



| | |
|--------------|---|
| Title | クロス表のカテゴリー統合 |
| Author(s) | 佐藤, 裕 |
| Citation | 年報人間科学. 1990, 11, p. 1-15 |
| Version Type | VoR |
| URL | https://doi.org/10.18910/12626 |
| rights | |
| Note | |

The University of Osaka Institutional Knowledge Archive : OUKA

<https://ir.library.osaka-u.ac.jp/>

The University of Osaka

大阪大学人間科学部 [1990年 3月]

『年報人間科学』第11号 1頁—15頁

クロス表のカテゴリー統合

佐 藤 裕

クロス表のカテゴリー統合

1. 本稿の目的

我々が実際に調査データなどを分析する際に、クロス表を構成する変数のカテゴリーを統合し、より小さなクロス表に加工したいという欲求を持つ場合がある。これは、主に次のような理由によるものである。

まず第一に、全ケース数が少なすぎたり、カテゴリーが多すぎたりして1つのセルの度数が少なくなりすぎる場合である。クロス表の分析でよく用いられるカイ2乗独立性検定はセルの期待度数があまり小さいと適用できないため、カテゴリーの統合によってこれを回避するのである。

第二に、カテゴリーが多すぎてクロス表が見にくく、その表が表わしていることが読みとりにくい場合である。これは特に順序尺度以上の変数をクロス表にするに行なわれる。例えば年齢は多くのデータでは実数で与えられているが、そのままでは、あまりに大きな表になるため、例えば10歳ごとのカテゴリーにまとめて表を作る。これもまたカテゴリーの統合である。

それでは、このような場合にはどのような基準でカテゴリー統合が行なわれているだろうか。実際には統合したい変数の度数分布を見て、なるべくそれぞれのカテゴリーの度数が均等になるようにしたり、縦または横のパーセント(条件付き確率)を見て、それが似ているカテゴリーを統合したり、あるいは先の年齢のように等間隔に区切ってカテゴリーを構成したりしている。

しかし、そのような方法でのカテゴリー統合は本当にデータの特徴をよく表わすことになるだろうか。カテゴリーの統合は元のクロス表の持っていた情報を確実に減少させるが、その失われた情報の中に本当に重要な情報はなかったと言えるのだろうか。これらはカテゴリー統合の統計的妥当性の問題である。

カテゴリーを統合する際の1つの基準は、この統合によって失われる情報量の評価から与えられるだろう。統合の割合に較べて情報量の損失が少ないとき、そのクロス表が持つ情報をより凝縮したような新しいクロス表が得られたと言えるのである。

また、こういった基準を明確にし、いくつかの統合の仕方を比較し、最も適切な統合パターンを選択することができるならば、クロス表の構造を明確にすることもできるだろう。例えば、年齢と何か別の変数のクロス表の場合、年齢の最適なカテゴリー化の選択は、どの年齢で大きな変化があるか、ということを見いだすことにもなるだろう。

このような観点からのクロス表のカテゴリー統合については、すでに坂元慶行氏の業績がある(1)。坂元氏はAIC(赤池情報量基準)を用いてカテゴリー統合を行ったモデルの評価を行い、また最適な統合法を選択するコンピュータプログラムも開発している。

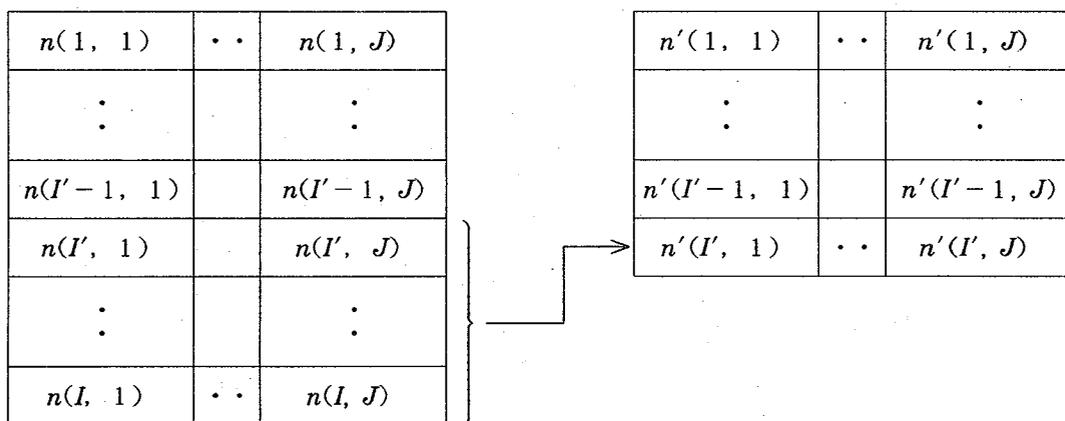
本稿は、坂元氏の方法をもとにして、実際のデータ分析を行う際により容易に利用できる方法と、坂元氏の方法を2変数のカテゴリーを同時に統合した場合へと拡張するモデルを提案するものである。

2. 1変数の統合

1) モデルの設定

$I \times J$ 表の I' 行から I 行までを1つのカテゴリーとしてまとめ、新たに $I' \times J$ 表を作ること考えてみよう(図1)。まとめる箇所が隣接しているのは表記の便宜のためであり、カテゴリーの順序は以後の議論に無関係であるので、離れたカテゴリーの統合も同様に扱える。

図1



カテゴリーを1つにまとめることができるための条件は、まとめるそれぞれのカテゴリーで、列の変数についての条件付き確率が等しい、すなわち、

$$p(j|i_1) = p(j|i_2) \quad (j=1, 2, \dots, J, I' \leq i_1, i_2 \leq I) \quad (2.1)$$

であるとする。例えば行を年齢、列をある意見とするクロス表の場合、意見の回答分布が等しい年齢カテゴリーは統合できる、と考えるのである。これはごく自然な条件であると思われる。この条件付き確率はセルの度数を用いて表わすと、

$$p(j|i) = \frac{n(i, j)}{n(i, \cdot)}$$

であるから、式(2.1)は、

$$\frac{n(i_1, j)}{n(i_1, \cdot)} = \frac{n(i_2, j)}{n(i_2, \cdot)} \quad (j=1, 2, \dots, J, I' \leq i_1, i_2 \leq I)$$

と書ける。したがって、

$$\frac{n(i_1, j)}{n(i_2, j)} = \frac{n(i_1, \cdot)}{n(i_2, \cdot)} \quad (j=1, 2, \dots, J, I' \leq i_1, i_2 \leq I) \quad (2.2)$$

ところが、式(2.2)の右辺は j に無関係であるので、左辺はすべての j について等しい。すなわち、

$$\frac{n(i_1, j_1)}{n(i_2, j_1)} = \frac{n(i_1, j_2)}{n(i_2, j_2)} \quad (j_1, j_2=1, 2, \dots, J, I' \leq i_1, i_2 \leq I) \quad (2.3)$$

である。また逆に、 $n(i, \cdot)$ は $n(i, j)$ をすべての j について加えたものだから、式(2.3)が成り立つならば式(2.2)も成り立つ。したがって式(2.3)は式(2.1)の必要充分条件である。式(2.3)は、統合してできたセルの度数を、各列に共通な比によって配分したものが、もとのセルの度数と等しい、ということの意味する。

カテゴリー統合の妥当性は、統合してできたクロス表から、式(2.3)の条件によってもとのクロス表を推定したときの適合度で評価できる。このときの推定モデルは以下のように表わすことができる。

もとのクロス表のセルの期待度数を $m(i, j)$ とすると、

$$m(i, j) = \begin{cases} \frac{a(i, j)}{n} & (i < I') \\ \frac{a(I', j) \cdot b(i)}{n} & (i \geq I') \end{cases} \quad (2.4)$$

である。

ここで、 $a(i, j)$ は統合してできたクロス表の i 行 j 列のセルの確率、 $b(i)$ は統合したセルをもとのセルに配分する比率である。ただし、

$$\sum_{i=1}^{I'} \sum_{j=1}^J a(i, j) = 1, \quad \sum_{i=I'}^I b(i) = 1 \quad (2.5)$$

である(2)。

2) 最尤推定値

モデルの適合度の評価のためには、まずそのモデルのパラメータと期待度数の最尤推定値を求めなくてはならない。

パラメータ $\{a(i, j), b(i)\}$ において、実現度数 $\{n(i, j)\}$ が見られる確率は、

$$\begin{aligned} & M(\{n(i, j)\} | \{a(i, j), b(i)\}) \\ &= \frac{n!}{\prod_{i=1}^{I'} \prod_{j=1}^J n(i, j)!} \prod_{i=1}^{I'} \prod_{j=1}^J a(i, j)^{n(i, j)} \prod_{i=I'}^I \prod_{j=1}^J a(i, j) b(i)^{n(i, j)} \end{aligned}$$

したがって、 $\{a(i, j), b(i)\}$ についての対数尤度は、パラメータに無関係な定数項を無視すると、

$$\begin{aligned}
& l(\{a(i, j), b(i)\}) \\
&= \sum_{i=1}^{I-1} \sum_{j=1}^J n(i, j) \log a(i, j) + \sum_{i=I}^I \sum_{j=1}^J n(i, j) \log a(I, j) + \sum_{i=I}^I \sum_{j=1}^J n(i, j) \log b(i) \\
&= \sum_{j=1}^J \left(\sum_{i=1}^{I-1} n(i, j) \log a(i, j) + \sum_{i=I}^I n(i, j) \log a(I, j) \right) + \sum_{i=I}^I \sum_{j=1}^J n(i, j) \log b(i) \quad (2.6)
\end{aligned}$$

ここで、

$$n'(i, j) = \begin{cases} n(i, j) & (i < I) \\ \sum_{i=I}^I n(i, j) & (i = I) \end{cases}$$

とすると、式(2.6)の右辺の第1項は、

$$\sum_{i=1}^{I-1} \sum_{j=1}^J n'(i, j) \log a(i, j)$$

となる。

式(2.6)を最大化する $a(i, j)$ を求めるためには、式(2.5)の制約より

$$a(I, j) = 1 - \sum_{i=1}^{I-1} a(i, j)$$

を式(2.6)に代入して $a(i, j)$ で偏微分し、これを0にするような $a(i, j)$ を求めればよい(3)。その結果、 $a(i, j)$ の最尤推定値 $\hat{a}(i, j)$ は、

$$\begin{aligned}
\hat{a}(i, j) &= \frac{n'(i, j)}{n} \\
&= \begin{cases} \frac{n(i, j)}{n} & (i < I) \\ \frac{\sum_{i=I}^I n(i, j)}{n} & (i = I) \end{cases}
\end{aligned}$$

である。

また同様に、 $b(i)$ の最尤推定値 $\hat{b}(i)$ は、

$$\hat{b}(i) = \frac{n(i, \cdot)}{\sum_{i=I}^I n(i, \cdot)}$$

よって、期待度数 $m(i, j)$ の最尤推定値は

$$\hat{m}(i, j) = \begin{cases} n(i, j) & (i < I) \\ \frac{\sum_{i=I}^I n(i, j) \cdot n(i, \cdot)}{\sum_{i=I}^I n(i, \cdot)} & (i \geq I) \end{cases}$$

カテゴリーを統合し、そのときの条件確率が等しいというモデルでの $m(i, j)$ の最尤推定値は、統合しないセルについては実現度数と等しく、統合したセルは統合してできたセルの度数を行の周辺度数によって比例配分したものである。

図2 aの2行目と3行目を統合するというモデルを考えると、その時の期待度数は図2 bのようになる。例えば、2行1列のセルの期待度数は、 $(100+0) \times 100 / (100+400) = 20$ である。

図2 a

| | | | |
|-----|-----|-----|------|
| 400 | 50 | 50 | 500 |
| 0 | 50 | 50 | 100 |
| 100 | 100 | 200 | 400 |
| 500 | 200 | 300 | 1000 |

図2 b

| | | | |
|-----|-----|-----|------|
| 400 | 50 | 50 | 500 |
| 20 | 30 | 50 | 100 |
| 80 | 120 | 200 | 400 |
| 500 | 200 | 300 | 1000 |

3) モデルの検定

統合の妥当性を検定するためには、尤度比統計量 L^2 を使うことができる。 L^2 は一般に以下の式で与えられる。

$$L^2 = 2 \sum_{i=1}^I \sum_{j=1}^J n(i, j) \log \frac{n(i, j)}{m(i, j)}$$

ところが、 $i < I'$ の範囲では、 $n(i, j) = m(i, j)$ であり、 Σ の右はすべて0になるので、カテゴリー統合の妥当性を検定する L^2_1 は次の式で求められる。

$$L^2_1 = 2 \sum_{i=I'}^I \sum_{j=1}^J n(i, j) \log \frac{n(i, j) \cdot \sum_{i=I'}^I n(i, \cdot)}{\sum_{i=I'}^I n(i, j) \cdot n(i, \cdot)}$$

この式は、統合する部分を1つのクロス表と見なし、独立性の検定を行なった場合の式と同じである。

自由度を求めるにはセルの数から (パラメータ数 + 1) を引けばよい。 $a(i, j)$ は式 (2.5) の制約より、 $I'J - 1$ 個、 $b(i)$ も同じく $I - I'$ 個であるから、自由度 df_1 は、

$$\begin{aligned} df_1 &= IJ - (IJ - 1 + I - I' + 1) \\ &= (I - I')(J - 1) \end{aligned}$$

$I - I'$ は結合する部分の行数 - 1 と等しいので、自由度もまた結合する部分を1つのクロス表と見なして独立性の検定を行なった場合と等しい。

以上のことから、カテゴリー統合の妥当性を検定するには、統合する部分を1つのクロス表と見なして独立性の検定を行なえばよいということがわかる。

また、統合してできたクロス表の独立モデルの尤度比統計量を L^2_2 とすると、

$$L^2_2 = 2 \sum_{j=1}^J \left(\sum_{i=1}^{I-1} n(i, j) \log \frac{n \cdot n(i, j)}{n(i, \cdot) n(\cdot, j)} + \sum_{i=I}^I n(i, j) \cdot \log \frac{n \cdot \sum_{i=I}^I n(i, j)}{\sum_{i=I}^I n(i, \cdot) n(\cdot, j)} \right)$$

これに、先ほどの L^2_1 を加えると、

$$\begin{aligned} L^2_1 + L^2_2 &= 2 \sum_{j=1}^J \left\{ \sum_{i=1}^{I-1} n(i, j) \log \frac{n \cdot n(i, j)}{n(i, \cdot) n(\cdot, j)} + \right. \\ &\quad \left. \sum_{i=I}^I n(i, j) \log \left(\frac{n \cdot \sum_{i=I}^I n(i, j)}{\sum_{i=I}^I n(i, \cdot) n(\cdot, j)} \cdot \frac{n(i, j) \cdot \sum_{i=I}^I n(i, \cdot)}{\sum_{i=I}^I n(i, j) \cdot n(i, \cdot)} \right) \right\} \\ &= 2 \sum_{j=1}^J \left(\sum_{i=1}^{I-1} n(i, j) \log \frac{n \cdot n(i, j)}{n(i, \cdot) n(\cdot, j)} + \sum_{i=I}^I n(i, j) \log \frac{n \cdot n(i, j)}{n(i, \cdot) n(\cdot, j)} \right) \\ &= 2 \sum_{i=1}^I \sum_{j=1}^J n(i, j) \log \frac{n \cdot n(i, j)}{n(i, \cdot) n(\cdot, j)} \end{aligned}$$

となり、これはもとのクロス表の独立モデルの尤度比統計量と等しい。すなわち、もとのクロス表の独立性の尤度比統計量を L^2_3 とすると、

$$L^2_3 = L^2_1 + L^2_2 \quad (2.7)$$

という式が成り立っているのである。また、自由度についても、

$$\begin{aligned} df_1 + df_2 &= (I-I')(J-1) + (I-1)(J-1) \\ &= (I-1)(J-1) \\ &= df_3 \end{aligned} \quad (2.8)$$

と、同様の関係が成り立つ。

すなわち、カテゴリー統合の妥当性を検定するには、もとのクロス表の L^2 と統合してできたクロス表の L^2 の差を自由度の差で検定すればよい、ということである。

これまでの議論は 1ヶ所のみ統合のみを扱っていたが、2ヶ所以上を統合する場合もこれをそのまま拡張することができる。統合の妥当性の L^2 は統合したそれぞれの箇所を 1つのクロス表と見なした独立性の L^2 をすべて足し合わせたものになり、また、式 (2.7)、(2.8) も 1ヶ所のみ統合と同様成り立っている。

3. 2変数の統合

(1) モデル設定と最尤推定値

次に、行と列の両方を統合する場合について考えてみよう。 $I \times J$ 表の I' 行目から I 行目と J' 列目から J 列目を統合するとすると、統合してできたクロス表からもとのクロス表を推定するモデルは、1変数の場合と同様にもとのクロス表のセルの期待度数を $m(i, j)$ とすると、

$$m(i, j) = \begin{cases} \frac{a(i, j)}{n} & (i < I', j < J') \\ \frac{a(I', j) \cdot b(i)}{n} & (i \geq I', j < J') \\ \frac{a(i, J') \cdot c(j)}{n} & (i < I', j \geq J') \\ \frac{a(I', J') \cdot b(i) \cdot c(j)}{n} & (i \geq I', j \geq J') \end{cases}$$

ここで、 $a(i, j)$ は統合してできたクロス表の i 行 j 列のセルの確率、 $b(i)$ 、 $c(j)$ は統合したセルをもとのセルに配分する比率である。ただし、

$$\sum_{i=1}^{I'} \sum_{j=1}^{J'} a(i, j) = 1, \quad \sum_{i=1}^{I'} b(i) = 1, \quad \sum_{j=1}^{J'} c(j) = 1 \quad (3.1)$$

$\{a(i, j), b(i), c(j)\}$ についての対数尤度は、パラメータに無関係な定数項を無視すると、

$$\begin{aligned} & l(\{a(i, j), b(i), c(j)\}) \\ &= \sum_{i=1}^{I'-1} \sum_{j=1}^{J'-1} n(i, j) \log a(i, j) \\ &+ \sum_{i=I'}^{I'} \sum_{j=1}^{J'-1} n(i, j) \log a(I', j) + \sum_{i=1}^{I'-1} \sum_{j=J'}^{J'} n(i, j) \log b(i) \\ &+ \sum_{i=1}^{I'-1} \sum_{j=J'}^{J'} n(i, j) \log a(i, J') + \sum_{i=I'}^{I'} \sum_{j=J'}^{J'} n(i, j) \log c(j) \\ &+ \sum_{i=I'}^{I'} \sum_{j=1}^{J'-1} n(i, j) \log a(I', J') + \sum_{i=1}^{I'-1} \sum_{j=J'}^{J'} n(i, j) \log b(i) \\ &+ \sum_{i=I'}^{I'} \sum_{j=J'}^{J'} n(i, j) \log c(j) \end{aligned} \quad (3.2)$$

ここで、

$$n'(i, j) = \begin{cases} n(i, j) & (i < I', j < J') \\ \sum_{i=I'}^{I'} n(i, j) & (i = I', j < J') \\ \sum_{j=J'}^{J'} n(i, j) & (i < I', j = J') \\ \sum_{i=1}^{I'} \sum_{j=J'}^{J'} n(i, j) & (i = I', j = J') \end{cases}$$

とし、式 (3.2) から $a(i, j)$ を含む項のみを抜き出すと、

$$\sum_{i=1}^{I'} \sum_{j=1}^{J'} n'(i, j) \log a(i, j)$$

となるので、1変数の場合と同様に、(3.1) の制約式を代入し、それぞれのパラメータで偏微分して $a(i, j)$ の最尤推定値を求めると、

$$\hat{a}(i, j) = \frac{n'(i, j)}{n}$$

また、 $b(i)$ を含む項は、

$$\sum_{i=I'}^I \sum_{j=1}^J n(i, j) \log b(i)$$

と整理できるので、(3. 1) の制約式から $b(i)$ の最尤推定値は

$$\hat{b}(i) = \frac{n(i, \cdot)}{\sum_{i=I'}^I n(i, \cdot)}$$

$c(j)$ の最尤推定値も同様に、

$$\hat{c}(j) = \frac{n(\cdot, j)}{\sum_{j=J'}^J n(\cdot, j)}$$

となる。

したがって、 $m(i, j)$ の最尤推定値は、

$$\hat{m}(i, j) = \begin{cases} n(i, j) & (i < I', j < J') \\ \frac{\sum_{i=I'}^I n(i, j) \cdot n(i, \cdot)}{\sum_{i=I'}^I n(i, \cdot)} & (i \geq I', j < J') \\ \frac{\sum_{j=J'}^J n(i, j) \cdot n(\cdot, j)}{\sum_{j=J'}^J n(\cdot, j)} & (i < I', j \geq J') \\ \frac{\sum_{i=I'}^I \sum_{j=J'}^J n(i, j) \cdot n(i, \cdot) \cdot n(\cdot, j)}{\sum_{i=I'}^I n(i, \cdot) \cdot \sum_{j=J'}^J n(\cdot, j)} & (i \geq I', j \geq J') \end{cases}$$

図2 aのクロス表の2行目と3行目、2列目と3列目をそれぞれ統合したときの期待度数は図3ようになる。例えば2行3列のセルの場合であれば、 $(50+50+100+200) \times 100 \times 200 / (100+400) / (200+300) = 32$ である。

図3

| | | | |
|-----|-----|-----|------|
| 400 | 40 | 60 | 500 |
| 20 | 32 | 48 | 100 |
| 80 | 128 | 192 | 400 |
| 500 | 200 | 300 | 1000 |

(2) モデルの検定

統合の妥当性を検定する L^2_1 は、やはり行も列も統合しない範囲は無視できるので、

$$L^2_1 = 2 \sum_{i=I'}^I \sum_{j=1}^{J'} n(i, j) \log \frac{n(i, j) \cdot \sum_{i=I'}^I n(i, \cdot)}{\sum_{i=I'}^I n(i, j) \cdot n(i, \cdot)}$$

$$+ 2 \sum_{i=1}^{I'} \sum_{j=J'}^J n(i, j) \log \frac{n(i, j) \cdot \sum_{j=J'}^J n(\cdot, j)}{\sum_{j=J'}^J n(i, j) \cdot n(\cdot, j)}$$

$$+ 2 \sum_{i=I'}^I \sum_{j=J'}^J n(i, j) \log \frac{n(i, j) \cdot \sum_{i=I'}^I n(i, \cdot) \cdot \sum_{j=J'}^J n(\cdot, j)}{\sum_{i=I'}^I \sum_{j=J'}^J n(i, j) \cdot n(i, \cdot) \cdot n(\cdot, j)}$$

これに、統合してできたクロス表の独立性 L^2

$$L^2_2 = 2 \sum_{j=1}^{J'} \sum_{i=1}^{I'} n(i, j) \log \frac{n \cdot n(i, j)}{n(i, \cdot) n(\cdot, j)}$$

$$+ 2 \sum_{i=I'}^I \sum_{j=1}^{J'} n(i, j) \log \frac{n \cdot \sum_{i=I'}^I n(i, j)}{\sum_{i=I'}^I n(i, \cdot) \cdot n(\cdot, j)}$$

$$+ 2 \sum_{i=1}^{I'} \sum_{j=J'}^J n(i, j) \log \frac{n \cdot \sum_{j=J'}^J n(i, j)}{\sum_{j=J'}^J n(\cdot, j) \cdot n(i, \cdot)}$$

$$+ 2 \sum_{i=I'}^I \sum_{j=J'}^J n(i, j) \log \frac{n \cdot \sum_{i=I'}^I \sum_{j=J'}^J n(i, j)}{\sum_{i=I'}^I n(i, \cdot) \cdot \sum_{j=J'}^J n(\cdot, j)}$$

を加えると、

$$L^2_1 + L^2_2 = 2 \sum_{i=1}^I \sum_{j=1}^J n(i, j) \log \frac{n \cdot n(i, j)}{n(i, \cdot) n(\cdot, j)}$$

となり、やはりもとのクロス表の独立性 L^2 と等しい。

自由度は、

$$df_1 = IJ - (I'J' - 1 + I - I' + J - J' + 1)$$

$$= IJ - I'J' - I + I' - J + J'$$

$$df_2 = (I' - 1)(J' - 1)$$

$$= I'J' - I' - J' + 1$$

$$df_3 = (I - 1)(J - 1)$$

$$= IJ - I - J + 1$$

であるので、やはり

$$df_1 + df_2 = df_3$$

が成り立っている。

2変数の統合の場合は、このように、もとのクロス表と統合してできたクロス表の独立性 L^2 を比較することによって間接的に妥当性を検定する方法が便利である。

ここまでで紹介した方法は、特別な計算プログラムを必要とせず、SPSSやSASなどの汎用統計パッケージを利用してだれでもすぐに計算することができる。特に、統合する前後のクロス表の独立性 L^2 を比較する方法は、1変数を統合する場合も両方の変数を統合する場合も利用でき、しかも汎用統計パッケージの出力を比較してカイ2乗分布表を調べるだけなので、非常に手軽で有効な方法である。

4. 統合法の選択

カテゴリーの統合を行なう場合に、いくつかの方法が考えられる場合がある。例えば年齢をカテゴリー化する場合に、10歳ごとに区切るのか、あるいはなるべくそれぞれのカテゴリーに含まれる人数が等しくなるようにするのか、といった場合である。このようないくつかの統合法から1つを選択するといった場合にも、今まで述べた方法が適用できる。

複数の統合法を比較するための基準としては、AIC (赤池情報量基準) が便利である。AICは次の式で求めることができる(4)。

$$AIC = L^2 - 2df$$

L^2 も df も今まで述べた方法で計算できるので、AICの計算も簡単である。こうしてすべての統合法について計算し、AICが最も小さい統合法が、もとのクロス表の持つ情報を最も凝縮して表わした最適なクロス表である。

クロス表の統合法は、すべての可能性を含めると、比較的小さなクロス表でも膨大な数になる。例えば5カテゴリーを統合する組合せは51通りだが、これが6カテゴリーになると202通りへと急増する。さらに行と列の両方を統合する組合せを考えると単純に計算してこれの2乗、 5×5 だと2601通りにのぼる。このような膨大な計算を行なうには当然コンピュータの利用が考えられる。コンピュータを用いればこの程度の計算なら楽々とこなしてくれるので機械的に最適なクロス表を選択するには有効である。

しかし、もっと大きなクロス表の場合は、組合せがあまりにも多くなりすぎ、またクロス表を構成する変数の性質や分析の目的によってはすべての組合せを計算する必要がない場合もある。例えば変数が順序尺度を構成していると考えられる場合は隣合うカテゴリー以外の統合は無意味であり、これを除くと、考えられる統合の組合せは大幅に減少する(5)。

従来の独立性の検定は、行または列のすべてのカテゴリーを1つにまとめた場合の妥当性の検討と等しい。様々な統合法はもとのクロス表と、すべてを1カテゴリーにまとめた場合との

中間に位置づけられる。コンピュータの利用によって最適なクロス表を選択することは、2つの変数が独立であるかそうでないか、という両極端の中間にいくつものレベルを設定し、どのレベルが最も適当かということを求めることであるとも言える。

5. 例題

最後に、前節で述べたような最適なクロス表の選択法を、実際のデータに適用して考えてみよう。表4は85年SSM調査報告書からのデータである(6)。これは40歳以上の既婚女性のキャリアタイプ別の年齢を表わしている。この表をそのまま見ると結婚退職や中年退職は年齢がやや低く、就業継続は年齢が高いように見えるが、それほどはっきりしたものではない。そこでまずキャリアタイプによる年齢構成の違いをはっきりさせるために、キャリアタイプのみを統合を行ってみると、AICが最も小さい統合法は表5のようになった。これを見ると、就業継続と不就業が類似した年齢構成を持ち、中年退職、再就職、結婚退職の3つも年齢構成が類似している。前者は比較的高い年齢が多く、後者は若い年齢層に多い。逆に、年齢によるキャリアタイプの違いに着目して年齢カテゴリーの統合を行なったのが表6である。この場合、年齢は順序尺度を構成すると考えて隣合うカテゴリーのみを統合した。表5と同様に、就業継続と不就業は高い年齢層ほど比率が多くなっており、中年退職は高い年齢層ほど比率が低くなっているが、再就職は年齢による差がほとんどなく、結婚退職も高い年齢層ほど比率が低くなっているものの、それほど大きな差ではない。最後に年齢、キャリアタイプ両方を統合したのが表7である。このときの最適な統合法は、それぞれの変数を別々に統合した場合と必ずしも一致しないことが注目される。この表から読み取れることは、就業継続が最も高い年齢層に集中し、ついで不就業も高い年齢層に偏っている。最も若い年齢層に集中しているのは中年退職であり、再就職と結婚退職はやや若い年齢層に偏っている。

表4と表7を見比べると、表4では細かい数値に紛れて読み取りにくかった傾向が、表4では非常にはっきりと現われてきている。このように、最適な統合法を求めることは複雑な構造を持つクロス表を解読する場合にも有効な方法であると言えよう。

表 4

| | | 年齢 | | | | | | |
|-------------|------|---------------|---------------|---------------|---------------|--------------|--------------|-----|
| 実数/比率 (%) | | 40~44歳 | 45~49歳 | 50~54歳 | 55~59歳 | 60~64歳 | 65~69歳 | 合計 |
| キャリア タイプ | 就業継続 | 5 (6.8) | 9 (12.2) | 19 (25.7) | 9 (12.2) | 11 (14.9) | 21 (28.4) | 74 |
| | 中年退職 | 10 (18.9) | 17 (32.1) | 11 (20.8) | 10 (18.9) | 2 (3.8) | 3 (5.7) | 53 |
| | 再就職 | 14 (14.3) | 27 (27.6) | 21 (21.4) | 13 (13.3) | 12 (12.2) | 11 (11.2) | 98 |
| | 結婚退職 | 78 (21.3) | 83 (22.6) | 70 (19.1) | 66 (18.0) | 34 (9.3) | 36 (9.8) | 367 |
| | 不就業 | 16 (14.5) | 17 (15.5) | 22 (20.0) | 21 (19.1) | 16 (14.5) | 18 (16.4) | 110 |
| | 合計 | 123 (17.5) | 153 (21.8) | 143 (20.4) | 119 (17.0) | 75 (10.7) | 89 (12.7) | 702 |

表 5

| | | 年齢 | | | | | | |
|-------------|------|---------------|---------------|---------------|---------------|--------------|--------------|-----|
| 実数/比率 (%) | | 40~44歳 | 45~49歳 | 50~54歳 | 55~59歳 | 60~64歳 | 65~69歳 | 合計 |
| キャリア タイプ | 就業継続 | 21 (11.4) | 26 (14.1) | 41 (22.3) | 30 (16.3) | 27 (14.7) | 39 (21.2) | 184 |
| | 不就業 | 102 (19.7) | 127 (24.5) | 102 (19.7) | 89 (17.2) | 48 (9.3) | 50 (9.7) | 518 |
| | 中年退職 | | | | | | | |
| | 再就職 | | | | | | | |
| | 結婚退職 | | | | | | | |
| | 合計 | 123 (17.5) | 153 (21.8) | 143 (20.4) | 119 (17.0) | 75 (10.7) | 89 (12.7) | 702 |

AIC=-12,329

表 6

| | | キャリアタイプ | | | | | |
|-----------|--------|--------------|-------------|--------------|---------------|---------------|-----|
| 実数/比率 (%) | | 就職継続 | 中年退職 | 再就職 | 結婚退職 | 不就業 | 合計 |
| 年齢 | 40~49歳 | 14 (5.1) | 27 (9.8) | 41 (14.9) | 161 (58.3) | 33 (12.0) | 276 |
| | 50~59歳 | 28 (10.7) | 21 (8.0) | 34 (13.0) | 136 (51.9) | 43 (16.4) | 262 |
| | 60~69歳 | 32 (19.5) | 5 (3.0) | 23 (14.0) | 70 (42.7) | 34 (20.7) | 164 |
| | 合計 | 74 (10.5) | 53 (7.5) | 98 (14.0) | 367 (52.3) | 110 (15.7) | 702 |

AIC=-13,954

表7

| 実数/比率(%) | | 年齢 | | | 合計 |
|-------------|------|---------------|---------------|---------------|-----|
| | | 40~49歳 | 50~59歳 | 60~69歳 | |
| キャリア タイプ | 就業継続 | 14 (18.9) | 28 (37.8) | 32 (43.2) | 74 |
| | 中年退職 | 27 (50.9) | 21 (39.6) | 5 (9.4) | 53 |
| | 再就職 | 202 (43.4) | 170 (36.6) | 93 (20.0) | 465 |
| | 結婚退職 | 33 (30.0) | 43 (39.1) | 34 (30.9) | 110 |
| | 不就業 | 276 (39.3) | 262 (37.3) | 164 (23.4) | 702 |
| | 合計 | | | | |

AIC=-17,041

注

- (1) 坂元慶行、『カテゴリカルデータのモデル分析』、共立出版、1985
- (2) 坂元慶行氏は同じ条件から、条件付き確率をパラメータとしたモデルをつくり、以下と同様の結論を得ている。本稿のモデルと坂元氏のモデルは同等であるが、後で述べる2変数の統合への拡張を容易にするためにこのモデルを採用した。
坂元慶行、同書、p31
- (3) この部分の計算については、以下を参照のこと。
坂元慶行、同書、p8
- (4) 正確には、最もパラメータ数の多いモデル(飽和モデル)のAICとの差である。
坂元慶行、同書、p44
- (5) 上田尚一氏は、最も適当な2つのカテゴリー統合を繰り返して、徐々にカテゴリー数を減らしていく階層的手法を提唱している。分析の目的によってはこれを用いることもできるだろう。一般には階層的手法での解と総当たりでの最適解は一致しない。以下の例題の計算は総当たり法で計算した。
上田尚一、『データ解析の方法』、朝倉書店、1982、p99
- (6) 岡本英雄、「女性の職業経歴」、『1985年社会階層と社会移動全国調査報告書 第4巻』
1985年社会階層と社会移動全国調査委員会、1989、pp61-74