

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:

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

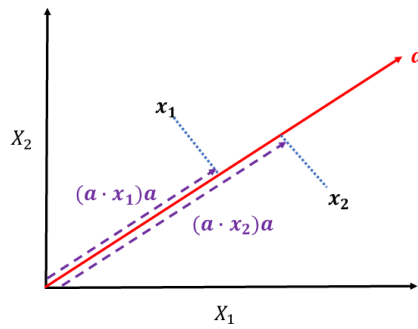


図 1: 観測データの直交射影 ($p = 2, n = 2$ の場合)

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)

第 1 主成分の解

ベクトル a の解

- 最適化問題

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

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

$$f(a) = a^T X^T X a = a^T \lambda a = \lambda$$

- 第 1 主成分負荷量: ベクトル 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

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

- 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

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:

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

第2主成分の解

第2主成分

- Gram 行列 $\tilde{X}^T \tilde{X}$ の固有ベクトル \mathbf{a}_1 の固有値は 0
- Gram 行列 $\tilde{X}^T \tilde{X}$ 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 $X^T X$ の第 k 固有値 λ_k に対応する固有ベクトル \mathbf{a}_k