

カテゴリカルデータの分析 — シフト係数とその応用

Analysis of categorical data — shift coefficient and its applications

松井敬*

Takashi Matsui

Email: matsui@dokkyo.ac.jp

本論文は1標本カテゴリカルデータのカテゴリ間のシフトの評価にかかわる問題を扱う。統計量として RCN (相対カテゴリ番号、シフト係数) を用い、まずこの統計量の数学的背景を簡略に述べた。次に、サンプリング方式の手続きにしたがって以下の2つの場合を扱った。一つは小標本の場合で、permutation test をこの統計量 RCN に適用し、カテゴリ反応を調べる手続きを述べた。もう一つは大標本の場合で、統計量 RCN の漸近分布を求め、これを利用してカテゴリ反応を分析している。さらに、RCN の有効な場面を探るために、permutation distribution の考えを用い、精密な形でカイ二乗適合度検定との検出力比較を行った。比較はカテゴリ確率に対し、slippage 型配置と step 型配置について行った。その結果、slippage 型では適合度検定に、step 型ではシフト係数にもとづく方法に有意性が見られたが、このことはシフト係数にもとづく方法が特にカテゴリカルシフトの判定に対して有効であることを示している。

A purpose of this paper is to extend the procedure for evaluating the shift of categorical response probabilities based on the RCN or shift coefficient. First, mathematical background of the statistic RCN is introduced briefly. Then, two separate methods are proposed according to the sampling procedures. One is for small sample case and the permutation test is applied to the statistic RCN with given data set. Another one is for large sample case and the asymptotic distribution of RCN is used to analyze the dataset from categorical response. Further, power comparisons between the test based on RCN and the χ^2 -goodness of fit (GOF) test are performed to make the efficacy of the proposed test clear. Comparisons are made for slippage and step configurations of parameters. For slippage configurations, power of the GOF test is superior to the test based on RCN, while for step configurations, power of the RCN test exceed the GOF test. These results suggest the significance of RCN test for finding the shift of categorical data.

*: 獨協大学名誉教授、獨協大学情報学研究所客員研究員

1. はじめに

順位の尺度あるいは間隔、比の尺度からの区分によって得られる表1の形式のカテゴリカルデータに対し、評価結果の分布—あるいは評価結果の左右への広がり—を測り全体の傾向を把握するための客観的な基準をつくることは興味あるテーマである。

表1. カテゴリカルデータ

カテゴリ	1	2	...	k	計
観測度数	n_1	n_2	...	n_k	n

カテゴリカルデータの場合には位置母数を持つ密度関数の場合のように、母数の左右へのシフトという形で傾向を捉えにくい点に問題がある。このため、多くのアンケート調査の結果分析ではカテゴリを合併したりして、カテゴリ度数の細部に踏み込まずにカテゴリ間の優劣の評価を下すといったことが行われている。このことは、データの有効利用といった点で問題がある。

筆者がこの問題にかかわった動機は次のような事例からである。それは、臨床試験における医業への医師の判定で、 n 人の被験者をたとえば（悪化、不変、やや改善、改善）と判定し、 $k=4$ の表2の形の結果を得たときに関連している。いま、AとBの2つのグループに対し次の（1）ないし（2）のような結果が得られたとする。

表2 説明例

(1)	悪化	不変	やや改善	改善	計
A	10	15	25	20	70
B	15	10	25	20	70
(2)	悪化	不変	やや改善	改善	
A	10	15	15	30	70
B	10	15	30	15	70

説明例にあげた（1）、（2）ともに全体として「改善している」が45、「改善していない」を25としている。ただし、（1）では改善していないグループの、（2）では改善しているグループの度数が異なっている。すなわち、改善、非改善の程度の度合いが異なった事例となっている。このような場合、（1）、（2）の2つのグループA、Bともに上位2カテゴリを合併して判断され、全体としては改善されたと説明することが多いと思われる。ところが、2つの例ともに明らかにカテゴリ間の差異はあり、これらがあまり考慮されないことになる。このことはアンケート調査の結果の分析についてもいえることで、このようなカテゴリ間の違いを考慮に入れた指標は作れないかということで、筆者らは第2節で述べるような考えにしたがって、ある考察結果を提案した(Choi and Matsui [1])。

新たな指標の作成にあたってはカテゴリカルデ

ータの度数分布を累積分布の差異に転化し、累積分布のズレを評価する量を考えた（後で述べるNMEA）。さらに、これにもとづいてシフト係数(RCN)を構成しカテゴリカルデータの解析のための統計量とした(Masui [5])。

シフト係数はカテゴリカルデータの度数の左右へのシフトを示す量としてピアソンの相関係数に類似の役割を果たすほかに、カテゴリカルデータの生起にかかわる各カテゴリ（セル）の生起確率（母数）に対する適合度を検定するためにも用いられる。

シフト係数の応用についてはこれまでも論文の中で述べてきたが、理論に偏しすぎた嫌いがあった。実際に応用しにくい面があった。そこで、本論文の目的は次のように述べられる。カテゴリ度数の総数を n としたときに

1. n が比較的小さい場合に permutation test（並べかえ検定）の考えを用いてカテゴリカルシフトの状況¹を明らかにするプロセスを示すこと。

2. n が比較的大きい場合に、シフト係数の正規近似を用いて仮説への適合性や左右へのシフトの是非を判定するプロセスを示すこと。

これらの目的のために、第2節では提案した統計量の由来と考え方、そして算法を簡潔に示した。第3～5節ではカテゴリの数 k が2, 3および4の場合の各論を上記目的にそって示している。なお、上記目的の中にもあげた並べかえ検定(Edgington [2], Ernst [3] そして Good [4] など)は、統計データの分析に大変有効であり、本論文の中ではこの考え方と方法についても詳しく説明している。

ところで、カテゴリの各項目を考慮に入れてカテゴリの傾向性を測ろうとする統計的方法はほとんど見られず、これが上述のようにいまだにカテゴリの合併という形で調査結果の分析などが行なわれていることにつながっている。ただし、 χ^2 -適合度検定についてはカテゴリ確率に対しシフトを想定した帰無仮説を導入できればここで提案している方法のカウンターパートとなる。そこで、第6章では与えられたカテゴリカル母数の仮説に対し、カテゴリカルデータの適合度を示す χ^2 -適合度検定とシフト係数による検定の結果との間の検出力の比較を行ない、その結果を示した。

2. 統計量とその背景

Choi and Matsui [1], Matsui and Choi [5] はカテゴリカル変数の左右へのシフトについて極端数(NMEA)という概念を使って考察している。本節ではカテゴリカルデータの度数の左右へのシ

¹シフトの状況とはそれぞれの目的に応じて得られるデータの区分に関連して、たとえば「改善している」と見てよいか、「賛成」あるいは「反対」、「わかりやすい」などのカテゴリ変数への方向性を示すことを意味している。

フトをどのような形でとらえ、計量化しているかという点を要約し説明しておく。
いま、表1に与えた形の k 個のカテゴリを持つあるカテゴリカル配列を

$$A(k) = (n_1, n_2, \dots, n_k), \quad n_i \geq 0, \quad \sum_{i=1}^k n_i = n, \quad (1)$$

とおく。これは、表1のカテゴリの度数分布に対応する表現である。このとき、 $A(k)$ とは異なる配列

$$B(k) = (m_1, m_2, \dots, m_k), \quad m_i \geq 0, \quad \sum_{i=1}^k m_i = n, \quad (2)$$

に対し、次のように定義する。

定義1： MEA および NMEA (MEA 数)

2つのカテゴリカル配列 $A(k)$ と $B(k)$ に対し、次の条件が満たされるとき、 $B(k)$ は $A(k)$ よりもより極端な配列 (more extreme arrangement) であるとよぶ。

$$\sum_{j=1}^i m_j \leq \sum_{j=1}^i n_j, \quad (i = 1, 2, \dots, k-1). \quad (3)$$

このとき配列 $B(k)$ は配列 $A(k)$ の MEA であるとよぶ。

さらに、配列 $A(k) = (n_1, n_2, \dots, n_k)$ に対し、この配列よりもより極端である配列の集合を考え、それらすべての配列の数を、配列 $A(k)$ についての NMEA (number of more extreme arrangements、MEA 数、右極端数) とよび、

$$NMEA(k; n; n_1, n_2, \dots, n_k) \quad (4)$$

と書くことにする。これは、カテゴリ数 k 、配列 $A(k)$ によって決まる数で、以後、必要に応じて $NMEA(A(k))$ と書くことにする。

上の定義の2つの配列間の不等号のかかわる式(3)で、その向きを反対にした場合には $A(k)$ は $B(k)$ よりも極端な配列となるが、配列 $B(k)$ を配列 $A(k)$ の LEA (less extreme arrangement) とよび、その総数を NLEA (左極端数) とよぶ。なお、2つの極端数の間には、配列の要素につて、次の関係がある。

$$NLEA(k; n; n_1, n_2, \dots, n_k) = NMEA(k; n; n_k, n_{k-1}, \dots, n_1). \quad (5)$$

極端数を図で示すと次のようになる。図1で左側の度数の配列 A、B の違いをカテゴリの度数の累積で示したのが右側の累積度数分布である。コロモゴロフ-スミルノフ検定では累積度数分布の最大差を考慮しているが、極端数では図で配列 A

の下側に来る配列 (例えば B) を MEA とよび、それらすべての配列の総数を NMEA とよんでいる。

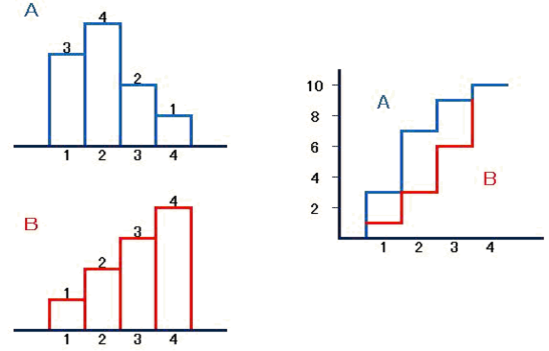


図1 2つの配列 A と B の関係

大きさ n の観測値を大きさ k の度数の配列 (n_1, n_2, \dots, n_k) に配分する方法の総数 $N(k, n)$ は

$$N(k, n) = \binom{n+k-1}{k-1} \quad (6)$$

となる。(1) にあげた配列 $A(k)$ で配列の要素 n_i を和 $\sum_{i=1}^k n_i = n$ となるように変化させた総数が $N(k, n)$ である。図1で見ると、原点 $(0,0)$ から点 $(k-1, n)$ にいたる各格子点を通る道筋 (path) の総数となっている。後の例で見るように、これら全配列の中で NMEA ないし NLEA は、カテゴリ配列の右ないし左へのシフトを考えるための有効な数量表現となっている。NMEA および NLEA とともに最大値は $N(k, n) + 1$ 、最小値は1である。(6)式で示されるカテゴリ総数の大きさは n と k に依存し、一般にはかなり大きな値となる。そこで本論文に示すような、 n が小さい場合 (小標本) と大きい場合 (大標本) のそれぞれの場合を整理した手法が望まれることになる。

ここで、配列 $A(k)$ の NMEA と NLEA の差をとることによって得られる CN (category number、カテゴリ番号) を (7)式のように定義する。これは与えられた配列 $A(k)$ に対する固有の番号で、カテゴリ度数の左右へのシフトの状況を示す値となっている。さらに、 $CN(A(k))$ を配列の総数 (6)で調整し、相対カテゴリ番号 (RCN、relative category number) を定義する。これを配列 $A(k)$ のシフト係数とよぶことにする。

定義2： $CN(A(k))$ 、 $RCN(A(k))$

$$CN(A(k)) = NMEA(A(k)) - NLEA(A(k)). \quad (7)$$

$$RCN(A(k)) = \frac{CN(A(k))}{N(k, n) - 1}. \quad (8)$$

RCN は任意のカテゴリ配列 $A(k)$ に与えられ

たカテゴリ計測の数量化システムで、 $RCN(A(k))$ の値は -1 から $+1$ までの値をとり、配列 $A(k)$ が左右対称の時にはゼロとなっている。その意味で、与えられたカテゴリ配列に対し、相関係数と似た働きを示しているが、後で述べるようにカテゴリ配列に与える確率との関係で適合度検定に見られるような仮説検定も容易に行うことができる側面を持っている。

定義 1 および定義 2 にあげた $NMEA(A(k))$, $NLEA(A(k))$ および $RCN(A(k))$ の計算式はカテゴリ数 k に依存し、配列ごとに表現される。これらの式は以下の節で個別に示すことにする。なお、 k の値がふえると式を構成する項数が多いという意味でかなり面倒である。この点について、これらの統計量を求めるための次の漸化式が存在する。

$$NMEA(k; n; n_1, n_2, \dots, n_k) = \sum_{i=0}^{n_1} NMEA(k-1; n; n_1+n_2-i, n_3, \dots, n_k). \quad (9)$$

この漸化式を用いると (4)、(5) および (8) 式などの数値計算も容易である。なお、R 言語を用いた MEA 数計算のためのスクリプトの例を補遺に記載しておいた。

カテゴリカルデータと確率モデル

さて、ここでカテゴリカルデータの生成に関し、次の確率モデルを導入しておく。これは、大きさ n のデータが k 個のカテゴリに配分される確率を p_1, p_2, \dots, p_k とした多項分布モデルである。

表 3 カテゴリカルデータと確率モデル

カテゴリ	1	2	...	k	計
観測度数	n_1	n_2	...	n_k	n
期待確率	p_1	p_2	...	p_k	1

表 3 で与えられた確率モデルの下で、上表によって示される配列 $A(k) = (n_1, n_2, \dots, n_k)$ のカテゴリ番号 $CN(A(k))$ の平均および分散は次のように求められる。

$$E[CN(A(k))] = \sum_{n_1+\dots+n_k=n} CN(A(k)) M(n_1, n_2, \dots, n_k) \quad (10)$$

$$V[CN(A(k))] = \sum_{n_1+\dots+n_k=n} \{CN(A(k)) - E[CN(A(k))]\}^2 \times M(n_1, n_2, \dots, n_k) \quad (11)$$

ここで $M(n_1, n_2, \dots, n_k)$ は表 3 のモデルの下で観測値の列 (n_1, n_2, \dots, n_k) を生成する多項分布確率で、次の形で表される。

$$M(n_1, n_2, \dots, n_k) = \frac{n!}{n_1! n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}. \quad (12)$$

この関係から与えられた確率モデルのもとで、RCN の平均、分散などの特性量を評価することができる。カテゴリ数 k ごとの個別の結果は以下の節で順次示している。

3. カテゴリ数が 2 の場合

本節ではカテゴリの数が 2 ($k=2$) の場合を述べる。

調査や観察における大きさ n の観測結果で、カテゴリが「はい、いいえ」、「すぎ、きらい」あるいは「男性、女性」のように二分されて表現される場合は数多く存在する。この場合の確率モデルが表 3 の $k=2$ で表現されているものである。

カテゴリ数が $k=2$ でサンプルの大きさが n 、カテゴリ度数が $A(2) = (n_1, n_2)$ の場合の極端数を $NMEA(A(2))$ とおくと、定義から

$$NMEA(A(2)) = n_1 + 1, \quad (13)$$

と表される。また、配列 $A(2) = (n_1, n_2)$, $(n_1 + n_2 = n)$ の総数 $N(2, n)$ は

$$N(2, n) = n + 1, \quad (14)$$

で、シフト係数 $RCN(A(2))$ は

$$RCN(A(2)) = \frac{n_2 - n_1}{n} = r_2 - r_1 = 1 - 2r_1, \quad (15)$$

となる。ここで $r_1 = n_1/n$, $r_2 = n_2/n$ である。ところで (12) 式に示した多項分布はこの場合は通常の二項分布と同じ形式で、配列 $A(2)$ のカテゴリ度数が得られる確率は次の形で表される。

$$\begin{aligned} \Pr(n_1, n_2) &= \frac{n!}{n_1! n_2!} p_1^{n_1} p_2^{n_2} \\ &= \frac{n!}{n_1! (n - n_1)!} p_1^{n_1} (1 - p_1)^{n - n_1}, \\ n_1 &= 0, 1, 2, \dots, n. \end{aligned} \quad (16)$$

したがって、 $k=2$ の場合のシフト係数 RCN にもとづく統計的推測は与えられたカテゴリ度数に対する二項確率を計算することによって進められる。 $RCN(A(2))$ の平均と分散は (10)、(11) 式によって次の通りである。

$$E(RCN(A(2))) = 1 - 2p_1, \quad (17)$$

0.04206 である。

$$V(RCN(A(2))) = \frac{4p_1(1-p_1)}{n}. \quad (18)$$

十分大きい n に対しては、通常の二項分布の正規分布による近似を RCN の分布に適用できて、

$$RCN(A(2)) \sim N(1 - 2p_1, \frac{4p_1(1-p_1)}{n}) \quad (19)$$

である。

シフト係数 $RCN(A(2))$ を用いてカテゴリ確率 p_1 についての信頼期間を求めたり、仮説検定を行うことも通常の二項分布の比率の検定の場合と同様に進められる。 n が小さい時には二項確率を用いた精密計算によって計算し、処理を進めることができる。また、 n が大きい場合には正規近似を用いて実行できる。

実際の進め方の例を以下に示しておく。この場合、帰無仮説 $H_0: p_1 = 1/2$ のもとで (18)式から

$$RCN(A(2)) \sim N(0, 1/n) \quad (20)$$

したがって、次の関係を使って検定を行うことができる。

$$Z = \sqrt{n}RCN(A(2)) \sim N(0, 1) \quad (21)$$

例 1 n が小さい場合

ある学生母集団からランダムに選んだ学生 15 人に喫煙の有無を聞いた。「すう」 $n_1 = 5$ 名、「すわない」 $n_2 = 10$ 名であった ($A(2) = (5, 10)$)。

シフト係数：

$RCN(A(2)) = (10 - 5)/15 = 1/3$ (右にシフト)

帰無仮説を $H_0: p_1 = 1/2$ とすると、 $RCN(A(2))$ の分布は二項分布 $B(15, 1/2)$ なので、この分布から精密計算ができて、

$\Pr(RCN(A(2)) \geq 1/3) = 0.1509$ (P-値) であることが分かる。

例 2 n が大きい場合

ある施策に対する賛否で、賛成が「男性」 $n_1 = 573$ 名、「女性」 $n_2 = 516$ 名であった ($A(2) = (573, 516)$, $n = 1089$)。

シフト係数：

$RCN(A(2)) = (516 - 573)/1089 = -0.0523$ (左にシフト)

帰無仮説を $H_0: p_1 = 1/2$ とし、正規近似を用いると (20)式から、

$Z = \sqrt{n}RCN(A(2)) = -1.7273$.
この値に対する P-値は正規分布確率から

4. カテゴリ数が 3 の場合

シフト係数を用いたカテゴリカルデータの分析がより有効な局面は、カテゴリの数 k が 3 以上の場合である。本節では RCN の分布の形状分析から始め、これにもとづく統計的推測の問題を並べかえ検定 (permutation test) を用いた精密な場合と、サンプルサイズ n を大とした漸近的な場合に分けて扱うことにする。

3つのカテゴリを持つ大きさ n のカテゴリカルデータが与えられているとする。各カテゴリの度数を $A(3) = (n_1, n_2, n_3)$, ($n_1 + n_2 + n_3 = n$) とおき、モデルとして生起確率 (p_1, p_2, p_3) を持つ表 3 の形式を想定するものとする。

このとき、 A の MEA 数は次の形で表される (Matsui and Choi [5], Matsui [6])。

$$NMEA(A(3)) = \frac{1}{2}(n_1 + 1)(n_1 + 2n_2 + 2). \quad (22)$$

また $A(3)$ のすべての配列の個数は (6)式によって

$$N(3, n) = \binom{n+2}{2} = \frac{1}{2}(n+2)(n+1) \quad (23)$$

であって、シフト係数 $RCN(A(3))$ は

$$RCN(A(3)) = \frac{(n_3 - n_1)(n + n_2 + 3)}{n(n+3)}. \quad (24)$$

である。

ここでまずシフト係数の例をあげ、その意味を考察してみよう。表 4 にあげたのは、大学卒業時の取得相単位数と成績評価 A、B、C それぞれの取得単位数を示している²。A 評価に 3 点、B 評価に 2 点、C 評価に 1 点を与え重み付平均を求めたのが GPA (grade point average) で、成績評価の指数としてよく用いられている。これを表中の GPA の欄にあげている。本論文で扱っているシフト係数を (24) 式によって計算したのが RCN の欄である。

RCN の値はカテゴリカル度数の分布にしたがって $-1 \sim +1$ の間の値をとり、カテゴリ度数の右ないし左へのシフトを示す指数になっている。表は成績上位者で RCN の値がプラスだが、各配列の度数を左右で入れ替えると値はマイナスとなる。また、 RCN と GPA の大きさにしたがって順序を与えたが、個別に観察すると RCN がシフ

²表 3 および表 4 は学部学生から得られた実際のデータで、成績上位者からの一部である。

ト（カテゴリ間の度数の変化）により敏感な指数となっていることが分かる。

表4 成績評価とシフト係数 ($k = 3$)

	評価と単位			単位計	順位		順位
	C	B	A		RCN	GPA	
1	0	10	132	142	0.9937	2.930	1
2	0	22	118	140	0.9725	2.843	2
3	0	36	162	198	0.9647	2.818	3
4	8	16	114	138	0.8553	2.768	7
5	8	20	120	148	0.8570	2.757	5
6	8	24	104	136	0.8278	2.706	9
7	4	34	102	140	0.8664	2.700	4
8	8	24	100	132	0.8209	2.697	10
9	6	32	106	144	0.8456	2.694	8
10	4	32	92	128	0.8554	2.688	6

シフト係数を用いた分析には n の値が小さい場合には次にあげる permutation distribution を用いた考え方を適用し、 $k = 2$ の場合と同様に進めるのが有効である。具体的には例を含め次節で説明する。

ここではまず、与えられた確率モデルのもとで RCN の分布について考えてみよう。

カテゴリ数 $k = 3$ の場合には、与えられたサンプルサイズ n に対し、カテゴリの配列 $A(3) = (n_1, n_2, n_3)$ の総数は (23) 式によって $N(3, n) = (n+1)(n+2)/2$ 個存在する。そのすべての配列に対し $RCN(A(3))$ の値が対応していてこれらの総体が、与えられた確率モデルのもとで RCN の分布を規定している。これを permutation distribution とよぶ。以下、本論文ではこの分布を“全配列分布”とよんでおく。

$k = 3$ の場合の RCN の全配列分布は $(n+1)(n+2)/2$ 個ある $RCN(A(3))$ の離散点に対して与えられる確率の分布で、

配列 $A(3)$	シフト係数	配列の確率
(n_1, n_2, n_3)	$RCN(A(3))$ の値	$\Pr(n_1, n_2, n_3)$

の対応によって計算される。 $\Pr(n_1, n_2, n_3)$ は配列 (n_1, n_2, n_3) に対する多項分布確率

$$\Pr(n_1, n_2, n_3) = \frac{n!}{n_1! n_2! n_3!} p_1^{n_1} p_2^{n_2} p_3^{n_3} \quad (25)$$

である。

RCN の分布の平均、分散の値はこの全配列分布によって任意の確率モデルに対し求められるが、その精密な計算式は(10), (11) によっても求められ、 $k = 3$ の場合の平均は次のように表される。計算式は k が大きくなるにつれかなり複雑な式となっている。

$$E(RCN(A(3))) = \frac{(p_3 - p_1)(n + 3 + (n - 1)p_2)}{n + 3} \quad (26)$$

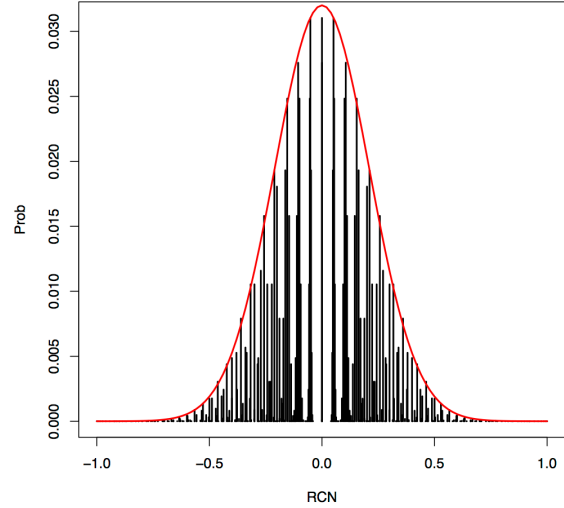


図2 $n = 25, (p_1, p_2, p_3) = (1/3, 1/3, 1/3)$

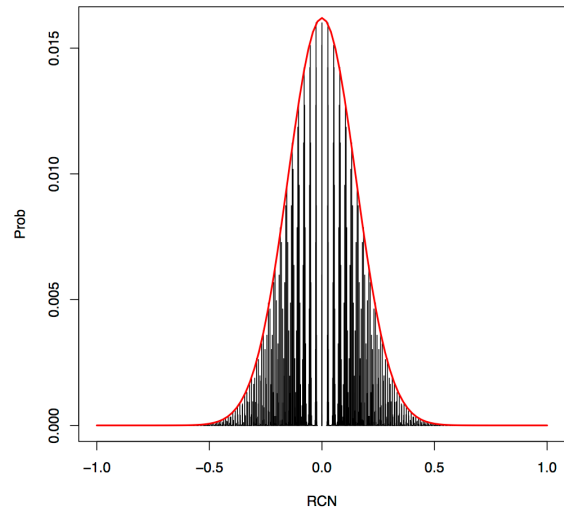


図3 $n = 50, (p_1, p_2, p_3) = (1/3, 1/3, 1/3)$

さて、このようにして得られる $(n+1)(n+2)/2$ 個の離散点に対し与えられる RCN の確率分布（全配列分布）の例を、異なる n と異なる確率モデルに対し図2、図3および図4に与えた。直感的に正規分布風に分布していることが分かるが、実際二項分布の場合と同様正規分布への良い近似を与えてくれている。

図2は母数モデルが等確率 $p_1 = p_2 = p_3$ で、 $n = 25$ の場合の全配列分布とそれに対応する平均と分散を持つ正規分布を重ねて描いた。図3についても同様に、 $n = 50$ の場合である。

確率モデルが等確率でない場合であっても、 n

が大きければかなり良い正規分布への近似が得られる。たとえば、確率モデルが $(p_1, p_2, p_3) = (0.2, 0.3, 0.5)$ で、 $n = 40$ の場合を示したのが図4である。

与えられた確率モデルと n のもとで RCN の平均と分散は一意に定まるので、 n が大きい場合には中心極限定理によって正規分布を近似させることができるのである。

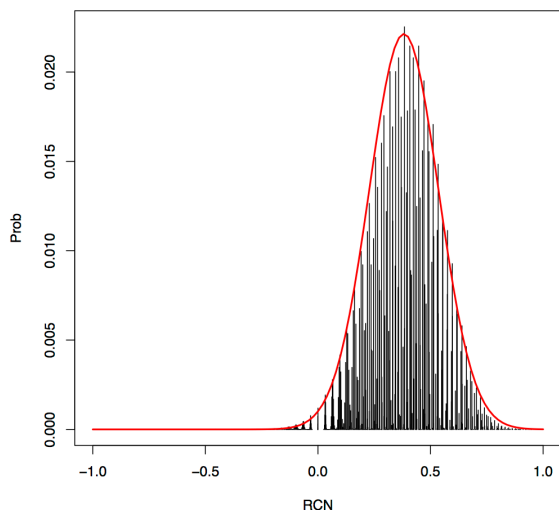


図4 $n = 40$, $(p_1, p_2, p_3) = (0.2, 0.3, 0.5)$

カテゴリ間の一様性の検定では帰無仮説として $H_0: p_1 = p_2 = p_3 = p_0$ (= 一定、等確率モデル) を採用する。この等確率モデルのもとで、 RCN の全配列分布によって得られる上側%点と、対応する平均、分散を与えてあてはめた正規分布による値をあげたのが表5である。これによって、 RCN の分布の正規性や、近似の度合いの程度を見ることができる。

表5 RCN の全配列分布と正規近似による上側%点 ($k = 3$)

n	RCN 分布		正規近似	
	1%	5%	1%	5%
10	0.7538	0.5385	0.7422	0.5247
30	0.4444	0.3212	0.4492	0.3176
50	0.3540	0.2528	0.3518	0.2487
100	0.2509	0.1773	0.2510	0.1774

RCN の取りうる値は離散点であって、上側確率を正確に %点に対応させることはできない。そこで表にあげた値を例によって示すと次の通りである。

$k = 3$, $n = 50$ の H_0 のもとでの分布から、たとえば RCN の上側 5% 値は 0.2528 と得られる。これは

$$\Pr(RCN(A(3)) \geq 0.2528) \leq 0.05$$

であって、 RCN の値が全配列分布の上側の 5% を超えない最小の値となっている。他のパーセント点についても同様の考え方でパーセント値を求める。したがってこの %点を用いた検定では、検定結果が保守的（与えられた有意水準で棄却されにくい）となることが考えられる。P-値は全配列分布から容易に求められるので、P-値を記しておくことが薦められる。

サンプルサイズ n が比較的小さい場合に全配列分布を用いた分析例は、第5節で $k = 4$ の場合で考察している。

大標本の場合

これまで見てきたように、 RCN の正規分布への近似はかなり良好で、 n が大きい（漸近的な）場合には、次の手続きによる統計的推測の方法が有効である。

(24)式にあげたシフト係数 $RCN(A(3))$ は、 n が大きい場合には次のように表される。これを $RCN(A(3))_{Asy}$ と書くことにする。

$$\begin{aligned} RCN(A(3))_{Asy} &= r_1^2 - r_3^2 - 2r_1 + 2r_3 \\ &= (1 - r_1)^2 - (1 - r_3)^2, \quad (r_i = n_i/n, i = 1, 2, 3). \end{aligned} \quad (27)$$

さらに、期待値の定義にしたがって計算すると $RCN(A(3))_{Asy}$ の平均と分散が次のように得られる。

$$E(RCN(A(3))_{Asy}) = (p_3 - p_1)(1 + p_2). \quad (28)$$

$$\begin{aligned} V(RCN(A(3))_{Asy}) &= \frac{1}{n} (6p_1^3 + 6p_3^3 - 2p_1^2p_3 - 2p_1p_3^2 \\ &\quad - 12p_1^2 - 12p_3^2 + 8p_1p_3 + 4p_1 + 4p_3) \end{aligned} \quad (29)$$

したがって n が十分大きい場合には表3の形式の任意の確率モデルに対し、シフト係数 RCN が、(28)、(29) 式の平均と分散を持つ正規分布にしたがうものとして分析を進めることができる。

なお、特に帰無仮説 $H_0: p_1 = p_2 = p_3$ のときには次が成立する。

$$E(RCN(A(3))_{Asy}) = 0, \quad (30)$$

$$V(RCN(A(3))_{Asy}) = \frac{321}{27n}. \quad (31)$$

したがって、一般によく見られるカテゴリ間の一様性（等確率）の想定のもとでの統計的検定には検定量

$$Z = \frac{RCN(A(3))_{Asy}}{\sqrt{\frac{32}{27n}}} \sim N(0,1) \quad (32)$$

を用いるのが有効である。

5. カテゴリ数が 4 の場合

本節ではカテゴリ数が 4 の場合を述べる。NMEA、RCN の式はかなり面倒な形となり式 (33)、(34) のように示される。数値的な計算では漸化式 (9) を用いることができ、補遺にあげた R 言語による例にしたがえば数値計算上は見かけほど面倒なものではない。

$$NMEA(A(4)) = \frac{1}{6}(n_1 + 1)\{(6n_2 + 3n_1 + 6)n_3 + 3n_2^2 + 3(n_1 + 3)n_2 + (n_1 + 2)(n_1 + 3)\} \quad (33)$$

$$CN(A(4)) = \frac{1}{6}\{n_4^3 - n_1^3 + 3(n_2 + n_3 + 2)(n_4^2 - n_1^2) + (6n_2n_3 + 9n_2 + 9n_3 + 11)(n_4 - n_1) + 3(n_3^2 + n_3)(n_4 + 1) - 3(n_2^2 + n_2)(n_1 + 1)\} \quad (34)$$

また、全配列の数は

$$N(4, n) = (n + 3)(n + 2)(n + 1)/6$$

なので、シフト係数 $RCN(4)$ は次のように書ける。

$$RCN(A(4)) = \frac{CN(A(4))}{N(4, n) - 1}$$

$RCN(A(4))$ の具体的な値を見るために、 $k = 3$ の時と同様ここでも成績評価に適用した例を示しておく。表 6 は成績評価を 4 区分 (AA、A、B、C) としたときの成績評価結果 (単位数) を示している。GPA (AA を 4、A を 3、B を 2、C を 1 とスコア化した時の重み付平均) の結果とあわせ作表している。GPA と RCN の値の大きさによる順位もあわせ示しているが、RCN が両端にある度数の影響を強く受け右ないし左へのシフトの状況を示していることが分かる。

表 6 R 成績評価とシフト係数 ($k = 4$)

	評価と単位数				単位計	RCN GPA		順位	
	C	B	A	AA				RCN	GPA
1	6	12	40	70	128	0.7701	3.359	2	1
2	2	18	44	70	134	0.8204	3.358	1	2
3	10	10	37	73	130	0.7098	3.331	8	3
4	5	8	65	64	142	0.7491	3.324	4	4
5	6	14	43	67	130	0.7490	3.315	5	5
6	7	16	46	69	138	0.7251	3.283	6	6
7	4	23	46	69	142	0.7497	3.268	3	7
8	5	15	59	57	136	0.7014	3.235	9	8
9	4	22	44	60	130	0.7222	3.231	7	9
10	2	31	42	55	130	0.6914	3.154	10	10

次に、サンプルサイズ n が比較的小さい場合と大きい場合についてのシフト係数による対応を述べておく。

小標本と全配列分布

本論文の目的の一つにシフト係数を用いたカテゴリカルデータの分析を、サンプルサイズが比較的小さい中で行いたいということがある。このために前節でも全配列分布と permutation test (並べかえ検定) を用いた考え方を示してきた。本節では事例を通しこの考え方と方法を説明したい。

例：教育法の比較

PC を用いたある科目の教育法の開発で、新しい方法 A と、従来からの方法 B とを 7 人の被験者に評価してもらい、次の結果を得た。評価は値の大きいほうが良い評価である。

表 7 教育法の評価と人数

教育法	1	2	3	4	計
A	0	0	3	4	7
B	0	3	3	1	7

このような比較実験はよくおこなわれているようだが、実際には統計的解析に手を焼く問題である。シフト係数を用いた並べかえ検定による手続きは次の通りである。

カテゴリ数 $k = 4$ で $n = 7$ なので、可能な配列総数は (6) 式から $N(4, 7) = 120$ である。このすべての配列について RCN の値と帰無仮説 $H_0: p_1 = p_2 = p_3 = p_4$ の下で (12) 式の多項分布確率を計算する。これによって表 8 の形式の結果がえられる。表 8 は作成される 120 の配列全体の表の一部で、RCN の値について降順に、累積確率をあわせ製表してある。

表 8 Permutation test のための表

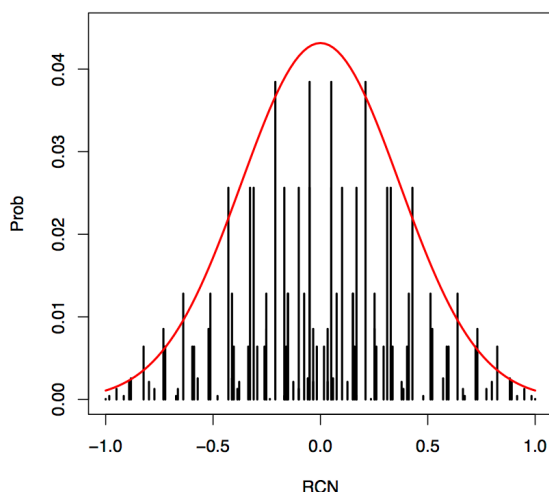
No	n1	n2	n3	n4	NMEA	NLEA	CR	RCN	Prob(Null)	CumProb
1	0	0	0	7	1	120	119	1.0000	0.00006	0.00006
2	0	0	1	6	2	119	117	0.9832	0.00043	0.00049
3	0	0	2	5	3	116	113	0.9496	0.00128	0.00177
4	0	1	0	6	3	112	109	0.9160	0.00043	0.00220
5	0	0	3	4	4	110	106	0.8908	0.00214	0.00433
6	0	1	1	5	5	110	105	0.8824	0.00256	0.00690
7	0	1	2	4	7	105	98	0.8235	0.00641	0.01331
8	0	0	4	3	5	100	95	0.7983	0.00214	0.01544
9	0	2	0	5	6	98	92	0.7731	0.00128	0.01672
10	0	1	3	3	9	96	87	0.7311	0.00854	0.02527
11	0	2	1	4	9	95	86	0.7227	0.00641	0.03168
12	1	0	0	6	4	84	80	0.6723	0.00043	0.03210
13	0	0	5	2	6	85	79	0.6639	0.00128	0.03339
14	0	2	2	3	12	88	76	0.6387	0.01282	0.04620
15	1	0	1	5	7	83	76	0.6387	0.00256	0.04877
16	0	1	4	2	11	82	71	0.5966	0.00641	0.05518
17	0	3	0	4	10	80	70	0.5882	0.00214	0.05731
18	1	0	2	4	10	80	70	0.5882	0.00641	0.06372

第5節にあげた式から、表7にあげた教育法Aの配列 $A(4)$ についてのシフト係数は $RCN(A(4)) = 0.8908$ 、教育法Bの配列 $B(4)$ については $RCN(B(4)) = 0.2521$ であり、配列 $A(4)$ の右へのシフト（教育法への高い評価）は大きい。さらに、仮説 $H_0: p_1 = p_2 = p_3 = p_4$ としたときの RCN についての全配列分布は図5³に示したように得られるが（表8はこの一部分で

ある）、この分布から、

$$\Pr(RCA(A(4)) \geq 0.8908) = 0.0043$$

すなわち、配列 $A(4)$ の上側 P-値は 0.0043 である。同様に、配列 $B(4)$ の P-値として 0.2521 が得られる。この結果から、有意水準を 5% としたとき帰無仮説 H_0 に対し、配列 $A(4)$ は有意に右にシフトしているといえるが、配列 $B(4)$ についてはそのようには言えないということが分かる。

図5 全配列分布 ($k = 4, n = 7$)

大標本の場合

n が十分大きい時、RCN の値は固有の平均と分散を持つので中心極限定理により正規分布によって近似できる。図6に、 $n = 40$ 、 $p_1 = p_2 = p_3 = p_4$ の場合、図7に、 $n = 50$ 、 $p_1 = 0.1$ 、 $p_2 = 0.2$ 、 $p_3 = 0.3$ 、 $p_4 = 0.4$ の場合の RCN の全配列分布とそれに対応する正規分布のグラフをあげた。また、表9には幾つかの n に対し、RCN の全配列分布から得られた上側 1%点、5%点と対応する正規分布の %点をあげておいた。

n が十分大きいときには計算は大変煩雑なものになるが、帰無仮説 $H_0: p_1 = p_2 = p_3 = p_4$ のもとで、 $CN(A(4))$ の平均と分散について次式が得られている。

$$E(CN(A(4))_{Asy}) = 0 \quad (35)$$

$$\begin{aligned} V(CN(A(4))_{Asy}) \\ = \frac{1}{1536} n(51n^4 + 522n^3 + 1961n^2 + 3138n + 2008) \end{aligned} \quad (36)$$

³図に示した RCN の全配列分布の平均は 0、標準偏差は 0.3683 である。なお、この分析と直接の関わりはないが、 n が小さい時の正規近似への近似の参考のために同じ平均と標準偏差を持つ正規分布を一縦軸側のスケールを調整した上で一当てはめて描いてある。

これから n が十分大きいとき

$$E(RCN(A(4))_{Asy}) = 0 \quad (37)$$

$$V(RCN(A(4))_{Asy}) = \frac{153}{128n} \quad (38)$$

したがって、 $k = 4$ の場合に RCN によって、カテゴリ間の一様性（等確率）の検定を行うには

$$Z = \frac{RCN(A(4))_{Asy}}{\sqrt{\frac{153}{128n}}} \quad (39)$$

が有効であることが分かる。

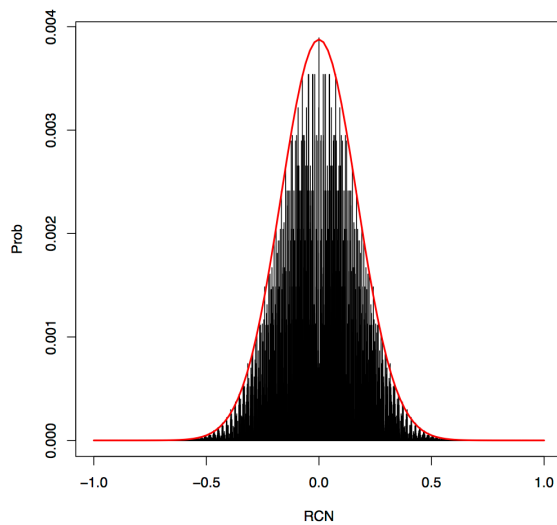


図6 全配列分布 ($p_1 = p_2 = p_3 = p_4$)

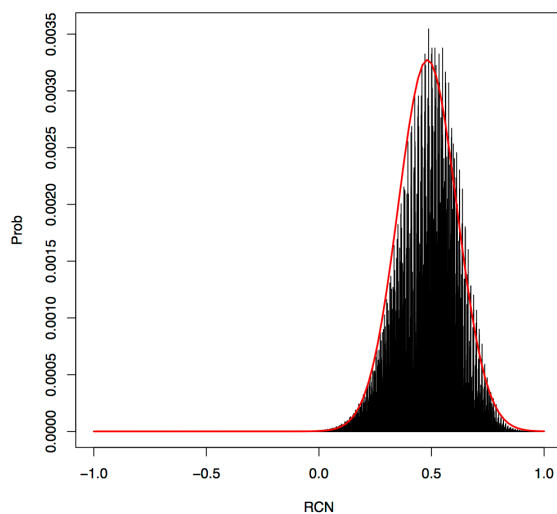


図7 全配列分布 ($p_1 = 0.1, p_2 = 0.2, p_3 = 0.3, p_4 = 0.4$)

表9 RCN の全配列分布と正規近似による上側%点 ($k = 4$)

n	RCN 分布		正規近似	
	1%	5%	1%	5%
10	0.7053	0.5333	0.7397	0.5230
30	0.4475	0.3178	0.4510	0.3189
50	0.3513	0.2509	0.3534	0.2499
100	0.2513	0.1785	0.2512	0.1783

例 表2への適用

ここで、 $k = 4$ の場合の例として表2の説明例にシフト係数を適用してみる。シフト係数 RCN の全配列分布は $N(4,70) = 62,196$ の配列から構成され、それから平均はゼロ、標準偏差は 0.12904 と得られる。シフト係数は補遺にあげた R 言語を用いて計算し、(39)式の正規近似による標準偏差にもこの値を適用した。全配列分布から得られる P-値、正規近似の結果を要約したのが表 10 である。

表10 説明例

(1)	全配列分布		正規近似	
	RCN	P-値	Z-値	P-値
A	0.2788	0.0152	2.1606	0.0154
B	0.1845	0.0762	1.4298	0.0764
(2)	RCN		正規近似	
	RCN	P-値	Z-値	P-値
A	0.4015	0.0008	3.1112	0.0009
B	0.1937	0.0673	1.5014	0.0666

まず、P-値を比較して分かるように、全配列分布による精密な結果に対し、正規近似が十分な結果をもたらしていることが分かる。

この例では表2のところでも述べたように「改善」の2つのカテゴリを合併し、改善の有無を帰無仮説 $H_0: p = 1/2$ の二項検定で行うと、

$$Z = (45 - 35) / \sqrt{70/4} = 2.3905$$

で P-値は 0.0084 と高度に有意となる。

シフト係数による解釈ではカテゴリ確率が一樣であるとの帰無仮説に対し、有意水準を 5% とした判定では (1) - Aは右にシフト、(1) - Bではシフトは認められない。また、(2) - Aは右にシフト、(2) - Bはシフトは認められないとなる。

(1) - B、(2) - Bともにシフトの傾向は微妙（統計用語にこのような言い方はないが）ではあるが、シフトは認められないという結論になっている。このように、シフト係数による検定結果はカテゴリ度数の差異を反映したものとなっている。

6. 適合度検定との比較

カテゴリカルデータの分析によく用いられる方法に適合度検定 (GOF, goodness of fit test) がある。これは表3のモデルのもとで、データにもとづいて母数 (p_1, p_2, \dots, p_k) について与えられた仮説への適合性を検定する方法である。GOF では表3のモデルに対し、 χ^2 -統計量

$$\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} \quad (40)$$

が自由度 $k-1$ のカイ二乗分布にしたがうことを用いて検定を行っている。この方法は χ^2 -統計量が近似的にカイ二乗分布にしたがうことを利用しているが、この GOF もこれまでに述べた全配列分布を利用して精密な検定を行うことができる。

そこで、本節では表3に与えたモデルに対し、RCN と GOF にもとづく分布のシフトないしは適合性の検定法を全配列分布を用い、検定の検出力を計算し比較検討してみたい。

比較は $k=3$ と $k=4$ について行ったが、帰無仮説

$$k=3 \Rightarrow H_0: p_1 = p_2 = p_3$$

$$k=4 \Rightarrow H_0: p_1 = p_2 = p_3 = p_4$$

に対し、次の2つの母数形式 (parameter configuration) を考える。 d はいずれもパラメータ間のギャップの大きさである。

Slippage configuration

$$k=3 \Rightarrow p_1 = p_2 = p_0, p_3 = p_0 + d \\ \Rightarrow p_1 = p_2 = \frac{1-d}{3}, p_3 = \frac{1+2d}{3}.$$

$$k=4 \Rightarrow p_1 = p_2 = p_3 = p_0, p_4 = p_0 + d \\ \Rightarrow p_1 = p_2 = p_3 = \frac{1-d}{4}, p_4 = \frac{1+3d}{4}.$$

Step configuration

$$k=3 \Rightarrow p_1 = p_0, p_2 = p_1 + d, p_3 = p_2 + d \\ \Rightarrow p_1 = \frac{1}{3} - d, p_2 = \frac{1}{3}, p_3 = \frac{1}{3} + d.$$

$$k=4 \Rightarrow p_1 = p_0, p_2 = p_1 + d, p_3 = p_2 + d, \\ p_4 = p_3 + d \\ \Rightarrow p_1 = \frac{1-6d}{4}, p_2 = \frac{1-2d}{4}, \\ p_3 = \frac{1+2d}{4}, p_4 = \frac{1+6d}{4}.$$

以下、 $k=3, n=50$ として検出力計算と比較

の過程を説明しよう。

1° $n_1 + n_2 + n_3 = n (= 50)$ となるすべての配列 (n_1, n_2, n_3) の総数は (6) 式によって $N(k, n) = N(3, 50) = \binom{52}{2} = 1,326$ 通りである。これらすべての場合について統計量 $RCN(A(3))$ を (24) 式によって計算する。

2° RCN の分布はセル確率 (p_1, p_2, p_3) によって配列 (n_1, n_2, n_3) ごとに多項分布によって確率を計算できるので、1,326 通りのすべての配列について計算した結果から図3に示すような全配列分布が得られる。

3° ステップ2°で $H_0: p_1 = p_2 = p_3 (= 1/3)$ として計算したのが null distribution (仮説の下での統計量の分布) である。この分布で RCN の値 $(-1 \sim +1)$ の降順に並べ、対応する確率の累積確率を計算しておけば null 分布における上側 (あるいは下側) のパーセント点を求めることができる (実際の検定では P-値によって対応するほうが簡単)。

4° $k=3, n=50$ の null 分布から、たとえば RCN の上側 2.5% 値は 0.29811 と得られる。これは

$$\Pr(RCN \geq 0.29811) \leq 0.025$$

であって、RCN の値が全配列分布の上側の 2.5% を超えない最小の値となっている。他のパーセント点についても同様の考え方でパーセント値を求める。このケースでは次の結果が得られる。

パーセント	RCN の限界点	実際の確率
0.5	0.3857	0.0050
1.0	0.3540	0.0093
2.5	0.2981	0.0249
5.0	0.2528	0.0483

RCN の分布は離散分布なので、実際の確率とした欄は対応するパーセント点を超えない内輪の累積確率の値である。

5° 検出力の計算

上にあげた slippage configuration の場合はセル確率が $p_1 = p_2 = (1-d)/3, p_3 = (1+2d)/3$ のもとで、 d の値を 0.05, 0.10, 0.15, ... と変えながら、それぞれの場合の全配列分布を形成する配列 A ごとに多項分布確率を計算する。そこでステップ4° の null の場合のパーセント点と合わせ $\Pr(RCN(A) \geq 0.29811 | d)$ とすればそれぞれの有意水準に対応する検出力を求めることができる。⁴

6° χ^2 -適合度 (GOF) 検定との比較

⁴ パーセント点の決め方を真のパーセント値よりも小さい値で決めているため (パーセント点の値は大きくなっていて) 検出力の値は小さめの値となっている。

GOF 検定についても RCN の場合と同様に進め、全配列分布の立場から検出力の計算と比較を行ってみる。

配列 (n_1, n_2, n_3) に対し

$$\chi^2 = \frac{3}{n}((n_1 - n/3)^2 + (n_2 - n/3)^2 + (n_3 - n/3)^2)$$

を $N(3, 50) = 1,326$ 個のすべての配列について計算し、その値と対応する配列の H_0 のもとでの多項分布確率を合わせて得られるのが χ^2 の null distribution である。null 分布のもとでパーセント点は次のように得られる。

%	GOF の 限界点	実際の 確率	χ^2 の 近似値
0.5	10.720	0.00444	10.597
1.0	9.160	0.00977	9.210
2.5	7.720	0.02357	7.378
5.0	6.280	0.04591	5.991

RCN の場合についてステップ 4°で述べたことは、この場合の GOF χ^2 -検定についても同様にいえる。GOF χ^2 -検定ではこの null 分布を自由度 2 のカイ二乗分布で近似して適用しているが、その場合のカイ二乗によるパーセント点を上の表にのせている。これは全配列分布を用いた精密な分布に対する χ^2 による近似ということである。

7°セル確率が $p_1 = p_2 = (1-d)/3$, $p_3 = (1+2d)/3$ のもとで多項確率と χ^2 の値を計算して得られるのが slippage configuration のもとでの全配列分布である。

d の値を変化させてそれぞれの場合の検出力を求め、これを総合して得られるのが検出力関数である。なお、検定は両側検定である。帰無仮説は $H_0: p_1 = p_2 = p_3$, 対立仮説は帰無仮説の否定である。

以上を要約すると、RCN、GOF の両統計量に対し

- ・帰無仮説のもとで全配列分布を求め、与えられた有意水準に対する棄却限界値を求める。

- ・slippage (あるいは step) configuration のもとで母数間のギャップ d の値を変化させながら全配列分布を求め、 H_0 の下での棄却限界値を超えるものの確率を求める (これがそのギャップに対する検出力である)。

計算は $k=3$ で $n=50$ と $n=100$, $k=4$ で $n=50$ と $n=100$ の場合に、slippage と step の両 configuration について行った。結果の一部を図 8～図 11 に示す。

slippage configuration については GOF が検出力が高い結果を示している。それに対し、step configuration では RCN が高い検出力を示して

いる。これは RCN が成績評価の例で見られるように、左右へのシフトにより敏感な統計量となっているためであると考えられる。

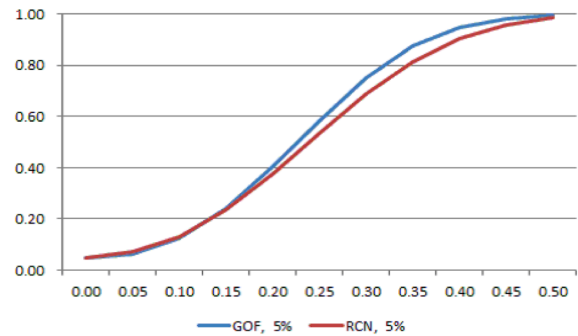


図 8 $k=3$, $n=50$, Slippage configuration

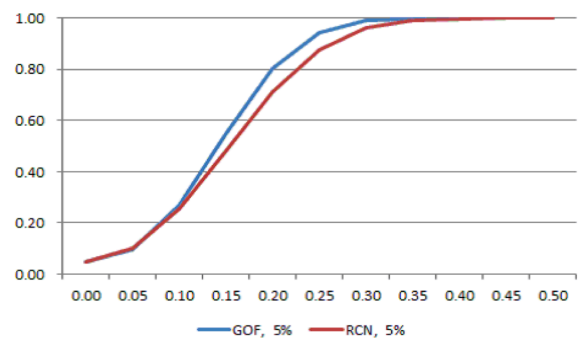


図 9 $k=4$, $n=100$, Slippage configuration

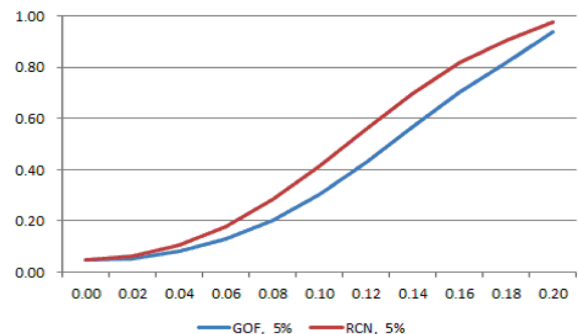


図 10 $k=3$, $n=50$, Slippage configuration

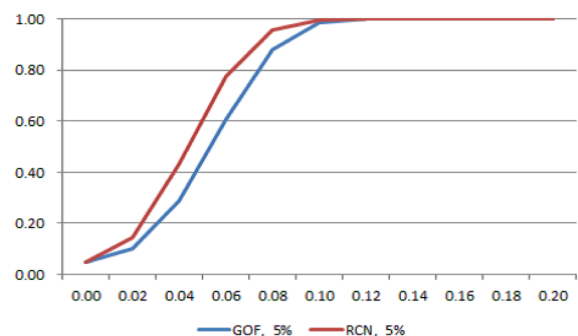


図 11 $k=4$, $n=100$, Slippage configuration

7. おわりに (考察)

「はじめに」で述べたように個々のカテゴリ度数を考慮したカテゴリカルデータの分析は大変面倒であるが、シフト係数はその点で優れていると考えられる。この理由の一つはカテゴリ番号 (CN) の生成にかかわっていて、この方法が配列 $A = (n_1, n_2, \dots, n_k)$, $\sum_{i=1}^k n_i = n$ を形づくる $N(k, n)$ 個の全配列に固有の番号を与えるシステムになっているためである。さらに、異なる配列が同じ番号を持つことはあるもののその重複率は小さく、定義が示すようにカテゴリカルシフトをうまく表現するように統計量が形成されている (Matsui[6])。

第5節の事例で示した並べかえ検定の手続きは一見面倒なようであるが、近似的に何々分布にしたがうという現存する多くの統計的手法の枠を取り払ってもう少し多用されてもよいと考えられる。それは、Ernst [3] も述べているように、計算性能が向上し全配列分布 (permutation distribution) を比較的簡単に作り出すことができること、それによってより精密な計算が可能となっていることなどによっている。この点からすると、この手法は今後もっと積極的に用いられるべき方法であると考えられる。そのためには正規近似を用いた既存の方法 (たとえば、適合度検定) について精密検定法のアプリケーションを準備しておくといえよう。実際、2 x 2 分割表には Fisher の精密検定の方法は準備されている。

論文の中で $k \geq 3$ の場合の大標本の説明では帰無仮説として等確率の下での考え方だけを示したが、これは任意の確率モデルに対しても有効である。ただその場合の分散の評価式が項数が多いという意味で膨大となり面倒なので、 n が比較的小さい場合に全配列分布を用いた方法で処理するよう記述しておいた。

本稿で述べた方法は小標本については $k \geq 5$ についてもそのまま適用される。 $k = 4$ の場合の CN の表現を (34) 式にあげたが、 $k \geq 5$ の場合の対応する式は項数が多いという意味で大変煩雑になる。ところがこの点は漸化式 (9) 式によって回避され、補遺にあげた R 言語によるスクリプトコード ($k = 2 \sim 5$ の場合をあげているが、 $k \geq 6$ も同じである) などを使って RCN の値は容易に計算できる。したがって、小標本の場合には全配列分布を使った取り組みを行うとよいであろう。

$k \geq 5$ の場合の大標本法については漸近分散の式が得られていない。そのため、この点については今後の課題として残されることになる。

補遺

R 言語によるスクリプトの例

NMEA の計算 ($k = 2 \sim 5$ までの計算)

```
#nmea_a.r
#NMEA recursive formula
#Cases for k = 2, 3, 4 and 5
#-----
#NMEA2
nmea2 <- function(n, n1) {
  s2 <- numeric(length=n+1)
  for (i in 0:n1) {
    s2[i+1] <- 1
  }
  sum(s2)
}
#NMEA3
nmea3 <- function(n, n1, n2) {
  s3 <- numeric(length=n+1)
  for (i in 0:n1) {
    s3[i+1] <- nmea2(n, n1 + n2 - i)
  }
  sum(s3)
}
#NMEA4
nmea4 <- function(n, n1, n2, n3) {
  s4 <- numeric(length=n+1)
  for (i in 0:n1) {
    s4[i+1] <- nmea3(n, n1 + n2 - i, n3)
  }
  sum(s4)
}
#NMEA5
nmea5 <- function(n, n1, n2, n3, n4) {
  s5 <- numeric(length=n+1)
  for (i in 0:n1) {
    s5[i+1] <- nmea4(n, n1 + n2 - i, n3, n4)
  }
  sum(s5)
}
#-----
```

RCN の計算

```
#rcn_a.r
#calculation of RCN for k = 2, 3, 4 and 5
#Use "nmea_a.r"
#-----
#RCN2
rcn2 <- function(n, n1) {
  rshift <- nmea2(n, n1)
  lshift <- nmea2(n, n-n1)
  cn2 <- lshift - rshift
  cn2/n
}
#RCN3
rcn3 <- function(n, n1, n2) {
  rshift <- nmea3(n, n1, n2)
  lshift <- nmea3(n, n-n1-n2, n2)
  cn3 <- lshift - rshift
  cn3/(n*(n+3)/2)
}
#RCN4
```

```
rcn4 <- function(n, n1, n2, n3) {
  rshift <- nmea4(n, n1, n2, n3)
  lshift <- nmea4(n, n-n1-n2-n3, n3, n2)
  cn4 <- lshift - rshift
  cn4/(choose(n+3,3)-1)
}
#RCN5
rcn5 <- function(n, n1, n2, n3, n4) {
  rshift <- nmea5(n, n1, n2, n3, n4)
  lshift <- nmea5(n, n-n1-n2-n3-n4, n4, n3, n2)
  cn5 <- lshift - rshift
  cn5/(choose(n+4, 4)-1)
}
#-----
```

実行：

```
> source("nmea_a.r")
> source("rcn_a.r")
> nmea3(30, 5, 10)
[1] 81 <--- (5, 10, 15) の NMEA
> nmea4(50, 20, 15, 10)
[1] 13216 <--- (50, 15, 10, 5) の NMEA
> rcn3(30, 5, 10)
[1] 0.4343 <--- (5, 10, 15) の RCN
> rcn4(50, 20, 15, 10)
[1] -0.4869 <--- (20, 15, 10, 5) の RCN
```

参考文献

- [1] Choi, K. and Matsui, T., Permutation test for probability shift of a discrete distribution, 情報科学研究, 第 10 号、pp1-7, 1992.
- [2] Edgington, E.S., Randomization tests, Third Ed., Marcel Dekker, Inc., 1995.
- [3] Ernst, M.D., Permutation methods: A basis for exact inference, Statistical Science, Vol.19, No.4, 2004.
- [4] Good, P., Permutation tests, Second Ed., Springer, 2000.
- [5] Matsui, T. and Choi, K., On some properties of the number of more extreme arrangements, 獨協経済、第 62 号、pp29-39, 1996.
- [6] Matsui, T., Testing a shift of categorical response probabilities for the associated probability model, 獨協経済、第 78 号、pp5-14, 2004.

(2013 年 9 月 30 日受付)
(2013 年 12 月 18 日採録)