

## On the number of alternatives and correlation.

|       |   |
|-------|---|
| メタデータ | 言語: jpn<br>出版者:<br>公開日: 2017-10-02<br>キーワード (Ja):<br>キーワード (En):<br>作成者:<br>メールアドレス:<br>所属: |
| URL   | <a href="http://hdl.handle.net/2297/5116">http://hdl.handle.net/2297/5116</a>               |

## 選択肢数と相関係数について<sup>1)</sup>

岡本安晴

### 1 はじめに

質問紙などにおいては、一般にデータは離散的である。これは、解答者に要求される反応型が言語的であるという利点がある。例えば、お菓子の甘さの評定のときに、100点満点で評定する場合と、甘い—甘くないの2件法で評定する場合とを比べてみればよい。2件法の場合、データの分析において2項目間の相関を求めるのにファイ係数が用いられることがある。背後に想定される要因が連続型のときに、それらの連続型の要因間の相関係数の推定値として用いられたファイ係数には大きな偏りの生じることがある (cf. 岡本 (1987))。質問紙等における選択肢数については、2件法の他、「甘い—どちらでもない—辛い」の3件法など種々の多件法がある。選択肢数が多い場合には、等間隔に数値を割り当てて連続型のデータとして処理されることがあるが、本稿では、この場合における算出された相関係数の偏りについての検討を試みる。

### 2 モデル

2変数  $x$  と  $y$  が相関係数  $\rho$  の2次元正規分布

$$N(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}\right) \quad (1)$$

に従い、かつ、 $(k+1)$  個の値  $C_{x,i} (i=0,1,\dots, k)$  により  $x$  の値域が

$$C_{x,i-1} < x \leq C_{x,i}$$

のようにカテゴリーに分けられているとする。但し、 $C_{x,0} = -\infty$ 、 $C_{x,k} = +\infty$ 。

確率変数  $X$  を次のようにとる。

$$X = i \quad \text{if } C_{x,i-1} < x \leq C_{x,i}$$

$y$  についても同様に

---

1) 計算はすべてパーソナルコンピュータ PC9801vm2 (NEC) で行い、XYプロッタ DXY-980 (Roland DG) により作図した。プログラムは Turbo Pascal (Borland) で書かれている。

$$-\infty = C_{y,0} < C_{y,1} < \dots < C_{y,k} = +\infty$$

なる  $(k+1)$  個の値  $C_{y,i}$  により  $k$  個のカテゴリーに分け、そのカテゴリーの値をとる確率変数を  $Y$  とおく。

このとき、 $N$  個のデータ  $(X, Y)$  に対して相関係数  $r_{XY}$  は次式により算出される。

$$r_{XY} = \frac{\frac{1}{N} \sum (X - \bar{X})(Y - \bar{Y})}{\sqrt{\frac{1}{N} \sum (X - \bar{X})^2} \sqrt{\frac{1}{N} \sum (Y - \bar{Y})^2}}$$

但し、

$$\bar{X} = \frac{1}{N} \sum X, \quad \bar{Y} = \frac{1}{N} \sum Y$$

$N \rightarrow \infty$  すると、

$$\lim_{N \rightarrow \infty} r_{XY} = \frac{(\sum_{i,j} ij P_{xy,ij}) - (\sum_i iP_{x,i})(\sum_j jP_{y,j})}{\sqrt{(\sum_i i^2 P_{x,i}) - (\sum_i iP_{x,i})^2} \sqrt{(\sum_j j^2 P_{y,j}) - (\sum_j jP_{y,j})^2}}$$

ここに、

$$P_{x,i} = \lim_{N \rightarrow \infty} (X=i \text{ であるデータの比率})$$

$$= \int_{C_{x,i-1}}^{C_{x,i}} \int_{-\infty}^{\infty} N(x,y) dy dx$$

$$P_{y,j} = \lim_{N \rightarrow \infty} (Y=j \text{ であるデータの比率})$$

$$= \int_{C_{y,j-1}}^{C_{y,j}} \int_{-\infty}^{\infty} N(x,y) dx dy$$

$$P_{xy,ij} = \lim_{N \rightarrow \infty} (X=i \text{ かつ } Y=j \text{ であるデータの比率})$$

$$= \int_{C_{x,i-1}}^{C_{x,i}} \int_{C_{y,j-1}}^{C_{y,j}} N(x,y) dx dy$$

### 3 結 果

$k=2, 4, 8$  の場合の  $\lim_{N \rightarrow \infty} r_{XY}$  を  $\rho$  の関数として描いたものが図-1 及び 2 である。 $C_{x,i}, C_{y,j}$  として表-1 に示す値を用いたものが図-1 であり、表-2 の値に対応するものが図-2 である。標準正規分布における  $C_{x,i}, C_{y,j}$  の位置については図-3 を参考にして頂きたい。

図-1 (表-1) は、データが中央のカテゴリーを中心にして対称に分布している場合であるが、選択肢数が 8 の場合は  $\lim_{N \rightarrow \infty} r_{XY}$  は  $\rho$  に十分近いと思われる。選択肢数が 2 であっても (ファイ係数の場合)  $\rho$  との関係はほぼ線型であるとみなせるが、因子分析等における相関係数の推定値として用いた場合には communality の低下が予想される。図-2 は表-2 に示されるようなデータのカテゴリー上の分布に偏りがある場合である。選択肢数が 2 の場合には  $\rho$  が 1 に近づくと  $\lim_{N \rightarrow \infty} r_{XY}$  はほぼ一定となっている。又、 $\rho$  の値が 0 から 1 の範囲の全域において、 $\lim_{N \rightarrow \infty} r_{XY}$  の値は  $\rho$  に比べて相対的に随分と小さい。即ち、岡本(1987)においても示されたように、このような場合にファイ係数を用いることは極めて危険である。選択肢数が 4 以上のときは、 $\rho$  との関係はほぼ線型であり  $\rho$  とともに増加している。これは、表-2 のような偏りの場合は、選択肢数が 4 以上のときデータの分布が 2 つ以上のカテゴリーに散らばっているため図-1 の選択肢数が 2 の場合程度の性質の良さがみられたものと思われる。

以上、図-1 と 2 の結果だけから判断すれば、選択肢数としては 8 程度あればカテゴリー値から直接算出された  $r_{XY}$  を相関係数の推定値として用いてもよいように思われるが、データのカテゴリー上の分布には注意する必要がある。少なくとも 2 つ以上のカテゴリーにデータが分布しているかどうか調べて、もしデータの大部分が 1 つのカテゴリーに集中している場合にはカテゴリー値から直接  $r_{XY}$  を算出して相関係数の推定値とすることは避けねばならない。

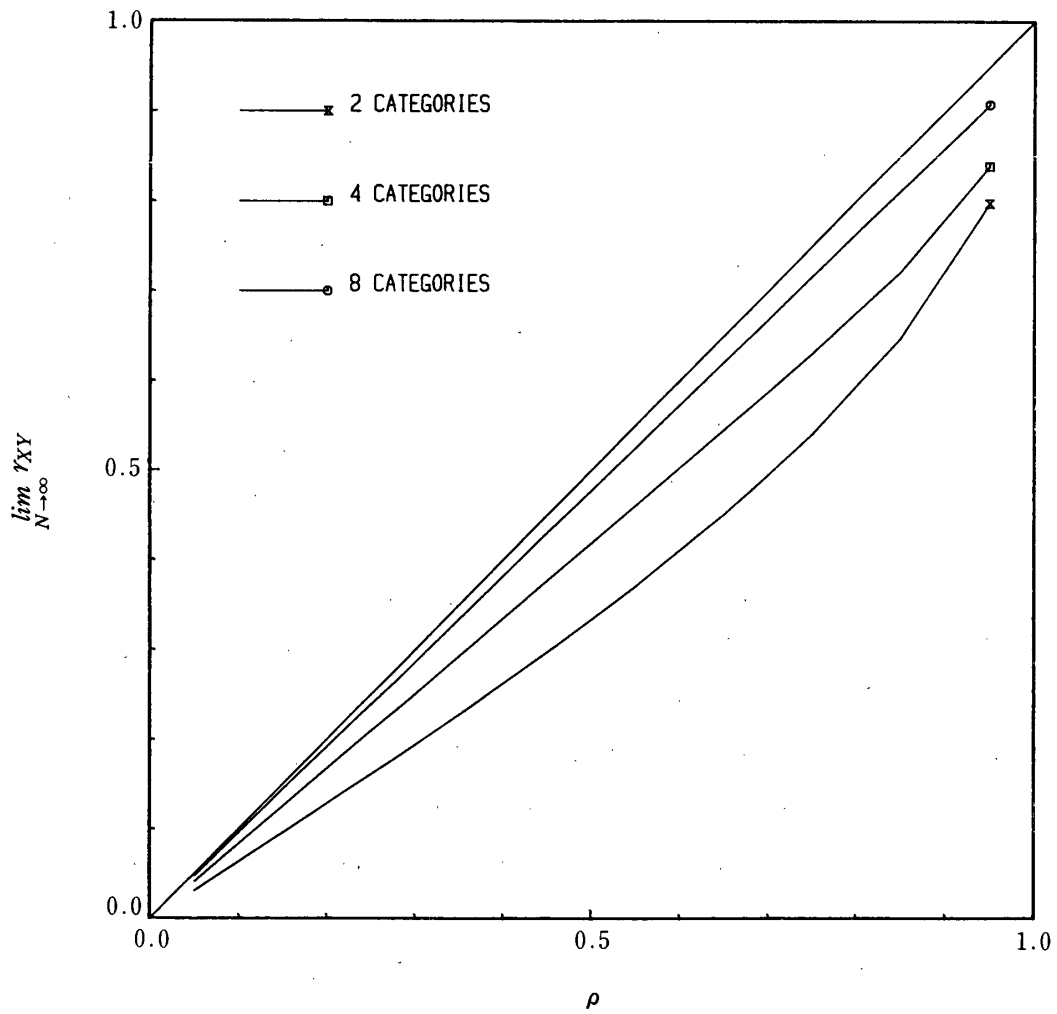


図-1 (1)式における $\rho$ の関数としての $\lim_{N \rightarrow \infty} r_{XY}$ のグラフ。各選択枝数 $k$  ( $k=2, 4, 8$ ) に対して $C_{\alpha,i}$  ( $\alpha=x, y; i=1, \dots, k-1$ ) は表-1に示された値が設定されている。

|       |            | $C_{\alpha,1}$ | $C_{\alpha,2}$ | $C_{\alpha,3}$ | $C_{\alpha,4}$ | $C_{\alpha,5}$ | $C_{\alpha,6}$ | $C_{\alpha,7}$ |
|-------|------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| $k=2$ | $\alpha=x$ | 0.0            |                |                |                |                |                |                |
|       | $\alpha=y$ | 0.0            |                |                |                |                |                |                |
| $k=4$ | $\alpha=x$ | -1.5           | 0.0            | 1.5            |                |                |                |                |
|       | $\alpha=y$ | -1.5           | 0.0            | 1.5            |                |                |                |                |
| $k=8$ | $\alpha=x$ | -2.25          | -1.5           | -0.75          | 0.0            | 0.75           | 1.5            | 2.25           |
|       | $\alpha=y$ | -2.25          | -1.5           | -0.75          | 0.0            | 0.75           | 1.5            | 2.25           |

表-1 カテゴリーの分割点 $C_{\alpha,i}$  ( $\alpha=x, y; i=1, \dots, k-1$ ) の値。分布が中央のカテゴリーを中心として対称になるように値が設定されている。

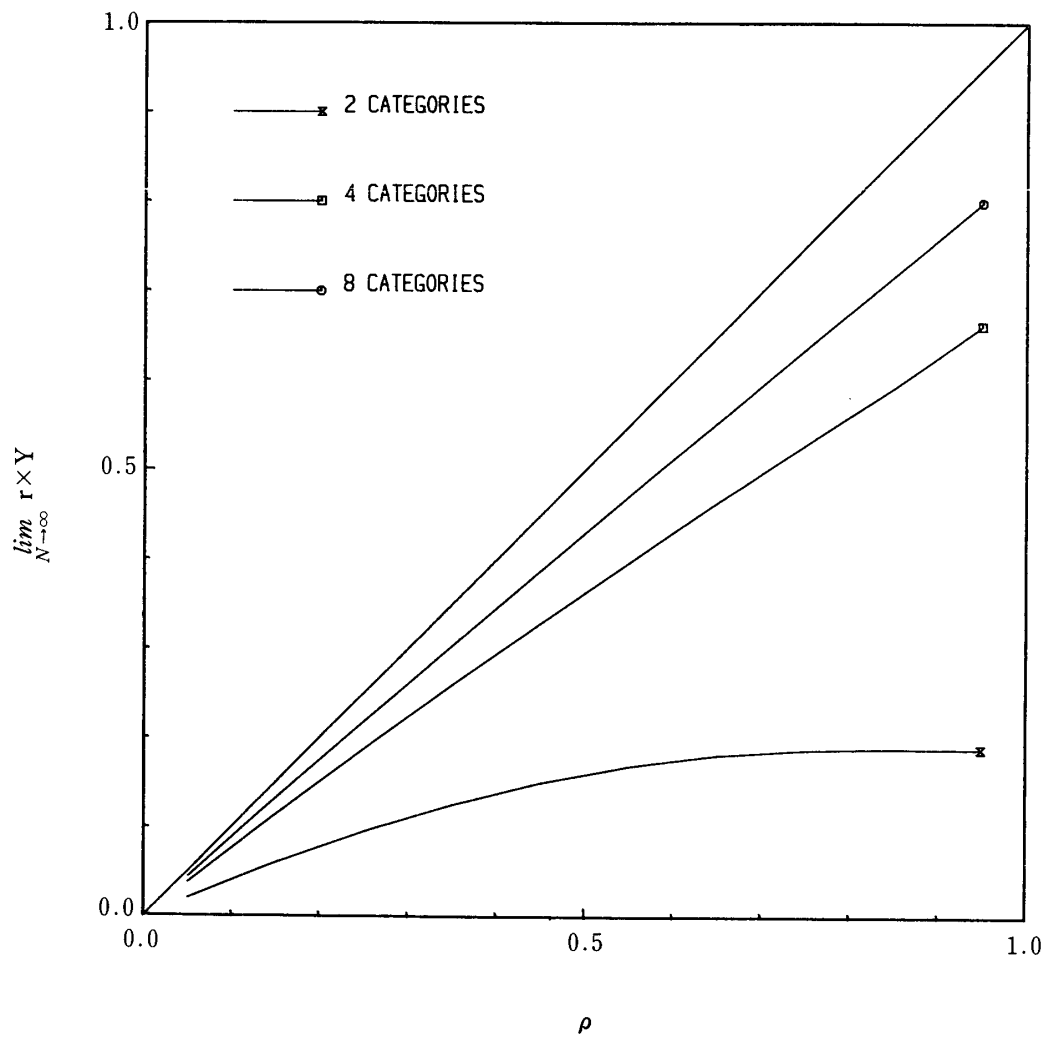


図-2 (1)式における $\rho$ の関数としての $\lim_{N \rightarrow \infty} r_{XY}$ のグラフ。各選択枝数 $k$  ( $k=2, 4, 8$ ) に対して $C_{\alpha,i}$  ( $\alpha=x, y; i=1, \dots, k-1$ ) は表-2に示された値が設定されている。

|       |            | $C_{\alpha,1}$ | $C_{\alpha,2}$ | $C_{\alpha,3}$ | $C_{\alpha,4}$ | $C_{\alpha,5}$ | $C_{\alpha,6}$ | $C_{\alpha,7}$ |
|-------|------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| $k=2$ | $\alpha=x$ | 1.0            |                |                |                |                |                |                |
|       | $\alpha=y$ | -1.0           |                |                |                |                |                |                |
| $k=4$ | $\alpha=x$ | 0.0            | 1.0            | 2.0            |                |                |                |                |
|       | $\alpha=y$ | -2.0           | -1.0           | 0.0            |                |                |                |                |
| $k=8$ | $\alpha=x$ | -0.5           | 0.0            | 0.5            | 1.0            | 1.5            | 2.0            | 2.5            |
|       | $\alpha=y$ | -2.5           | -2.0           | -1.5           | -1.0           | -0.5           | 0.0            | 0.5            |

表-2 カテゴリーの分割点 $C_{\alpha,i}$  ( $\alpha=x, y; i=1, \dots, k-1$ ) の値。カテゴリ上の分布に偏りが生じるように値が設定されている。

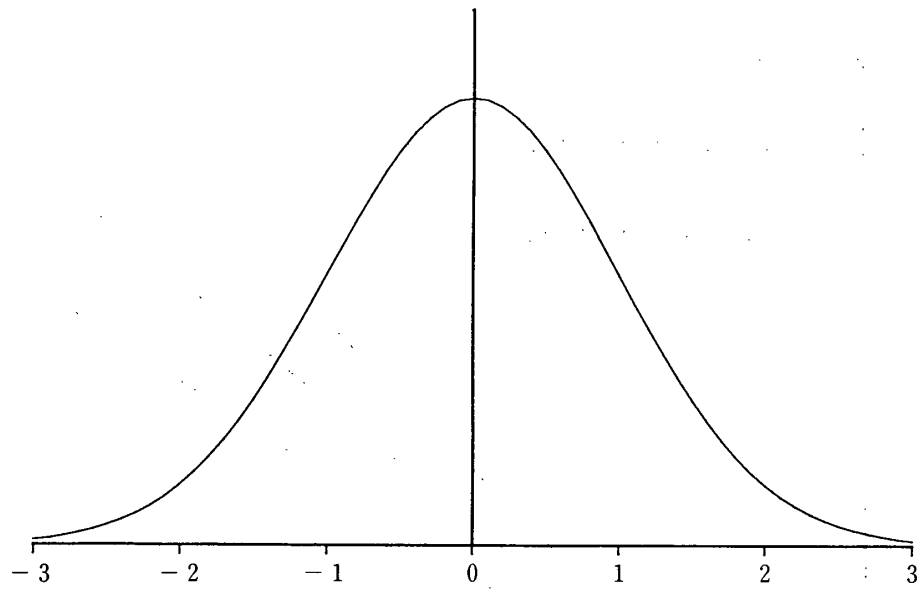


図-3 標準正規分布のグラフ。

引用文献

岡本安晴 1987 ファイ係数について 金沢大学文学部論集行動科学篇 第7号