

数理統計 補助資料 ～主成分分析～

2023年度2学期: 月曜1限, 水曜3限
担当教員: 石垣 司

多変量データの表記法

- データ i : $x_i = [x_{i1}, \dots, x_{iP}]^T$ ($i = 1, \dots, N$)
 N : データ数, P : 変数の数(データの次元)
- 全データ: $X = [x_1, \dots, x_N]^T$

ID	変数1	変数2	...	変数 P	
1	x_{11}	x_{12}	...	x_{1P}	x_1^T
2	x_{21}	x_{21}	...	x_{2P}	x_2^T
⋮	⋮	⋮	⋮	⋮	⋮
i	x_{i1}	x_{i2}	...	x_{iP}	x_i^T
⋮	⋮	⋮	⋮	⋮	⋮
N	x_{N1}	x_{N2}	...	x_{NP}	x_N^T



主成分分析

多変量データの傾向を説明する合成指標を作成し、その傾向を可視化や要約するための手法

- 次元削減法(次元圧縮, 次元縮約)の一種
- 多変量データの傾向を代表する低次元の主成分を見つける

多変量データ

ID	変数1	変数2	...	変数 P
1	x_{11}	x_{12}	...	x_{1P}
2	x_{21}	x_{21}	...	x_{2P}
⋮	⋮	⋮	⋮	⋮
i	x_{i1}	x_{i2}	...	x_{iP}
⋮	⋮	⋮	⋮	⋮
N	x_{N1}	x_{N2}	...	x_{NP}

合成指標
(主成分)

第1主成分

⋮

第M主成分

$M < P$

多変量データをよく説明できる
少数の軸を作り出す



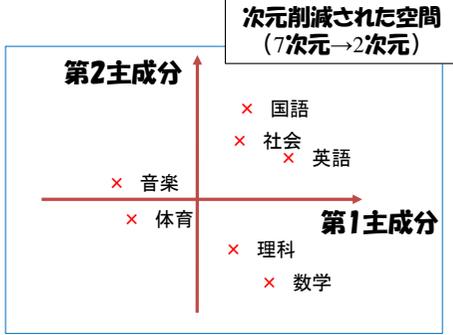
Karl Pearson 1857-1936

主成分分析による多変量データの可視化

多変量データの傾向を直感的に把握したい

- 例: 300人分の中学校のテスト(仮想データ)

生徒 No.	国語	数学	理科	社会	英語	音楽	体育
1	83	60	55	81	90	50	93
2	70	80	78	80	55	44	59
3	50	90	95	70	80	80	49
4	60	44	44	99	78	73	30
5	57	80	80	50	67	64	59
6	55	65	70	65	67	30	70
7	80	73	66	46	55	58	88
8	98	40	50	88	99	93	54
9	55	77	88	40	89	88	97
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮



可視化の場合は2次元に次元削減

主成分の内容を結果から解釈

- 第1主成分: 受験用学力
- 第2主成分: 文系・理系能力

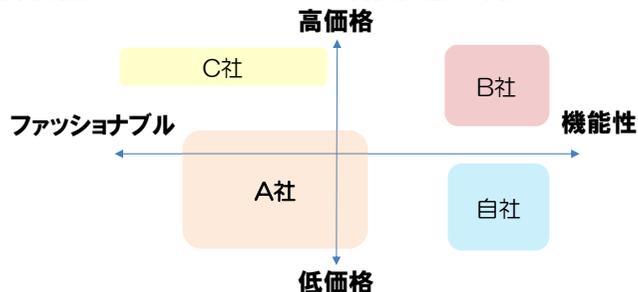
マーケティングにおける応用例

商品・サービスの差別化 ⇒ 優位性獲得のポジション

- 顧客ニーズの理解と価値の提供
- そのポジションを実現するマーケティング・ミックスの策定

知覚マップ (←本日は知覚マップをデータから作成)

- 消費者の目線での各ブランドの市場内での位置づけを可視化
- 新製品のポジショニングや既存製品の再ポジショニングに利用



5

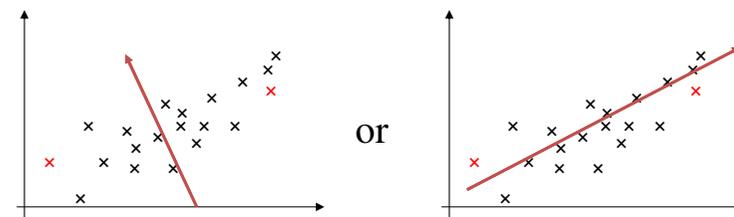
次元削減

高次元データを低次元データへ変換する方法

- 本日の講義内容は可視化を目的とした2次元への次元削減を主に説明するが、3次元以上の空間への次元削減も可能

主成分分析は高次元のデータを低次元空間へ射影する次元削減の手法

- 下図の2次元から1次元への次元削減はどちらが良い？
⇒ 情報の損失を最小化する次元削減が望ましい



6

主成分分析の指標

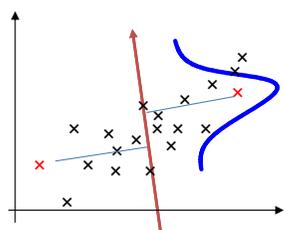
次元削減による情報損失は必至

- その上で、次元削減後もデータ間の違いが分かりやすい情報損失が少ない次元削減が望ましい

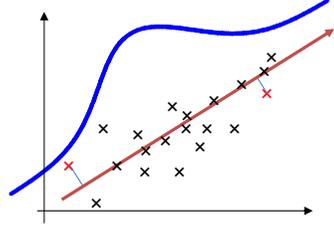
主成分分析での情報損失の測り方

⇒ 次元削減後の低次元空間でのデータの分散

低次元空間での分散が小さい



低次元空間での分散が大きい

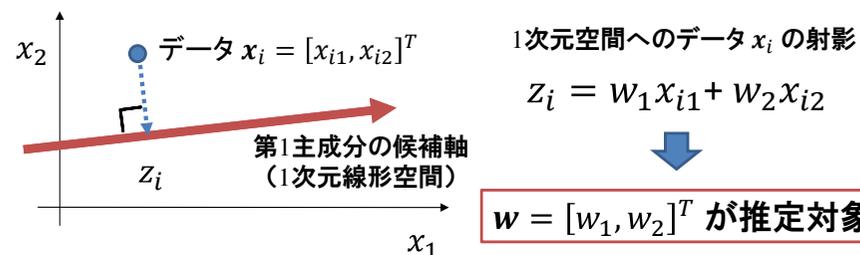


7

主成分の定義

第1主成分：次元削減後の分散を最大化する射影軸

- 第 i 主成分：第 $1, \dots, i-1$ 主成分と直交する空間で分散を最大化する射影軸
- z_i ：データ x_i を第1主成分の候補軸へ射影した点



1次元空間へのデータ x_i の射影
 $z_i = w_1 x_{i1} + w_2 x_{i2}$

$w = [w_1, w_2]^T$ が推定対象

第1主成分の求め方

- $\{z_1, \dots, z_N\}$ の分散を最大化する w_1, w_2 を推定

8

第1主成分の求め方 #1

$\{z_1, \dots, z_N\}$ の標本分散を最大化する w を決める

Notation

- 第1主成分の重み: $w = [w_1, \dots, w_P]^T$
- データ i の第1主成分: $z_i = w^T x_i$
- 変数 p の標本平均: $\bar{x}_p = \frac{1}{N} \sum_{i=1}^N x_{ip}$
- データの標本分散共分散行列: $S_x = \begin{bmatrix} S_{11} & \cdots & S_{1P} \\ \vdots & \ddots & \vdots \\ S_{P1} & \cdots & S_{PP} \end{bmatrix}$
- i 番目と j 番目の変数の標本共分散: S_{ij}

9

第1主成分の求め方 #2

$\{z_1, \dots, z_N\}$ の標本分散の最大化

- $\{z_1, \dots, z_N\}$ の標本平均 \bar{z} と標本分散 S_z check!

$$\bar{z} = \frac{1}{N} \sum_{i=1}^N z_i = \sum_{p=1}^P w_p \bar{x}_p$$

$$S_z = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})^2 = w^T S_x w$$

制約条件: $\|w\|^2 = w^T w = 1$

- w の要素を大きくすると S_z も大きくなっていくために制約条件を導入。知りたいのは射影軸の方向のみ

10

第1主成分の求め方 #3

制約条件付き2次形式の最大化問題

- 目的関数: $w^T S_x w$,
- 制約条件: $w^T w = 1$
- この最適化は固有値分解: $S_x w = \lambda w$ に帰結
 - 固有値: $\lambda_1 > \lambda_2 > \dots > \lambda_P \geq 0$
 - 対応する固有ベクトル: $\hat{w}_1, \hat{w}_2, \dots, \hat{w}_P$
 - ⇒ 最大固有値に対応する固有ベクトルが分散 S_z を最大化
- 最適解: $w^* = \hat{w}_1 = [\hat{w}_{11}, \dots, \hat{w}_{1P}]^T$

前回の講義内容(分散共分散行列の固有値分解)を思い出して、この結果をイメージしてみよう

11

第2以降の主成分以降の求め方

求め方: $\begin{cases} \text{目的関数: } w^T S_x w \\ \text{制約条件: } w^T w = 1 \text{ and } w^T \hat{w}_1 = 0 \end{cases}$

第2主成分の定義は第1主成分と直交する射影軸

- 実は、最適解は前頁で求めた $w_2^* = \hat{w}_2 = [\hat{w}_{21}, \dots, \hat{w}_{2P}]^T$

累積寄与率: $\frac{\lambda_1 + \dots + \lambda_M}{\lambda_1 + \lambda_2 + \dots + \lambda_P} \quad (M < P)$

- 多変量データがもつ情報(分散)を、第 M 主成分まででどのくらい保持できているかの指標

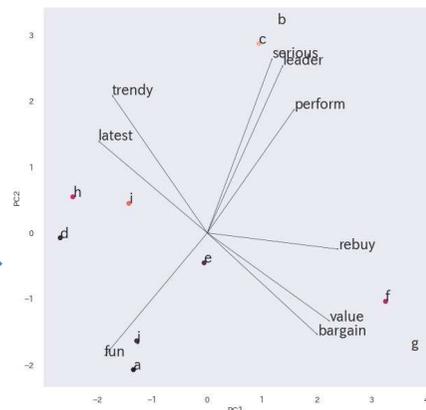
12

主成分分析を用いた知覚マップ #1

例: Consumer Brand Rating Data (simulated data)

- 10種のコーヒーブランドの模擬調査データ
 - 100人が各ブランドについて9の観点から評価(1点~10点)
- ブランド: $a \sim j$
 観点: perform, leader, latest, fun, serious, bargain, value, trendy, rebuy

ID	perform	leader	...	rebuy
a	x_{11}	x_{12}	...	x_{1P}
b	x_{21}	x_{22}	...	x_{2P}
\vdots	\vdots	\vdots	\ddots	\vdots
j	x_{j1}	x_{j2}	...	x_{jP}



C. Chapman, E.M. Feit, "R for Marketing Research and Analytics", Springer 2015

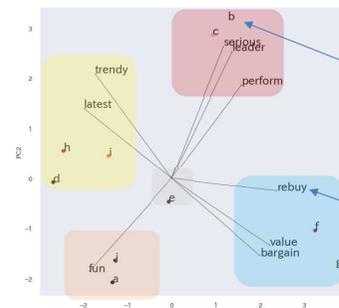
主成分分析を用いた知覚マップ #2

例: Consumer Brand Rating Data (simulated data)

- ブランド i の第1主成分(主成分得点): $z_i^{(1)} = \hat{w}_1^T x_i$
- ブランド i の第2主成分(主成分得点): $z_i^{(2)} = \hat{w}_2^T x_i$
- 変数 j と第1&第2主成分の相関係数(因子負荷量): $(\sqrt{\lambda_1} \hat{w}_{1j}, \sqrt{\lambda_2} \hat{w}_{2j})$

ブランド i の
 圧縮次元上の座標
 $(z_i^{(1)}, z_i^{(2)})$

変数 j の圧縮次元上
 の因子負荷量
 $(\sqrt{\lambda_1} \hat{w}_{1j}, \sqrt{\lambda_2} \hat{w}_{2j})$



ブランド b の
 圧縮次元上の座標
 $(z_b^{(1)}, z_b^{(2)})$

変数 "rebuy" の
 圧縮次元上の因子負荷量
 $(\sqrt{\lambda_1} \hat{w}_{1, rebuy}, \sqrt{\lambda_2} \hat{w}_{2, rebuy})$

C. Chapman, E.M. Feit, "R for Marketing Research and Analytics", Springer 2015

演習問題

次の主成分分析の結果から読み取れる各ブランドの特徴についてまとめなさい。また、ブランド a から j の内一つを自由に選択し、そのブランドが取り得る戦略について議論しなさい。ただし、各ブランドを販売する企業の社会的意義、経営戦略、経営状況、強みや弱みなどの状況は自由に設定してよい

