

主成分分析

基本的な考え方

村田 昇

2020.11.03

講義の予定

- 第1日: 主成分分析の考え方
- 第2日: 分析の評価と視覚化

主成分分析の考え方

主成分分析

- PCA (Principal Component Analysis)
- 多数の変量のもつ情報の分析・視覚化:
 - 変量を効率的に縮約して少数の特徴量を構成する
 - 特徴量に關与する変量間の關係を明らかにする

分析の枠組み

- X_1, \dots, X_p : 変数
- Z_1, \dots, Z_d : 特徴量 ($d \leq p$)
- 変数と特徴量の關係: (線形結合)

$$Z_k = a_{1k}X_1 + \dots + a_{pk}X_p \quad (k = 1, \dots, d)$$

- 特徴量は定数倍の任意性があるので以下を仮定:

$$\|a_k\|^2 = \sum_{j=1}^p a_{jk}^2 = 1$$

主成分分析の用語

- 特徴量 Z_k :
第 k 主成分得点 (principal component score)
または
第 k 主成分
- 係数ベクトル a_k :
第 k 主成分負荷量 (principal component loading)
または
第 k 主成分方向 (principal component direction)

分析の目的

- 目的:
主成分得点 Z_1, \dots, Z_d が変数 X_1, \dots, X_p の情報を効率よく反映するように主成分負荷量 $\mathbf{a}_1, \dots, \mathbf{a}_d$ を観測データから **うまく** 決定する
- 分析の方針: (以下は同値)
 - データの情報を最も保持する変量の **線形結合を構成**
 - データの情報を最も反映する **座標軸を探索**
- 教師なし学習 の代表的手法の 1 つ:
 - 次元縮約: 入力をできるだけ少ない変数で表現
 - 特徴抽出: 情報処理に重要な特性を変数に凝集

第1主成分の計算

記号の準備

- 変数: x_1, \dots, x_p (p 次元)
- 観測データ: n 個の (x_1, \dots, x_p) の組

$$\{(x_{i1}, \dots, x_{ip})\}_{i=1}^n$$

- $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$:
 i 番目の観測データ (p 次元空間内の 1 点)
- $\mathbf{a} = (a_1, \dots, a_p)^\top$:
長さ 1 の p 次元ベクトル

係数ベクトルによる射影

- データ \mathbf{x}_i の \mathbf{a} 方向成分の長さ:

$$\mathbf{a}^\top \mathbf{x}_i \quad (\text{スカラー})$$

- 方向ベクトル \mathbf{a} をもつ直線上への点 \mathbf{x}_i の直交射影

$$(\mathbf{a}^\top \mathbf{x}_i) \mathbf{a} \quad (\text{スカラー} \times \text{ベクトル})$$

幾何学的描像

ベクトル \mathbf{a} の選択の指針

- 線形結合での見方
ベクトル \mathbf{a} を **うまく** 選んで観測データ $\mathbf{x}_1, \dots, \mathbf{x}_n$ の情報を最も保持する 1 変量データを構成:

$$\mathbf{a}^\top \mathbf{x}_1, \mathbf{a}^\top \mathbf{x}_2, \dots, \mathbf{a}^\top \mathbf{x}_n$$

- 座標軸での見方
観測データの **ばらつき** を最も反映するベクトル \mathbf{a} を選択:

$$\arg \max_{\mathbf{a}} \sum_{i=1}^n (\mathbf{a}^\top \mathbf{x}_i - \mathbf{a}^\top \bar{\mathbf{x}})^2, \quad \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i,$$

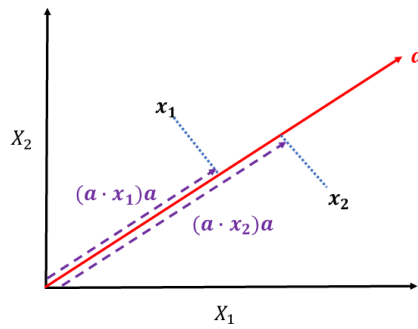


図 1: 観測データの直交射影 ($p = 2, n = 2$ の場合)

ベクトル a の最適化

- 最適化問題

制約条件 $\|a\| = 1$ の下で以下の関数を最大化せよ:

$$f(a) = \sum_{i=1}^n (a^T x_i - a^T \bar{x})^2$$

- この最大化問題は必ず解をもつ:
 - $f(a)$ は連続関数
 - 集合 $\{a \in \mathbb{R}^p : \|a\| = 1\}$ はコンパクト (有界閉集合)

演習

問題

- 以下の問に答えなさい.
 - 評価関数 $f(a)$ を以下の中心化したデータ行列で表しなさい.

$$X = \begin{pmatrix} x_1^T - \bar{x}^T \\ \vdots \\ x_n^T - \bar{x}^T \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 & \cdots & x_{1p} - \bar{x}_p \\ \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & \cdots & x_{np} - \bar{x}_p \end{pmatrix}$$

- 上の結果を用いて次の最適化問題の解の条件を求めなさい.

$$\text{maximize } f(a) \quad \text{s.t.} \quad a^T a = 1$$

第1主成分の解

ベクトル a の解

- 最適化問題

$$\text{maximize } f(a) = a^T X^T X a \quad \text{s.t.} \quad a^T a = 1$$

- 固有値問題

$f(a)$ の極大値を与える a は $X^T X$ の固有ベクトルとなる

$$X^T X a = \lambda a$$

第1主成分

- 求める \mathbf{a} は行列 $X^T X$ の最大固有ベクトル (長さ 1)
- このとき $f(\mathbf{a})$ は行列 $X^T X$ の最大固有値

$$f(\mathbf{a}) = \mathbf{a}^T X^T X \mathbf{a} = \mathbf{a}^T \lambda \mathbf{a} = \lambda$$

- 第1主成分負荷量: ベクトル \mathbf{a}
- 第1主成分得点:

$$z_{i1} = a_1 x_{i1} + \cdots + a_p x_{ip} \quad (i = 1, \dots, n)$$

Gram 行列の性質

Gram 行列の固有値

- $X^T X$ は非負定値対称行列
- $X^T X$ の固有値は 0 以上の実数
 - 固有値を重複を許して降順に並べる

$$\lambda_1 \geq \cdots \geq \lambda_p \quad (\geq 0)$$

- 固有値 λ_j に対する固有ベクトルを \mathbf{a}_j (長さ 1) とする

$$\|\mathbf{a}_j\| = 1 \quad (j = 1, \dots, p)$$

Gram 行列のスペクトル分解

- $\mathbf{a}_1, \dots, \mathbf{a}_p$ は互いに直交 するようとることができる

$$j \neq k \quad \Rightarrow \quad \mathbf{a}_j^T \mathbf{a}_k = 0$$

- 行列 $X^T X$ (非負値正定対称行列) のスペクトル分解:

$$\begin{aligned} X^T X &= \lambda_1 \mathbf{a}_1 \mathbf{a}_1^T + \lambda_2 \mathbf{a}_2 \mathbf{a}_2^T + \cdots + \lambda_p \mathbf{a}_p \mathbf{a}_p^T \\ &= \sum_{k=1}^p \lambda_k \mathbf{a}_k \mathbf{a}_k^T \end{aligned}$$

固有値と固有ベクトルによる行列の表現

演習

問題

- 以下の問に答えなさい.
 - Gram 行列のスペクトル分解において λ_j と \mathbf{a}_j が固有値・固有ベクトルとなることを確かめなさい.

$$X^T X = \sum_{k=1}^p \lambda_k \mathbf{a}_k \mathbf{a}_k^T$$

- 以下の行列を用いて Gram 行列のスペクトル分解を書き直しなさい.

$$A = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_p^T \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & \lambda_p \end{pmatrix}$$

第2主成分以降の計算

第2主成分の考え方

- 第1主成分:
 - 主成分負荷量: ベクトル \mathbf{a}_1
 - 主成分得点: $\mathbf{a}_1^T \mathbf{x}_i$ ($i = 1, \dots, n$)
- 第1主成分負荷量に関してデータが有する情報:

$$(\mathbf{a}_1^T \mathbf{x}_i) \mathbf{a}_1 \quad (i = 1, \dots, n)$$

- 第1主成分を取り除いた観測データ: (分析対象)

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i - (\mathbf{a}_1^T \mathbf{x}_i) \mathbf{a}_1 \quad (i = 1, \dots, n)$$

第2主成分の最適化

- 最適化問題
制約条件 $\|\mathbf{a}\| = 1$ の下で以下の関数を最大化せよ:

$$\tilde{f}(\mathbf{a}) = \sum_{i=1}^n (\mathbf{a}^T \tilde{\mathbf{x}}_i - \mathbf{a}^T \bar{\tilde{\mathbf{x}}})^2 \quad \text{ただし} \quad \bar{\tilde{\mathbf{x}}} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i$$

演習

問題

- 以下の問に答えなさい.
 - 以下の中心化したデータ行列を X と \mathbf{a}_1 で表しなさい.

$$\tilde{X} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top - \bar{\mathbf{x}}^\top \\ \vdots \\ \tilde{\mathbf{x}}_n^\top - \bar{\mathbf{x}}^\top \end{pmatrix}$$

– 上の結果を用いて次の最適化問題の解を求めなさい.

$$\text{maximize } \tilde{f}(\mathbf{a}) \quad \text{s.t. } \mathbf{a}^\top \mathbf{a} = 1$$

第2主成分の解

第2主成分

- Gram 行列 $\tilde{X}^\top \tilde{X}$ の固有ベクトル \mathbf{a}_1 の固有値は 0
- Gram 行列 $\tilde{X}^\top \tilde{X}$ の最大固有値は λ_2
- 解は第2固有値 λ_2 に対応する固有ベクトル \mathbf{a}_2
- 以下同様に第 k 主成分負荷量は $X^\top X$ の第 k 固有値 λ_k に対応する固有ベクトル \mathbf{a}_k

解析の事例

データセットについて

- 総務省統計局より取得した都道府県別の社会生活統計指標の一部
 - 総務省 <https://www.e-stat.go.jp/SG1/estat/List.do?bid=000001083999&cyclo=0>
 - データ https://noboru-murata.github.io/multivariate-analysis/data/japan_social.csv
- * Pref: 都道府県名
- * Forest: 森林面積割合 (%) 2014 年
- * Agri: 就業者 1 人当たり農業産出額 (販売農家) (万円) 2014 年
- * Ratio: 全国総人口に占める人口割合 (%) 2015 年
- * Land: 土地生産性 (耕地面積 1 ヘクタール当たり) (万円) 2014 年
- * Goods: 商業年間商品販売額 [卸売業 + 小売業] (事業所当たり) (百万円) 2013 年

社会生活統計指標の分析

- データ (の一部) の内容

	Forest	Agri	Ratio	Land	Goods
Hokkaido	67.9	1150.6	4.23	96.8	283.3
Aomori	63.8	444.7	1.03	186.0	183.0
Iwate	74.9	334.3	1.01	155.2	179.4
Miyagi	55.9	299.9	1.84	125.3	365.9
Akita	70.5	268.7	0.81	98.5	153.3
Yamagata	68.7	396.3	0.88	174.1	157.5
Fukushima	67.9	236.4	1.51	127.1	184.5
Ibaraki	31.0	479.0	2.30	249.1	204.9
Tochigi	53.2	402.6	1.55	199.6	204.3
Gumma	63.8	530.6	1.55	321.6	270.0
Saitama	31.9	324.7	5.72	247.0	244.7
Chiba	30.4	565.5	4.90	326.1	219.7

Tokyo	34.8	268.5	10.63	404.7	1062.6
Kanagawa	38.8	322.8	7.18	396.4	246.1
Niigata	63.5	308.6	1.81	141.9	205.5
Toyama	56.6	276.1	0.84	98.5	192.4
Ishikawa	66.0	271.3	0.91	112.0	222.9
Fukui	73.9	216.1	0.62	98.5	167.3

- データの散布図

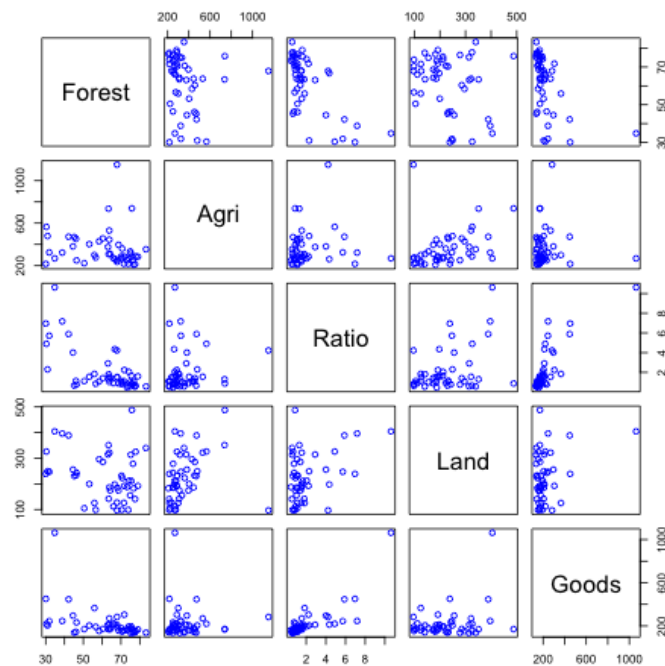


図 2: 散布図

- データの箱ひげ図
- 主成分負荷量を計算 (正規化後)

	PC1	PC2	PC3	PC4	PC5
Forest	-0.4871498	0.1045813	-0.45748795	0.6859649	-0.26815060
Agri	0.1339190	0.8115056	0.47912767	0.3045447	0.03483694
Ratio	0.5851294	-0.1511042	0.04467249	0.1640953	-0.77837539
Land	0.3547649	0.4851374	-0.74167904	-0.2897485	0.06885892
Goods	0.5258481	-0.2689436	-0.09517368	0.5708093	0.56238052

- 主成分方向から読み取れること:
 - 第 1: 人の多さに関する成分 (正の向きほど人が多い)
 - 第 2: 農業生産力に関する成分 (正の向きほど高い)
- 主成分得点の表示

次週の予定

- 第 1 日: 主成分分析の考え方
- 第 2 日: 分析の評価と視覚化

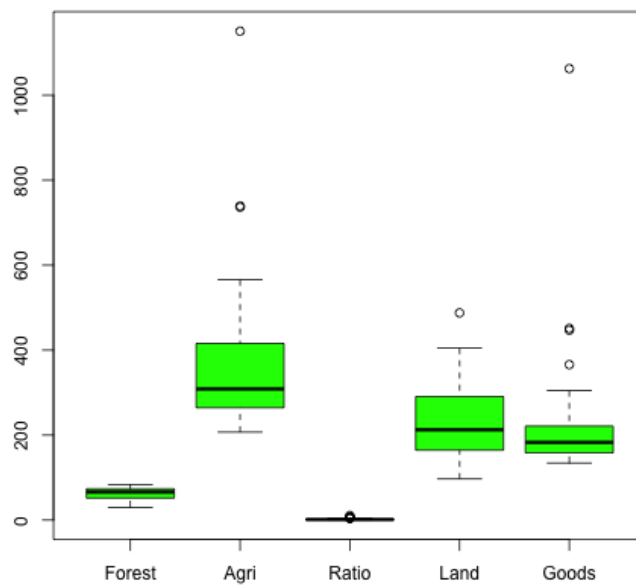


図 3: 箱ひげ図

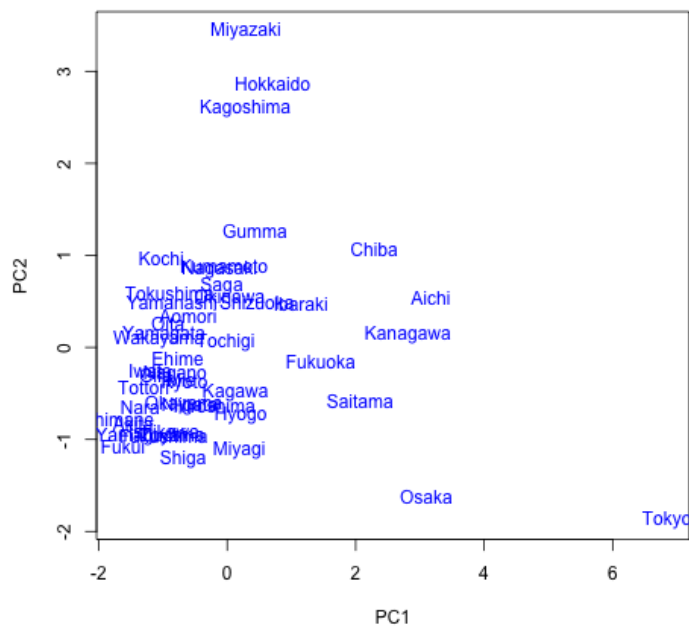


図 4: 主成分得点による散布図