

オッズ比の高次元への拡張とその性質

矢野 環

オッズ比 (OR) は、近年効果量としての意義も認められ、よく利用されている。またメタ分析の手法も確立している。この 2×2 で定義されている OR を高次元に拡張し、その OR が望ましい不変性を持つように定義する。これまで拡張の提案はあるが、不変性が十分でない。OR はその分散の推定を元に、通常のオッズ比と同じくメタ分析を行うこともできる。拡張されたオッズ比が 1 となるのは、行 (列) が線形従属な時である。実用上、 2×2 の場合とほぼ同様の処理ができ、関連する検定よりも良い基準とできる。

1. 緒言

オッズ比 (Odds ratio, 以下 OR) は従来から疫学・医学・薬学などで良く利用されていた。近年では、効果量としての価値も認められるようになり、英語教育等でも注目されている。

OR を高次元の行列に対して定義するには、 2×2 小行列群の OR の集合を考えるのが主流であり [1, 2, 5]、その可視化の工夫もされ [6, 7, 8]、代数幾何学的考察もある [9]。一個の量としては、Goodman-Kruskal's γ から $OR_G = (1 + \gamma) / (1 - \gamma)$ とする [3, 4] などの提案がある。しかし、それらの量は本来の OR の満たす性質の一部は持っているものの、重要な不変性が失われているなどの欠点もある。このため、 2×2 以外の複数の行列が与えられた時、その群の異質性や同質性を判断しようとする、解りやすい不変量がなかった。

もし OR と同様な不変性を持つ量を定めることができるならば、一群の行列が与えられた場合、その量の Mantel-Haenszel 型 estimate や、メタ分析の手法での推定値が、その群の性質を表現することになるであろう。

以下において、正方行列の場合を主体として、OR の高次元への拡張を提案し、その性質を調べ、実際の適用例を解説する。また全般的理解の為に、

交代群の不変式の一般論から始める。

2. 交代群の不変式

V を体 K 上の n 次元ベクトル空間とし、その座標関数の集合 $\{x_1, x_2, \dots, x_n\}$ を X とする。 X には、 n 次対称群 S_n が添え字の置換として作用する。即ち、置換 $\sigma \in S_n$ に対して、 X 上の置換 $x_i \mapsto x_{\sigma(i)}$ が引き起こされる。

この S_n の X 上への作用は、自然に多項式環 $K[X]$ 上に誘導される。その S_n の作用での不変式環は、対称式の集合 S である。 S_n の元は偶置換と奇置換にわかれ、偶置換全体の集合 A_n は S_n の index 2 の (正規) 部分群をなす。置換は符号 $\text{sign}(\sigma)$ をもち、これは群準同型を与える。

$$\text{sign} : S_n \rightarrow \{\pm 1\}$$

偶置換は符号 1, 奇置換は符号 -1 である。 A_n の不変式であって、 S_n の不変式ではないものを交代式とよぶ。 S と交代式の全体は再び環をなし、半不変式の環と呼ばれる。

交代式 Δ が存在し、半不変式 $= S + \Delta S$ となることが知られている。 K の標数が 2 でない場合は、 Δ として基本差積をとることができる。

$$\Delta = \prod_{i < j} (x_i - x_j)$$

右辺を展開し、係数が 1 である単項式の和を Δ_p ,

係数が-1となる単項式の和を Δn とすれば、

$$\Delta = \Delta p - \Delta n$$

となる。 K の標数が2の場合は、 Δ は対称式であり、半不変式 $= S + \Delta p S = S + \Delta n S$ となる。

以下 K の標数 $\neq 2$ とする。

次に、 $n \times n$ 行列 $M=(x_{ij})$ を考える。第二添え字に置換を含む単項式

$$m(\sigma) = \prod_{i=1}^n x_{i\sigma(i)}$$

は、 i 行 $\sigma(i)$ 列成分の積であり、各行各列から一か所ずつ採択している。この第二添え字を S_n 全体に及ぼした和、

$$\text{perm}(M) = \sum_{\sigma \in S_n} m(\sigma)$$

は、 M のpermanentと呼ばれ、グラフ理論で重要である。また、置換の符号をつけた和

$$\det(M) = \sum_{\sigma \in S_n} \text{sign}(\sigma) m(\sigma)$$

は、 M の行列式(determinant)である。

定義1 detp, detn

$$\text{detp}(M) = \sum_{\sigma \in A_n} m(\sigma),$$

$$\text{detn}(M) = \sum_{\sigma \in S_n \setminus A_n} m(\sigma),$$

と定義する。

命題1 detp, detn いずれも A_n 不変式である。

$$\det(M) = \text{detp}(M) - \text{detn}(M),$$

$$\text{perm}(M) = \text{detp}(M) + \text{detn}(M),$$

となる。

M^σ を奇置換 σ に従い M の行を入れ替えた行列とすれば(列でも同様)次式が成立する。

$$\text{detp}(M^\sigma) = \text{detn}(M), \text{detn}(M^\sigma) = \text{detp}(M).$$

3. 高次元 OR と類似不変量

体 K 上の $n \times n$ 行列全体を $M(n, K)$ と記す。以下、定義と命題を列挙する。命題の証明は容易であるので略す。

定義2 OR

$$\text{Dom} = \{M \in M(n, K); \text{detp}(M) \neq 0, \text{detn}(M) \neq 0\}$$

と置く。

$$\text{OR}: \text{Dom} \rightarrow K \setminus \{0\}$$

$$\text{OR}(M) = \text{detp}(M) / \text{detn}(M)$$

により定義する。

命題2 ORは次の性質を満たす。

1. $\text{OR}(cM) = \text{OR}(M)$, $0 \neq c \in K$.

2. M^T を M の転置行列とすれば

$$\text{OR}(M^T) = \text{OR}(M).$$

3. M' を、 M の2つの行、あるいは2つの列を入れ替えた行列とすれば $\text{OR}(M') = 1 / \text{OR}(M)$.

4. $\text{OR}^{-1}(1) = \text{Dom} \cap \{M; \det(M) = 0\}$.

また、Yule's Q の拡張が定義できる。 $Q(M)$ はGoodman-Kruskal's γ の類似である。

定義3 Q

$$Q(M) = (\text{OR}(M) - 1) / (\text{OR}(M) + 1)$$

と定義する。 OR の定義より次式が成立する。

$$Q(M) = \det(M) / \text{perm}(M).$$

拡張された OR は、加法・乗法と0でない除法について閉じている K の部分集合上に定義できる。例えば、0以上の実数 R_+ においても定義できる。それが通常の OR の拡張となっている。

OR が1となるのは、 $\det(M) = 0$ となる場合、つまり、 M がmulti-colinear(M の行、また列が一次従属)な時である。

以下、 $K = R_+$ とする。 $-1 < Q(M) < 1$ である。

ファイ係数の拡張は次の式で与えられる。下記において、colSums(), rowSums()はそれぞれ列和と行和ベクトル、prod(v)はベクトルvの成分すべての積を意味する。また、sqrt()は平方根である。

定義4 Φ

$$\Phi(M) = \frac{\det(M)}{\sqrt{\text{prod}(\text{colSums}(M), \text{rowSums}(M))}}$$

対数オッズ比も同様に定義する。

定義5 logOR

$\log\text{OR}: \text{Dom} \rightarrow R$ を次式で定める。

$$\log\text{OR}(M) = \log(\text{OR}(M))$$

命題3 関数 $Q, \Phi, \log\text{OR}$ のいずれかを T とする。

このとき T は次の性質を満たす。

1. $T(cM) = T(M)$, $0 \neq c \in R_+$.

2. M^T を M の転置行列とすれば

$$T(M^T) = T(M).$$

3. M' を、 M の2つの行、あるいは2つの列を入れ替えた行列とすれば

$$T(M') = -T(M).$$

従って、

$$T(M^\sigma) = \text{sign}(\sigma)T(M), \sigma \in S_n.$$

4. $T^{-1}(0) = \text{Dom} \cap \{M: \det(M)=0\}$.

命題 3.1 によれば、 T は（そして $\log(OR_G)$ も） $n^2 - 1$ 次元射影空間上の関数となる。さらに、行番号の置換と列番号の置換として、 $S_n \times S_n$ が Dom 上に作用しており、命題 3.3 により T は $\text{sign} \times \text{sign}$ の符号変化を受ける。その作用により、 T は $An \times An$ 不変である。

命題 4

M'' を、 M の一つの行（又は一つの列）を c 倍 ($c \neq 0$) した行列とする。 $T=OR, Q, \log OR$ について $T(M'') = T(M)$ が成立する。

本来の OR は勿論命題 4 を満たすので、その拡張である関数も同様であることが好ましい。

ファイ係数はこの性質を満たさず、拡張である Φ でも一般に $\Phi(M'') \neq \Phi(M)$ となる。

Agresti の OR_G は命題 4 を満たさず、命題 2.3 も満たさない。従って、 $\log(OR_G)$ は命題 3.3 を満たさない。これは、 OR_G また Goodman-Kruskal's γ が、元々順序因子のクロス集計を考慮して定義されていることから生じている。

4. $\log OR$ の性質。同質性

4.1 同質性と推定値

複数の行列が与えられた時、その同質性もしくは異質性を判断する必要が生じることがある。例えば、複数の治験例での OR を統合して、その全体のシステムでの OR の推定値を求めたい場合などである。以下行列は非負整数成分とする。

$$G = \{M_1, M_2, \dots, M_k; M_i \text{ in } \text{Dom}\}$$

とする。このとき、 G の同質性や異質性を判定し、同質ならば不変量を統合する手法として、Mantel-Haenszel 検定とその estimate、replicated G-test、又メタ分析での各種モデルが知られている。検定の場合、異質と判定する p の閾値は .05 に拘らず、分散分析の交互作用と同じく、分野によって定めるのが適切であろう。

Mantel-Haenszel estimate は、行列によって成分の値の大きさに偏りがあっても、多少緩和して加重平均する簡易な目安であり、高次元の OR に直接的に拡張可能である。ここで $\text{sum}(M)$ は行列

M の成分の和を意味する。

定義 6 Mantel-Haenszel type estimate.

$$MHe(G) = \frac{\sum_i \det p(M_i) / \text{sum}(M_i)^{n-1}}{\sum_i \det n(M_i) / \text{sum}(M_i)^{n-1}}$$

メタ分析での OR 統合の手法を適用するには、 $\log OR$ の分散（標準誤差 se の 2 乗相当であるが分散とよぶ習慣である）の推定値が必要となる。 M が 2×2 行列の場合は、各成分の逆数和である次式がよく知られている。

$$v(M) = 1/m_{11} + 1/m_{12} + 1/m_{21} + 1/m_{22}$$

高次元での v 相当を理論的に求めることも可能ではあるが、シミュレーションで容易に推定することができるので、実用上問題ない。そこで、 $(\log OR(M_i), v(M_i))_{i=1, \dots, n}$ をメタ分析に引き渡せば、モデルに従った推定値を得る。その結果を通行の Forest 図によって表現する (R の metfor パッケージ等) ことも可能である。

4.2 実例

簡単な例を挙げる。ここここでは後述する近傍作成によって得た 3 つの行列を用いる。

$$Mg = \{M_1, M_2, M_3\},$$

$$M_1 = \begin{bmatrix} 2 & 5 & 8 \\ 7 & 8 & 5 \\ 11 & 7 & 12 \end{bmatrix}, \quad M_2 = \begin{bmatrix} 4 & 3 & 8 \\ 7 & 7 & 6 \\ 9 & 10 & 11 \end{bmatrix}, \quad M_3 = \begin{bmatrix} 3 & 4 & 8 \\ 7 & 6 & 7 \\ 10 & 10 & 10 \end{bmatrix}.$$

Mg に対する、replicated G-test（以下 replG）と、Cochran-Mantel-Haenszel 検定（以下 CMH）の結果は下記の通りである。実際はより多くの情報を出力するが、ここでは略記した。末尾の vib は、後述の近傍シミュレーションによる分散 (se) 推定値で

表 1 replG, CMH の結果

	G	df	p	VG	OR	logOR	vib
1	5.306	4	.257	.194	.719	-.329	0.980
2	2.220	4	.695	.130	1.056	1.055	1.078
3	2.020	4	.732	.125	1.052	1.050	1.080
total	9.546	12	.656		.927	-.076	
pooled	6.661	4	.155	.131	.947	-.055	1.019
heterog	2.885	8	.941				

Cochran-Mantel-Haenszel test
 $M^2 = 6.528, df = 4, p\text{-value} = 0.163$

ある。また、VG は、G 値から求めたクラメール V である（通常の V とほぼ等しい）。pooled は、the pooled matrix であり、3つの行列の単純和である。

CMHの結果は、replGのpooled matrixに対する結果と近いものである。replGのheterog（異質性検定 heterogeneity G-value であり、 $p <$ 閾値の場合に異質と判定する）からして、同質性が保証され、MHe=0.927が有効である。その値をtotalのORの箇所に記入している。

一方、メタ分析（Random-effects model. Restricted maximum-likelihood estimator）の結果は図2の通りである。なお、数値はORとその信頼区間として与えているが、棒グラフのスケールはlogORに従って配置している。OR=1の位置に点線を付している。その前後の目盛りは $1/\sqrt{2}$, $\sqrt{2}$ である。対数尺なので、一定比が等間隔となる。

Meta分析のOR統合値は0.945（信頼区間はその右[0.680,1.312]）であり、replGにおけるOR(pooled)=0.947やMHe(G)=0.927と近い値であるが、いつもそうなるわけではなく、かなり異なることもある。ここでは信頼区間が1を跨ぐ。

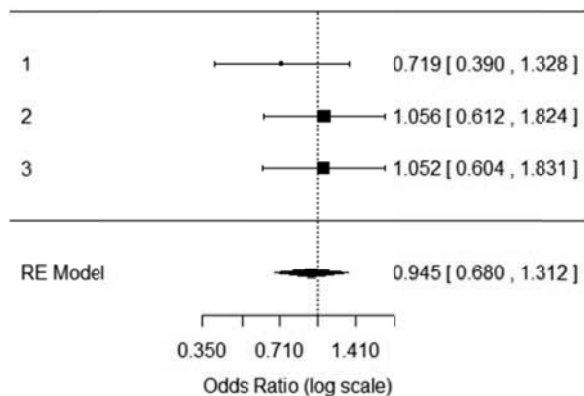


図2 meta分析の結果のForest図

いずれにせよ、これらの数値が整合していることがわかる。

5. 計算の手続

5.1 全体の手続

この手法を、大量の行列を含むGに適用することは無く、せいぜい20個以内であろう。また、サイズnも実用上7まで、あるいは5までと思えばよい。Likert scale 7段階の2つの変数のクロ

スを取ったとして7×7である。しかし、両端{1, 2}{6, 7}を統合すれば5×5でよい。

従って、detやpermの計算におけるオーバーヘッドはそれほど深刻ではない。必要なdetp, detnは直接再帰的に計算することにすればよい。もしも大規模ならば、detとpermの既存の高速近似計算からdetp, detnを求めることになると思われる。

Rを用いた素朴なdetp, detnの計算は付録に掲げる。通常の計算では有効数字は15桁程度であり加算でも桁落ちが心配されるが、Rでは多倍長が簡単に行えるので、detp, detnを例えば30桁の整数として中途計算を行い、除算してORを求める段階で通常の数値にすればよい。その実例も付録に含めた。

メタ分析のForest図の出力には、簡単なwrapperを設計すればよい。このため、replGのスク립トにおいて、行列の近傍を生成し、logORの分散を求めるようにして、その結果をmetafor::rmaに引き渡す。

5.2 近傍データの作成とその性質

以下、Mの成分は非負整数とする。

行列Mの近傍(neighborhood)の生成にあたっては、まずMをクロス表と見て、Utable(M)によって、クロスを取った元の2つの質的変数を再現する。それらの長さは $N = \text{sum}(M)$ である。この質的変数の組からN個のサンプリングをnrep回行い、各々をtableすればMに近い行列nrep個を得る。そのnrep個の行列に対してlogORを求め、標本分布をみて、分散 $v(M)$ を求めることができる。状況確認にはnrep=100で行い、本番ではnrep=1000, 2000程度に取ればよい。

Mの成分に0が含まれる場合、この方法ではサンプルのその成分は常に0であることに注意されたい。必要があれば、強制的に1を入れる。

この近傍データは十分によく分布するであろうか。またこのときlogORはどのように分布しているであろう。これらを確認するために、専用のスク립トを作成した。

そのスク립トでは、行列空間を n^2 次元のベクトル空間とみて、元の行列Mと近傍の行列の群 $G(M)$ を配置し、主成分分析を行う。そのPC1-PC2主成分スコア散布図に、95%許容楕円(信頼楕円)を付加し、かつPC1への射影の分布のdensityとrugを配置する。さらにlogORの分

布の density を平均を 0 にそろえて上書きする。標準的な凡例を左上に、また右上にデータフレーム名と nrep を記載する。

図 3 は ee の近傍での 1000 個の分布である。

$$ee = \begin{bmatrix} 20 & 3 & 3 \\ 3 & 20 & 3 \\ 3 & 3 & 20 \end{bmatrix}$$

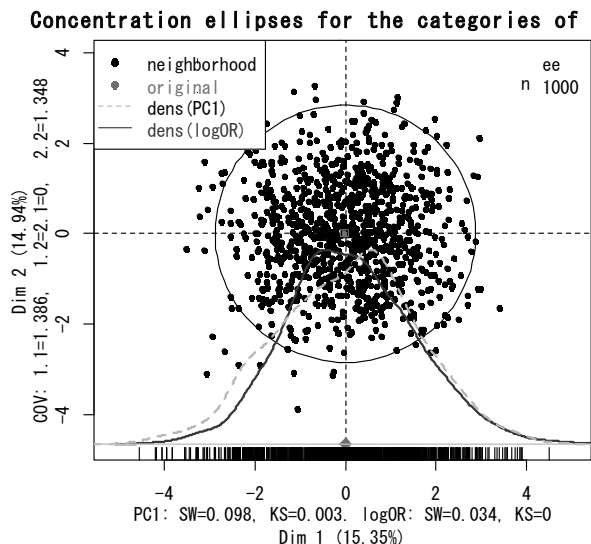


図 3 ee の近傍行列と logOR の分布

元の ee は実際中央（原点）に来ている。まさに近傍に相応しく分布しているのがわかる。

図 3 の灰色の破線が PC1 方向への射影の density であり、実線が logOR の分布の density である。この場合はよく重なっている。logOR の分布は尖度が大きくなることもある。

この ee の近傍 10 個に対してメタ分析を適用した結果の一例は図 4 のとおりである。この例の場合、MHe=15.935 であった。また、replG の異質

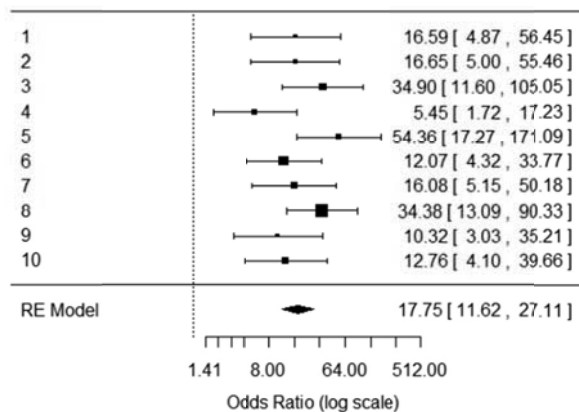


図 4 ee の近傍 10 個のメタ分析

性検定は $p=0.646$ となり同質と判断する。

CMH : $M^2=686.66$, $df=4$, $p\text{-value} < 2.2e-16$ となるがこれは $G(\text{pooled})=658.648$, $df=4$ と同程度の内容であり、役に立つ情報とはいえない。

6. 矩形行列への拡張

もとより、本稿の結果を $r \times c$ 行列 ($r < c$) に拡張することは可能である。すべての $r \times r$ 小行列 $N=cCr$ 個、または代表として連続する r 列からなる $c-r+1$ 個の $r \times r$ 小行列に対する OR の集合を対象とすればよい。この場合、余次元 1 の代数的集合 L を除いた $r \times c$ 行列全体の集合からの写像となる。

$$OR : M(r,c; K) \setminus L \rightarrow K^N \text{ 又は } K^{c-r+1}$$

特に $r=2$ の場合は、通行の方式と一致する。

さらにこの OR の集合を一つの代表値にするには MHe あるいは、メタ分析による推定値を用いればよい。

一例として、Agresti (2013) p.395, Table10.5 の 4×6 行列を取り上げる。その列 ABCDEF の連続する 4 列づつを取った 3 つの 4×4 行列の OR とメタ分析による統合は図 5 の通りである。

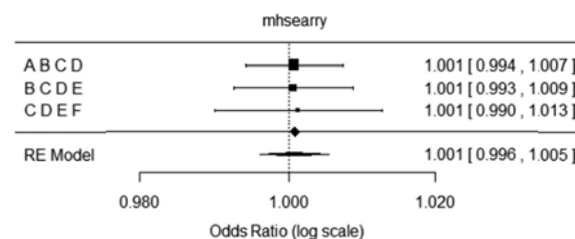


図 5 Cross-Classification of Mental Health Status and Socioeconomic Status.[1] Table10.5.

7. 結言

高次元の OR を、十分な不変性を持つように定義した。CMH が役に立つ情報を提供しないのに比較して、効果量として十分な機能を果たすと認められる。今後より精密な性質の検討を行う。

参考文献

- [1] Agresti A.(2013). *Categorical Data Analysis*. 3rd edition. Wiley, New York.
- [2] Agresti A.(1984). *Analysis of Ordinal Categorical Data*. Wiley, New York.
- [3] Agresti A.(1980). Generalized odds ratios for ordinal data. *Biometrics* 36:59-67.
- [4] Edwardes MD & E.Baltzan (2000). The generalization of the odds ratio, risk ratio and risk difference to $r \times k$ tables. *Stat Med.* 19(14):1901-14.
- [5] Subramanyam K.(1988). Analysis of Odds Ratios in $2 \times n$ Ordinal Contingency Tables. *J. Mult. Anal.* 27:478-493.
- [6] De Rooij M. and Anderson, C.J.(2007). Visualizing, summarizing and comparing odds ratio structures, *Methodology* 3:139-48.
- [7] D'Ambra, L., Camminatiello I. and Sarnacchiaro P. (2013) Generalized log odds ratio analysis for the association in two-way contingency table, *Electric Book "Advances in Latent Variables"*, Eds. Brentari E., Carpita M., Vitae Pensiero, Milan, Italy, ISBN978 88 343 25568.
- [8] Sarnacchiaro P., Gallo M. and D' Ambra L.(2014). Three-way decomposition of weighted log-odds ratio for customer satisfaction analysis, *Innovation and Society 2013 Conference, IES 2013. Procedia Economics and Finance* 17:30-38.
- [9] Slovkoic A.B. and Fienberg S.E.(2010). Algebraic geometry of 2×2 contingency tables, "Algebraic and Geometric Methods in Statistics", Eds. Gibilisco P., Riccomago E., Rogantin M.P., Wynn H.P., Cambridge University Press, Cambridge. pp.63-82.

附録 A.

A.1. detp, detn の計算

パッケージ Rmpfr によって多倍長計算を行う。関数 mpfr (k, 100) は、k を 100bit, $2^{100} \approx 10^{30.1}$ で取り扱うという意味であり、整数で 30 桁が確保される。出力 pn により、

```
pn[1]=detp, pn[2]=detn を受け取った後で
as.numeric(pn[1]/pn[2])
```

として、通常の数に変換する。

```
library(Rmpfr) # 事前に外部で実行しておく
rcsvm.det <- function(x) # x: matrix
{
  n <- nrow(x)
  pn <- c(mpfr(0,100),mpfr(0,100))
  if (n==2) {
    pn[1] <- mpfr(x[1,1],100)*mpfr(x[2,2],100)
    pn[2] <- mpfr(x[1,2],100)*mpfr(x[2,1],100)
    return(pn)
  } else {
    indx <- rep(c(1,2),n)
    for (cnt in 1:n){
      pn[1] <- pn[1] + mpfr(x[1,cnt],100)*
        Recall(x[-1, -cnt])[indx[cnt]]
      pn[2] <- pn[2] + mpfr(x[1,cnt],100)*
        Recall(x[-1, -cnt])[indx[cnt+1]]
    }
    return(pn)
  }
}
```

ここで Recall は R における再帰的呼び出しであり、これによって、全体の関数名を変更しても、コーディングを変更する必要が無い。上記の関数は行列式の小行列展開そのものを素朴にコーディングしたものであり、n が増加すると時間がかかる。通常の数で実行する場合、上記のコードから mpfr(,100) の文字を削除すればよい。

A.2. Untable について

Untable は、クロス集計表から元のデータを復元する。例えば下記の行列 mat に対して DescTools::Untable(mat) とすれば右のとおり 2 因子変数の data.frame が生成される。

```
> mat
```

```
  A B
```

```
a 1 3
```

```
b 2 4
```

```
> table(Untable(mat))
```

```
  Var2
Var1 A B
```

```
 a 1 3
```

```
 b 2 4
```

> Untable(mat)		
	Var1	Var2
1	a	A
2	b	A
3	b	A
4	a	B
5	a	B
6	a	B
7	b	B
8	b	B
9	b	B
10	b	B

A.3. Agresti の Generalized Odds Ratio. OR_G

行列 M (r 行 c 列) に対し

$$OR_G = P/Q$$

$$P = \sum_{i < j} \sum_{s < t} m_{is} m_{jt}, \text{ \# concordant pairs}$$

$$Q = \sum_{i < j} \sum_{s > t} m_{is} m_{jt}, \text{ \# discordant pairs}$$

と定義される。P, Q はいずれも $rC_2 \cdot cC_2$ 個の和からなる。

$x \leftarrow \text{DescTools::ConDisPairs}(M)$ と置けば $OR_G = x\$C/x\D である。

A.4. Forest 図実例

実例として、Agresti[1] の Table 8.9, 8.21, 11.19, 12.3 の行列群、また 13.10 を 3 通りに見た場合、さらに Agresti[3] の p.61 の例を図 A に提示する。

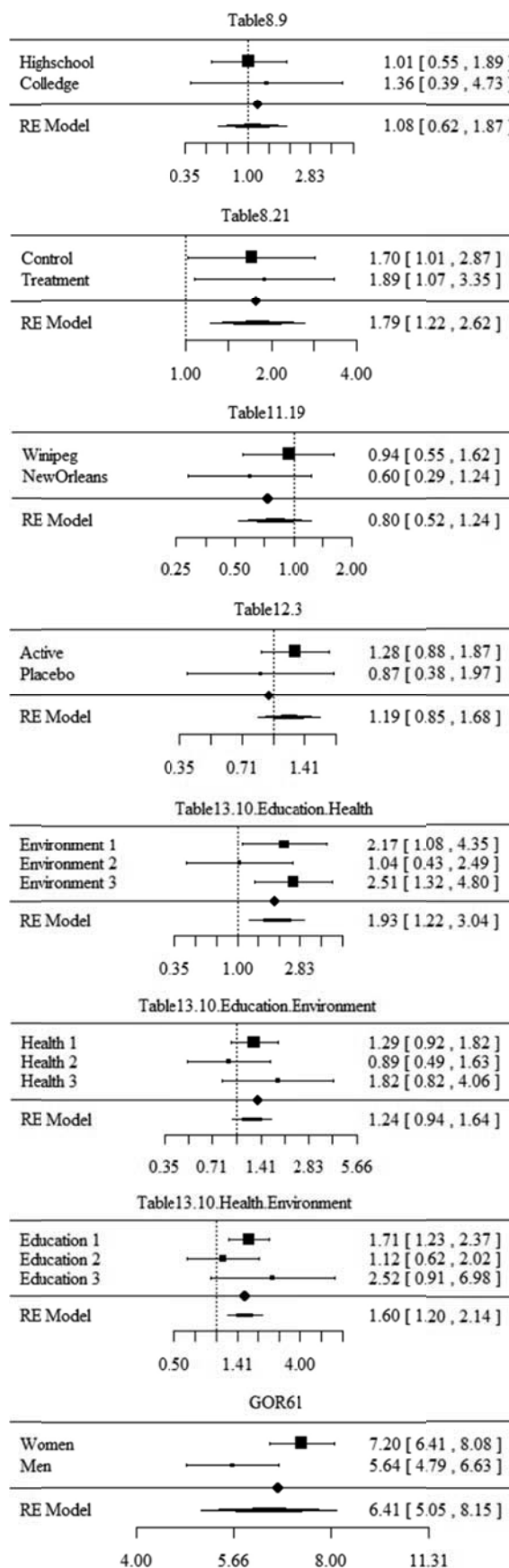


図 A Agresti の書籍・論文にある実例