

# Independent Component Biplot

Su Jin Lee<sup>a</sup> · Yong-Seok Choi<sup>a,1</sup>

<sup>a</sup>Department of Statistics, Pusan National University

(Received October 11, 2013; Revised January 13, 2014; Accepted January 17, 2014)

---

## Abstract

Biplot is a useful graphical method to simultaneously explore the rows and columns of a two-way data matrix. In particular, principal component factor biplot is a graphical method to describe the interrelationship among many variables in terms of a few underlying but unobservable random variables called factors. If we consider the unobservable variables (which are mutually independent and also non-Gaussian), we can apply the independent component analysis decomposing a mixture of non-Gaussian in its independent components. In this case, if we apply the principal component factor analysis, we cannot clearly describe the interrelationship among many variables. Therefore, in this study, we apply the independent component analysis of Jutten and Herault (1991) decomposing a mixture of non-Gaussian in its independent components. We suggest an independent component biplot to interpret the independent component analysis graphically.

Keywords: Biplot, independent component analysis, non-Gaussian.

---

## 1. 서론

행렬도(biplot)는 이원표 자료행렬(two-way data matrix)의 행과 열을 한 그림에 동시에 나타내는 탐색적 방법으로, 복잡한 다변량 분석 결과를 보다 쉽게 파악할 수 있는 장점이 있다. Gabriel (1971)에 의해서 주로 개발되었고, 지금까지 행렬도에 대한 연구는 다양한 분야에서 활발하게 진행되어 왔다. 특히 그 중에서 Choi와 Shin (2013, Chapter 2)가 비정칙값분해(singular value decomposition)를 활용하여 제안한 주성분인자 행렬도(principal component factor biplot; PCFB)는 인자분석(factor analysis)을 통해서 변수들 간의 상호의존 구조를 탐색하기 위한 시각적 도구이다.

자료에 따라 잠재된 변수들이 독립(independent)이고 비가우시안(non-Gaussian) 분포를 가진다는 사전 정보가 있을 때, 주성분법을 이용한 인자분석을 적용하면 원래 변수들의 상호 관계를 잘못 해석하는 경우도 있다. 따라서 Jutten과 Herault (1991)는 인자분석의 개념을 확장한 독립성분분석(independent component analysis)을 제안하였다.

따라서 본 연구에서는 독립성분분석을 응용하여 원래 변수들 간의 상호 관계를 기하학적으로 살펴볼 수 있는 시각적 도구인 독립성분 행렬도(independent component biplot; ICB)를 제안하려 한다. 2절에서는 주성분인자분석과 주성분인자 행렬도, 독립성분분석과 독립성분 행렬도의 기초 이론을 설명하려 한다. 3절에서는 독립성분 행렬도의 활용 사례를 보이고, 주성분인자 행렬도의 결과와 비교하고자 한다.

---

<sup>1</sup>Corresponding author: Professor, Department of Statistics, Pusan National University, 2, Busandaehak-ro 63beon-gil, Geumjeong-Gu, Busan 609-735, Korea. E-mail: yschoi@pusan.ac.kr

## 2. 독립성분 행렬도

### 2.1. 주성분인자분석과 주성분인자 행렬도

이 절에서는 Choi와 Shin (2013, Chapter 2), Choi와 Jung (2003, Chapter 3)을 참고로 하여 기존의 주성분인자 행렬도를 소개하기로 하자. 일반적으로 인자분석은 변수들 간의 상호의존 구조를 나타내는 공분산행렬에서 공통인자(common factor)를 추출하고, 이들을 해석하여 원래 변수들이 나타내는 복잡한 구조를 쉽게 파악하기 위한 자료축약(data reduction)기법이다.

크기가  $p \times 1$ 인 확률벡터  $\mathbf{x} = (x_1, \dots, x_p)^t$ 가 평균벡터  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^t$ 와 공분산행렬  $\boldsymbol{\Sigma}$ 를 가질 때, 인자분석 모형은

$$\mathbf{x} - \boldsymbol{\mu} = \mathbf{L}\mathbf{f} + \boldsymbol{\varepsilon}$$

와 같다. 여기서  $\mathbf{f} = (f_1, \dots, f_m)^t$ 는 관측되지 않는(unobservable) 크기가  $m \times 1$ 인 공통인자벡터이고,  $\mathbf{L} = (l_{ij})$ ,  $i = 1, \dots, p; j = 1, \dots, m$ 은 크기가  $p \times m$ 인 인자적재행렬(factor loading matrix)이고,  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_p)^t$ 은 크기가  $p \times 1$ 인 특수인자(specific factor) 벡터이다.

일반적으로 인자분석에서는 통계적 추론을 위해서 다음의 몇 가지 가정을 필요로 한다. 먼저 공통인자의 평균은 0이고 분산은 1이며, 특수인자의 평균은 0이고 분산은  $\psi_i$ ,  $i = 1, \dots, p$ 라 가정한다. 그리고 공통인자와 특수인자는 서로 독립이라 가정한다. 이때 공분산행렬  $\boldsymbol{\Sigma}$ 는 다음의 공통인자분해(common factor decomposition)

$$\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^t + \boldsymbol{\Psi}$$

가 성립함을 알 수 있고, 여기서  $\boldsymbol{\Psi} = \text{diag}(\psi_1, \dots, \psi_p)$ 는 크기가  $p \times p$ 인 특수인자의 분산을 대각원소로 가지는 특수분산행렬(specific variance matrix)이다.

인자분석의 목적은 인자적재행렬  $\mathbf{L}$ 과 특수분산행렬  $\boldsymbol{\Psi}$ 를 추정하는 것인데, 이를 위해 가장 널리 이용되는 공분산행렬의 고유값과 고유벡터를 활용하는 주성분법을 이용하는 주성분인자분석(principal component factor analysis)을 살펴보도록 한다. 인자적재행렬은 공분산행렬의 분해로 얻어질 수 있으므로, 공분산행렬  $\boldsymbol{\Sigma}$ 의 스펙트럼분해(spectral decomposition)를 고려하면

$$\boldsymbol{\Sigma} = \mathbf{V}\boldsymbol{\Lambda}_{\lambda^2}\mathbf{V}^t \quad (2.1)$$

와 같고, 여기서  $\boldsymbol{\Lambda}_{\lambda^2} = \text{diag}(\lambda_1^2, \dots, \lambda_p^2)$ 는 공분산행렬  $\boldsymbol{\Sigma}$ 의 고유값  $\lambda_i^2$ 을 대각원소로 갖는 크기가  $p \times p$ 인 대각행렬이고,  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_p)$ 는 공분산행렬  $\boldsymbol{\Sigma}$ 의 고유벡터  $\mathbf{v}_i$ 를 열벡터로 가지는 크기가  $p \times p$ 인 직교행렬이다. 식 (2.1)의 고유값들 중에서 마지막  $p - m$ 개가 처음  $m$ 개에 비해 상대적으로 무시할 수 있을 만큼 충분히 작다면, 공분산행렬을 근사적으로  $m$ 개의 인자들로 나타낼 수 있다. 즉, 주성분인자분석에서 인자적재행렬의 추정치는

$$\mathbf{L}_{(m)} = (\lambda_1 \mathbf{v}_1, \lambda_2 \mathbf{v}_2, \dots, \lambda_m \mathbf{v}_m) = \mathbf{V}_{(m)} \boldsymbol{\Lambda}_{\lambda(m)}$$

같이 나타나고, 특수분산행렬의 추정치는 다음과 같다.

$$\hat{\boldsymbol{\Psi}} = \boldsymbol{\Sigma} - \mathbf{L}_{(m)} \mathbf{L}_{(m)}^t.$$

일반적으로 주성분인자분석의 목적은 인자분석 모형의 모수인 인자적재행렬을 추정하는 것이지만 추가적으로 공통인자의 추정치인 주성분인자 점수(principal component factor score)를 구할 수 있다. 이

는 인자분석 모형의 타당성을 검토하는데 사용될 뿐만 아니라, 다른 통계분석에서 새로운 자료로 사용되기도 한다. 크기가  $m \times 1$ 인 주성분인자 점수는

$$\hat{\mathbf{f}} = \left( \mathbf{x}^t \mathbf{v}_1, \mathbf{x}^t \mathbf{v}_2, \dots, \mathbf{x}^t \mathbf{v}_m \right)^t = \mathbf{V}_{(m)}^t \mathbf{x}$$

와 같이 나타난다.

주성분인자 행렬도는 주성분인자분석에서 서로 상관이 높은 변수들이 공통인자에서 비슷한 인자적재값을 가지는지를 보이는 시각적 도구로, 변수와 공통인자의 공분산을 나타내는 인자적재행렬을 좌표점 행렬로 이용한다. 이를 주성분인자분석에서  $\mathbf{L}$ 과  $\Psi$ 를 추정할 때 이용되는 스펙트럼분해 대신에 행렬도에서 주로 활용하는 대수적 알고리즘인 비정칙값분해를 이용한다.

이를 위해 크기가  $p \times 1$ 인 확률벡터  $\mathbf{x}$ 가  $n$ 개의 관찰값을 가지는 경우 크기가  $n \times p$ 인 자료행렬  $\mathbf{X} = (x_{ij}), i = 1, \dots, n; j = 1, \dots, p$ 를 고려할 수 있다. 여기서  $\mathbf{X}$ 의 각 열을 나타내는 변수들의 평균  $\bar{x}_{.j} = \sum_{i=1}^n x_{ij}/n$ 을 뺀 중심화(centering) 자료행렬  $\mathbf{X}_c = (x_{ij} - \bar{x}_{.j}), i = 1, \dots, n; j = 1, \dots, p$ 의 비정칙값분해

$$\mathbf{X}_c = \sum_{k=1}^p \mathbf{u}_k \lambda_k \mathbf{v}_k^t = \mathbf{U} \mathbf{\Lambda}_\lambda \mathbf{V}^t$$

를 고려한다. 여기서  $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_p)^t$ 는 직교열벡터  $\mathbf{u}_k = (u_{1k}, \dots, u_{nk})^t$ 를 갖는 크기가  $n \times p$ 인 직교행렬이고,  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_p)^t$ 는 직교열벡터  $\mathbf{v}_k = (v_{1k}, \dots, v_{nk})^t$ 를 갖는 크기가  $p \times p$ 인 직교행렬이다. 그리고  $\mathbf{\Lambda}_\lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ 는  $\lambda_1 \geq \dots \geq \lambda_p > 0$ 인 비정칙값을 대각원소로 갖는 크기가  $p \times p$ 인 대각행렬이다. 표본공분산행렬  $\mathbf{S}_{n-1}$ 과 중심화 자료행렬  $\mathbf{X}_c$  사이의 관계는

$$\mathbf{S}_{n-1} = (n-1)^{-1} \mathbf{X}_c^t \mathbf{X}_c$$

와 같이 나타나고, 표본공분산행렬  $\mathbf{S}_{n-1}$ 의 스펙트럼분해는

$$\mathbf{S}_{n-1} = (n-1)^{-1} \mathbf{V} \mathbf{\Lambda}_\lambda^2 \mathbf{V}^t = (n-1)^{-1} \sum_{k=1}^p \lambda_k^2 \mathbf{v}_k \mathbf{v}_k^t$$

와 같이 나타나게 된다. 이때 고유벡터  $\mathbf{v}_k$ 와 비정칙값벡터  $\mathbf{v}_k$ 는 동일함을 알 수 있다. 즉, 인자적재행렬을 구할 때 대수적 계산상으로  $\mathbf{X}_c$ 의 공분산행렬  $\mathbf{S}_{n-1}$ 을 구해서 스펙트럼분해를 하는 것보다  $\mathbf{X}_c$ 의 비정칙값분해를 이용하는 것이 훨씬 편리함을 알 수 있다. 따라서 자료행렬  $\mathbf{X}_c$ 의 비정칙값분해를 통해서 구한 직교행렬  $\mathbf{V}$ 와 대각행렬  $\mathbf{\Lambda}_\lambda$ 로 추정된 인자적재행렬은

$$\hat{\mathbf{L}}_{(m)} = (n-1)^{-\frac{1}{2}} \mathbf{V}_{(m)} \mathbf{\Lambda}_{\lambda(m)} = (n-1)^{-\frac{1}{2}} \left( \lambda_1 \mathbf{v}_1, \lambda_2 \mathbf{v}_2, \dots, \lambda_m \mathbf{v}_m \right)$$

과 같고, 이를  $m$ 차원의 주성분인자 행렬도의 좌표점 행렬(coordinate points matrix)이라 한다. 원자료에 대한  $m$ 차원 주성분인자 행렬도의 설명력이라 할 수 있는 근사적합도(goodness-of-fit of the approximation)는 주성분인자분석에서 몇 개의 공통인자를 결정하는 문제와 직접적으로 연관되어 있으므로,  $m$ 개의 고유값이 전체 고유값의 합에서 차지하는 비율

$$\text{fit}_{\text{PCF}} = 1 - \frac{\left\| \mathbf{S}_{n-1} - \hat{\mathbf{L}}_{(m)} \hat{\mathbf{L}}_{(m)}^t \right\|^2}{\left\| \mathbf{S}_{n-1} \right\|^2} = \frac{\lambda_1^2 + \dots + \lambda_m^2}{\sum_{k=1}^p \lambda_k^2}$$

로 나타낸다. 여기서  $\lambda_k^2$ 은 중심화 자료행렬  $\mathbf{X}_c$ 의 비정칙값의 제곱이고, 표본공분산행렬  $\mathbf{S}_{n-1}$ 의 고유값과 일치한다.

## 2.2. 독립성분분석과 독립성분 행렬도

이 절에서는 Hyvarinen과 Oja (2000), Izenman (2008)를 참고로 하여 독립성분분석의 기초이론과 알고리즘을 요약하고, 독립성분분석을 기하학적으로 해석하기 위한 시각적 도구인 독립성분 행렬도를 제안하려 한다.

독립성분분석은 고차원 자료에서 잠재된 독립변수를 추출하고 변수들 간의 상호 관계를 파악하는 다변량 통계기법이다. 일반적으로 독립성분분석은 비가우시안 신호들의 혼합에서 독립성분 신호들을 찾는 문제를 해결하기 위해 적용되어 왔다.

관측되는 크기  $p \times 1$ 인 확률벡터  $\mathbf{x} = (x_1, \dots, x_p)^t$ 의 독립성분분석 모형은

$$\mathbf{x} = \mathbf{A}\mathbf{s} \quad (2.2)$$

와 같다. 여기서  $\mathbf{A} = (a_{ij}), i = 1, \dots, p; j = 1, \dots, m$ 는 크기가  $p \times m$ 인 혼합행렬(mixing matrix)이고,  $\mathbf{s} = (s_1, \dots, s_m)^t$ 는 관측되지 않는 크기가  $m \times 1$ 인 독립성분(independent component)벡터이다. 그리고 독립성분을 찾기 위해서 식 (2.2)를

$$\mathbf{s} = \mathbf{W}\mathbf{x} \quad (2.3)$$

로 변형할 수 있는데, 여기서  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m)^t$ 는 크기가  $m \times p$ 인 분리행렬(separating matrix)이다. 독립성분분석의 목적은 혼합행렬  $\mathbf{A}$ 와 독립성분  $\mathbf{s}$ 를 추정하는 것인데, 이를 위해 다음과 같은 가정을 필요로 한다. 먼저 독립성분은 통계적으로 독립이고 비가우시안 분포를 따른다고 가정한다.

식 (2.2)와 식 (2.3)을 이용해서 임의로 추정된 하나의 독립성분 점수(independent component score)  $\hat{\mathbf{s}}$ 을 고려하면

$$\hat{\mathbf{s}} = \hat{\mathbf{W}}\mathbf{x} = \hat{\mathbf{W}}\mathbf{A}\mathbf{s}$$

와 같다. 즉, 추정된 독립성분 점수  $\hat{\mathbf{s}}$ 은 원래 독립성분  $\mathbf{s}$ 의 선형결합으로 표현 된다. 일반적으로 독립인 확률변수의 합은 각각의 확률변수보다 가우시안에 가까우므로, 독립성분 점수  $\hat{\mathbf{s}}$ 은 독립성분  $\mathbf{s}$ 보다 가우시안에 가깝고, 독립성분 점수  $\hat{\mathbf{s}}$ 가 최소한의 가우시안성을 가질 때 독립성분  $\mathbf{s}$ 와 같게 된다. 그러므로 독립성분 점수  $\hat{\mathbf{s}}$ 의 비가우시안성을 최대로 하는 분리행렬  $\hat{\mathbf{W}}$ 을 구하면 독립성분을 추정할 수 있다.

여러 가지 비가우시안성을 측정하는 척도 중에서 독립성분분석에서는 계산상의 편의를 위해서 엔트로피(entropy)를 변형한 네젠트로피(negentropy)를 이용하게 된다.  $\mathbf{s}$ 의 네젠트로피는

$$J(\mathbf{s}) = H(\mathbf{s}_{gauss}) - H(\mathbf{s})$$

와 같고 여기서  $\mathbf{s}_{gauss}$ 는  $\mathbf{s}$ 와 같은 공분산을 가지는 가우시안 확률변수이고,  $f(\mathbf{s})$ 가  $\mathbf{s}$ 의 확률밀도함수일 때  $H(\mathbf{s}) = -\int f(\mathbf{s})\log f(\mathbf{s})d\mathbf{s}$ 가 엔트로피이다. 네젠트로피는  $\mathbf{s}$ 가 가우시안 변수이면 0이 되고, 비가우시안 변수이면 0이 아니다. 그러나 실제 네젠트로피를 구하는 것은 계산과정이 복잡하므로 근사한 네젠트로피

$$J(\mathbf{s}) \approx [E\{G(\mathbf{s})\} - E\{G(\mathbf{z})\}]^2 \quad (2.4)$$

를 이용하게 된다. 여기서  $\mathbf{z}$ 는 표준정규분포를 따르는 확률변수이고,  $G$ 는 비이차(non-quadratic)함수이다. 주로 독립성분분석에서 이용하는 비이차함수는

$$G(\mathbf{s}) = \frac{1}{\alpha} \log \cosh(\alpha \mathbf{s}), \quad 1 \leq \alpha \leq 2 \quad \text{또는} \quad G(\mathbf{s}) = -e^{-\frac{\mathbf{s}^2}{2}}$$

이다.

독립성분의 비가우시안성을 최대화하는 방법인 Hyvarinen과 Oja (1997), Hyvarinen (1999)가 제안한 고정점 반복법인 FastICA 알고리즘을 살펴보도록 한다. 일반적으로 FastICA 알고리즘을 적용하기 전에 확률벡터  $\mathbf{x}$ 의 평균이 0이 되도록 하는 중심화과정과 공분산행렬이 단위행렬이 되도록 하는 백색화(whitening) 과정을 시행한다. 확률벡터  $\mathbf{x} = (x_1, \dots, x_p)^t$ 가 평균벡터  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^t$ 와 공분산행렬  $\boldsymbol{\Sigma}$ 를 가질 때, 중심화과정을 시행한  $\mathbf{x}_c$ 는

$$\mathbf{x}_c = \mathbf{x} - \boldsymbol{\mu}$$

와 같고 백색화과정은 식 (2.1)의 스펙트럼분해를 이용하여 다음과 같이 나타난다.

$$\mathbf{x}_w = \boldsymbol{\Lambda}_\lambda \mathbf{V}^t \mathbf{x}_c.$$

앞에서 설명한 것과 같이 FastICA 알고리즘의 목적은 식 (2.3)의  $\mathbf{W}\mathbf{x}$ 의 비가우시안성을 최대화하는 분리행렬  $\mathbf{W}$ 를 찾는 것이다. 알고리즘의 설명을 위해  $m = 1$ 인 단일독립성분  $s = \mathbf{w}^t \mathbf{x}$ 를 고려하면,  $E\{\mathbf{w}^t \mathbf{x}\} = \|\mathbf{w}\|^2 = 1$ 의 조건을 만족하고 식 (2.4)를 최대화하는  $\mathbf{w}$ 를 찾는 것이다. 라그랑지 승수법(Lagrange multiplier method)을 이용하여

$$F(\mathbf{w}) = E\{G(\mathbf{w}^t \mathbf{x})\} - \frac{\lambda}{2} (\|\mathbf{w}\|^2 - 1) \quad (2.5)$$

와 같이 나타낼 수 있고, 여기서  $\lambda$ 는 라그랑지 승수이다. 식 (2.5)를 최대화하기 위해 뉴턴-랩슨법(Newton-Raphson method)을 이용하여

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \left( \frac{\partial^2 F(\mathbf{w}_k)}{\partial \mathbf{w}_k^2} \right)^{-1} \left( \frac{\partial F \mathbf{w}_k}{\partial \mathbf{w}_k} \right)$$

로 나타낼 수 있다. 이를 정리하면

$$\mathbf{w}_{k+1} \leftarrow E(\mathbf{x}g(\mathbf{w}_k^t \mathbf{x})) - \mathbf{w}_k E(g'(\mathbf{w}_k^t \mathbf{x}))$$

로 나타낼 수 있다. 여기서  $g = \partial G / \partial \mathbf{w}$ 는 비이차함수  $G$ 의 1차 도함수이고,  $g' = \partial^2 G / \partial^2 \mathbf{w}$ 는 2차 도함수이다. 이때 수렴하는 값을 찾을 때까지 위의 과정을 반복하면  $\mathbf{w}$ 를 찾을 수 있다.

크기가  $p \times 1$ 인 확률벡터  $\mathbf{x}$ 가  $n$ 개의 관찰값을 가질 때, 크기가  $n \times p$ 인 자료행렬  $\mathbf{X} = (x_{ij})$ ,  $i = 1, \dots, n; j = 1, \dots, p$ 의 독립성분분석 모형은

$$\mathbf{X} = \mathbf{S}\mathbf{A}^t \quad (2.6)$$

와 같고, 여기서  $\mathbf{S} = (\mathbf{s}_1, \dots, \mathbf{s}_n)^t$ 는 크기가  $n \times m$ 인 독립성분 행렬이고  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_p)^t$ 는 크기가  $p \times m$ 인 혼합행렬이다.

행렬도 관점에서 식 (2.6)을 살펴보면, 자료행렬  $\mathbf{X}$ 의  $(i, j)$ 번째의 원소를  $n$ 개의 표시점  $\mathbf{s}_1, \dots, \mathbf{s}_n$ 과  $p$ 개의 표시점  $\mathbf{a}_1, \dots, \mathbf{a}_p$ 의 내적인

$$x_{ij} = \mathbf{s}_i^t \mathbf{a}_j, \quad i = 1, \dots, n; j = 1, \dots, p$$

와 같이 나타낼 수 있다. 즉, 식 (2.6)은 크기가 각각  $n \times m$ 과  $p \times m$ 인  $\mathbf{S} = (\mathbf{s}_1, \dots, \mathbf{s}_n)^t$ 와  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_p)^t$ 에 의해서

$$\mathbf{X} = \mathbf{S}\mathbf{A}^t = \begin{pmatrix} \mathbf{s}_1^t \\ \vdots \\ \mathbf{s}_n^t \end{pmatrix} (\mathbf{a}_1, \dots, \mathbf{a}_p)$$

와 같이 인자분해(factorization)로 재표현되는 것을 알 수 있다. 이를 만족하는 표시점  $\mathbf{s}_1, \dots, \mathbf{s}_n$ 과  $\mathbf{a}_1, \dots, \mathbf{a}_p$ 에 의해서 자료행렬  $\mathbf{X}$ 의  $n$ 개의 행과  $p$ 개의 열에 대한 그림을 그릴 수 있는데, 이를 독립성분 행렬도라 제안한다. 여기서  $\mathbf{s}_1, \dots, \mathbf{s}_n$ 와  $\mathbf{a}_1, \dots, \mathbf{a}_p$ 는 행렬도에서 행 좌표점 벡터와 열 좌표점 벡터로 볼 수 있다.

좌표점 행렬  $\mathbf{S}$ 와  $\mathbf{A}$ 에서 임의의  $m$ 개의 열을 고려한 부행렬  $\mathbf{S}_{(m)}$ 과  $\mathbf{A}_{(m)}$ 에 의해서

$$\mathbf{X}_{(m)} = \mathbf{S}_{(m)} \mathbf{A}_{(m)}^t$$

와 같이 표현될 때, 자료행렬  $\mathbf{X}$ 가 행렬  $\mathbf{X}_{(m)}$ 에 의해 근사적으로 일치한다면  $\mathbf{X}_{(m)}$ 의 행렬도는  $\mathbf{X}$ 의  $m$ 차원 독립성분 행렬도를 제공하게 된다.

이때  $m$ 개의 열을 고려하는 방법은 2.1절에서 원자료에 대한  $m$ 차원 주성분인자 행렬도의 설명력으로 정의한 근사적합도를 응용하여 제안한다. 크기가  $p \times 1$ 인  $\mathbf{a}_j$ ,  $j = 1, \dots, m$ 를 열벡터로 가지는 혼합행렬  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_m)$ 를 고려한다. 이 때 혼합행렬  $\mathbf{A}$ 가 열벡터  $\mathbf{a}_j$ ,  $j = 1, \dots, m$ 에 의해 얼마나 근사적으로 일치하는지를 알 수 있는 근사정도값을

$$\text{fit}_j = 1 - \frac{\|\mathbf{A} - \mathbf{A}^{(j)}\|^2}{\|\mathbf{A}\|^2}, \quad j = 1, \dots, m$$

이라 제안한다. 여기서  $\mathbf{A}^{(j)} = (\mathbf{0}, \dots, \mathbf{0}, \mathbf{a}_j, \mathbf{0}, \dots, \mathbf{0})$ ,  $j = 1, \dots, m$ 는  $\mathbf{A}$ 의  $j$ 번째 열벡터  $\mathbf{a}_j$ 만 남기고 나머지 열벡터를  $\mathbf{0}$ 으로 바꾼 크기가  $p \times m$ 인 행렬이다. 이 때 구한  $m$ 개의 근사정도값  $\text{fit}_1, \text{fit}_2, \dots, \text{fit}_m$ 을 큰 것부터 크기순으로 배열하여

$$\text{fit}_{(1)} \geq \text{fit}_{(2)} \geq \dots \geq \text{fit}_{(m)}$$

이라 한다. 이에 대응하는  $\mathbf{A}$ 의 열벡터를 각각  $\mathbf{a}_{(1)}, \mathbf{a}_{(2)}, \dots, \mathbf{a}_{(m)}$ 이라 하면 이들에 의한 크기가  $p \times r$ 인 부행렬  $\mathbf{A}_{(m)}$ 은

$$\mathbf{A}_{(m)} = (\mathbf{a}_{(1)}, \mathbf{a}_{(2)}, \dots, \mathbf{a}_{(m)})$$

와 같다. 이때  $\mathbf{A}_{(m)}$ 을  $m$ 차원 독립성분 행렬도의 좌표점 행렬이라 제안한다. 끝으로 원자료에 대한  $m$ 차원 독립성분 행렬도의 설명력으로 정의할 수 있는 근사적합도를 다음과 같이 제안한다.

$$\text{fit}_C = \text{fit}_{(1)} + \text{fit}_{(2)} + \dots + \text{fit}_{(m)}.$$

### 3. 활용 사례

Lathauwer 등 (2000)의 심전도 자료는 5분 동안 500Hz의 주파수를 사용하는 8개의 채널을 임신부의 피부에 부착하여 측정된 것이다. Ch.1부터 Ch.5는 아기와 가까운 임신부의 복부(abdominal) 쪽에 부착하고, Ch.6부터 Ch.8은 임신부의 심장에 가까운 가슴(chest) 쪽에 각각 부착하였다. Figure 3.1에는 각 채널에서 0.002초 단위로 5분 동안 측정된 자료를, 수평축에는 시간으로 수직축에는 측정된 2500개의 값으로 그린 8개의 채널에 대한 그림이 정리되어 있다.

심전도 자료의 주성분인자 행렬도를 통하여 변수들 간의 관계를 살펴볼 수도 있지만, 일반적으로 심전도 자료에서는 비가우시안 신호들의 혼합에서 독립성분 신호들을 찾는 문제가 발생하므로 보다 고차원통계를 적용하여 변수들의 관계를 살펴보는 방법이 요구된다. 그러므로 본 연구에서는 심전도 자료의 독립성분 행렬도를 주성분인자 행렬도와 비교하여 변수들 간의 관계를 살펴본다.



Figure 3.1. Plot for electrocardiography of 8 channel

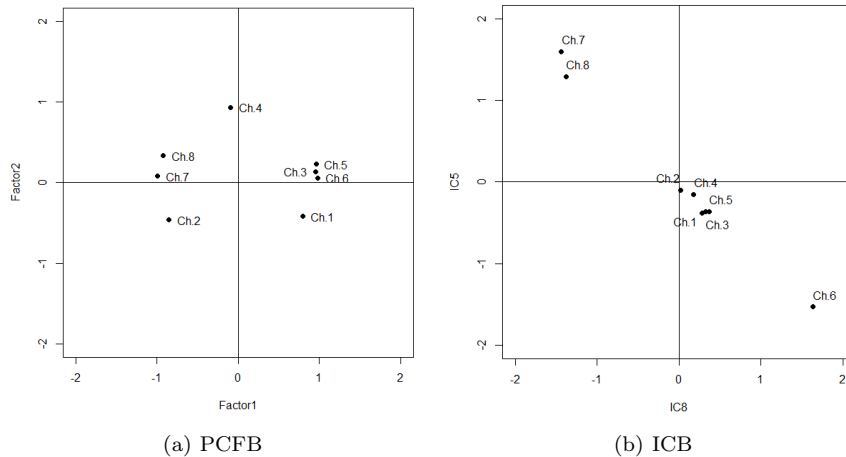


Figure 3.2. Biplots for electrocardiography: PCFB(Principal component factor biplot), ICB(Independent component biplot)

먼저 Figure 3.2에 (a)는 심전도 자료에 대한 주성분인자 행렬도를 나타낸 것이고, Factor1과 Factor2는 행렬도의 관점에서 제1축과 제2축을 나타내고, (b)는 심전도 자료에 대한 독립성분 행렬도를 나타낸 것이고, IC8과 IC5는 행렬도의 관점에서 제1축과 제2축을 나타낸다. 이 행렬도의 좌표점값과 근사적합도는 Table 3.1에 나타나 있으며, 특히 근사적합도는 주성분인자 행렬도는 92.97%, 독립성분 행렬도는 82.11%로 원래 변수들의 관계를 충분히 설명한다고 할 수 있다. 먼저 주성분인자 행렬도를 살펴보면, 제1축에 대해 왼쪽의 음의 방향에 놓인 Ch.2, Ch.7, Ch.8과 오른쪽의 양의 방향에 놓인 Ch.1, Ch.3, Ch.5, Ch.6은 이질적인 성격을 나타냄을 알 수 있다. 제2축의 설명력이 제1축에 비해 다소 낮지만 고려해서 살펴보면 Ch.2와 Ch.7, Ch.8은 서로 다른 방향에 놓여있고, Ch.1과 Ch.3, Ch.5, Ch.6도 서로 다른 방향에 놓여있음을 알 수 있다. 반면에 독립성분 행렬도를 살펴보면, 제1축에 대해 왼쪽의 음의 방향에 놓인 Ch.7, Ch.8과 오른쪽의 양의 방향에 놓인 Ch.1, Ch.3, Ch.4, Ch.5, Ch.6은 이질적인 성격을 나타냄을 알 수 있다.

**Table 3.1.** Coordinates and goodness-of-fits of the approximation for two biplots.

Channel	PCFB		ICB	
	Factor1	Factor2	IC8	IC5
Ch.1	-0.795	0.418	0.285	-0.383
Ch.2	0.850	0.463	0.022	-0.101
Ch.3	-0.957	-0.139	0.323	-0.364
Ch.4	0.089	-0.931	0.181	-0.150
Ch.5	-0.963	-0.228	0.371	-0.360
Ch.6	-0.978	-0.059	1.635	-1.530
Ch.7	0.990	-0.081	-1.435	1.596
Ch.8	0.920	-0.334	-1.382	1.292
Goodness-of-fits(%)	74.86	92.97	68.26	82.11

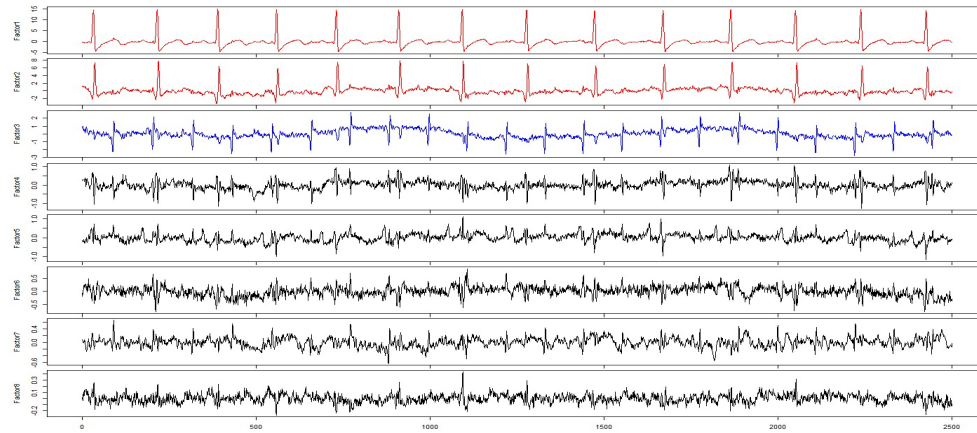
대체로 원자료의 진폭(amplitude)이 Ch.1에서 Ch.8로 갈수록 증가하는데, 이는 임신부의 복부 쪽에서 가슴 쪽으로 자료를 측정해서이다. 높은 진폭을 가지며 거의 동일한 파형을 가지는 Ch.7과 Ch.8은 임신부의 심박동에 관련된 채널이고, 나머지 채널들은 상대적으로 낮은 진폭을 가져 태아의 심박동에 관련된 채널이다. Ch.7과 Ch.8이 서로 동질적인 성격을 나타낸다는 해석은 두 행렬도 모두에서 얻을 수 있지만, 나머지 Ch.1부터 Ch.6까지의 경우 두 행렬도의 해석이 상이함을 알 수 있다. 주성분인자 행렬도는 Ch.1과 Ch.2와 Ch.3, Ch.5, Ch.6이 이질적인 성격을 나타냄을 알 수 있었지만, 독립성분 행렬도는 Ch.1, Ch.3, Ch.4, Ch.5, Ch.6이 동질적인 성격을 나타냄을 알 수 있었다. 즉 채널들 간의 관계를 주성분인자 행렬도보다 독립성분 행렬도가 더 명확하게 나타냄을 알 수 있다.

끝으로 주성분인자 점수와 독립성분 점수를 이용하여 차이를 확인하여 보자. Figure 3.3에 (a)는 주성분인자 점수 그림을 나타낸 것인데, Factor1과 Factor2가 대체로 유사한 파형을 가지며 임신부의 심박동을 나타내고 있고, Factor3이 태아의 심박동을 나타냄을 알 수 있다. 그리고 나머지 Factor4부터 Factor8까지는 뚜렷한 특징이 나타나지 않아 오차를 나타내고 있다. Figure 3.3에 (b)는 독립성분 점수 그림을 나타낸 것인데, IC2, IC4, IC5, IC8은 임신부의 심박동을 나타내고 있고, IC3과 IC6은 태아의 심박동을 나타냄을 알 수 있다. 그리고 IC1은 호흡과 관련된 성분을 나타내고 있고, IC7은 오차를 나타낸다고 할 수 있다. 즉 독립성분 점수가 주성분인자 점수보다 더 명확하게 잠재된 변수들을 찾아냄을 확인할 수 있다.

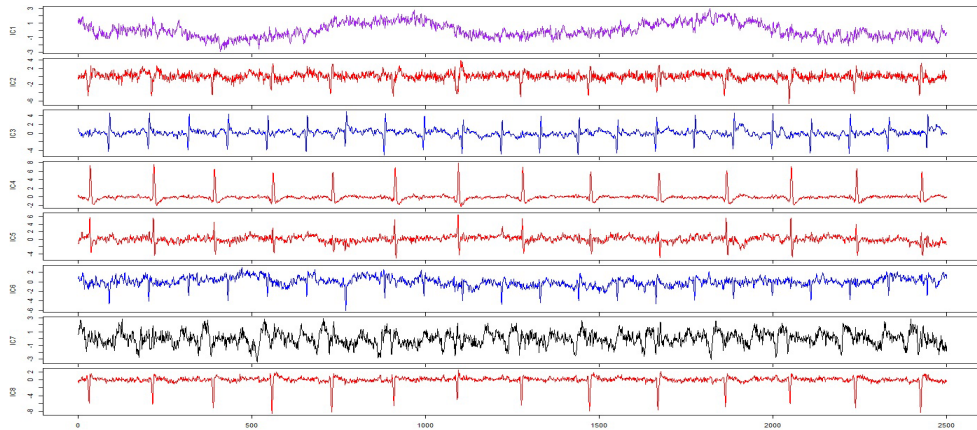
#### 4. 결론

주성분인자 행렬도는 다변량 자료에서 주성분인자분석을 통해 원래 변수들 간의 상호의존 구조를 파악하기 위한 탐색적 방법으로, 변수들 간의 복잡한 구조를 시각적으로 표현하여 한눈에 파악할 수 있는 장점이 있다. 하지만 다변량 자료가 비가우시안 신호들의 혼합에서 독립성분 신호들을 찾는 문제를 가지고 있는 경우, 잘못된 해석 또는 불분명한 해석을 하게 된다. 따라서 본 연구에서는 이 문제를 해결하기 위해 주로 적용되는 독립성분분석을 기하학적으로 살펴볼 수 있는 시각적 도구인 독립성분 행렬도를 제안하였다. 그리고 이를 활용하여 Lathauwer 등 (2000)의 임신부에게 8개 채널을 부착하여 5분 동안 측정된 심전도 자료에서 독립성분 행렬도의 예를 보이고, 주성분인자 행렬도와와의 비교를 통해 그 차이점을 설명하였다. 결과적으로 비가우시안 신호들의 혼합에서 독립성분 신호들을 찾는 경우에는 주성분인자 행렬도보다는 독립성분 행렬도가 변수들의 성격을 더욱 명확하게 나타내었고, 보다 나은 해석을 가능하게 해주었다. 그리고 이러한 행렬도간의 차이를 인자점수와 독립성분점수를 활용하여 비교한 결과에서도 유사한 해석을 할 수 있었다. 추후 독립성분 행렬도의 유용성과 일반성을 위해서는 보다 다양한 성





(a) Principal component factor score plot



(b) Independent component score plot

**Figure 3.3.** Score plot for electrocardiography

공사례를 제공해야 할 것이다.

## References

- Choi, Y. S. and Jung, K. M. (2003). *Method and Application of Multivariate Analysis by using SAS*, Free Academy, Seoul.
- Choi, Y. S. and Shin, S. M. (2013). *Understanding Biplots Analysis Using R*, Free Academy, Gyeonggi-Do.
- Gabriel, K. R. (1971). The biplot graphics display of matrices with applications to principal analysis, *Biometrika*, **58**, 453–467.
- Hyvarinen, A. (1999). Fast and robust fixed-point algorithms for independent component analysis, *IEEE Transactions on Neural Networks*, **10**, 626–634.
- Hyvarinen, A. and Oja, E. (1997). A fast fixed-point algorithm for independent component analysis, *Neural Computation*, **9**, 1483–1492.

- Hyvarinen, A. and Oja, E. (2000). Independent component analysis: Algorithms and application, *Neural Networks*, **13**, 411–430.
- Izenman, J. A. (2008). *Modern Multivariate Statistical Techniques*, Springer, New York.
- Jutten, C. and Herault, J. (1991). Blind Separation of Sources, part I: An Adaptive Algorithm Based on Neuromimetic Architecture, *Signal Processing*, **24**, 1–10.
- Lathauwer, L. De, De Moor, B. and Vandewalle, J. (2000). Fetal electrocardiogram extraction by blind source subspace separation, *IEEE Transactions on Biomedical Engineering*, **47**, 567–572.

## 독립성분 행렬도

이수진<sup>a</sup> · 최용석<sup>a,1</sup>

<sup>a</sup>부산대학교 통계학과

(2013년 10월 11일 접수, 2014년 1월 13일 수정, 2014년 1월 17일 채택)

---

### 요약

행렬도(biplot)는 이원표 자료행렬(two-way data matrix)의 행과 열을 한 그림에 동시에 나타내는 탐색적 방법으로, 복잡한 다변량 분석 결과를 보다 쉽게 파악할 수 있는 장점이 있다. 특히 주성분인자 행렬도(principal component factor biplot; PCFB)는 인자분석을 통해서 변수들 간의 상호의존 구조를 탐색하기 위한 시각적 도구이다. 자료에 따라 잠재된 변수들이 독립(independent)이고 비가우시안(non-Gaussian) 분포를 가진다는 사전 정보가 있을 때, Jutten과 Herault (1991)가 제안한 독립성분분석(independent component analysis)을 이용한다. 이 경우 주성분법을 이용한 인자분석을 적용하면 원래 변수들의 상호 관계를 잘못 해석할 수도 있다. 따라서 본 논문에서는 자료에 따라 잠재된 변수들이 독립이고 비가우시안 분포를 가진다는 사전 정보가 있을 때, 독립성분분석을 응용하여 원래 변수들 간의 상호 관계를 기하학적으로 살펴볼 수 있는 시각적 도구인 독립성분 행렬도(independent component biplot; ICB)를 제안하려 한다.

주요용어: 독립성분분석, 비가우시안, 행렬도.

---

<sup>1</sup>교신저자: (609-735) 부산광역시 금정구 부산대학교로 63번길 2, 부산대학교 통계학과, 교수.  
E-mail: yschoi@pusan.ac.kr