

クレジット:

UTokyo Online Education 統計データ解析Ⅱ 2018 小池祐太

ライセンス:

利用者は、本講義資料を、教育的な目的に限ってページ単位で利用することができます。特に記載のない限り、本講義資料はページ単位でクリエイティブ・コモンズ 表示-非営利-改変禁止 ライセンスの下に提供されています。

<http://creativecommons.org/licenses/by-nc-nd/4.0/>

本講義資料内には、東京大学が第三者より許諾を得て利用している画像等や、各種ライセンスによって提供されている画像等が含まれています。個々の画像等を本講義資料から切り離して利用することはできません。個々の画像等の利用については、それぞれの権利者の定めるところに従ってください。



統計データ解析 (II) 第 10 回

小池祐太

2018 年 6 月 14 日

1 重回帰分析: 交互作用モデル・変数の非線形変換

2 主成分分析

- 目的
- Rでの実行
- 計算法
- 分析の評価
 - 寄与率
 - バイプロット
- 補足: 計算の詳細の説明
 - 主成分の計算 ($d = 1$ の場合)
 - 主成分の計算 ($d = 2$ の場合)
 - 寄与率の計算

重回帰分析: 交互作用モデル・変数の非線形変換

- 目的変数を Y , 説明変数を X_1, \dots, X_p で表す
- 前述のように, 説明変数として $X_j X_k$ (交差項と呼ばれる) や $\log X_j$ といったものを新たに加えることで, 目的変数 Y と説明変数 X_1, \dots, X_p の非線形な関係をモデル化することができる
- R においてはこのような回帰式の柔軟なモデル化を可能にするためのモデル式の記述法が実装されている
- 実行例 `cross.r`

主成分分析

- **主成分分析 (principal component analysis, PCA)**

- ▶ 多数の変数/データが与えられたとき、変数/データたちのもつ情報を効率的に縮約して少数の特徴量を構成することで、変数/データ間の関係を明らかにするための分析法
- ▶ イメージ: 株価指数. ただし、通常の株価指数とは異なり、純粋にデータの情報のみを使って指数を作成する
- ▶ 「教師なし学習」の代表的手法の1つ

- 主成分分析では、特徴量は変数/データたちの線形結合として構成
- 幾何学的にいうと、観測データの含まれる p 次元空間にうまく座標軸を設定することにより、その座標軸上にデータのもつ情報が最大限反映されるようにすることが目的

- 数式で表すと, 与えられた変数を X_1, \dots, X_p としたとき, d を p 以下の正の整数とし (通常 p より小さくとる), X_1, \dots, X_p の線形結合として表される d 個の変数

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

を, もとの変数のもつ情報を最大限保持しつつ適切に構成することが目的

- ここで, Z_k と Z_k の (0 でない) 定数倍は互いに同じ情報量をもつので, そのような定数倍の任意性をなくすため, ベクトル $\mathbf{a}_k := (a_{1k}, \dots, a_{pk})^\top$ の長さが 1 となるようにする. すなわち,

$$\|\mathbf{a}_k\|^2 := \sum_{j=1}^p a_{jk}^2 = 1$$

と仮定

主成分分析

- k 番目の特徴量 Z_k を第 k **主成分 (得点) (principal component (score))** と呼ぶ
- Z_k の係数ベクトル \mathbf{a}_k を第 k **主成分方向 (principal component direction)** または第 k **主成分負荷量 (principal component loading)** と呼ぶ
- 主成分分析の目的
 - ▶ 主成分 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 X_1, \dots, X_p たちの観測データから「うまく」決定する

主成分分析: Rでの実行

- Rには主成分分析を実行するための関数として、関数 `prcomp()` および関数 `princomp()` が用意されている
- 前者と後者には計算法に若干の違いがあり、一般には前者の方が数値計算の観点からみると優れている(後者はS言語との互換性を重視した実装となっている)
- 従って以下では関数 `prcomp()` を利用する
- 実行例 `pca-simulate.r`

主成分分析: 計算法

- まず $d = 1$ の場合を考える
- 組 (X_1, \dots, X_p) に対する n 個の観測データ $\{(x_{i1}, \dots, x_{ip})\}_{i=1}^n$ が与えられているとする
- i 番目の観測データに対応する p 次元ベクトルを $\mathbf{x}_i := (x_{i1}, \dots, x_{ip})^\top$ とする
- **目的** 長さ 1 の p 次元ベクトル $\mathbf{a} = (a_1, \dots, a_p)^\top$ を「うまく」選んで、観測データ $\mathbf{x}_1, \dots, \mathbf{x}_n$ のもつ情報を最大限保持するように 1 変量データ $\mathbf{a} \cdot \mathbf{x}_1, \dots, \mathbf{a} \cdot \mathbf{x}_n$ を構成する
 - ▶ 各 \mathbf{x}_i は p 次元空間内の点だとみなせるが、このとき $(\mathbf{a} \cdot \mathbf{x}_i)\mathbf{a}$ はベクトル \mathbf{a} で張られる部分空間 (直線) への点 \mathbf{x}_i の直交射影に一致する
 - ▶ すなわち、 $\mathbf{a} \cdot \mathbf{x}_i$ はベクトル \mathbf{x}_i の \mathbf{a} -方向成分であると解釈できる

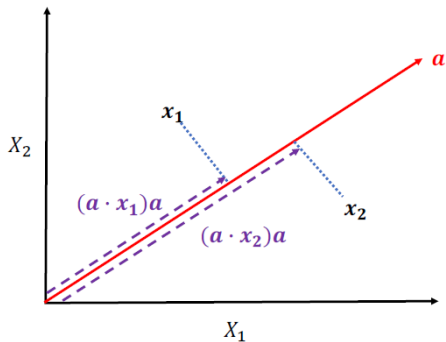


図 1: ベクトル a への観測データの直交射影 ($p = 2, n = 2$ の場合)

- このような幾何学的考察から、ベクトル \mathbf{a} の適切な選び方の指針としては以下の考え方が自然:
 - ▶ 構成した特徴量がもとのデータのばらつきを最大限反映するように、 $\mathbf{x}_1, \dots, \mathbf{x}_n$ たちのばらつきが最も大きい方向 \mathbf{a} を選ぶ。すなわち、

$$\sum_{i=1}^n (\mathbf{a} \cdot \mathbf{x}_i - \mathbf{a} \cdot \bar{\mathbf{x}})^2$$

を最大化するように \mathbf{a} を選ぶ。ここに、

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

であり、従って $\mathbf{a} \cdot \bar{\mathbf{x}}$ は $\mathbf{a} \cdot \mathbf{x}_1, \dots, \mathbf{a} \cdot \mathbf{x}_n$ の平均に対応する。

主成分分析: 計算法

- 以上をまとめると, 関数

$$f(\mathbf{a}) = \sum_{i=1}^n (\mathbf{a} \cdot \mathbf{x}_i - \mathbf{a} \cdot \bar{\mathbf{x}})^2$$

を制約条件 $\|\mathbf{a}\| = 1$ の下で最大化するように, ベクトル \mathbf{a} を選ぶのが方針となる

主成分分析: 計算法

- そのようなベクトル \mathbf{a} は存在して次の性質をもつ (詳細な導出は配布資料参照. 時間があれば後で説明):
 - ▶ $f(\mathbf{a})$ は行列 $\mathbf{X}^T \mathbf{X}$ の固有値であり, \mathbf{a} はこの固有値に対する固有ベクトルである.

ただし, $n \times p$ 行列 \mathbf{X} を

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

で定義する ($\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$)

主成分分析: 計算法

- 従って, 求めるべき \mathbf{a} は行列 $\mathbf{X}^T \mathbf{X}$ の最大固有値に対する固有ベクトルで長さ 1 のものであり, このとき $f(\mathbf{a})$ はその最大固有値に一致する
- このようにして求めたベクトル \mathbf{a} が **第 1 主成分負荷量** となる. また, 観測データから計算される **第 1 主成分** たちは

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

となる

- 実行例 `pca-simulate.r`

主成分分析: 計算法

- 次に $d \geq 2$ の場合を考える
- まず記号を準備する
 - ▶ $\mathbf{X}^T \mathbf{X}$ は非負定値対称行列だから, その固有値はすべて 0 以上の実数である
 - ▶ そこで, $\mathbf{X}^T \mathbf{X}$ の固有値を (重複を許して) 降順に並べて $\lambda_1 \geq \dots \geq \lambda_p (\geq 0)$ と書くことにする
 - ▶ さらに, 各 $j = 1, \dots, p$ について \mathbf{a}_j を λ_j に対する固有ベクトルとする
 - ▶ このとき, $\mathbf{a}_1, \dots, \mathbf{a}_p$ はそれぞれ長さ 1 かつ互いに直交するようにとることができる
 - ▶ すなわち,

$$\|\mathbf{a}_j\| = 1 \quad (j = 1, \dots, p), \quad j \neq k \Rightarrow \mathbf{a}_j \cdot \mathbf{a}_k = 0 \quad (1)$$

が成り立つと仮定してよい

主成分分析: 計算法

- 1つめの特徴量は前節で構成した第1主成分を用いる
- 前節の議論より, ベクトル \mathbf{a}_1 は第1主成分方向に対応する
- 第1主成分方向に関してデータが有する情報はベクトル $(\mathbf{a}_1 \cdot \mathbf{x}_i)\mathbf{a}_1$ ($i = 1, \dots, n$) にすべて縮約されているので, 第1主成分 $\mathbf{a}_1 \cdot \mathbf{x}_i$ ($i = 1, \dots, n$) にすべて含まれている
- 従って, 2つめの特徴量を構成する指針として, 観測データから第1主成分方向の成分を取り除いたデータ

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

に対して, 前節と同様の考え方でこれらのデータたちのばらつきが最も大きい方向 \mathbf{a} を求めて, 特徴量 $\mathbf{a} \cdot \mathbf{x}_i$ ($i = 1, \dots, n$) を構成するのが自然である

主成分分析: 計算法

- すなわち,

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

を制約条件 $\|\mathbf{a}\| = 1$ の下で最大化するような \mathbf{a} を選べばよい

- $\Rightarrow \mathbf{X}^T \mathbf{X}$ の 2 番目に大きい固有値 λ_2 に対応する固有ベクトル \mathbf{a}_2 がそのようなベクトル \mathbf{a} を与えることを示すことができる
 - ▶ 詳細な計算は配布資料参照. 時間があれば後で説明
- 同様の議論を繰り返すことによって, 第 k 主成分負荷量は $\mathbf{X}^T \mathbf{X}$ の k 番目に大きい固有値 λ_k に対応する固有ベクトル \mathbf{a}_k に取ればよいことがわかる
- 実行例 `pca.r`

分析の評価: 寄与率

- 構成した主成分が元のデータがもっていた情報をどの程度保持しているかを評価する方法の1つとして、回帰分析の場合と同様に、その主成分のばらつき (分散) がもとのデータのばらつき (分散) をどの程度説明できているかを評価する方法がある
- 第 k 主成分に対しては、これは

$$\frac{\frac{1}{n-1} \sum_{i=1}^n (\mathbf{a}_k \cdot \mathbf{x}_i - \mathbf{a}_k \cdot \bar{\mathbf{x}})^2}{\frac{1}{n-1} \sum_{i=1}^n \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}$$

を計算することに相当する

分析の評価: 寄与率

- この量を第 k 主成分の**寄与率 (proportion of variance)** と呼ぶ
- 計算すると,

$$\sum_{i=1}^n (\mathbf{a}_k \cdot \mathbf{x}_i - \mathbf{a}_k \cdot \bar{\mathbf{x}})^2 = \lambda_k, \quad \sum_{i=1}^n \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 = \sum_{l=1}^p \lambda_l$$

が成り立つことがわかる (詳細な計算は配布資料参照. 時間があれば後で説明)

- 以上より, 第 k 主成分の寄与率は,

$$\frac{\lambda_k}{\sum_{l=1}^p \lambda_l}$$

で計算される

分析の評価: 寄与率

- 第1主成分から第 k 主成分までを特徴量として用いた際に説明できるデータのばらつき割合は第 k 主成分までの**累積寄与率 (cumulative proportion)** と呼ばれ, 第1主成分の寄与率から第 k 主成分の寄与率までの総和として計算される:

$$\frac{\sum_{l=1}^k \lambda_l}{\sum_{l=1}^p \lambda_l}.$$

- 累積寄与率はいくつの主成分を用いるべきかの基準として用いられる. 一般に, 累積寄与率が80%程度の主成分を使って分析を行うことが多い
- Rでは, 寄与率および累積寄与率は, 関数 `prcomp()` のアウトプットに関数 `summary()` を適用することで計算できる
- 実行例 `pca-summary.r`

分析の評価: バイプロット

- 関連がある2枚の散布図を1つの画面に表示する散布図を**バイプロット (biplot)** という
- 主成分分析では、得られた主成分の意味を解釈するために、主成分方向の散布図と主成分の散布図を対応づけて分析を進める場合が多い
- より具体的には、2つの主成分方向 $\mathbf{a}_k = (a_{1k}, \dots, a_{pk})^\top$ と $\mathbf{a}_l = (a_{1l}, \dots, a_{pl})^\top$ に着目する
 - ▶ 主成分方向の散布図とは点 $\{(a_{jk}, a_{jl})\}_{j=1}^p$ の散布図であり、各変数が着目した主成分方向にどれだけの変動成分をもつかを図示している (主成分分析ではこれらの点を対応するベクトルで描画することが多い)
 - ▶ 一方で、主成分の散布図とは点 $\{(\mathbf{a}_k \cdot (\mathbf{x}_i - \bar{\mathbf{x}}), \mathbf{a}_l \cdot (\mathbf{x}_i - \bar{\mathbf{x}}))\}_{i=1}^n$ の散布図であり、各サンプルが着目した主成分方向にどれだけの変動成分をもつかを図示している

分析の評価: バイプロット

- Rでは, 関数 `prcomp()` のアウトプットに関数 `biplot()` を適用することでバイプロットを描画できる
- 実行例 `biplot2.r`

主成分の計算 ($d = 1$ の場合)

- 第 1 主成分負荷量は, 関数

$$f(\mathbf{a}) = \sum_{i=1}^n (\mathbf{a} \cdot \mathbf{x}_i - \mathbf{a} \cdot \bar{\mathbf{x}})^2$$

を制約条件 $\|\mathbf{a}\| = 1$ の下で最大化するようなベクトル \mathbf{a} として選ぶのであった

- $f(\mathbf{a})$ は連続関数であり, また集合 $\{\mathbf{a} \in \mathbb{R}^p : \|\mathbf{a}\| = 1\}$ はコンパクト (有界閉集合) であるから, この最大化問題は解をもつ

- 更に, Lagrange の乗数法から, 求めるべき解は, Lagrange 関数

$$L(\mathbf{a}, \lambda) = f(\mathbf{a}) + \lambda(1 - \|\mathbf{a}\|^2)$$

の勾配を 0 にするベクトルである

- いま,

$$f(\mathbf{a}) = \sum_{i=1}^n \left(\sum_{j=1}^p a_j (x_{ij} - \bar{x}_j) \right)^2$$

と書けるから, 各 $j = 1, \dots, p$ について

$$\begin{aligned} \frac{\partial L}{\partial a_j}(\mathbf{a}, \lambda) &= 2 \sum_{i=1}^n \left(\sum_{k=1}^p a_k (x_{ik} - \bar{x}_k) \right) (x_{ij} - \bar{x}_j) - 2\lambda a_j \\ &= 2 \sum_{k=1}^p \left(\sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \right) a_k - 2\lambda a_j \end{aligned} \quad (2)$$

が成り立つ

- 従って, (2) の右辺第 1 項は p 次元ベクトル $\mathbf{X}^\top \mathbf{X} \mathbf{a}$ の第 j 成分に等しいことがわかる
- 以上より, 求めるべき \mathbf{a} は, 方程式

$$\mathbf{X}^\top \mathbf{X} \mathbf{a} = \lambda \mathbf{a} \quad (3)$$

の解となることがわかる

- 特に, λ は p 次正方形行列 $\mathbf{X}^\top \mathbf{X}$ の固有値であり, \mathbf{a} は λ に対する固有ベクトルとなる
- また, (3) の両辺に左から \mathbf{a}^\top をかけると,

$$\mathbf{a}^\top \mathbf{X}^\top \mathbf{X} \mathbf{a} = \lambda \mathbf{a}^\top \mathbf{a} = \lambda \|\mathbf{a}\|^2 = \lambda$$

を得る

主成分の計算 ($d = 1$ の場合)

- ここで,

$$\mathbf{X}\mathbf{a} = ((\mathbf{x}_1 - \bar{\mathbf{x}}) \cdot \mathbf{a}, \dots, (\mathbf{x}_n - \bar{\mathbf{x}}) \cdot \mathbf{a})^\top$$

であることに注意すれば,

$$\mathbf{a}^\top \mathbf{X}^\top \mathbf{X} \mathbf{a} = (\mathbf{X}\mathbf{a})^\top \mathbf{X}\mathbf{a} = \|\mathbf{X}\mathbf{a}\|^2 = f(\mathbf{a})$$

を得る

- このことから, $f(\mathbf{a})$ は行列 $\mathbf{X}^\top \mathbf{X}$ の固有値となる
- 要約すると, 次のようになる:
 - ▶ $f(\mathbf{a})$ は行列 $\mathbf{X}^\top \mathbf{X}$ の固有値であり, \mathbf{a} はこの固有値に対する固有ベクトルである.

主成分の計算 ($d = 2$ の場合)

- 第 2 主成分負荷量は、関数

$$\sum_{i=1}^n |\tilde{\mathbf{x}}_i - \bar{\mathbf{x}} - (\mathbf{a}_1 \cdot \bar{\mathbf{x}})\mathbf{a}_1|^2$$

を制約条件 $\|\mathbf{a}\| = 1$ の下で最大化するようなベクトル \mathbf{a} として選ぶのであった

- 前と同様に、行列 $\mathbf{X}_{(-1)}$ を

$$\mathbf{X}_{(-1)} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top - \bar{\mathbf{x}}^\top - (\mathbf{a}_1 \cdot \bar{\mathbf{x}})\mathbf{a}_1^\top \\ \vdots \\ \tilde{\mathbf{x}}_n^\top - \bar{\mathbf{x}}^\top - (\mathbf{a}_1 \cdot \bar{\mathbf{x}})\mathbf{a}_1^\top \end{pmatrix}$$

で定義すれば、 \mathbf{a} は行列 $\mathbf{X}_{(-1)}^\top \mathbf{X}_{(-1)}$ の最大固有値に対する固有ベクトルとなることがわかる

UTokyo Online Education 統計データ解析 II 2018 小池祐太 CC BY-NC-ND

主成分の計算 ($d = 2$ の場合)

- ここで, 行列 $\mathbf{X}_{(-1)}$ は以下のように書けることに注意する (E_p は p 次単位行列):

$$\mathbf{X}_{(-1)} = \mathbf{X} - \mathbf{X}\mathbf{a}_1\mathbf{a}_1^\top = \mathbf{X}(E_p - \mathbf{a}_1\mathbf{a}_1^\top).$$

- 従って,

$$\mathbf{X}_{(-1)}^\top \mathbf{X}_{(-1)} = (E_p - \mathbf{a}_1\mathbf{a}_1^\top) \mathbf{X}^\top \mathbf{X} (E_p - \mathbf{a}_1\mathbf{a}_1^\top)$$

であるから, 条件 (1) より

$$\mathbf{X}_{(-1)}^\top \mathbf{X}_{(-1)} \mathbf{a}_1 = \mathbf{0}, \quad \mathbf{X}_{(-1)}^\top \mathbf{X}_{(-1)} \mathbf{a}_j = \lambda_j \mathbf{a}_j \quad (j = 2, \dots, p)$$

が成り立つ

主成分の計算 ($d = 2$ の場合)

- これは $0, \lambda_2, \dots, \lambda_p$ が $\mathbf{X}_{(-1)}^\top \mathbf{X}_{(-1)}$ の固有値であり, $\mathbf{a}_1, \dots, \mathbf{a}_p$ がそれぞれに対する固有ベクトルであることを意味する
- 従って, $\mathbf{X}_{(-1)}^\top \mathbf{X}_{(-1)}$ の最大固有値は λ_2 であり, 求めるべきベクトルは \mathbf{a}_2 である

寄与率の計算

- $\mathbf{a}_k \cdot \mathbf{x}_i - \mathbf{a}_k \cdot \bar{\mathbf{x}}$ が (列) ベクトル $\mathbf{X}\mathbf{a}$ の第 i 成分に対応することと, 列ベクトル \mathbf{v} に対して $\|\mathbf{v}\|^2 = \mathbf{v}^\top \mathbf{v}$ が成り立つことに注意すれば,

$$\begin{aligned} \sum_{i=1}^n (\mathbf{a}_k \cdot \mathbf{x}_i - \mathbf{a}_k \cdot \bar{\mathbf{x}})^2 &= \|\mathbf{X}\mathbf{a}_k\|^2 = \mathbf{a}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{a}_k \\ &= \lambda_k \mathbf{a}_k^\top \mathbf{a}_k = \lambda_k \|\mathbf{a}_k\|^2 = \lambda_k \end{aligned}$$

を得る

寄与率の計算

- さらに, $\|\mathbf{x}_i - \bar{\mathbf{x}}\|^2$ が行列 $\mathbf{X}\mathbf{X}^\top$ の第 i 対角成分であることと, 直交行列 A を $A = (\mathbf{a}_1, \dots, \mathbf{a}_p)$ で定義すれば,

$$A^\top \mathbf{X}\mathbf{X}A = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix}$$

と対角化されることに注意すれば, 以下が成り立つ:

寄与率の計算

$$\begin{aligned}\sum_{i=1}^n \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 &= \text{tr}(\mathbf{X}\mathbf{X}^\top) = \text{tr}(\mathbf{X}^\top\mathbf{X}) = \text{tr}(\mathbf{X}^\top\mathbf{X}E_p) \\ &= \text{tr}(\mathbf{X}^\top\mathbf{X}AA^\top) = \text{tr}(A^\top\mathbf{X}^\top\mathbf{X}A) = \sum_{l=1}^p \lambda_l\end{aligned}$$

- ここで、積 BC および CB が定義されるような行列 B, C に対して $\text{tr}(BC) = \text{tr}(CB)$ が成り立つことを用いた