

WordMiner™における多次元データ解析 — 設計指針と主な特徴の紹介 —

WordMiner活用セミナー
2006年9月7日
於: 日本電子計算株式会社

テキスト・マイニング研究会
<http://wordminer.comquest.co.jp/>
統計数理研究所
大隅 昇
ohsumi@ss.ij4u.or.jp

All rights reserved. Copyright by Noboru Ohsumi, ISM Professor Emeritus.

はじめに: 本日のトークの情報源

- 本日配付のテキスト資料「テキスト型データのマイニング」
- テキスト・マイニング研究会ホームページから参照可能
<http://wordminer.comquest.co.jp/>
- WordMiner Tipsの中の「技術解説」から種々のペーパーが参照、ダウンロード可能
<http://wordminer.comquest.co.jp/wmtips/analysis.html>
- その他, 各種の文献・資料(邦文, 欧文など)
<http://wordminer.comquest.co.jp/biblio/index.html>
- 出版物発刊の予告
 - 「テキスト型データのマイニング: WordMinerによる事例解析」(ナカニシヤ出版)
 - 「調査方法論研究」(仮題), 翻訳書(朝倉書店)
Survey Methodology, by R.M. Groves and others, John Wiley.

2

本日のトークの内容

- 現状のテキスト・マイニング(TM)をどう考えるか?
- WordMinerの設計指針
- 多次元データ解析の概要
 - 多次元データ解析の分析手順
 - 対応分析法・数量化法III類
 - クラスタ化法(ハイブリッド法)⇒今回は概念のみ
 - 有意性テスト他 ⇒配付資料・テキストの「補足資料」参照
- 分析対象とするデータとデータ表の構造
 - データをどう考えるか, どのようなデータを扱うのか
 - データ表の特徴(どのようなデータ表を扱うか, その理由)
- とくに, 扱うデータとデータ表に集中して話す.
- そのため, 対応分析法とデータ表の関係をミニチュア・データで確認する. ⇒別資料とした.

3

お断り

- ある程度, 統計的データ解析の基礎知識を前提として話す(How-to的ではない).
 - 例: 統計値(平均値, 分散, 標準偏差, 相関係数, ...)
 - 例: 質的データとは何か? ↔ 量的データとは何か?
 - 例: クロス表とは何か?
 - 例: 質的データを計量化するとは何か, あるいは数量化とは何か?
- その他の基礎知識
 - クロス表の独立性検定, ピアソンのカイ二乗統計量
 - 統計的検定の基礎概念
 - 確率分布(例: 超幾何分布, 二項分布, 正規分布など)
- かなり圧縮化して話すので, 不明箇所があれば, トーク中でも躊躇せずに質問していただく.
- その他, 時間があれば, あるいはQ&Aとして対応.
- 別資料を用意したので, ここでは抜粋・要約して述べる.

4

テキスト・マイニングの何が問題か？

- 正しい理解が徹底しないこと。
- 依然としてテキスト・マイニングへの過剰期待があること。
- (WordMinerに限らず)ソフトを十分に使いこなせていない。
- 客観的に見直す必要があること(適用可能性の再評価)。
- ソフト提供者の側としては、ノウハウの可能な範囲の情報開示。
- 用いる方法論の内容の透明化(暗箱化の回避)。
- 利用者の理解とリテラシー向上が必要(誤用・濫用の回避)。
- テキスト・マイニング研究会(TM研)は、微速ではあるが少しでもWordMinerへの理解を徹底したいと考えての普及活動の一環として進めている。

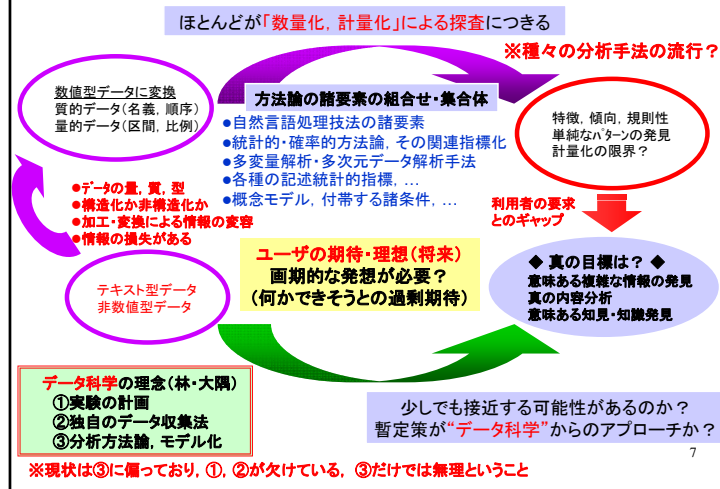
5

定性情報・質的情報への関心の高まり

- 定性情報の取得、あるいは定性調査におけるデータ収集方式(data collection mode)には様々な方法がある。
- 調査環境の変化、とくに電子的調査情報取得手法(CASIC, CADAC)の研究の進歩がある。
- テキスト型データの取得が容易となった(電子的取得)
 - インターネット調査, グループインタビュー(GI), フォーカス・グループ(FG)などの電子化(オンラインFG: OFGなど)
 - コール・センター, コンタクト・センターなどでのデータ収集
 - ITスキル, 技術改善が優先・重視されている
 - データ解析の本質が軽視, ソフトウェア依存となっている
- 様々な分野におけるテキスト型データを含む質的・定性的情報の分析への要求の高まりがある
 - 介護・福祉研究, 看護学研究, 人事管理, 企業組織研究, ...
 - 定量的な選択肢型の情報収集だけでは適切に対応しきれない

6

現状のテキスト・マイニングは？



7

何に注目すべきか？

- 現実には「レシピ・ライク or how-to-do」に対応できない場面が多々ある。
- 対象・現象別にテイラード方式(tailored design)をとること。
- 要は「うまくできた」だけでなく、「失敗例」に注目すべき。
⇒「なぜ、うまくできないのか？」を再検討する
- どこに分析の困難性があるかの情報の確認と開示。
- 適用可能性の範囲(出来ること, 出来ないこと)を明示。
- 「データ」とは何か, それをどう考えるか?
- 現象解明の分析対象とするデータの吟味・議論の重要性。
- 現象解明の客観性・科学性の担保は?
- 「知識」として活用するには? (「知識」から「知恵」への転換)

8

とくに「調査」における3つの見方(データ科学3原則)

- ① まず、質問をいかに「設計」するか (designの問題)
- ② つぎに、いかにデータを「集める」のか (data collection mode)
そして標本(サンプル)をいかに「抽出する」のか
- ③ 以上を吟味のうえで「分析」を行う (analyzing), モデル化の適否を考察する (model-like approach)

※分析に合った対象の選択と利用手法の相性をどう勘案、見極める
※WordMinerは③の一部を支援するツールにすぎないこと



- テキスト型データ (textual data) の解析では、これに加えて、...
- 質的データ・定性情報の 計量化・定量化の方法 をどう考えるか
- テキスト型データ以外の型の データの併用 (選択肢型質問, 属性情報のようなコード化データ, 質的変数) をどう利用可能とするか

9

WordMiner™の設計指針

- テキスト型データ だけを分析対象としたデータ解析ソフトウェアではないこと。
- 基本的には 社会調査型のデータの処理 を想定していること。
- 非構造的なテキスト型データ, 例えば 自由回答質問 (open-ended questions, free answer) や 自由記述型のデータの処理 に適している。
- 選択肢型質問 (closed questions), 属性情報 (デモグラフィック要因) と自由回答・自由回記述との併用分析が可能 (質的変数 として利用)。
- 談話形式 (discourse, narrative form), エスノグラフィー的 (ethnographic) な環境で取得したデータセットも可。
- 負荷を軽くし、価格を抑えるためにデータベース機能はなし。
- 日本語処理機能を極力軽装化したこと (分かち書き処理のみ, 細かい形態素解析までは行わない)。

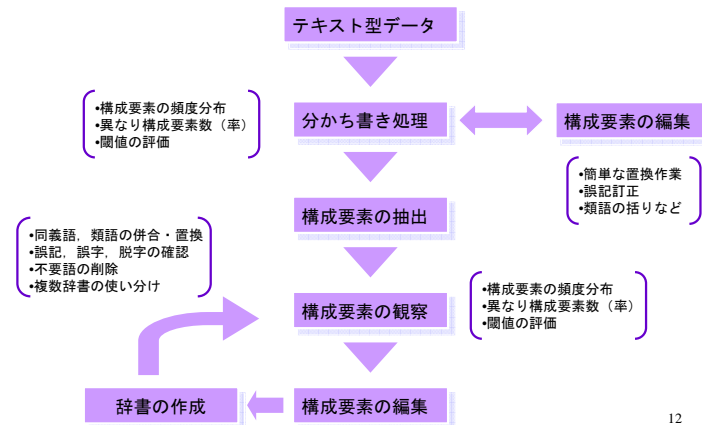
10

多次元データ解析の概要

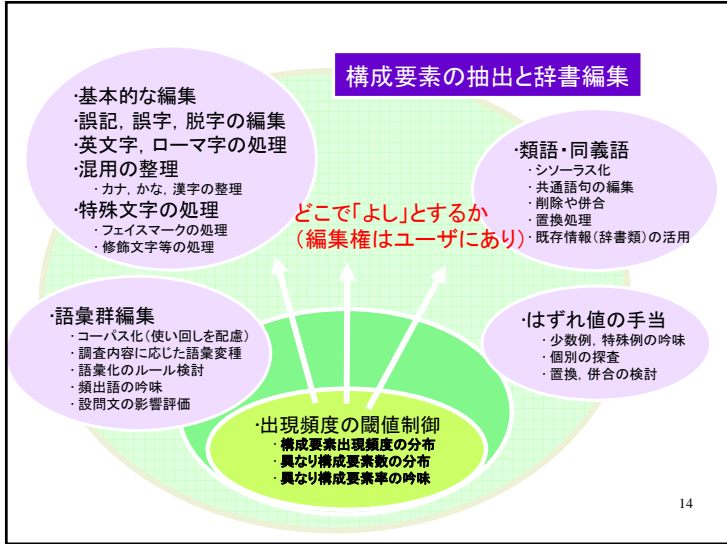
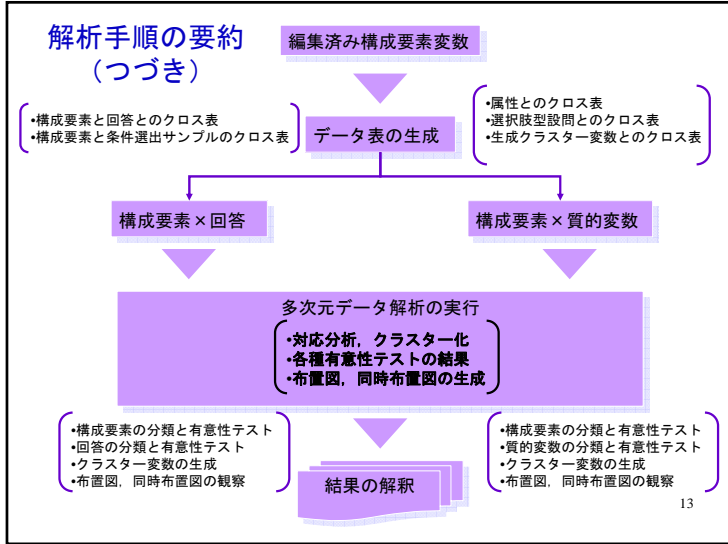
- 多次元データ解析の分析手順はきわめて簡略化されている。
 - 基本的には 以下の2手法だけ からなり立っている。
 - 対応分析法・数量化法III類
 - クラスター化法 (階層的分類と非階層的分類のハイブリッド法)
- 個々の 手法・技法の使い方に固有の特徴 があるが、ここでは省略する (テキストに一部記述されている)。
 - 各種有意性テスト (構成要素頻度, カイ二乗距離によるテスト, 他)
- 種々の データ表 に対応する。
- とくに、テキスト型データと 質的変数 (例: 選択肢型質問, 属性情報など) との併用分析。
- 適用手法をなるべく明示化し 「何を行っているか」 を開示していること。
- 複雑だ、難しいと思いきまぬこと。同時に、無理解・節操のない使い方も困る (どんなデータに適用できるのか)。

11

多次元データ解析手順の要約



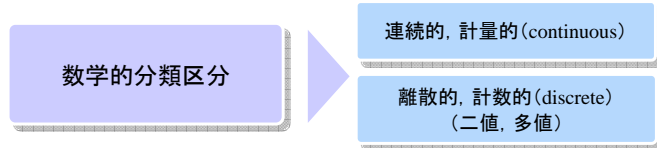
12



- ### 構成要素の抽出と辞書編集の考え方
- 出現頻度の頻度分布と閾値制御
 - ・ 構成要素数の制御と出現頻度分布の観察
 - ・ 異なり構成要素数の分布の観察
 - ・ 異なり構成要素率の吟味(サンプル数, 総構成要素数などの比較)
 - 語彙群(≒コーパス, コーポラ)の作成はユーザが必要に応じて対応する.
 - 類語・同義語(シソーラス)の整理についても同様に対応する.
 - 基本は, 以下を使い分ける(⇒作成辞書の流用, 相互利用).
 - ・ 分かち書き処理情報から最小限の編集(ゴミ除去程度)を行う
 - ・ 分析対象に応じた固有辞書の編集・作成
 - ・ 作成した辞書を他の課題(プロジェクト)で使う, 再利用
 - 構成要素の吟味・編集は考えるほど際限ない作業となる.
 - 複雑な日本語の特徴を知る, どこまで関与(日本語とは何か).

- ### データをどう考えるか, どんなデータを扱うのか
- 数学的分類
 - ・ 連続的変数か離散変数か
 - ・ 統計学のテキストなどにある分類
 - 「尺度(scale)」による分類
 - ・ 質的データ(名義尺度, 順序尺度)
 - ・ 量的データ(区間尺度, 比例尺度)
 - 両者を勘案して使い分けること
 - 定性調査においては尺度分類の方が説明しやすい
 - とくにテキスト型データは質的データであると考えること
- 16

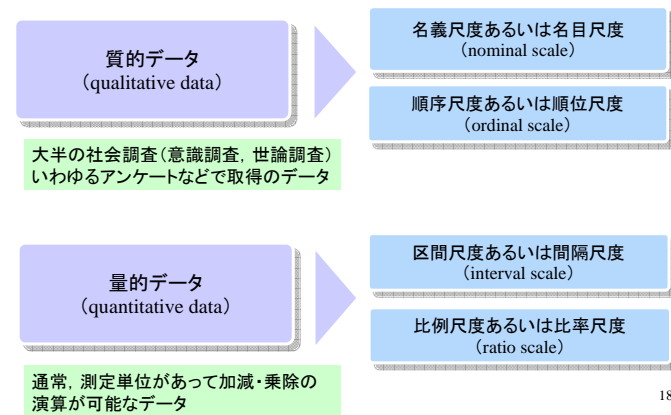
数学的分類



- いわゆる, 一般の統計学, 数理統計学に基づく議論は概ねこれで済む。
- 社会調査**(意識調査, 世論調査, アンケート)のデータや**テキスト型データ**を扱うにはこれだけでは扱いにくい。
- 一般に**質的研究**と言われる分野のデータも同様である。

17

尺度による分類



18

尺度による分類と数学的分類の要約[表1]

尺度による分類

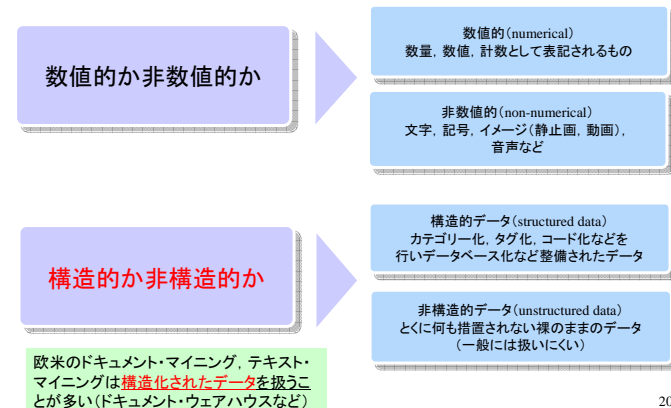
	質的データ		量的データ	
	名義尺度	順序尺度	区間尺度	比例尺度
連続量	(この組み合わせは考えられない)	音の強さの段階的区分 色度, 光沢度	温度 (°C) 硬度 比重	単位を持つ測定値 データの大部分 (長さ, 重さなど)
	機械名 作業者名 工場名 原産地名, など	段階的評価の成績データ 調査票の選択式質問における選択肢 ('満足」「やや満足」「満足でない」) など	TVのチャンネル 体育館の利用日数 車の故障台数, など	車の走行台数 都市内人口 参加者数 家の戸数
離散量	性別 (男、女) 「あり, なし」 (有, 無) スイッチの状態 (「入, 切」) など	物体の大きさ (大きい, 小さい) 濃度 (濃い, 薄い) 硬さ (硬い, 柔らかい) など	旅行経験の有無 (回数を考慮に入れば多値データとなる)	瓶入りと缶入りのジュース単価 (二値の分類区分で層化)

数学的分類

こうした分類区分を目安とすればよい。調査設計, データ取得時から意識すること。

19

テキスト型データの解析では以下の区分も重要



20

数量化法Ⅲ類と対応分析法(⇒テキスト, 3~5ページ)

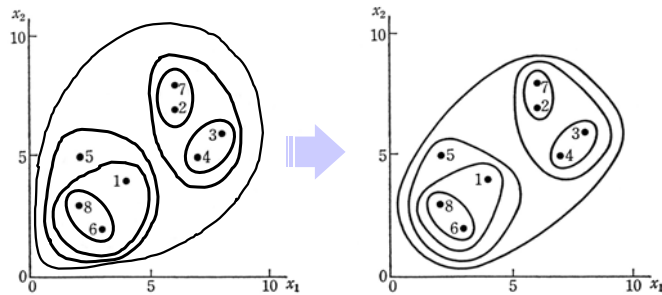
- 手法は同等の手法であり, 質的データの数量化の方法である。
- 数量化法Ⅲ類(quantification method, type III)
 - 鮎戸弘氏の命名による俗称
 - 林知己夫氏により提唱された手法(1955~1956年頃)
 - 正式には「パターン分類(法)」という(この年代に稀有な発想)
 - 数量化法(quantification methods)の一つ
- 対応分析法(Analyse Factorielle des Correspondances)
 - ベンゼクリ氏(J.-P. Benzécri)により提唱(1962年頃)
 - 正式にはAFC(Analyse Factorielle des Correspondances)
 - CA: Correspondence Analysisとして英語圏で紹介された
 - 対応分析(法)と命名(大隅・林他⇒国内で初めて紹介, M.Rouxを招聘)
 - コレスポネンス分析, コレスポネンス・アナリシスなど(多分, その本質的な意味が分からなかったのでこんな名称が出た)
- その他, 同等あるいは類似の手法が多数ある(テキスト参照)²¹

クラスター化法(自動分類法)

- 分類手法を大別すると, 階層的分類法と非階層的分類法とがある。
- WordMinerでは, 階層的分類法と非階層的分類法をハイブリッドして利用する, いわゆるハイブリッド法となっている。
- 階層的分類法としてウォード法(Ward's method)を用いる。
- 非階層的分類法うちの分割最適化型の一つであるk-平均法(k-means method)を用いる。
- クラスターの初期化に階層的分類法を用い(ボトムアップ), データの再配置・更新に(トップダウン)非階層的分類法を用いて調整を行う。
- 特徴として, 大量データの分類が可能なこと, はずれ値の影響が緩和されることなどがある。
- クラスター化による類型化で, 新たな質的変数が生成される(クラスター変数), これを使った二次分析が可能となる。

22

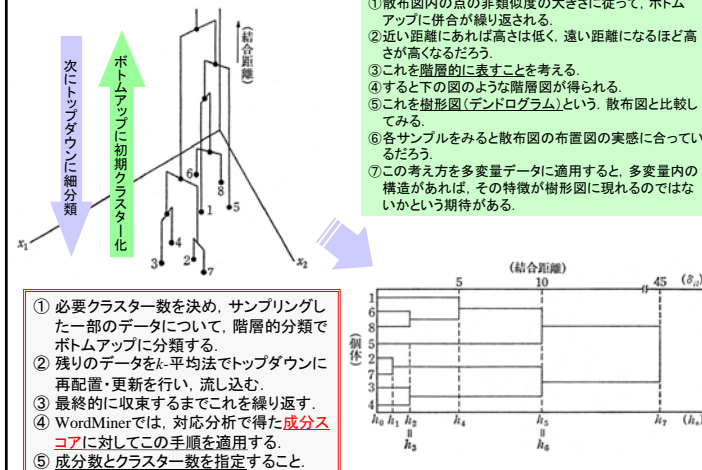
階層的分類法の考え方(クラスター生成の例)



2変量の人工データ, その散布図で観察
この関係を操作手順に従って, 模式図としてみると, ...

23

ハイブリッド法の考え方(模式図)



質的データの数量化とは？

- そもそも定性情報は**量的処理操作が難しい**、あるいはできないことがある、よって**数量化(quantification)**が必要である。
- テキスト型データの場合は、各種の**事前処理**も必要となる。
 - 例: 分かち書き処理
 - 例: 語彙・類語(コーポラ, シソーラス)などの情報整理
- 大まかには“次元データ解析に適した形式の**データを生成するまで**”の操作、次に不定形・非構造的な“**質的情報を数量化する**”までの操作、さらにその後の“**詳細な解析を行うための解析手順**”とがある。
- さらに**遡った重要な手当て**として、例えば調査であれば「**調査計画, データ取得方式, 調査方式(モード), 調査票**」の**効果的な設計**が必須。
- 調査以外の他の分野(研究, 実務)でも、ほぼ同様の対応が必要なのは。

25

WordMinerの考え方

- テキスト型データや選択肢型質問から得た**質的データをまず計量化・数量化する**。
- 質的データの**数量化の操作手順の一環として対応分析法(数量化法Ⅲ類)**を用いる。
- さらに、**量的データ化した数量化スコア(成分スコア)を用いて、クラスター化**他の処理を行う。
- 生成した**クラスター変数**を新たな質的変数として利用した分析を行うことの可能。
- この考え方は多くの他のテキスト・マイニング・ツールでも同様である。
- それを表だって主張しないで、あたかも“**新しい方法論がある**”かのような説明やキャッチコピーが多いのは問題。

26

数量化とは 一簡単な例示による確認 [例3]

- ある調査における取得データを取り上げる。

質問A: あなたは、いま住んでいるまちが気に入っていますか。(一つ選ぶ)

1. たいへん気に入っている
2. まあ気に入っている
3. あまり気に入っていない
4. 気に入っていない

質問B: あなたが住んでいる地区は、都市としては、**緑(みどり)**が多いと感じますか、**それとも少ないと感じますか**。(一つ選ぶ)

1. かなり多い
2. 多いほうである
3. ふつう
4. 少ない
5. 少ないほうである

※尺度の分類によればいずれも「**順序尺度, 名目尺度**」である

27

ある調査データの一部[(回答者) × (項目)](表7)

		WordMinerはどちらのタイプも扱える									
		選択肢(テキスト)のまま					選択肢をコード化				
サンプル番号	地区番号	回答コード	いつ頃から現在地に住んでいますか。	近(の)所(の)や公園(の)どのくらいありますか。	あなたは、いま住んでいるまちが気に入っていますか。	住んでいる地区は、都市としては、緑(みどり)が多いと感じますか。(選択肢)	(1)住んでいる地区は、都市としては、緑(みどり)が少ないと感じますか。(選択肢)	あなたは、住んでいる地区は、都市としては、緑(みどり)が多いと感じますか。(コード)	住んでいる地区は、都市としては、緑(みどり)が多いと感じますか。(コード)	緑(みどり)が多いと感じますか。(コード)	その補填や説明は、多いで行けず、分らないと答えていますか。
30	35	1	56	1	1	1	1	1	1	1	1
29	35	1	57	2	2	2	2	2	2	2	2
27	35	1	45	3	3	3	3	3	3	3	3
20	35	1	56	3	3	3	3	3	3	3	3
25	35	1	53	1	1	1	1	1	1	1	1
23	35	1	42	2	2	2	2	2	2	2	2
22	35	1	54	1	1	1	1	1	1	1	1
19	35	1	45	2	2	2	2	2	2	2	2
17	35	1	47	4	4	4	4	4	4	4	4
15	35	1	54	2	2	2	2	2	2	2	2
14	35	1	56	2	2	2	2	2	2	2	2
13	35	1	42	3	3	3	3	3	3	3	3
12	35	1	56	4	4	4	4	4	4	4	4
8	35	1	54	2	2	2	2	2	2	2	2
7	35	1	54	2	2	2	2	2	2	2	2
6	35	1	45	2	2	2	2	2	2	2	2
2	35	1	57	3	3	3	3	3	3	3	3
1	35	1	44	5	5	5	5	5	5	5	5
4	35	1	42	3	3	3	3	3	3	3	3
11	35	1	46	2	2	2	2	2	2	2	2
30	30	1	54	4	4	4	4	4	4	4	4
28	30	1	49	3	3	3	3	3	3	3	3
26	30	1	54	4	4	4	4	4	4	4	4
27	30	1	54	1	1	1	1	1	1	1	1
20	30	1	57	4	4	4	4	4	4	4	4
25	30	1	45	2	2	2	2	2	2	2	2
23	30	1	37	5	5	5	5	5	5	5	5
22	30	1	54	4	4	4	4	4	4	4	4
21	30	1	37	5	5	5	5	5	5	5	5
19	30	1	37	1	1	1	1	1	1	1	1

確認: 原データ, 環境意識調査データ

28

「質問」の特徴と留意点

- いわゆる「**質的データ**」である
- とくにこの2問は「**順序尺度**」である(**名義尺度**であって選択肢の意味に**序列の関係**がある)

- さてここで、以下の問題を考える。

Q: **このようなデータに対して以下の操作は可能だろうか**

- ① **四則演算**を行うこと、例えば平均値や標準偏差を求めることが可能か。
- ② 主成分分析、因子分析など、原則として「**量的データ**」を**対象とした手法**は適用できるのか。
- ③ クロス表を作成し**比率データを観察**するのはなぜか。

29

これらの解答は、...

- ① 四則演算を行うこと、例えば平均値や標準偏差を求めることが可能か。

答え: 形式的には可能でも、**演算の意味**があるかは保証されない。

- ② 主成分分析、因子分析など、原則として「**量的データ**」を対象とした手法は適用できるのか。

答え: これも形式的適用はあり得るが、実はいろいろと**問題**がある、とくに因子分析の利用には細心の注意が必要。

- ③ クロス表を作成し**比率データを観察**するのはなぜか。

答え: 質的データの情報のもっとも簡単な「**計量化・数量化**」の方法であるから(対応分析法に密接に関係)。

30

対応分析法が扱うデータ表

- 対応分析法・数量化法III類の原理から考えると、そこで扱えるデータ表の形式にはかなり**自由度がある**.
 - 一般に対応分析法・数量化法III類で扱うデータ表の形式
 - WordMinerで扱うデータ表の形式
- 扱うデータ表の**相互の関係を理解**することが肝要である。
- これについて概略を述べる。後でテキストも良く読んでいただきたい。
- 数量化法III類と対応分析法が扱うデータ表は**一見すると異なる**ようにみえる(ここに誤解がある)。実は**同じであったり、ある種の變形したデータ表である**にすぎないこと。

31

対応分析法で扱うデータ表の特徴

- 原則として「**二元のデータ表(クロス表型)**」を基本
- 二値の応答型データ(「yes」「no」型、0-1型)である場合
- 「二元のデータ表」の特徴(条件)は、
 - データ表の各要素(各セル内の値)が**非負の数値**
 - 行または列の「**プロフィール(相対比率)**」が**意味のある**データ
 - データ表の行または列の「**比率パターン(比率の分布)**」が**意味を持つ**データ表
- 以上を満たせば、**ほぼ、どのようなデータ表でも利用**できる。
- 換言するとさまざまなデータ表の**變形**が用意できる。
- WordMinerを利用するため、**テキスト型データあるいはそのままなせるデータをこれに適合させる工夫**が必要。

32

例えば, ...

- 通常の **二元クロス表** (分割表) を基本形式として考える.
- (0, 1) 型データ行列 (二元クロス表の特別な場合に相当)
- **多重クロス表・パート表** (「多元クロス表」ではないことに注意)
- 多くの統計表 (数値が非負の集約化データで上の条件を満たすとき).
- この「二元のデータ表」をどのように作成するか.
- つまり **データ収集法** (data collection mode) と **取得計画・調査計画** (design) の問題である.
- 何でも質的データ (⇒ 自由回答・自由記述, テキスト型データ) であればよいとはならない. 意図・目的を持って取得すること.

33

項目*I*と項目*J*の(2元)クロス表(表19)

$$F = (f_{ij}) \quad (f_{ij} \geq 0, i \in I, j \in J)$$

		項目 <i>J</i>						
		1	2	...	<i>j</i>	...	<i>n</i>	行和
表側 項目 <i>I</i>	選択肢							
	1	f_{11}	f_{12}	...	f_{1j}	...	f_{1n}	f_{1+}
	2	f_{21}	f_{22}	...	f_{2j}	...	f_{2n}	f_{2+}
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
	<i>i</i>	f_{i1}	f_{i2}	...	f_{ij}	...	f_{in}	f_{i+}
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
<i>m</i>	f_{m1}	f_{m2}	...	f_{mj}	...	f_{mn}	f_{m+}	
列和	f_{+1}	f_{+2}	...	f_{+j}	...	f_{+n}	f_{++}	

34

例をみる(テキスト, 93ページ, 例3)

- (表7)にあるような調査データがある. これは「(回答・サンプル) × (多数項目・変量)の**多変量構造のデータ表**である.
- この中の2項目(2つの質問AとB)を指定して**クロス表(表6)**を生成する(テキスト, 9ページ参照).
- これに対応分析法を適用すれば, 2つの質問AとBの関係(対応)が測れるはずである, と考える.
- このクロス表の「**表側(ひょうそく)**」と「**表頭(ひょうとう)**」に何を置くかで, 様々な変形があり得る.

35

クロス表の例(表6): (項目A × 項目B)のクロス表

度数 列%	1. かなり多い	2. 多いほう	3. ふつう	4. 少ない	4. 少ないほう	行和 行%
1. たいへん気に入っている	166 54.43 31.68	239 27.19 45.61	86 18.49 16.41	7 10.14 1.34	26 11.40 4.96	524 26.93
2. まあ気に入っている	131 42.95 10.61	598 68.03 48.42	324 69.68 26.23	36 52.17 2.91	146 64.04 11.82	1,235 63.46
3. あまり気に入っていない	6 1.97 3.49	40 4.55 23.26	55 11.83 31.98	20 28.99 11.63	51 22.37 29.65	172 8.84
4. 気に入っていない	2 0.66 13.33	2 0.23 13.33	0 0.00 0.00	6 8.70 40.00	5 2.19 33.33	15 0.77
列和	305 15.67	879 45.17	465 23.90	69 3.55	228 11.72	1,946

確認: 環境意識調査(1973 cases)

36

対応分析法で扱う典型的なデータ表の例

例番号	対応する表	行または表側項目	列または表頭項目	テキスト該当ページ
例1	表2, 表3(*)	サンプル	銘柄	90ページ
例2	表4, 表5(*)	サンプル	「好む」清涼飲料水の銘柄	91,92ページ
例3	表7(*)	サンプル(回答者)	質問項目A, B	94ページ
	表6	質問項目A	質問項目B	93ページ
例4	表8(*)、表9(*)	サンプル(回答者)	質問項目A, B, C, ...	95ページ
	表12	質問項目A	質問項目B	96ページ
例5	表15(*)	サンプル(回答者)	質問項目I, J	99ページ
	表16	質問項目I	質問項目J	100ページ

いずれもWordMinerで対応処理が可能である。

演習:

どのようにデータ表を作ってインポートすればよいだろうか。
上の(*)印を付けた表の形式はすべて対応する。

37

データ表の相互の関連(重要)

- 数量化法III類は, (0, 1)型データとすることが多い。
 - (サンプル) × (もの, 項目, カテゴリー)のデータ表(例:表2~5)
 - アイテム・カテゴリー型のデータ表(例:表11の右側, 表17など)
 - こうしたデータ表しか扱えないと思われる節がある(誤解)
- 対応分析法ではより一般的に2元のデータ表を扱う。
 - (回答・サンプル) × (多変量項目)型から出発(例:表3,5,7,8,15など)
 - アイテム・カテゴリー型(インジケータ行列)(例:表11の右側, 表17,91ページの表4など)
 - 多重クロス表・パート表(Burt's tables, Burt's matrix)(例:表13,14,18, ~5,90ページの表2,93ページの表6など)
 - この他の二元表タイプ
 - これらのデータ表の間の関係が重要
 - 実は同じことを考えている場合が多い(テキストに若干記述あり)

38

参考:データ表の相互の関係を例4で確認

- データ表の関係を例でみる(94ページ, 例4:自治体の意識調査)
- 元となる「(回答・サンプル) × (項目)型」データ表(表8)
- 質的データ(名義尺度, 順序尺度)が多いことに注意
- コード変数, 文字変数(テキスト型データ)の表記が混在
- ある2項目を切り出し(指定すると)表9を得る
- それをコード化すると(しなくてもよいが)表10となる
- 表10をアイテム・カテゴリー型(インジケータ行列)に展開(表11)
- 表10(表9)からクロス表を生成(表12)
- 表11(アイテム・カテゴリー型)から行列演算(行列の積)で, 表13の多重クロス表(パート表:Burt's table)を生成
- 表12のクロス表は表13の多重クロス表内のブロック行列となる
- 以上の関係は多数項目(多変量)となっても同様になる
- 演習問題2とその「補足」に要約した(テキスト, 135ページから)

39

多重クロス表の例(表13を表11から生成する)

質問	質問	質問 A					質問 B				
		守っている	まあ守っている	あまり守っていない	守っていない	無回答	お寺詣りをよくする	たまにお寺詣りをする	あまりお寺詣りをしない	お寺詣りをしない	無回答
問 A	守っている	106	0	0	0	0	41	26	22	15	2
	まあ守っている	0	167	0	0	0	25	67	45	30	0
	あまり守っていない	0	0	84	0	0	6	13	34	31	0
	守っていない	0	0	0	41	0	1	6	7	27	0
	無回答	0	0	0	0	15	1	4	1	2	7
問 B	お寺詣りをよくする	41	25	6	1	1	74	0	0	0	0
	たまにお寺詣りをする	26	67	13	6	4	0	116	0	0	0
	あまりお寺詣りをしない	22	45	34	7	1	0	0	109	0	0
	お寺詣りをしない	15	30	31	27	2	0	0	0	105	0
	無回答	2	0	0	0	7	0	0	0	0	9

- ① 表12のクロス表がブロック行列で入っている
- ② 対角ブロック行列がそれぞれの項目の周辺度数分布
- ③ 当然, 対称行列となっている

質問Aと質問Bのクロス表

40

テキスト型データはなぜ質的データか？

- テキスト型データは、質的データにどのように変換できるのかを考える。
- テキスト型データがまず **分析処理単位** に分けられていなければならない。
- つまり何らかの形で **扱い単位** を決めること⇒ WordMinerではこれを「**構成要素 (fragments)**」という。
- 構成要素はどんな単位であってもよい(ゴミも入っている)。
- WordMinerの標準装備の機能として **分かち書き処理機能** がある、これを行うと区切りのない日本語テキスト型データが分かち書き処理され **構成要素の単位** に分解される。
- よって何らかの「**処理単位 = 構成要素**」が作ればよいので **他の分かち書きツール** を使ったデータを用いても良い。
 - 例: 事例紹介、樋口耕一氏のペーパー (KH Coder) (テキスト, 123ページから)
 - 例: 茶釜などの **形態素解析ソフト** を使う (WordMinerでそのまま処理)

41

WordMinerで扱うデータ表形式 (表27)

表側項目: I	表頭項目: J
構成要素変数 (分かち書き, キーワード)	回答 (サンプル), 個体
構成要素変数 (分かち書き, キーワード)	質的変数 (選択肢型設問, 属性項目等)
構成要素変数 (分かち書き, キーワード)	クラスター変数 ※) クラスター・メンバーシップ情報から得られるクラスター変数は質的変数に変換して名義尺度データとして使う

- ① 構成要素変数, つまり抽出・編集した単語・語句と何をクロスさせるかで様々な変形がありうる。
- ② 所与のデータを項目 I と項目 J に対応させることで, 分析の範囲を拡張できる。

42

WordMinerにおけるデータ表のイメージ図 (図8)

		分かち書きで得られる構成要素 (単語, 語句, キーワード…)							
「回答者・サンプル」 あるいは 「質的変数・属性」	1	$w_1^{(1)}$	$w_2^{(1)}$...	$w_j^{(1)}$...	$w_k^{(1)}$...	
	2	$w_1^{(2)}$	$w_2^{(2)}$...	$w_j^{(2)}$...	$w_k^{(2)}$...	
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	
	i	$w_1^{(i)}$	$w_2^{(i)}$...	$w_j^{(i)}$	$w_l^{(i)}$...
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	
n	$w_1^{(n)}$	$w_2^{(n)}$...	$w_j^{(n)}$...				

- ① 「(回答・サンプル) × (構成要素)」のデータ表の場合
- ② 「(構成要素) × (質的変数)」のデータ表の場合
- ③ その他, 必要に応じて変数指定を行い, このデータ表に合わせる。

43

例: Web調査の質問 (テキスト, 121ページ)

問3. 次に、あなたと「インターネット」とのかかわりについてお伺いします。

3-1. あなたご自身にとって「インターネット」は、どのようなことがらに活用できると思いますか。どんなことでも結構ですので、以下になるべく具体的にご記入ください。

3-2. では、一般的に「インターネット」は、どのようなことがらに活用できると思いますか。なるべく、他にはないような活用法を、どんなことでも結構ですので、以下になるべく具体的にご記入ください。

44

「(回答・サンプル) × (構成要素)」のデータ表の例(表28)

サンプル	構成要素(キーワードを用いたとき)
1	為、しちべ、利用、家族、遊園地、公園、売べ物屋、情報収集、調査
2	あまり、セキュリティ、必要、ミーティング、世間話、仕事
3	新製品、スペック、価格、お店
4	役所、証明書発行、受け取り
5	旅行、計画、観光地、チェック、お店、情報収集
6	情報収集、調査、メール、座席予約、航空機、列車、オークション
7	地図検索、鉄道、乗り換え、検索、その他、時々、必要、情報検索
8	通信販売、申し込み、旅行、情報収集
9	情報ツール
10	あまり、ふつう、店舗、販売、商品、販売店、ショッピング、建築図面作成用、CADデータ、ダウンロード
11	自分、興味、事例、容易、公式、専門家、情報
12	日常生活、中、帰省時、飛行機、時刻表、育児、経験談、アドバイス、仕事、必要、情報、特定人物、活動、著書
13	情報収集
14	電話、手紙、かわり
15	仕事上、事、出張、際、ホテル、情報、等
16	パソコン、周辺機器、仕様、価格、懸賞、応募、ドライバ、ダウンロード、ゲーム
17	掲示板、一つ、場所、みんな、話
18	調べ物、ショッピング、オークション
19	情報、収集、自己、PR
20	ニュース、天気、行事情報、仕事、情報
21	映画、書籍、情報入手、収入検索、単語、等、検索、メール
22	専門的、事例、情報収集
23	メール、一番、仕事、不明瞭、確認、美術館、博物館、映画、その他、催し物、情報収集、たまに、オークション、お食事、電車、時刻表、経路
24	調べ物、ホームページ、サイト
25	友人、知人、連絡
26	百科事典
27	趣味、人、交流、勉強、場所、交通機関、時間
28	天気予報、道路状況、宿泊情報、等、行先、情報収集、辞書、新聞
29	不、知識、情報、時間、辞書、新聞、地図、最近、ネット、使用、事

確認: 原データ、環境意識調査データ

生成した「(回答・サンプル) × (構成要素)」のクロス表(表29)

サンプル	行和	H P	いろいろ	いろいろ	いろいろ	お店	その他	ときに	やり	やりとり	アーカイブ	イベント	インターネット	オークション	オンラインショッピング	ゲーム
1	25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	25	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0
3	17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	14	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0
7	14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8	13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9	13	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0
10	13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
11	13	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
12	12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
13	12	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0
14	12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
15	12	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
16	12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	11	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
18	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
20	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
22	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
24	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

このタイプのデータ表は行列の寸法が大きく、かつ要素の度数が非常に速になることが特徴
※固有値はきわめて小さい値となることが多い ⇒ 固有ベクトルの向きの動きにも注意

確認: WM環境意識調査データ(回答) × (構成要素)

「(構成要素) × (質的変数)」のデータ表の例(表30)

サンプル	性別	年齢区分	学年区分	表取特	職業	構成要素(ここでキーワード)
1	男性	4.35才~39才	性/4.35才~39才	既婚	営業職	為、しちべ、利用、家族、遊園地、公園、売べ物屋、情報収集、調査
2	男性	5.40才~44才	性/5.40才~44才	既婚	研究開発職	あまり、セキュリティ、必要、ミーティング、世間話、仕事
3	女性	5.40才~44才	性/5.40才~44才	既婚	主婦専業	新製品、スペック、価格、お店
4	男性	5.40才~44才	性/5.40才~44才	既婚	労務職	役所、証明書発行、受け取り
5	女性	2.25才~29才	性/2.25才~29才	既婚	主婦専業	旅行、計画、観光地、チェック、お店、情報収集
6	男性	5.40才~44才	性/5.40才~44才	既婚	研究開発職	情報収集、調査、メール、座席予約、航空機、列車、オークション
7	男性	4.35才~39才	性/4.35才~39才	既婚	無職・その他	地図検索、鉄道、乗り換え、検索、その他、時々、必要、情報検索
8	女性	2.25才~29才	性/2.25才~29才	既婚	主婦専業	通信販売、申し込み、旅行、情報収集
9	男性	3.30才~34才	性/3.30才~34才	既婚	自営業とその他	情報ツール
10	男性	6.45才~49才	性/6.45才~49才	既婚	専門職	あまり、ふつう、店舗、販売、商品、販売店、ショッピング、建築図面作成用、CADデータ、ダウンロード
11	男性	3.30才~34才	性/3.30才~34才	未婚	無職・その他	自分、興味、事例、容易、公式、専門家、情報
12	女性	3.30才~34才	性/3.30才~34才	既婚	専門職	日常生活、中、帰省時、飛行機、時刻表、育児、経験談、アドバイス、仕事
13	男性	9.60才~64才	性/9.60才~64才	既婚	無職・その他	必要、情報、特定人物、活動、著書
14	男性	8.55才~59才	性/8.55才~59才	既婚	管理職	電話、手紙、かわり
15	男性	7.50才~54才	性/7.50才~54才	既婚	販売・保安・サービス	パソコン、周辺機器、仕様、価格、懸賞、応募、ドライバ、ダウンロード
16	男性	5.40才~44才	性/5.40才~44才	既婚	営業職	ゲーム
17	男性	5.40才~44才	性/5.40才~44才	既婚	技術職	掲示板、一つ、場所、みんな、話
18	女性	3.30才~34才	性/3.30才~34才	既婚	パート・アルバイト	調べ物、ショッピング、オークション
19	女性	1.25才未満	性/1.25才未満	未婚	自由業	情報収集、自己、PR
20	女性	3.30才~34才	性/3.30才~34才	既婚	技術職	ニュース、天気、行事情報、仕事、情報

質的変数として年齢区分を指定
構成要素変数としてキーワード抽出で生成した変数を指定

生成した「(構成要素) × (年齢区分)」のクロス表(表31)

通番	行和	1.25才未満	2.25才~29才	3.30才~34才	4.35才~39才	5.40才~44才	6.45才~49才	7.50才~54才	8.55才~59才	9.60才~64才
117	情報	270	39	42	41	36	45	26	19	7
121	情報収集	130	11	19	21	27	20	14	10	2
109	検索	99	15	14	19	15	17	6	7	1
33	メール	95	12	13	11	19	17	4	10	5
46	検索	79	12	9	14	11	11	5	9	2
84	仕事	74	8	5	14	14	11	9	8	3
162	友人	60	7	6	9	12	9	9	4	1
145	車	58	6	11	10	8	5	6	5	1
149	入車	58	7	8	4	10	12	4	6	2
168	旅行	56	1	8	9	10	9	2	7	2
55	管理	55	11	8	11	9	7	3	3	0
91	書	54	11	10	16	5	1	5	2	1
99	自分	49	9	6	15	3	6	3	5	1
158	ショッピング	48	3	7	12	8	3	7	3	1
150	調べ物	46	3	11	9	9	7	2	2	0
170	通勤	46	4	5	6	6	9	5	3	3
94	映画	42	2	8	6	7	6	7	1	0
135	調べ物	40	8	0	13	4	10	2	2	0
110	収集	36	3	6	4	6	8	2	0	0
31	ホームページ	34	6	6	5	6	3	1	6	0
128	人	34	11	10	6	4	2	0	1	0
93	新聞	33	3	4	4	4	2	2	4	1
165	利用	33	1	4	7	4	8	5	1	2
16	コミュニケーション	29	7	4	5	8	3	1	1	0
78	購入	29	3	5	2	4	5	4	1	0
13	オークション	28	3	5	5	6	5	0	2	0
113	映画	28	3	5	5	4	5	1	2	0
153	応募	28	3	4	3	5	5	1	1	1

一般には質的変数の選択肢数が大きくはないので、要素内の度数が比較的小まるとなるのが特徴
※固有値の値も、比較的大きい値が出る⇒よって寄与率も目安となる

確認: WM環境意識調査データ(構成要素) × (質的変数)

補足情報:

- ここからは、(質問に備えての)補足情報として要約する。
- 記述の大半は配付資料(テキスト)に圧縮して書かれている。
- それ以上の情報を必要とするときは、関連書籍を参照する。
- とくに、多次元データ解析の諸方法については、このことへの理解にそれなりの手間と時間を要する。
- テキストやここに補足とした情報の中の「キーワード」となる言葉の感じ(こんなことを言っているらしい)を汲み取ることから始めること。

49

いわゆる「数量化」の原点は何か？

- 数量化法Ⅲ類の誕生の経緯(故林知己夫氏)
- その考え方・思想は何かを簡潔に言えば、...
 - 質的データに対して、数値はアブリアリに与えるべきではない。
 - 線形性(線形モデル)をその名目的な数値にそのまま想定はできない。
 - 「数量」は現象を説明するであろうデータに基づいて作られるもの(別で作るもの、「数量の作り方」を考えるべきこと)。
 - 新たな(なるべく線形となるような)座標空間(スコア)を作り出すこと。
- この発想は、そのまま定性情報である「テキスト型データ」に当てはまるであろう(選択肢はこの一つではないが)。
- つまりテキスト型データは一旦計量化・数量化した後に、さらなる分析に進むべきであるという選択肢があるだろう。
- 換言すると、生のまま数値(コード)として扱う処理には問題があるのではないか？

50

対応分析法(AFC)とは？

- 質的データ、とくに「調査型データ」は多くの場合そのまま計量化して使えない(ここは数量化法と似た発想)。
- 質的データの原点はクロス表型データ表にある。
- つまり、質的データ(名義尺度、順序尺度)情報を集約化したデータ表である。
- クロス表(分割表)を基礎情報と考えると、これは比率データで観察を行うように、視点が表側の側と表頭の側と2つの方向から分析できる。
- このとき、それぞれ比率データは多次元空間内に布置する多次元データと見なすことができる(⇒後述)。
- この視点から、クロス表型データ表の表側、表頭の対応関係を測ることができるのではないか。

51

古典的な手法: クロス表の独立性の検定

- 既存の方法論として、分割表(クロス表)の「独立性の検定」がある(例:「ピアソンのカイ二乗統計量」を用いる)。
- この発想は、以下のような考え方である。
- 表側と表頭の2つの項目*i, j*の間には「関係がない(独立)」という帰無仮説をたてる。
- つまり表側と表頭にある2つの項目は無関係という独立モデル($p_{ij} = p_{i+}p_{+j}$)。⇒テキスト、表19、図2などを参照
- これが統計的に棄却されれば帰無仮説を棄却、よって表側と表頭の2つの項目*i, j*の間には何らかの関係がないとはいえない(関係がありそうと言えらるだろう)とする検定法(隔靴搔痒)。

52

対応分析法の本質

- ベンゼクリはこれを別の視点から謎解きした。
- クロス表の表側と表頭のそれぞれの項目の相対比率データを考える(「**プロフィール**」と名付けた)。
- プロファイルのカイニ乗距離(加重付距離)を考え、これが近いものは近い位置にあるとする(加重化した比率データが似ているものは近いとする)。
- つまりプロフィールを多次元空間内に布置する多次元データと考え、その空間内での加重平均指標を作り次元の縮約を行う(主成分分析のような合成指標化を考える)。
- プロファイルは行と列との両方から観察できるから、**双対性**を考慮して分析を行う。
- 一見すると、**数量化法Ⅲ類**と異なる定式化のように見えるが実は**同等の手法**である。

53

ピアソンのカイニ乗統計量との関係

- 対応分析のアプローチでは、**ピアソンのカイニ乗統計量**が重要な役割を果たす。
- 定式化の結果として(そのようになるように定式化して)、**ピアソンのカイニ乗統計量**と**密接な関係**にある。
- 本来、クロス表は2つの項目間に何らかの意味があるとして観測(測定)したはずなのに、**独立性の検定の帰無仮説(独立モデル)**のような設定では情報の活用が十分ではない。
- よって、クロス表の行と列との2項目間の関連性を主成分分析型手法とすることで、固有値(=相関の情報に相当)の大きさを測ることを可能とした。
- つまり、2つの項目間の関連性と対応関係を計量的測れることになる(ここでも**質的データの計量化**となる)。

54

対応分析の仕組み: 人工データによる確認

- **二元のクロス表(型)**を基本のデータ表とする(表19)
- このデータ表から**プロフィール**を作る(表20, 図2)
 - 行のプロフィール(行の相対比率のデータ)⇒式(8)
 - 列のプロフィール(列の相対比率のデータ)⇒式(9)
- プロフィールをそれぞれ行あるいは列の**多次元空間内のデータ**と考え、この空間内での**次元縮約**を行う(加重平均による合成指標化)⇒つまり主成分(成分)を作る(図3)
- 要点は、
- **比率**を考えることで、行と列との**双方向から**データを観察することに注意する(双対性に関連する)(例えば式(19), (20))
- ここでいう「多次元データ, 多次元空間」とは何かが重要(図3)。

55

数値例で見るのが理解を容易にする

- 数式の誘導やその意味付けをある程度知ることは重要。
- 必要最小限の数理はテキストに記したのでこれを参照。
- ここでは何より、**対応分析法の仕組み**を知ること。
- これを実装するWordMinerの機能を理解すること。
- 例示したデータ表のうちから「例5」のレストラン評価を用いる。
- これは対応分析法の仕組みを理解するために簡略化した人工データである。
- まずこの例5の「意味, 内容」(データ表の意味)を理解する。
- それぞれの情報, データ表, 用語の確認などはテキストで確認のこと。

56

①「質問」と「選択肢」の確認

- まず用いる例を挙げる。[例5]

質問I: 次に挙げるレストランのうち、あなたがお気に入りのレストランはどれですか？

- | | | | |
|---------|---------|----------|---------|
| 1. さとみ | 2. パツハ | 3. ムガール | 4. いりふね |
| 5. コルシカ | 6. クラーク | 7. ロゴスキー | 8. きくみ |
| 9. ラ・マレ | 10. かりや | | |

質問J: その選択時の評価基準は次の3つのうちのどれでしょうか？

1. 味 2. 量 3. 工夫・サービス

サンプル数(回答者数)が $N=1,284$ (人)、2項目(I と J)の多変量データ構造のデータ表(表15)から、クロス表(表16)が得られる。ここで質問を「項目」と呼び、選ぶカテゴリーを「選択肢」と名付ける。この場合、項目 I の選択肢数は $m=10$ 、項目 J のそれは $n=3$ となる。

57

②クロス表の生成

- 表16のクロス表を作成する。ここでは「行に項目 I 」を「列に項目 J 」を充てた。(表19参照)
- クロス表は多次元データである(⇒後述、図2、図3)。
- このクロス表を以下の式で表す(式(1)、(2))。

$$\mathbf{F} = (f_{ij})_{m \times n} \quad (f_{ij} \geq 0, i \in I, j \in J) \quad \text{〈項目}I\text{と項目}J\text{のクロス表、ここで寸法は}m \times n\text{である〉}$$

$$I = \{1, 2, \dots, m\}, \quad J = \{1, 2, \dots, n\} \quad \text{〈項目}I\text{と項目}J\text{の選択肢} \\ \text{※項目}I\text{の選択肢数は}m\text{個、項目}J\text{の選択肢数は}n\text{個}$$

「(項目 I) × (項目 J)」のクロス表 ※表19を確認

58

項目 I と項目 J のクロス表(表19)

$$\mathbf{F} = (f_{ij})_{m \times n} \quad (f_{ij} \geq 0, i \in I, j \in J)$$

表頭

		(項目 I) × (項目 J)のクロス表							
		項目 J							
		選択肢	1	2	...	j	...	n	行和
表側	項目 I	1	f_{11}	f_{12}	...	f_{1j}	...	f_{1n}	f_{1+}
		2	f_{21}	f_{22}	...	f_{2j}	...	f_{2n}	f_{2+}
		⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
		i	f_{i1}	f_{i2}	...	f_{ij}	...	f_{in}	f_{i+}
		⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
m	f_{m1}	f_{m2}	...	f_{mj}	...	f_{mn}	f_{m+}		
列和		f_{+1}	f_{+2}	...	f_{+j}	...	f_{+n}	f_{++}	

59

③プロフィールを作る

- 行のプロフィール, つまり行の相対度数(相対確率)を求める。いわゆる「行和を1」と揃えた(行100%とした)表と思えばよい, これが表21である。[図2の左側の流れ]
- 列のプロフィール, 「列和を1」とした列の相対度数(相対確率)を求める。これが表22である。[図2の右側の流れ]

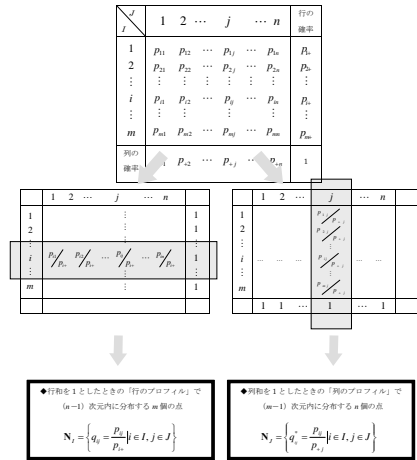
$$\mathbf{N}_I = \left\{ q_{ij} = \frac{p_{ij}}{p_{i+}} \mid i \in I, j \in J \right\} \quad \text{〈行のプロフィール〉 式(8)}$$

$$\mathbf{N}_J = \left\{ q_{ij}^* = \frac{p_{ij}}{p_{+j}} \mid i \in I, j \in J \right\} \quad \text{〈列のプロフィール〉 式(9)}$$

※式(8)、(9)と図2を確認

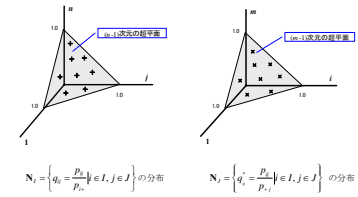
60

プロフィールの関係

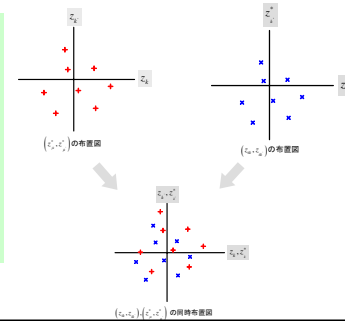


61

成分スコアの 布置をイメージ で示すと、...



- ① 重心座標系の空間内で次元縮約を行う
- ② データ表から得られた成分スコアを布置図とする(行成分スコア, 列成分スコア)
- ③ 必要に応じて, 行成分スコアと列成分スコアとの同時布置図とする



62

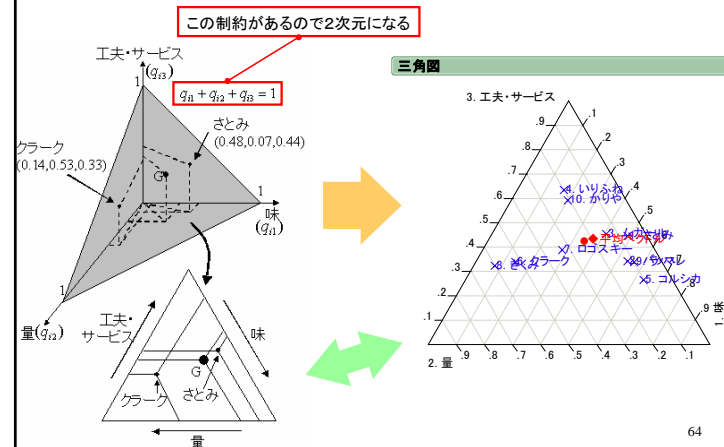
④ 多次元空間に布置することの意味(表21,22)

- (i) **行のプロフィール**とは「項目:評価基準の3つの選択肢」(=3次元空間内)に「項目:レストランの10の選択肢」が布置するデータ空間と考える(行の向きに行和=1と揃えたことに注意). [表21]
 - (ii) 同じく, **列のプロフィール**とは「項目:レストランの10の選択肢」(=10次元空間内)に「項目:評価基準の3つの選択肢」が布置するデータ空間と見ることできる(ここは列和=1と揃えた). [表22]
- 1) この一般的なクロス表(表19)から得られる行プロフィール, 列プロフィールの関係を図式化したものが図2と図3である.
 - 2) ここで「**行と列の両方向**」から見ていることに注意(行と列を転置しても情報は変わらない; **双対性**がある).

※図2と図3を確認

63

行のプロフィールの意味を図で確認(図4)



64

さらに, ...

- 1) 例題についてさらに「**行のプロフィール**」側からの観察を続ける。つまり, 図2, 図3の左側のパスを考える。
- 2) このとき, $n=3$ であるから行プロフィールを実際に「**3次元空間内**」に描いてみると図4の左側の図となるだろう。
- 3) ここで「**行和=1**」という制約があるから, 10のレストランの布置は, 実は $n-1=2$ (次元)の平面内に入る(自由度が1だけ減る)。
- 4) 実際に, 図4の左の図の網かけ部の**平面上**に分布する。
- 5) これをそのまま(点の布置関係を保持したまま)射影すると図4の右側の図(三角座標系の図=**三角図**)となるだろう。
- 6) この例は, **視認できる3次元(2次元)の説明**であるが, 多次元になって**次元数が上がっても考え方は同じ**である。

65

⑤次元数の縮約を行うこと

- 1) 考えるデータ布置の空間が異なるが多次元空間内での**次元縮約**を考えることには違いがない。
重心座標系 (barycentric coordinate system)
三角図: 三角座標系 (triangular coordinate system)
- 2) 高次元の空間に布置されるデータを**少数次元内に縮約**せねばならない。
- 2) **主成分分析と同じようなこと(加重和を作る)**が考えられるのか?
- 3) そのためのデータの構造は(作り方は)?
数式(10)~(14)のような変換を行ったデータを扱えばよい, またこの形でないと不都合を生じる。
- 4) こうする理由がある(例: **ピアソンのカイニ乗統計量**と関係)

※図3の布置図イメージを確認

66

⑥データ行列の分解(固有値問題他)

- 1) データ行列を作りその共分散行列の**固有値問題に帰着**する(あるいは元のデータ表の特異値分解:SVD)
- 2) データがある形であることを除けば多くの合成指標型手法(主成分分析など)と同じ解法となる。
- 3) 固有値, 固有ベクトルを求める
- 4) **固有値と寄与率**が情報縮約の程度を知る指標
- 5) **プロフィール(を加工したある形)の加重平均(=成分スコア)**を求めることに帰着[成分スコア=数値化スコア, 数値化得点]
- 6) 固有ベクトルが加重平均の式の係数に相当
(注: 式(10),(11)に固有ベクトルを加重とする一次結合式)
- 7) **成分スコア**(数値化スコア, 数値化得点)の算出
- 8) **行と列との双方向**から考えるから**成分スコアも2組ある**

67

⑦成分スコア, 固有値の性質を確認

- 1) 成分スコアは項目 I の選択肢 i ($i \in I$)と項目 J の選択肢 j ($j \in J$)のそれぞれに対して付与される。[表25, 図6参照]
- 2) 両者の成分スコアの関係が重要(とくに双対性)
- 3) 成分の数, つまり**固有値の数**はクロス表の行数(m)と列数(n)の少ない方から1を引いた数: $K = \min\{m, n\} - 1$ となる
 z_{ik} ($i \in I, k = 1, 2, \dots, K$) (選択肢 i に対する第 k 成分の成分スコア)
 z_{jk}^* ($j \in J, k = 1, 2, \dots, K$) (選択肢 j に対する第 k 成分の成分スコア)

※図3の布置図イメージを確認

68

成分スコアと元の確率行列の関係(表25)

		項目 J					成分スコア						
		1	2	...	j	...	n	1	2	...	k	...	K
項目	1	p_{11}	p_{12}	...	p_{1j}	...	p_{1n}	z_{11}	z_{12}	...	z_{1k}	...	z_{1K}
	2	p_{21}	p_{22}	...	p_{2j}	...	p_{2n}	z_{21}	z_{22}	...	z_{2k}	...	z_{2K}

	I	p_{i1}	p_{i2}	...	p_{ij}	...	p_{in}	z_{i1}	z_{i2}	...	z_{ik}	...	z_{iK}

成分スコア	m	p_{m1}	p_{m2}	...	p_{mj}	...	p_{mn}	z_{m1}	z_{m2}	...	z_{mk}	...	z_{mK}
	1	z_{11}^*	z_{12}^*	...	z_{1j}^*	...	z_{1n}^*	↑ 行の項目 I の選択肢の成分スコア					
	2	z_{21}^*	z_{22}^*	...	z_{2j}^*	...	z_{2n}^*	↑					
	↑					
	k	z_{k1}^*	z_{k2}^*	...	z_{kj}^*	...	z_{kn}^*	← 列の項目 J の選択肢の成分スコア					
...	←						
...	k'	$z_{k'1}^*$	$z_{k'2}^*$...	$z_{k'j}^*$...	$z_{k'n}^*$	←					
...	←						
...	K	z_{K1}^*	z_{K2}^*	...	z_{Kj}^*	...	z_{Kn}^*	←					

69

補足1: 固有値と寄与率

$0 \leq \lambda_k \leq 1$ ($k=1, 2, \dots, K; K = \min\{m, n\} - 1$) (第 k 成分の固有値)

$$tr(\mathbf{V}) - 1 = \sum_{k=1}^K \lambda_k \quad (K = \min\{m, n\} - 1) \quad \text{(固有値の総和)} \quad \text{式(24)}$$

(ここで, $tr(\mathbf{V})$ は行列 \mathbf{V} のトレース=対角要素の和, を示す)

$$V_k = \frac{\lambda_k}{\sum_{k=1}^K \lambda_k} \times 100(\%) \quad \left(\begin{array}{l} k=1, 2, \dots, K \\ K = \min\{m, n\} - 1 \end{array} \right) \quad \text{(第}k\text{成分の寄与率の式)} \quad \text{式(25)}$$

70

補足2: 固有値とピアソンのカイニ乗統計量の関係

固有値とカイニ乗統計量の関係

$$tr(\mathbf{V}) - 1 = \frac{\chi^2}{N} = \sum_{k=1}^K \lambda_k \quad (K = \min\{m, n\} - 1) \quad \text{式(26)}$$

ピアソンのカイニ乗統計量

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{N(p_{ij} - p_{i+} p_{+j})^2}{p_{i+} p_{+j}} = \sum_{i=1}^m \sum_{j=1}^n \frac{\left(f_{ij} - \frac{f_{i+} f_{+j}}{N} \right)^2}{\frac{f_{i+} f_{+j}}{N}} \quad \text{式(28)}$$

※この関係は対応分析法を考えるうえできわめて重要

71

参考: ピアソンの χ^2 統計量の一般型

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{(\text{実現度数}_{ij} - \text{期待度数}_{ij})^2}{\text{期待度数}_{ij}}$$

- ①クロス表の表側と表頭という2つの分布の一種の距離となっている(乖離度を測る)。
- ②ただし, 独立モデルのとき, もっとも小さくなるような距離。
- ③現実には, 表側と表頭との間の相関・関連を知りたい, つまり独立でない程度を知りたいはずである。
- ④ピアソンの χ^2 統計量では直接はこれを測れない。
- ⑤式(26)はこれを固有値により測ろうとしていることに注意。

72

⑧ 双対性について

- 2項目 I, J の各選択肢に付与の成分スコア間の関係に注目
- いわゆる「**双対性(duality)**」がある。(きわめて重要)

$$z_{ik} = \frac{1}{\sqrt{\lambda_k}} \sum_{j=1}^n \left(\frac{p_{ij}}{p_{i+}} \right) z_{jk}^* \quad (i \in I, k = 1, 2, \dots, K) \quad \text{式(19)}$$

(行の成分スコアは列のその
のプロファイルの加重平均)

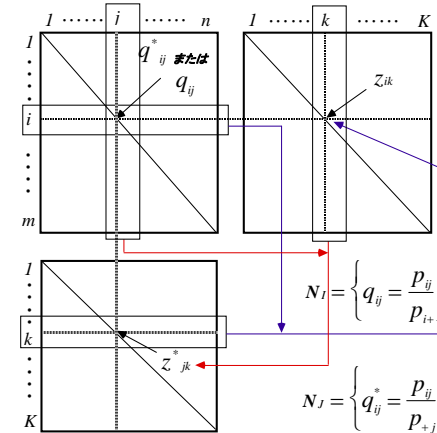
$$z_{jk}^* = \frac{1}{\sqrt{\lambda_k}} \sum_{i=1}^m \left(\frac{p_{ij}}{p_{+j}} \right) z_{ik} \quad (j \in J, k = 1, 2, \dots, K) \quad \text{式(20)}$$

(列の成分スコアは行のその
のプロファイルの加重平均)

※これら**成分スコアの関係は図6と表25のように考える。**

73

(i) 双対性の考え方(図6)



74

(ii) 成分スコアと元の確率行列の関係(表25)

		項目 J					成分スコア									
		1	2	...	j	...	n	1	2	...	k	...	k'	...	K	
項目	1	p_{11}	p_{12}	...	p_{1j}	...	p_{1n}	z_{11}	z_{12}	...	z_{1k}	...	$z_{1k'}$...	z_{1K}	
	2	p_{21}	p_{22}	...	p_{2j}	...	p_{2n}	z_{21}	z_{22}	...	z_{2k}	...	$z_{2k'}$...	z_{2K}	
	
	I	i	p_{i1}	p_{i2}	...	p_{ij}	...	p_{in}	z_{i1}	z_{i2}	...	z_{ik}	...	$z_{ik'}$...	z_{iK}
	
...	m	p_{m1}	p_{m2}	...	p_{mj}	...	p_{mn}	z_{m1}	z_{m2}	...	z_{mk}	...	$z_{mk'}$...	z_{mK}	
成分スコア	1	z_{11}^*	z_{12}^*	...	z_{1j}^*	...	z_{1n}^*	↑ 行の項目 I の選択肢の成分スコア ← 列の項目 J の選択肢の成分スコア								
	2	z_{21}^*	z_{22}^*	...	z_{2j}^*	...	z_{2n}^*									
									
	k	z_{k1}^*	z_{k2}^*	...	z_{kj}^*	...	z_{kn}^*									
									
...	k'	$z_{k'1}^*$	$z_{k'2}^*$...	$z_{k'j}^*$...	$z_{k'n}^*$									
...									
...	K	z_{K1}^*	z_{K2}^*	...	z_{Kj}^*	...	z_{Kn}^*									

75

⑨ 成分スコアの布置図と同時布置図

- 行の選択肢への成分スコア, 列の選択肢への成分スコアの**ドットプロット図(1次元)**や**散布図(布置図)**を描く。
- 同じ成分軸について行と列の成分スコアを重ねた図を**同時布置図**という。

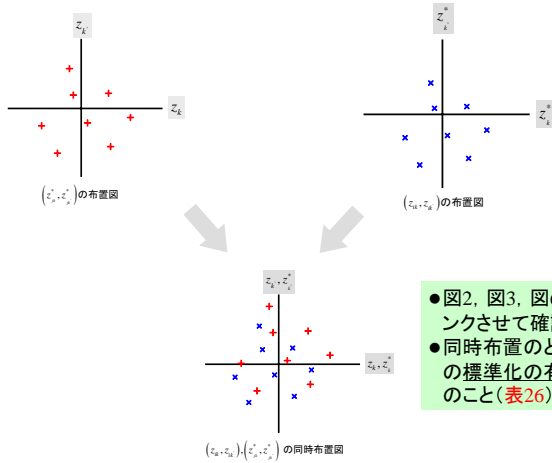
$$(z_{ik}, z_{ik'}) \quad \begin{cases} i = 1, 2, \dots, m \\ k, k' = 1, 2, \dots, K \\ K = \min\{m, n\} - 1 \end{cases} \quad \text{(行の選択肢への成分スコア) 式(21)}$$

$$(z_{jk}^*, z_{jk'}^*) \quad \begin{cases} i = 1, 2, \dots, m \\ k, k' = 1, 2, \dots, K \\ K = \min\{m, n\} - 1 \end{cases} \quad \text{(列の選択肢への成分スコア) 式(22)}$$

※図3の布置図イメージを確認

76

布置図と同時布置図 (図3)



- 図2, 図3, 図6, 表25をリンクさせて確認
- 同時布置のとき固有値の標準化の有無に注意のこと(表26)

「レストランと評価基準」の例で数値を確認

2項目への成分スコア

		成分スコア	
		第1成分スコア	第2成分スコア
項目と選択肢	成分	z_{1j}	z_{2j}
項目 J	さとみ	0.40067	-0.09077
	パツハ	0.39656	0.12200
	ムガール	0.19686	-0.08210
	いりふね	-0.20169	-0.40820
	コルシカ	0.54972	0.25857
	クラーク	-0.66717	0.25584
	ロゴスキー	-0.21980	0.10024
項目 J	さくみ	-0.85898	0.30915
	ラ・マレ	0.46355	0.11909
	かりや	-0.16472	-0.32610
	成分	z_{1j}^*	z_{2j}^*
項目 J	味	0.52347	0.17643
	量	-0.65787	0.25247
	工夫・サービス	-0.06055	-0.28561

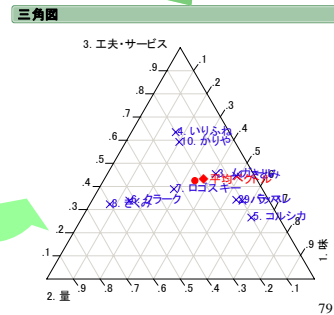
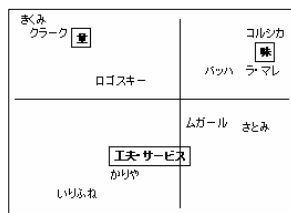
固有値と寄与率

主成分 k	固有値 λ_k	寄与率 (%)
1	0.19766	76.71
2	0.06002	23.29

※成分スコアは表25, 図6と対応させて確認のこと
 ①固有値の数は $K = \min\{m, n\} - 1 = 2$ となった。
 ②2成分に対する成分スコアが算出される。

※表23, 24に相当
 ※表25との対応に注意

得られた布置図と三角図の比較 (図4, 図5)



※図4と図5を比較のこと

元の2次元上の10のレストランの三角図布置が左の成分スコアとして再現されている(2成分スコアとして)

付録: 質的データの数量化の意味 (数値例)

- 数量化の意味をより理解するために簡単な数値例をみる。
- 対応分析法の理解をより深める。
- この例を通じて林の数量化法の考え方を知る。
- なるべく対応分析の機能と構造が見えるように作ってみる。
- 対応分析法と数量化法III類は“数理的には”同等であることを知る。
- 簡単な人工データを用意する(2例, いずれも人工データ)。
 - 例1: (サンプル:人) × (項目:銘柄)の対応分析
 - 例2: 「短文」による対応分析の機能確認
- ミニチュアの(構造を想定した)人工データを使って試用・体験することを推奨する。

例1:「(サンプル) × (銘柄)」のデータ表 (表32)

サンプル	銘柄	次年度調査の銘柄	一番好きな銘柄	性別	年齢区分
サンプル1	銘柄B, 銘柄E, 銘柄F	●銘柄E, ●銘柄F	◆銘柄B	▼男性	★30代
サンプル2	銘柄F	●銘柄F, ●銘柄B	◆銘柄F	▼男性	★40代
サンプル3	銘柄C, 銘柄F	●銘柄F	◆銘柄C	▼男性	★30代
サンプル4	銘柄B, 銘柄C, 銘柄E, 銘柄F	●銘柄C, ●銘柄B	◆銘柄E	▼男性	★30代
サンプル5	銘柄B, 銘柄C, 銘柄F	●銘柄B, ●銘柄C, ●銘柄F	◆銘柄C	▼男性	★30代
サンプル6	銘柄A, 銘柄B, 銘柄C, 銘柄E	●銘柄A, ●銘柄B	◆銘柄A	▼女性	★30代
サンプル7	銘柄A, 銘柄B, 銘柄D, 銘柄E	●銘柄D, ●銘柄E	◆銘柄B	▼女性	★20代
サンプル8	銘柄C, 銘柄F	●銘柄C, ●銘柄F	◆銘柄F	▼男性	★40代
サンプル9	銘柄A, 銘柄B, 銘柄E	●銘柄B, ●銘柄E	◆銘柄E	▼女性	★30代
サンプル10	銘柄A, 銘柄D, 銘柄E	●銘柄A, ●銘柄E	◆銘柄D	▼女性	★30代

- ①ここで「サンプル」と「銘柄」の2項目を選ぶ。
- ②前にみた例に合わせて項目Iをサンプル、項目Jを銘柄と対応させてみる。
- ③右側「次年度調査の銘柄」から「年齢区分」までは追加処理機能の説明(⇒テキスト).

81

①「(サンプル) × (銘柄)」のクロス表 (表33)

銘柄	銘柄 A	銘柄 B	銘柄 C	銘柄 D	銘柄 E	銘柄 F	行和
サンプル 1	0	1	0	0	1	1	3
サンプル 2	0	0	0	0	0	1	1
サンプル 3	0	0	1	0	0	1	2
サンプル 4	0	1	1	0	1	1	4
サンプル 5	0	1	1	0	0	1	3
サンプル 6	1	1	1	0	1	0	4
サンプル 7	1	1	0	1	0	0	4
サンプル 8	0	0	1	0	0	1	2
サンプル 9	1	1	0	0	1	0	3
サンプル 10	1	0	0	1	1	0	3
列和	4	6	5	2	6	6	29

項目「サンプル」の10の選択肢

項目「銘柄」の6選択肢

確認: (サンプル) × (銘柄)

82

②検討事項: 情報の見方

- 「(サンプル) × (銘柄)」のクロス表としたが、前にみた“2項目のクロス表”とは少しイメージが異なる。これはクロス表になっているのか。
- 「サンプル」を表側項目Iとし「銘柄」を表頭項目Jとし、「好きな銘柄」に回答(回答) = 1としたのでデータ表の要素として「度数 = 1」を充てたと考える。
- つまりこれもクロス表の特別なケースである。
- 林の数量化法III類では、これを(もの) × (項目・反応)の二値型データ表と考える。
- 同時に対応分析には「分布の同等性」という重要な性質がある(⇒後述、重要)。
- これが保持できるような二元表であればほとんど適用できる。

83

③固有値、寄与率と成分スコアの算出 (表34, 36)

k	固有値	寄与率 (%)	累積寄与率 (%)
1	0.6260	61.41	61.41
2	0.1877	18.41	79.82
3	0.1345	13.19	93.01
4	0.0452	4.43	97.45
5	0.0260	2.55	100.00

- ①固有値の数: $K = \min\{m, n\} - 1 = 6 - 1 = 5$ となる。
- ②人工データの例であり、かつある構造をもつように作ったので始めの固有値の寄与率が高い。
- ③始めの2成分で全情報の約8割を占める。
- ④固有ベクトルを使って成分スコアを求めると表36のようになった。
- ⑤ここで行(サンプル)と列(銘柄)の双方の成分スコアがある(図3、表25を参照)。この関係を良く理解すること。
※成分スコア⇒表36 ⇒表25も参照

84

④第1成分スコアの大きさで行、列を並べ替える(表35)

ID	銘柄D	銘柄A	銘柄E	銘柄B	銘柄C	銘柄F	行和	サンプルの 第1成分スコア
サンプル2	0	0	0	0	0	1	1	1.2969
サンプル3	0	0	0	0	1	1	2	1.1538
サンプル8	0	0	0	0	1	1	2	1.1538
サンプル5	0	0	0	1	1	1	3	0.7206
サンプル4	0	0	1	1	1	1	4	0.3785
サンプル1	0	0	1	1	0	1	3	0.1678
サンプル6	0	1	1	1	1	0	4	-0.2432
サンプル9	0	1	1	1	0	0	3	-0.6611
サンプル7	1	1	1	1	0	0	4	-0.9102
サンプル10	1	1	1	0	0	0	3	-1.1650
列和		4	6	6	5	6	29	
銘柄の 第1成分スコア	-1.3113	-0.9414	-0.5125	-0.1153	0.7997	1.0261		

- ①きれいに線形に並んでいる。この相関はどの程度？
- ②第2成分以上のスコアでも同じように並べ替えができる。

85

⑥要約:ここで数量化の意味を再考する

- 元のデータ表(クロス表, 表33)の行(サンプル)と列(銘柄)の各選択肢に新たな数量(成分スコア)を付与できた。
- それを成分別に観察すると、成分スコアの相関(係数)が固有値の正の平方根に相当する。
- つまり、選択肢に付与された成分スコアが数量(量的データ)として機能する。これが数量化の目標であった。
- 対応分析法では、プロフィールに注目したが、林の数量化法III類では、正にこれが付与した数量の相関を最大にするという最適化を行うことに相当する。
- データ表(クロス表)は多次元データであるが、これが少数の次元で説明できる(元の多次元の情報の損失がなるべく少ないような線形の加重平均を作った)。
- しかし、明らかに主成分分析のような量的データの分析とは異なる考え方である(質的データの数量化)。

87

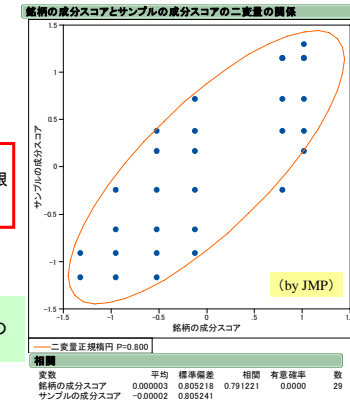
⑤第1成分スコアによる散布図(図9)

サンプルの 番号	銘柄の成分 スコア	サンプルの 成分スコア
1	-0.513	0.168
1	-0.115	0.168
2	1.026	0.168
2	1.026	1.297
3	0.800	1.154
3	1.026	1.154
4	-0.513	0.379
4	-0.115	0.379
4	0.800	0.379
4	1.026	0.379
5	-0.115	0.721
5	0.800	0.721
5	1.026	0.721
6	-0.941	-0.243
6	-0.513	-0.243
6	-0.115	-0.243
6	0.800	-0.243
7	-1.311	-0.910
7	-0.941	-0.910
7	-0.513	-0.910
7	-0.115	-0.910
8	0.800	1.154
8	1.026	1.154
9	-0.941	-0.661
9	-0.513	-0.661
9	-0.115	-0.661
10	-1.311	-1.165
10	-0.941	-1.165
10	-0.513	-1.165



- 相関係数: $r=0.7912$
- 固有値の正の平方根
 $\sqrt{\lambda_1} = 0.7912$

(Twin-map,
Dual-mapなどの
呼称あり)



- ①スコアの実寸で図を描くと縦軸、横軸の間隔(スコアの間隔)が等間隔でない。
- ②つまり数量化により元の選択肢(質的データ)が量的データに変換されている(“数量化”された)。
- ③成分スコアの相関係数: $r=0.7912$ となった。これが第1固有値の正の平方根に一致する。

例2:「短文」による対応分析法の機能の確認

- WordMinerにおける対応分析法の機能を知るには、
 - ①人工的に短い文章を作ってみる
 - ②雑誌、小説・エッセイ、その他書籍に登場する短文を集める
 - ③できるだけ、構造が分かるような文章を用意する
など、日頃から心がけるとよい。
- ここでは、ごく単純な文章を用意した。
 - (i)「私が文章を書く」を基本にいくつか変形文を作る。
 - (ii)「否定語」を入れる。
 - (iii)単語「文章」がない文をいくつか入れる。
- このとき、対応分析の結果に文章の特徴がどう表れるかを観察する。
- 実際のテキスト型データにはさらに複雑になるからきめ細かい探査が重要となる。

88

①データ表, 分かち書き, キーワード, 他

サンプル番号	(i)原文	(ii)分かち書きのみ	(iii)キーワード	(iv)分かち書き編集(助詞, 句読点削除)
SEQ	ID	原文	短文-分かち書き	短文-分かち書き-編集_all
[00000001]	1	私は、書かない。	私	私 書かない
[00000002]	2	私が書いた文章である。	私 文章	私 書いた 文章 ある
[00000003]	3	私に書いた文章です。	私 文章	私 書いた 文章 です
[00000004]	4	私の書けない文章だ。	私 文章	私 書けない 文章 だ
[00000005]	5	私が書いた文章である。	私 文章	私 書いた 文章 ある
[00000006]	6	私には書けない文章です。	私 文章	私 には 書けない 文章 です
[00000007]	7	私と書いた文章である。	私 文章	私 と 書いた 文章 ある
[00000008]	8	文章には書けない私のこと。	文章 私	文章 には 書けない 私
[00000009]	9	私が書く。	私	私 書く
[00000010]	10	私が書いた。	私	私 書いた
[00000011]	11	私と書いた。	私	私 書いた
[00000012]	12	私を書いた文章である。	私 文章	私 書いた 文章 ある

- ①図8のイメージ図を参照のこと。
- ②「分かち書き」そのものを用いるとき((ii)の欄)
- ③分かち書きから「助詞」を削除((iv)の欄)
- ④キーワード抽出の結果(動詞系「書く、書いた...」が除外されてしまう)((iii)の欄)

※これらと比べると「同じ原文を用いても加工で情報が変容」することが分かる。
つまり分析目的に応じた編集加工が必要となる(客観性を失わない範囲で)。

確認: WM短文の例

89

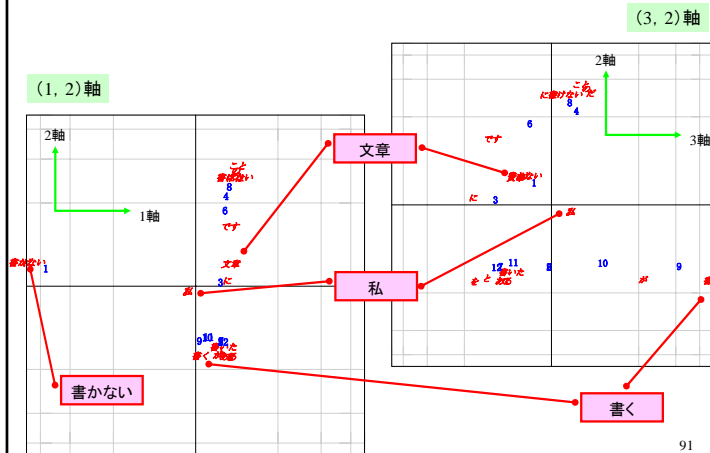
②「分かち書き」の構成要素の分布

構成要素	構成要素数	サンプル度数
私	12	12
文章	8	8
書いた	7	7
ある	4	4
が	4	4
で	4	4
書けない	3	3
です	2	2
と	2	2
には	2	2
の	2	2
こと	1	1
だ	1	1
に	1	1
は	1	1
を	1	1
書かない	1	1
書く	1	1

- ①利用頻度の高い語は「私、文章、書いた」など⇒共通に使われている語句
- ②利用頻度の少ない単語、とくに「書く、書かない」
- ③これらの単語が布置図のどこに位置するか
- ④元の語句の並び・語句の順序はどう現れるか

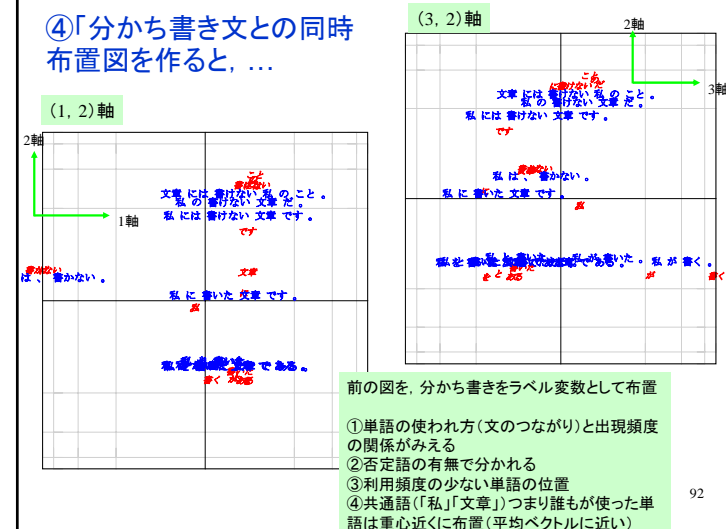
90

③同時布置図を作ると, ...



91

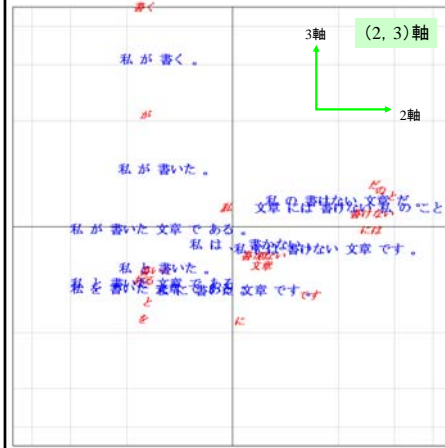
④「分かち書き文との同時布置図を作ると, ...



- 前の図を, 分かち書きをラベル変数として布置
- ①単語の使われ方(文のつながり)と出現頻度の関係がみえる
 - ②否定語の有無で分かれる
 - ③利用頻度の少ない単語の位置
 - ④共通語(「私」「文章」)つまり誰もが使った単語は重心近くに布置(平均ベクトルに近い)

92

⑤前の図の軸の向きを変えると視点も変わる



- ①文字情報が布置図内に入り
るので、軸の向きを変えても見える
情報に違いがある。
- ②固有値、寄与率などを考慮
しながら、高次成分の軸の組
み合わせを観察すること。
- ③併せて寄与度(絶対寄与度、
相対寄与度)の観察も行う。