# CA/MCAの数理 (その 2)

津田塾大学 数学・計算機科学研究所 藤本一男 kazuo.fujimoto2007@gmail.com

### 履歴

- ・2025/09/04 計量分析セミナー2025夏用に改訂
- ・2023/09/06 計量分析セミナー2023夏用に作成

# CA/MCAのリザルトを解釈するために

- ・生成される二つの空間
- ・分解される分散と座標軸
- ・SVDの結果から得られる行ポイント、列ポイントの座標軸
  - ・主座標と標準座標
- ・各ポイントの座標軸への寄与率(絶対的寄与率、 contribution)
- ・各ポイントと座標軸の相関。cos<sup>2</sup>

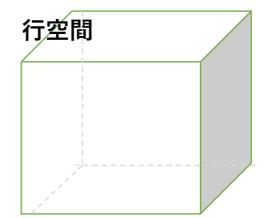
# 分析手順

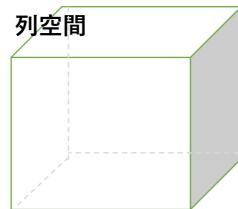
- ・分析対象のデータを前処理
  - ・度数分布表など、通常のデータの前処理を行う。
  - ・NAを"NA"など文字列に変換
  - CAやMCAのfunctionが求めている属性に変換
  - table、data.frame
  - factor
- CA/MCAを行う
- 慣性率を確認する。固有値の棒グラフ
- ・行マップ、列マップを描き、座標軸を解釈する
  - ・ここが最初の関門

### CA は 行空間と列空間を生成する

|        | col1 | ••• | coln | rowS<br>um |
|--------|------|-----|------|------------|
| rowl   |      |     |      |            |
| :      |      |     |      |            |
| rown   |      |     |      |            |
| colSum |      |     |      |            |







m x n 行列

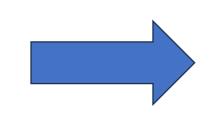
空間生成とは、座標軸の生成 を意味する dim ] ···. dimn.

> これらの座標軸は、データ 表の慣性(分散)が空間生成 に用いられた変数が分解され たもの。

> これらの慣性(分散)は、 行空間と列空間で同じ 値となる。(ここに、行分析 と列分析の相互浸透があらわ れている。cf カイ2乗距離を もちいているため。主成分 分析PCAでは、それはできな い。)

### CA**の三つの**result

|     | 強盗 | 詐欺 | 破壊 | 行和 |
|-----|----|----|----|----|
| オスロ |    |    |    |    |
| 中部  |    |    |    |    |
| 北部  |    |    |    |    |
| 列和  |    |    |    |    |



| 固有値 $\lambda$ | Dim1 | Dim 2 |
|---------------|------|-------|
| λ             |      |       |

| 行座標 F | Dim1 | Dim2 |
|-------|------|------|
| オスロ   |      |      |
| 中部    |      |      |
| 北部    |      |      |

| 列座標 G | Dim1 | Dim2 |
|-------|------|------|
| 強盗    |      |      |
| 詐欺    |      |      |
| 破壊    |      |      |

#### 集計データ

| 回答者 | 変数1 | 変数 2 | ••• | 変数n |
|-----|-----|------|-----|-----|
| 1   |     |      |     |     |
| 2   |     |      |     |     |
| 3   |     |      |     |     |
| :   |     |      |     |     |
| n   |     |      |     |     |

多重対応分析 MCAによる 空間生成 Active変数

選択肢回答変数

Active変 数 Pas sive 変数

#### 生成された座標軸

| 回答者   | 分散率 | 累積分散率 | 修正分散率    | 累積修正分散率 |
|-------|-----|-------|----------|---------|
| Dim.1 |     |       |          |         |
| Dim.2 |     | 割包    | _        |         |
| Dim.3 |     | 吉J c  | <b>=</b> |         |
| :     |     |       |          |         |
| Dim.n |     |       |          |         |

#### 個体空間座標值

| 回答者 | Dim.1 | Dim.2 ···.         |               | Dim.n |
|-----|-------|--------------------|---------------|-------|
| 1   |       |                    |               |       |
| 2   |       | <b>□□</b> +==      | . <del></del> |       |
| 3   |       | —— <mark>座標</mark> | 世             |       |
| :   |       |                    |               |       |
| n   |       |                    |               |       |

#### 変数空間座標値

| 変数  | Dim.1 | Dim.2             | •••      | Dim.n |
|-----|-------|-------------------|----------|-------|
| 変数1 |       |                   |          |       |
| 変数2 |       |                   |          |       |
| 変数3 |       | <mark>座標</mark> [ | す        |       |
| :   |       | <u> </u>          | <u> </u> |       |
| 変数n |       |                   |          |       |

#### 射影される変数座標値

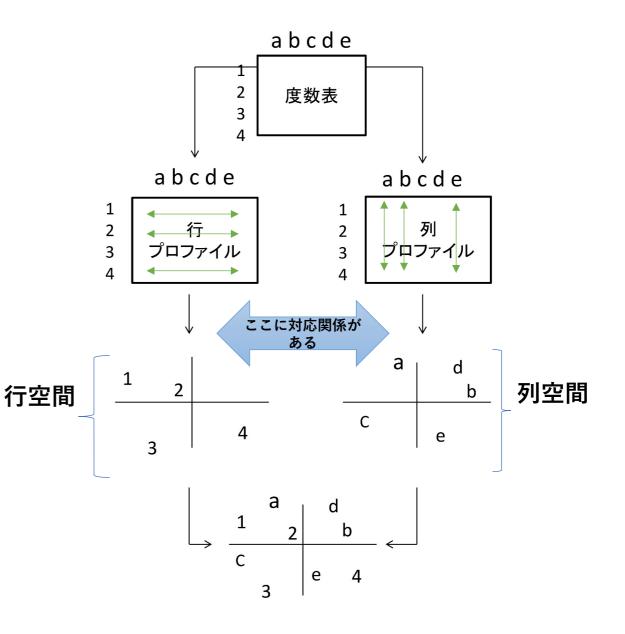
| Passive変数 | Dim.1 | Dim.2 ···. |   | Dim.n |
|-----------|-------|------------|---|-------|
| 1         |       |            |   |       |
| 2         |       | 座標值        | 直 |       |

# 主座標と標準座標

## 二つの空間と座標系

- 、二つの空間の重ね合わせ
- これは、それぞれのグラフを見てもらうのがよい。
- 数理的に追いたい人は、『対応分析の理論と実践』を 読んでください。

| マップ   | 行ポイント<br>の座標 | 列ポイント<br>の座標 | option            |
|-------|--------------|--------------|-------------------|
| 対称マップ | 主座標          | 主座標          | Symmetric         |
|       | 主座標          | 標準座標         | rowprinci<br>pal  |
|       | 標準座標         | 主座標          | colprintci<br>pal |
| 不可能!  | 標準座標         | 標準座標         | X                 |



『対応分析入門』p6、図1.2に加 筆

# 追加変数による空間解釈

# 行と列の関係

- ・各行ポイントは、列ポイント全てと結びついている。
- 各列ポイントは、行ポイント全てと結びついている。
- ・その関係は、transition formura(遷移公式、推移公式)として以下のようにかける。『対応分析の理論と実践』p246
  - F (行主座標) = D<sub>r</sub>-¹PΓ (Pは元表、Γは列標準座標)
  - ・G(列主座標) $= D_c^{-1}P^T\Phi$  ( $P^T$ は元表の転置。 $\Phi$ は行標準座標)
- ・標準座標は、平均O分散1にスケーリングされています(標準化)。

# 追加変数(サプリメンタリ・ポイント)

- ・元表の行和、列和は、各行、各列の質量(weight)と呼ばれる。 ・この質量がある行や列が、行空間、列空間を生成する。
- ・ところで、行ポイント、列ポイントは、遷移公式によって結びついている。
  - 例:
    - ・ MCAで変数空間で選択されたカテゴリの組み合わせによって、個体空間での個体の「位置」が 決まります。
- そこで、質量をもたないプロファイルを考えると、そのプロファイルは、 反対側の空間に座標をもつことができる。
- こうして、空間生成には寄与せずに、内部構造を分析するための変数を考えることができる。
  - 例えば:
    - 空間生成には、性別、年齢変数は用いずに、追加変数として生成された空間にplotする。
    - MCAでの構造化データ解析はこの仕組みを活用します。

### CAの入力表と二つのプロファイル

0.56 0.15 0.29

行分析 行プロファイル

|     | 強盗       | 詐欺       | 破壊       | 行和          |            |
|-----|----------|----------|----------|-------------|------------|
| オスロ | 395      | 245<br>6 | 175<br>8 | 4609        |            |
| 中部  | 147      | 153      | 916      | 1216        |            |
| 北部  | 694      | 327      | 134<br>7 | 2368        |            |
| 列和  | 123<br>6 | 293<br>6 | 402<br>1 | 8列分析<br>列プロ | f<br>コファイル |

このクロス表の変数、ノルウェイの都市名と 犯罪名は、Clausen1987=2015で使われて い

|       | 強盗   | 詐欺   | 破壊   | 行和 |
|-------|------|------|------|----|
| オスロ   | 0.09 | 0.53 | 0.38 | 1  |
| 中部    | 0.12 | 0.13 | 0.75 | 1  |
| 北部    | 0.29 | 0.14 | 0.57 | 1  |
| 平均行比率 | 0.15 | 0.36 | 0.49 | 1  |

|       | オスロ  | 中部   | 北部   | 列和 |
|-------|------|------|------|----|
| 強盗    | 0.32 | 0.12 | 0.56 | 1  |
| 詐欺    | 0.84 | 0.05 | 0.11 | 1  |
| 破壊    | 0.44 | 0.23 | 0.33 | 1  |
| 平均列比率 | 0.56 | 0.15 | 0.29 | 1  |

\*元表を転置して配置していある

る事例のもの。

### CAの入力表と二つのプロファイル

|       | 強盗   | 詐欺   | 破壊   | 行和 |
|-------|------|------|------|----|
| オスロ   | 0.09 | 0.53 | 0.38 | 1  |
| 中部    | 0.12 | 0.13 | 0.75 | 1  |
| 北部    | 0.29 | 0.14 | 0.57 | 1  |
| 平均行プロ | 0.15 | 0.36 | 0.49 | 1  |
| ファイル  |      |      |      |    |

行比率として要素計算したものが、 行プロファイル。 列和の行は、平均行プロファイルとなる。

|       | オスロ  | 中部   | 北部   | 列和 |
|-------|------|------|------|----|
| 強盗    | 0.32 | 0.12 | 0.56 | 1  |
| 詐欺    | 0.84 | 0.05 | 0.11 | 1  |
| 破壊    | 0.44 | 0.23 | 0.33 | 1  |
| 平均列比率 | 0.56 | 0.15 | 0.29 | 1  |

列比率として要素計算したものが、 列プロファイル。 行和の列は、平均列プロファイルとなる。

\*元表を転置して配置していある

# 行と列関係:対応関係

- ・行の各ポイントは、列のすべてのポイントに関連している。
- 列の各ポイントは、行のすべてのポイントに関連している。

・行と列の「対称性」。これを実現しているのが、点間距離をカイ2乗距離で定義していること。

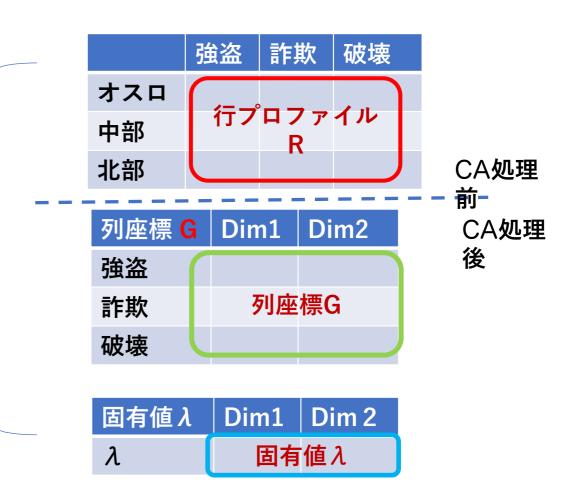
### 遷移公式: Transition Formula



| 行座標 F | Dim1 | Dim2 |
|-------|------|------|
| オスロ   |      |      |
| 中部    |      |      |
| 北部    |      |      |



 $f_{ik} = \frac{1}{\sqrt{\lambda_k}} \sum_{j} \left(\frac{p_{ij}}{r_i}\right) g_{ik}$ 



\*Greenacre1984:64 \*\*

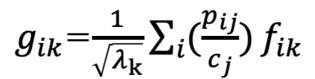
Greenacre2017=2020:108

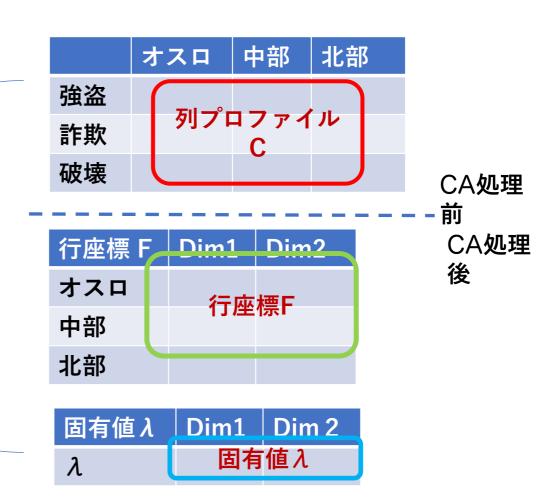
### 遷移公式: Transition Formula



| 列座標 G | Dim1 | Dim2 |
|-------|------|------|
| 強盗    |      |      |
| 詐欺    |      |      |
| 破壊    |      |      |







# 行列の演算形式:遷移公式

3行

3行

3列

行プロファイル R 3行3列 2列

列座標 G 3 行 2 列 対角行列 2 x 2

行座標 F 3 行 2 列

行も列も 形式は同じ

3列

列プロファイル C 3行3列 2列

行座標 F 3 行 2 列 対角行列 2 x 2

列座標 G 3 行 2 列

# 追加変数を計算できる仕組み

# 追加変数/カテゴリ/個体

- ・空間生成に関係するポイント
  - Active変数
- ・CAの時はあまり問題にならないが、MCAでは次のことが問題になる。
  - ・関連のない(まざると解釈不能になるカテゴリ)も、SVDを行えば数 学的には何らかの座標を獲得できる。
  - ・しかし、解釈不能….
- ・そこで、空間生成には、同質(まざっても解釈可能という意味)なカテゴリを用いる。Active変数
- ・しかし、その連関はなくても、関連を見たい変数はある。
  - 年齢、性別、年収…..

## 追加変数の射影というアイデア

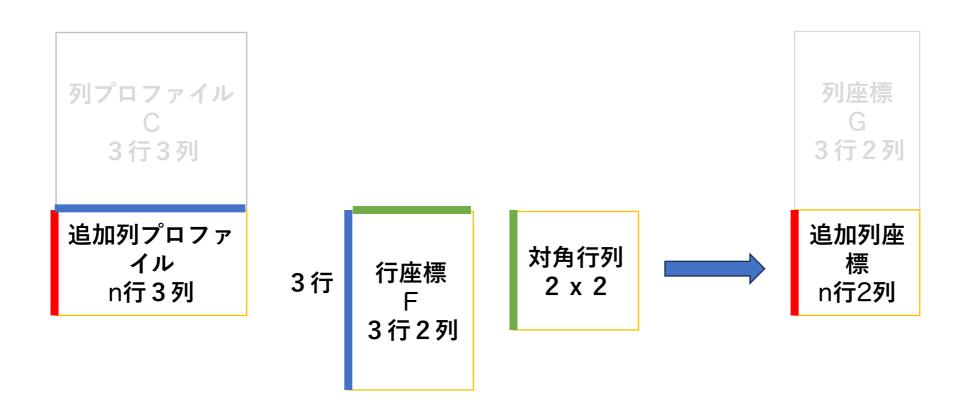
- ・空間成績を行う変数の選択
  - Active変数
- これで生成される二つの空間の構造を、例えば性別、年代、もしくはその合成変数で分析したいときには、それら(性別など)を追加変数(サプリメンタリ変数)として、射影する。
- ・射影するには、座標が取得される必要がある。
- ・この空間生成には寄与しないが、座標を取得できる方法を、 「追加変数」として得る。

### 追加変数の座標が計算できるしくみ

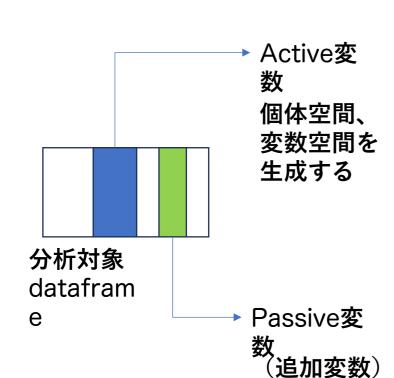


行列の掛け算では、右行列の列数と左行列の行数が同じであれば、計算可能。 つまり、列プロファイルに追加列プロファイルを加えると、列座標と固有値の 対角行列をかけて、追加行座標が得られることになる。 ここで、この追加列は、空間にはなんら影響を与えていない。

# 実際の計算では、追加変数分だけでよい



# Active変数と追加変数の関係



指示行列化(index 円和r<sup>(x)</sup>列和 C 標準化残差行列 junkカテゴリ指定 特異値分解 SVD

 $\begin{array}{c} \Phi = D_r^{-1/2} U \\ D_{\alpha} V^t & \Gamma = D_c^{-1/2} V \\ F = D_r^{-1/2} U D_{\alpha} = \Phi D_{\alpha} \\ G = D_c^{-1/2} V D_{\alpha} = \Gamma D_{\alpha} \end{array}$ 

指示行列化(index matrix)

matrix) 転置して行和を求め、行比率行列 (プロファイル)化(C<sub>sup</sub>)

Passive変数の座標 $G_{sup} = C_{sup} F D_{\lambda}^{-1/2}$ 

Active変数の列座標とPasssive変数(追加 変数)の列座標(ともに主座標)で散布図を描く。

# CA & MCA

#### CA & MCA

- ・CAは2変数、クロス集計表を入力とし、MCAは行が個体、列 が変数の表(集計表)を入力する。
- ・しかし、functionの内部では、変数は回答カテゴリに分解されて、0/1でコーディングされたindicator(指示)行列に展開されて、その表に対してCAが行われている。
  - ・もう一つのMCAはburt行列という2変数ごと総当たりクロス表に対するCA。今回は、指示行列版MCAしか使いません。『対応分析の理論と実践』第18章
- ・この指示行列には厳しい制約がある。
  - ・0/1で表示されるけれども、MA回答のような0/1展開ではない。

# 指示行列の重要性

- ・変数内は、複数のカテゴリに分割されている。
  - その中に1が立つものが必ず1つあること。
  - ・つまり、行和は、変数総数になる。
  - これがMA回答の0/1とは異なる部分。
- ・ではMA回答はどうコーディングするのか。
  - その変数内の回答を合計すると1になるように配分する。
  - ・選択肢が10個あって、3つ選ばれていたら、一つには1/3を配分。
  - 「いくつまで」「いくつでも」「いくつ」によって、コーディング方 法がかわるので、やっかい。
  - ・「参考資料」の「MAコーディングの問題」参照。

# この指示行列ルールがGDAでは重要

- ・平方和の分解で、
  - 全分散=群間分散+郡内分散 が成り立つ前提。
- ・MCAをやっていて実践的に直面する、ジャンクカテゴリの処理 に関係する。
  - speMCA カテゴリ選択MCA
  - CSA 個体選択MCA

ともに、全体の行和、列和は維持しており、全体のMCAとの比較を可能にする。(GDAのところで説明します。)