

## 応用統計学 2017 主成分分析に関する補足資料

2017年10月18日 (水)

清 智也 sei@mist.i.u-tokyo.ac.jp

<http://ur0.pw/yTzt>

先週の講義で主成分分析について少し混乱したので、整理しておく。

まず、データ行列を

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p) = \begin{pmatrix} \mathbf{x}^{(1)'} \\ \vdots \\ \mathbf{x}^{(n)'} \end{pmatrix} \in \mathbb{R}^{n \times p}$$

と記す。ただし各列はあらかじめ標準化されているものとする：

$$\bar{x}_i = \frac{1}{n} \mathbf{1}'_n \mathbf{x}_i = 0, \quad V(\mathbf{x}_i) = \frac{1}{n} (\mathbf{x}_i - \bar{x}_i \mathbf{1}_n)' (\mathbf{x}_i - \bar{x}_i \mathbf{1}_n) = 1, \quad i = 1, \dots, p.$$

先週の講義では、

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}' = \sum_{k=1}^p d_k \mathbf{u}_k \mathbf{v}_k'$$

と特異値分解し、右辺の第  $k$  項  $d_k \mathbf{u}_k \mathbf{v}_k'$  を第  $k$  主成分と呼んだ。しかし一般には変数  $\mathbf{x} \in \mathbb{R}^p$  に対して  $\mathbf{v}_k' \mathbf{x}$  のことを  $\mathbf{x}$  の第  $k$  主成分と呼ぶことが多いようである<sup>1</sup>。この用語の下で主成分に  $\mathbf{x} = \mathbf{x}^{(t)}$  を代入した値を**主成分得点**と呼ぶ。第  $k$  主成分の主成分得点は

$$\mathbf{v}_k' \mathbf{x}^{(t)} = (\mathbf{X} \mathbf{V})_{tk} = (\mathbf{U} \mathbf{D})_{tk} = d_k (\mathbf{u}_k)_t, \quad t = 1, \dots, n,$$

となる。これをすべての個体について並べてできるベクトルは  $d_k \mathbf{u}_k$  となる。

また、主成分得点と各変量の相関係数を**因子負荷量**<sup>2</sup>と呼ぶ。第  $k$  主成分の因子負荷量は

$$\begin{aligned} \text{Corr}(d_k \mathbf{u}_k, \mathbf{x}_i) &= \frac{\text{Cov}(d_k \mathbf{u}_k, \mathbf{x}_i)}{\sqrt{V(d_k \mathbf{u}_k) V(\mathbf{x}_i)}} \\ &= \frac{1}{\sqrt{n}} \mathbf{u}_k' \mathbf{x}_i \quad (\because \bar{u}_k = \bar{x}_i = 0, \mathbf{u}_k' \mathbf{u}_k = 1, V(\mathbf{x}_i) = 1) \\ &= \frac{1}{\sqrt{n}} d_k (\mathbf{v}_k)_i \end{aligned}$$

となる。これをすべての変量について並べてできるベクトルは

$$\frac{d_k}{\sqrt{n}} \mathbf{v}_k$$

となる。

<sup>1</sup>たとえば永田 靖, 棟近 雅彦「多変量解析法入門」, サイエンス社。

<sup>2</sup>主成分負荷量ともいう。

さて、主成分分析の説明では、共分散行列  $S = (1/n)X'X$  を出発点とすることが多い。すなわち、 $S$  の第  $k$  固有ベクトルを  $v_k$  とおくと、ベクトル  $x \in \mathbb{R}^p$  の第  $k$  主成分を  $v_k'x$  と定義するのである。この定義は上の定義と一致する。実際、特異値分解を用いると

$$S = \frac{1}{n}X'X = \frac{1}{n}VDU'UDV' = \frac{1}{n}VD^2V'.$$

となり、 $S$  のスペクトル分解が得られる。また  $S$  の固有値は  $\tilde{\lambda}_k := d_k^2/n$  と表され、 $\tilde{\lambda}_k$  を用いて因子負荷量を表すと

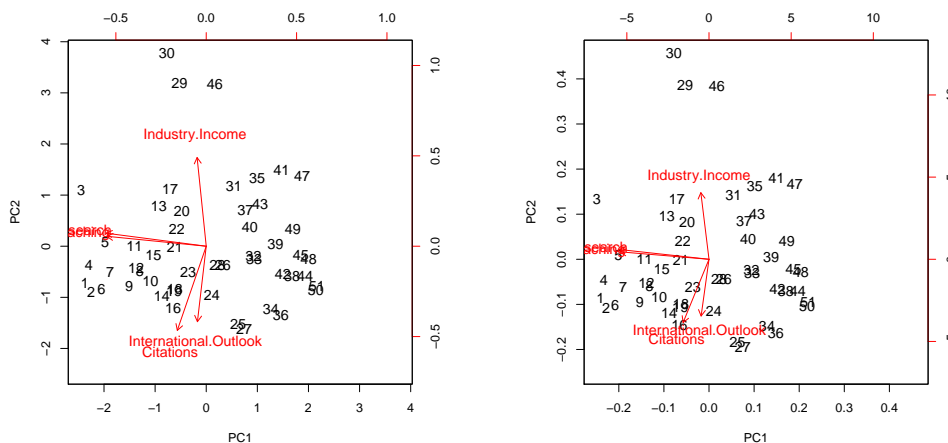
$$\frac{d_k}{\sqrt{n}}v_k = \sqrt{\tilde{\lambda}_k}v_k$$

となる。先週の講義で混乱したのは、この式の左辺と右辺を混同していたことが原因だった。

**バイプロットに関する補足** 先週の講義で説明したバイプロットは、第1, 第2主成分得点をプロットした散布図に、変数軸をベクトルとして描き加えたものであった (下図(a))。この方法は、幾何学的には個体空間  $\mathbb{R}^p$  において  $v_1, v_2$  の張る平面への直交射影を考えている。ただし**個体空間**とは各個体  $x^{(t)}$  が属す空間という意味である<sup>3</sup>。

これに対し、変量空間  $\mathbb{R}^n$  において  $u_1, u_2$  の張る平面への直交射影を考えることもできる (下図(b))。**変量空間**とは各変量  $x_i$  が属す空間という意味である。

後者の利点は、因子負荷量 (の定数倍) を図から読み取れることである。R 言語の関数 `biplot` では、後者がデフォルトで用いられる。また、これら以外の計量を考えることも可能である<sup>4</sup>。どの計量を選んでも、個体と変数の内積  $d_1u_{t1}v_{i1} + d_2u_{t2}v_{i2}$  は不変である。



(a) 個体空間におけるバイプロット  $(d_1u_1, d_2u_2)$  と  $(v_1, v_2)$   
`biplot(prcomp(X), scale=0)`  
 (b) 変量空間におけるバイプロット  $(u_1, u_2)$  と  $(d_1v_1, d_2v_2)$   
`biplot(prcomp(X))`

<sup>3</sup>柴田 里程「データ分析とデータサイエンス」, 近代科学社。

<sup>4</sup>K. R. Gabriel, The biplot graphic display of matrices with application to principal component analysis, *Biometrika*, 58 (3), 453-467.