによる方法では

え。以上のところまで $\phi(t)$ を計算しても $\sigma_0^2(x)$ に影響はない。実際においては、重要と考えられる最大の抽出区間に等しい長さ上でそれが直線と見做せるようになる点まで $\phi(t)$ をつなけばおそろく十分であろう。

仮定が正しいと仮定して、 t_0 までの断片のみを考える場合に含まれる誤差は

$$\sum_{v=1}^{\infty} \frac{B_{2v}}{(2v)! n^{2v}} \left[\phi^{(2v-1)}(1) - \phi^{(2v-1)}(t_0) \right]$$

で与えられる。

もし級数の形を用いるなら、幾つかの特殊の仮定が必要である。 $\phi(t)$ が $t=t_0$ から先で線型と仮定すれば、 $\phi'(1) = \phi'(t_0)/(1-t_0)$ で、高次の微係数は 0 である。これと別な仮定は、領域 $(0, 1-t_0)$ で $X(S)$ が平均値 μ_1 をもち、 $(t_0, 1)$ なる領域で $X(S)$ と独立とするものである。但し平均値は μ_2 。このとき、 $t > t_0$ に対して

$$\phi(t) = (1-t)(\mu_1 - a)(\mu_2 - a)$$

となる。観測値を観察してこの仮定の正しいことを認められれば $\phi(t)$ はトランセクトの最初と終りの平均値から直接決定することができる。この仮定が正確には成立たなくても、 $\phi(t)$ はこのずれにそう影響されない。例えば μ_1 および μ_2 が断片でなくて、それぞれ異なる

$$\mu_1 = \alpha_1 + \beta_1 t, \quad \mu_2 = \alpha_2 + \beta_2 t$$

なる傾向にたうとすれば、

$$\phi(t) = (1-t)(\alpha_1 - a)(\alpha_2 + \beta_2 - a) + \frac{(1-t)^2}{2} \left[\beta_1 (\alpha_2 + \beta_2 - a) - \beta_2 (\alpha_1 - a) \right] - \frac{(1-t)^3}{6} \beta_1 \beta_2$$

これから

$$\phi^{(1)}(1) = -(\alpha_1 - a)(\alpha_2 + \beta_2 - a)$$

$$\phi^{(2)}(1) = \beta_1 \beta_2$$

となる。平均値が領域 $(0, 1-t_0)$ および $(t_0, 1)$ で断片であるという仮定によつて

$$\phi^{(1)}(t) = -(\alpha_1 - \alpha)(\alpha_2 + \beta_2 - \alpha) - \frac{(1-t_0)}{2} [\beta_1(\alpha_2 + \beta_2 - \alpha) + \beta_2(\alpha_1 - \alpha)] + \frac{(1-t_0)^2}{4} \beta_1 \beta_2$$

となるから、平均値が妥当な標本にもとづくものでありまたも。をできるだけ大きくとるなら、この誤差は傾向が善しい（即ち β_1 および β_2 が大きい）場合を除いてそう大きくなるまいであろう。

3.2 計画のデータを求めるため、パイロット調査（ m 点の）を行なう場合には、 n より大きいかまたは小さい n の値について $\sigma_S^2(n)$ を計算したい。このとき標本平均値の分散は $m=n$ になる標本それぞれから計算されるべきである。

n が事實上 m より小さいとき例えばその $1/2$ もしくはそれ以下であれば、級数式がゆるやかに収斂するからグラフによる方法が最もよい。 n が m と同じオーダーかそれよりも大きいときには、 $\phi(t)$ の値と決定する場合の抽出変動が測定すべき差と違わなくなつてくるから、級数の形を用いなければならぬ。

$t=0$ で必要な微係数は *Geogory - Newton* の公式によつて決定できる。即ち

$$\phi^{(1)}(0) = m (\Delta \psi(0) - \frac{1}{2} \psi(0) \dots)$$

$$\phi^{(3)}(0) = m^3 (\Delta^3 \psi(0) - \frac{3}{2} \Delta^4 \psi(0) \dots)$$

とおけばよい。ここで $\Delta \psi(0) = \psi(\frac{1}{m}) - \psi(0)$

$$\Delta^2 \psi(0) = \psi(\frac{2}{m}) - 2\psi(\frac{1}{m}) + \psi(0) \text{ である。}$$

$\Delta^i \psi(0)$ の抽出分散は、 i とともに急速に増大するから、 $\phi^{(1)}(0)$ 、 $\phi^{(3)}(0)$ 等を推定するには、少数の項のみを用いる方がよい。また同じ理由で、 $\sigma_S^2(n)$ を実際に展開するときも項数は少ない方がよい。このことは級数展開が有用なのは、 n が十分大きくて、収斂が迅速な場合に限ることを示している。というのは他の場合には、善しい偏りが生ずるからである。実際においては多くの場合、分散の正確な推定値が要求されることはなく、かなり大きい誤差でも許容

できるものである。

3.3 観測値に周期的な傾向があるときには、周期が $\frac{1}{m}$ またはその倍数に近いと非常な偏りを生ずる原因となる。というのは、グラフによる方法において計算された $\psi(\frac{1}{m})$ の値を通る滑らかな曲線を引くことは妥当性を失うからである。

級数式も確實な結果を与える速度には収斂しない。もしデータの性質から周期の存在することか予想できれば、この困難は $\frac{1}{m}$ が周期の十分小さい部分となるようにすれば避けられる。

Finney は森林調査において生ずる予期されない周期変動の一例を示した。そうしてこの場合には簡単な説明は何も与えられないように思えた。この種の影響を避ける唯一の方法は、観測値の間隔をできるだけ小さくとり、少なくとも周期効果が見出されないことのないようにすることである。

3.4 $\chi(S)$ が 0 または 1 の値をとる特別な場合には $a=0$ とおくと方程式 (8) は非特異な形となる。この場合 χ の個数を \sqrt{u} とかくと

$$\psi(\frac{u}{m}) = \frac{\sqrt{u}}{m}$$

となる。

この場合には

$$0 \leq \phi(t) \leq (1-t)$$

であるから $t=1$ における勾配 (gradient) の限界を

$$-1 \leq \phi^{(1)}(1) \leq 0$$

のように定めることができる。

もしトランセクトの浮動の点の間に相関がないと仮定すれば $\phi^{(1)}(1) = -\rho_1, \rho_2$ である。但し ρ_1 と ρ_2 はトランセクトの始めと終りでの χ の割合である。

もし m が大きければ第一偏差 (first difference) はほとんど常視できるから

$$\sigma_S^2(n) \approx \frac{m_r}{12n^2} \dots \dots \dots (9)$$

が成立する。但し m_r は m 個の観測値における χ と θ の ラン (var)

の数である。

この近似式は π が非常に大きいときにのみ正しい。その理由は隣接点間の相関が高いから、ランの数は π が少し変化してもそれにうれてあまり大きく変らないからである。例えは観測値を交互にゆかして系列内のランの数を求めることができる。もしこれが π とあまり異ななければこの式を分散推定の簡便式として用いることができる。

4 実験結果

実験データは道生の被覆度の記録で、これは生きている植生 (living vegetation)、枯死した植生 (dead vegetation) 裸地の三つに分類されている。この植生は、New Zealand 南島 (South Island) の南 alps の分水嶺以東の山岳地帯にある土着のカヤ草原の一部からなっている。外景は Waikakariri 川流域において行なわれた。ここの土壌の大部分は、海拔高 7000 フィート以上に亘る峻険な硬砂岩質の山岳から導かれたものである。一世紀にわたって、自然の降塵によつて多数のブロックにわけられた。このカヤの草地に山羊が放牧されてきた。

草地解析 (Pasture analysis) の点分析法 point method が用いられ、長さ 10 チェーンおよびそれ以上のライントランセクトが恒久的に敷設された。解析器のは 2 インチ間隔で鉄のピンを植えた長さ 20 インチの金属製の棒である。毎年夏の同じ日に解析器をトランセクト上に定位置に置いて 各々のピンに触れる植物の数を記録した。これによつてトランセクト上の 2 インチごとの観測値の相関が得られた。

色々買なつた条件の 4 つのトランセクトを送った。これらはオ 4、7、8、11 トランセクトで、その簡潔な説明を Appendix にあけてある (省略)

我々は生きている (living) のと生していない場合、枯死 (dead) しているのと枯死していない場合および裸地または非観

表 1 ムに面あどに同じ地波の点が続いている点の数 ν_u

u	トランセクト 4			トランセクト 7			トランセクト 11			19418
	L	D	B	L	D	B	L	D	B	L
0	1570	201	2226	4713	343	1227	2955	127	676	3335
1	1362	58	1992	4315	84	933	2787	33	514	2750
2	1252	38	1895	4194	55	848	2718	24	456	2498
3	1169	29	1815	4116	43	802	2668	14	424	2340
4	1100	25	1745	4061	36	750	2635	16	441	2214
5	1046	27	1681	4015	31	717	2601	12	381	2151
6	974	20	1623	3991	33	694	2574	11	365	2077
7	919	16	1579	3974	27	668	2556	8	350	2042
8	872	19	1530	3936	21	647	2534	6	338	2025
9	829	13	1488	3924	24	630	2518	8	320	1999
10	797	11	1545	3909	18	626	2506	4	314	1974
11	773	9	1426	3888	18	615	2499	3	306	1962
12	762	11	1411	3869	18	595	2478	4	294	1968
13	743	15	1378	3855	17	586	2464	7	281	1974
14	733	11	1382	3856	21	581	2451	6	271	1981
15	724	8	1367	3847	17	581	2444	6	268	1979
16	724	9	1319	3832	18	561	2434	8	260	1966
17	730	4	1373	3812	18	550	2427	7	259	1936
18	728	12	1368	3809	17	557	2418	6	251	1927
19	719	12	1363	3788	14	544	2414	3	250	1915
20	718	11	1363	3784	13	530	2408	4	245	1936
21	709	14	1348	3786	21	535	2404	1	243	1951
22	704	13	1336	3788	24	522	2401	5	253	1941
23	693	11	1330	3779	17	524	2404	5	249	1945
24	694	11	1325	3778	21	522	2389	8	236	1942
25	696	10	1322	3767	25	513	2380	3	232	1944
26	700	10	1326	3755	22	501	2387	5	236	1958
27	703	14	1335	3745	16	492	2395	4	236	1945
28	-	-	-	3745	19	498	2398	4	240	1931
29	-	-	-	3746	14	483	2396	4	244	1918
30	-	-	-	3745	19	475	2398	9	237	1910
35	-	-	-	-	-	-	-	-	-	1875
40	-	-	-	-	-	-	-	-	-	1826
合計(点)	3997			6283			3758			6321
観測率(%)	3928	5.03	55.69	75.01	5.46	19.53	78.63	5.38	19.99	52.76

地の三つの場合を別々に考えることにする。生態学的にはこの最初
のものが最も重要である。我々は点Sでの地面の状態を生活してい
る植生があれば $\chi(S) = 1$ なければ $\chi(S) = 0$ として記述すること
にする。その他の場合も同様である。

表2 表 $\sigma_S^2(\pi)$ 一覧表

トランセクト 4 $n = 377$

$$V_L = \frac{10^4}{\pi^2} (46.75 - 0.1P^2 \frac{\pi^2}{\pi^2} \dots)$$

$$V_D = \frac{10^4}{\pi^2} (40.30 - 0.311 \frac{\pi^2}{\pi^2} \dots)$$

$$V_B = \frac{10^4}{\pi^2} (52.03 - 0.381 \frac{\pi^2}{\pi^2} \dots)$$

トランセクト 7 $n = 6283$

$$V_L = \frac{10^4}{\pi^2} (102.32 - 0.650 \frac{\pi^2}{\pi^2} \dots)$$

$$V_D = \frac{10^4}{\pi^2} (74.17 - 0.592 \frac{\pi^2}{\pi^2} \dots)$$

$$V_B = \frac{10^4}{\pi^2} (75.85 - 0.473 \frac{\pi^2}{\pi^2} \dots)$$

トランセクト 11 $n = 3758$

$$V_L = \frac{10^4}{\pi^2} (40.59 - 0.222 \frac{\pi^2}{\pi^2} \dots)$$

$$V_D = \frac{10^4}{\pi^2} (27.53 - 0.239 \frac{\pi^2}{\pi^2} \dots)$$

$$V_B = \frac{10^4}{\pi^2} (39.99 - 0.217 \frac{\pi^2}{\pi^2} \dots)$$

トランセクト 8 $n = 6321$

$$V_L = \frac{10^4}{\pi^2} (138.48 - 0.664 \frac{\pi^2}{\pi^2} \dots)$$

上にあげた場合に対する生活しているもの、枯死、裸地 (L, D, B)
の V_u の値を表1表に与えてある。

表1表をみると、約50インチ (即ち $u > 25$) 以上の間隔への
場合の枯死点 (dead point) の観測値に対する V_u の値 (無作為
変動を別にすると) 相関がない (即ち $(n-u)P^2$ 。但し P は枯死の
比率) 場合に予想されるものと大体等しい。生活および裸地帯の
割合はこの間隔では多少相関がある。しかしこれは大きい抽出間隔の
場合を除いて、直線と考え得るのに十分な理由づくり変化する。

このことは、ここに表示してない中 (七) の値で与えられる分散に
対する寄与を無視すること、同等である。

級数展開の係数を求めるため $\tau = 0$ における三次までの微分のみ
を用いた。これは抽出誤差が大きくなるのを避けるためである。

$\tau = 1$ では $\psi^{(1)}(1) = -P^2$ ととり、他のすべての導函数は0ととつた

表2表に生活、枯死、および裸地 (V_L, V_D, V_B) の割合 (%) の
分散を級数式で与える。トランセクト内の点の数は n だから、 n
点を抽出すると間隔は $2\pi/n$ インチとなる。

直接計算して求めた分散の値と観測を無作為に行なう何様な標本
抽出で期待される分散 $\sqrt{V_L}, \sqrt{V_D}, \sqrt{V_B}$ を比較のために表3表に示
す。間隔が狭いとき V_L と V_B は $\sqrt{V_L}$ と $\sqrt{V_B}$ よりはるかに小さい。
枯死した植生の割合は比較的小さく、生活している地帯は多少無作
為に点在するということのために V_D と $\sqrt{V_D}$ 間の差は小さい。一方
生活している植生では、特にトランセクト4でそうであるが、小さ
い裸地又は岩石地に無々と群生する傾向があるから、大体10~20
インチの狭い間隔では中 (七) の値が大きくなる。トランセクト8の
差が小さいのは地被がむしろより一葉であることを示している。

表3表には、10インチ $n = \pi/5$ の間隔に対する近似公式 (9) か
ら計算した分散の値をあげてある。この一致は満足なものである。
この例は表3.2で注意した問題の説明を与えるものである。近似
公式は、 $\psi^{(1)}(0)$ の推定に $\Delta\psi(0)$ のみを用い $\sigma_S^2(\pi)$ の級数でそれか
らあとのすべての項を無視する場合と殆んど同等である。10イン
チ間隔の場合、この方法で数値積分から分散の近似値を求めると、
表2表で与えたと同様に $\Delta^2\psi(0)$ および $\Delta^3\psi(0)$ を用いる級数
形よりよい結果が得られる。

理論の正しさおよび間接には模型の十分さをチェックするため、
それぞれトランセクト上の最初の20点中から無作為に選んだ1点
より出発し、20番目ごとの点をとつた10個の標本と40番目ご
との点をとつた20個の標本を作る。これらの標本平均値の分散は
(1) で与えられる $\sigma_S^2(\pi, \tau)$ の実験的な推定値を与える。

表3 相関のあるデータから抽出した系統的標本と無作為標本の分散の比較

間隔 (20/24寸)	10インチ	20インチ	40インチ	60インチ	80インチ
トランセクト 4					
V _L	0.604(0.543*)	1.746	6.088(5.852**)	11.156	16.223(9.223**)
V' _L	2.996	5.992	11.985	17.977	23.970
V _D	0.370(0.373*)	0.902	2.091(3.233**)	3.330	4.568(2.538**)
V' _D	0.571	1.143	2.286	3.429	4.572
V _B	0.586(0.610*)	1.959	6.219(2.059**)	11.118	16.017(19.519**)
V' _B	3.090	6.180	12.360	18.540	24.720
トランセクト 7					
V _L	0.400(0.420*)	1.141	2.942(6.200**)	4.991	7.020(9.120**)
V' _L	1.496	2.992	5.983	8.975	11.967
V _D	0.267(0.273*)	0.644	1.462(1.863**)	2.282	3.100(1.918**)
V' _D	0.409	0.819	1.637	2.456	3.270
V _B	0.341(0.310)	0.842	2.331(3.229**)	3.480	5.300(7.065**)
V' _B	1.257	2.514	5.029	7.543	10.058
トランセクト 11					
V _L	0.435(0.510*)	1.606	4.090(2.494**)	7.819	11.549(15.366**)
V' _L	2.235	4.470	8.940	13.410	17.881
V _D	0.243(0.277*)	0.633	1.503(1.769**)	2.444	3.386(4.776**)
V' _D	0.434	0.869	1.738	2.606	3.475
V _B	0.437(0.463*)	1.389	3.463(3.426**)	6.400	9.337(13.535**)
V' _B	1.917	3.833	7.667	11.500	15.334
トランセクト 8					
V _L	0.661(0.610*)	1.791	4.957(3.776**)	8.474	11.210(6.744**)
V' _L	1.969	3.939	7.878	11.817	15.755

* $\sigma_s^2(n) = nr / 12n^2$ から計算した近似推定値

** 抽出実験から求めた推定値

但し $n=20$ または $n=40$ で $n \cdot n = n$ である。しかし $\sigma_s^2(n, n)$ はまた $\sigma_s^2(n) - \sigma_s^2(n \cdot n)$ の推定値とも考えられるから、 $\sigma_s^2(n, n)$ の実験的推定値に、(9)で $n=n$ において与えられる、 $\sigma_s^2(n \cdot n)$ を加えれば $\sigma_s^2(n)$ の推定値を求めることができる。これらは表3のグラフ法(間隔40および80インチ)で求められた対応する分散の次に括弧書きで示してある。これらの比較は厳密でないがしかし実験的分散のグラフ法で得られた分散に対する比の平均値は1.0946となるから十分と思われる。

5 要約

有限母集団から抽出された系統的標本の平均値の分散を求める Madrow & Madrow (1944) の公式を連続的母集団に適用し、ある仮定のもとに単一標本から系統的標本の平均値の分散を求める方法を展開した。そうしてこれを生態学研究における植生の被覆度の推定に適用した。

R. M. Williams

参考文献

- Cochran, W. G. (1946) Relative accuracy of systematic and stratified random samples for certain class of populations. A.M.S. 17, 164-77
- Finney J. J. (1948) Random and systematic sampling in timber surveys. Forestry, 22, 64-99
- (1950) An example of periodic variation in forest sampling. Forestry, 23, 96-111
- Hazel, A. A. (1938) Sampling error in timber surveys. J. agric. Res. 57, 713-36.
- Jowett, G. H. (1952) The accuracy of systematic sampling from conveyor belts. Appl Statist 1, 50-9
- Madrow, W. G. (1949) On the theory of systematic

Sampling II A.M.S. 20, 333-54

Madow W.G. (1953) On the theory of systematic

Sampling III A.M.S. 24 101-6

Madow W.G. & Madow L.H. (1944) On the theory
of systematic sampling I A.M.S. 15, 1-24

Osborne J.G. (1942) Sampling errors of systematic
and random surveys of cover type
areas. T.A.S.A. 33 256-64

Quenouille M.H. (1949) Problems in plane

Sampling A.M.S. 20, 355-75

Yates F. (1948) Systematic sampling *Phil Trans*

A 241, 345-77