

K-Subspace clustering and its application in sparse component analysis

Zhaoshui He^{1,2} and Andrzej Cichocki¹

1-Laboratory for Advanced Brain Signal Processing RIKEN Brain Science Institute,
Wako-shi, Saitama 351-0198, Japan

2- School of Electronic and Information Engineering, South China University of Technology
Guangzhou, China, 510640

Abstract. The K-subspace clustering algorithm is established for sparse component analysis and overcome the difficulty that conventional SCA algorithms can not overcome. The conventional SCA algorithm can only perform single dominant SCA, can not perform multiple dominant SCA, but the proposed SCA algorithm based on K-subspace clustering can overcome this difficulty.

1 Introduction

Clustering is one of the most widely used techniques for exploratory data analysis. It can be used for signal compression, blind source separation, feature extraction, regularization in inverse problems, especially in the applications of sparse component analysis and undetermined blind source separation [1]-[5][8]-[11]. Usually the clustering plays a key role in sparse component analysis.

The conventional SCA is based on the single-dominant-component assumption that for each time, there is only one dominant component and others components are insignificant. So the scatter plot of observation signals shows clear line orientations geometrically. But for multiple dominant components SCA, for example, if there are two or more dominant components for each time, there will be no line orientations at the scatter plot of observation signals and the conventional SCA algorithms will fail. However Georgiev, Theis and Cichocki [1] pointed out that as long as the number of nonzero components in source signals is smaller than the number m of observation signals, SCA is solvable and gave the corresponding algorithm. But they did not consider the noise. In practice it is impossible to rigorously satisfy their assumption (the number of nonzero components in source signals is smaller than the number m of observation signals). So here we continue to discuss the “multiple dominant components SCA problem”. The K-subspace clustering approach is proposed to perform multiple dominant components SCA method in this paper.

2 Problem statement of K-subspace clustering

Give a $m \times n$ complex valued matrix:

$$\mathbf{B} = \begin{pmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{m1} & \cdots & b_{mn} \end{pmatrix} = (\mathbf{b}_1, \dots, \mathbf{b}_n) \in \mathbb{C}^{m \times n}, \quad (1)$$

where C denotes complex valued number. If $m \leq n$, assume that any a $m \times m$ submatrix of \mathbf{B} is invertible. Give an integer $k (< m)$ and a data set $\mathbf{X} = (\mathbf{x}(1), \dots, \mathbf{x}(T)) \in C^{m \times T}, (T \gg m)$. So \mathbf{B} has C_n^k different $m \times k$ submatrice and denote them as $\mathbf{B}_1, \mathbf{B}_2, \dots, \mathbf{B}_{C_n^k}$. The column vectors of each $m \times k$ submatrix of \mathbf{B}_i can span a linear space. Without the loss of generality, we suppose that \mathbf{B}_i is composed by the column vectors $\mathbf{b}_{i1}, \mathbf{b}_{i2}, \dots, \mathbf{b}_{ik}$ of matrix \mathbf{B} , where $\{i1, i2, \dots, ik\} \subset \{1, 2, \dots, n\}$. The column vectors $\mathbf{b}_{i1}, \mathbf{b}_{i2}, \dots, \mathbf{b}_{ik}$ span a linear space $span\{\mathbf{B}_i\} = span\{\mathbf{b}_{i1}, \mathbf{b}_{i2}, \dots, \mathbf{b}_{ik}\}$. So we can obtain C_n^k linear spaces using $m \times k$ submatrice $\mathbf{B}_1, \mathbf{B}_2, \dots, \mathbf{B}_{C_n^k}$. Indicate $K = C_n^k$. The K-subspace clustering problem can be described as follows: for the given data set $\mathbf{X} = (\mathbf{x}(1), \dots, \mathbf{x}(T))$, how to estimate a $m \times n$ matrix \mathbf{B} such that the distance sum $\sum_{t=1}^T \min_{i=1, \dots, C_n^k} d(\mathbf{x}(t), \mathbf{B}_i)$ is minimized, where $d(\mathbf{x}(t), \mathbf{B}_i)$ means the distance from point $\mathbf{x}(t)$ to the linear space spanned by the k column vectors of \mathbf{B}_i .

3 The distance formula from a point to a sub-linear-space

Consider a point $\mathbf{p} = (p_1, \dots, p_m)^T$ (see Figure 1), we attempt to calculate the distance $d(\mathbf{p}, \mathcal{L})$ from the point \mathbf{p} to linear space \mathcal{L} , whose basis matrix is

$$\mathbf{L} = \begin{pmatrix} l_{11} & \cdots & l_{1k} \\ \vdots & \ddots & \vdots \\ l_{m1} & \cdots & l_{mk} \end{pmatrix} = (\mathbf{l}_1, \dots, \mathbf{l}_k). \quad (2)$$

This problem can convert to search a point \mathbf{l}_* (in the space \mathcal{L}), which is closest to \mathbf{p} , i.e., we can deal with this problem by solving the following optimization problem:

$$\begin{cases} d^2(\mathbf{p}, \mathcal{L}) = \min_{\mathbf{l}} g(\mathbf{l}) = \min_{\mathbf{l}} \|\mathbf{p} - \mathbf{l}\|^2, \\ \text{subject to: } \mathbf{l} \in \mathcal{L}, \end{cases} \quad (3)$$

where $\mathbf{l} = (l_1, \dots, l_m)^T \in C^m$, $\mathbf{p} = (p_1, \dots, p_m)^T \in C^m$. And \mathbf{l} can be linearly represented:

$$\mathbf{l} = \alpha_1 \mathbf{l}_1 + \cdots + \alpha_k \mathbf{l}_k = \mathbf{L}\boldsymbol{\alpha}. \quad (4)$$

From equation (4), solving problem (3) can be described as

$$\begin{aligned} \min_{\boldsymbol{\alpha} \in C^k} g(\boldsymbol{\alpha}) &= \min_{\boldsymbol{\alpha} \in C^k} \|\mathbf{p} - \mathbf{L}\boldsymbol{\alpha}\|^2 = \min_{\boldsymbol{\alpha} \in C^k} \langle \mathbf{p} - \mathbf{L}\boldsymbol{\alpha}, \mathbf{p} - \mathbf{L}\boldsymbol{\alpha} \rangle \\ &= \min_{\boldsymbol{\alpha} \in C^k} (\mathbf{p} - \mathbf{L}\boldsymbol{\alpha})^H (\mathbf{p} - \mathbf{L}\boldsymbol{\alpha}) = \min_{\boldsymbol{\alpha} \in C^k} (\mathbf{p}^H - \boldsymbol{\alpha}^H \mathbf{L}^H) (\mathbf{p} - \mathbf{L}\boldsymbol{\alpha}), \end{aligned} \quad (5)$$

where $\langle \cdot, \cdot \rangle$ stands for inner product. For any real valued function $f(z)$ of a complex valued variable z the gradients with respect to the real and imaginary part are obtained by taking derivatives formally with respect to the conjugate quantities z^* , ignoring the nonconjugate occurrences of z [2][6][7], i.e.,

$$\frac{\partial f(z)}{\partial R(z)} + i \frac{\partial f(z)}{\partial I(z)} = 2 \frac{\partial f(z)}{\partial z^*}. \quad (6)$$

Therefore the derivative of cost function $g(\bullet)$ with respect to α^* is

$$\frac{\partial g(\bullet)}{\partial \alpha^*} = -\mathbf{L}^H \mathbf{p} + \mathbf{L}^H \mathbf{L} \alpha. \quad (7)$$

$$\text{Let } \partial g(\bullet) / \partial \alpha^* = 0, \text{ we get } \alpha = (\mathbf{L}^H \mathbf{L})^{-1} \mathbf{L}^H \mathbf{p}. \quad (8)$$

Substitute equation (8) and equation (4) into (3), we have

$$d(\mathbf{p}, \mathcal{L}) = \sqrt{\mathbf{p}^H \mathbf{p} - \mathbf{p}^H \mathbf{L} (\mathbf{L}^H \mathbf{L})^{-1} \mathbf{L}^H \mathbf{p}}. \quad (9)$$

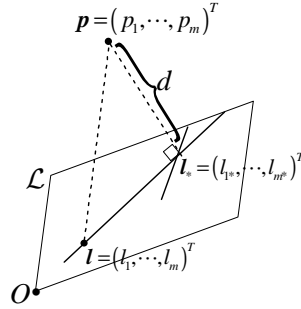


Fig.1 The distance from a point to a line

4 K-subspace clustering algorithm

K-subspace clustering algorithm (here $K = C_n^k$) can be outlined as follows:

- (1) Input the observation dataset $\mathbf{X} = (\mathbf{x}(1), \dots, \mathbf{x}(T)) \in C^{m \times T}$, ($T \gg m$). Initialize the basis matrix \mathbf{B} as $\hat{\mathbf{B}} \in C^{m \times K}$, normalize $\hat{\mathbf{B}}$ by $\hat{\mathbf{b}}_j = \hat{\mathbf{b}}_j / \|\hat{\mathbf{b}}_j\|$, $j = 1, \dots, n$. And denote its K ($K = C_n^k$) $m \times k$ submatrices as $\hat{\mathbf{B}}_1, \hat{\mathbf{B}}_2, \dots, \hat{\mathbf{B}}_K$.
- (2) **Partition stage:** calculate each distance $d(\mathbf{x}(t), \hat{\mathbf{B}}_i)$, $i = 1, \dots, K$ from the observation sample point $\mathbf{x}(t)$ to linear space $\text{span}\{\hat{\mathbf{B}}_i\}$ using the distance formula (9). The observation sample point $\mathbf{x}(t) \in \theta(\hat{\mathbf{B}}_i)$ if and only if $d(\mathbf{x}(t), \hat{\mathbf{B}}_i) = \min\{d(\mathbf{x}(t), \hat{\mathbf{B}}_j), j = 1, \dots, K\}$, where $\theta(\hat{\mathbf{B}}_i)$ is a data vector set. By this means, assign the T observation sample points $\mathbf{x}(1), \dots, \mathbf{x}(T)$ in observation matrix $\mathbf{X} = [\mathbf{x}(1), \dots, \mathbf{x}(T)]$ into K different clusters $\theta(\hat{\mathbf{B}}_i)$, $i = 1, \dots, K$.
- (3) Update the cluster centroids matrix stage: update the cluster centroids matrix $\hat{\mathbf{B}}$:
 - 3.1) Consider each cluster $\theta(\hat{\mathbf{B}}_i)$, assume it contains $T^{(i)}$ entries $\mathbf{x}^{(i)}(1), \dots, \mathbf{x}^{(i)}(T^{(i)})$, which compose a matrix $\mathbf{X}^{(i)} = [\mathbf{x}^{(i)}(1), \dots, \mathbf{x}^{(i)}(T^{(i)})]$. For the

symmetrical matrix $\frac{1}{T}(\mathbf{X}^{(i)})^H \mathbf{X}^{(i)} = \mathbf{V}^{(i)} \mathbf{D}^{(i)} (\mathbf{V}^{(i)})^H$, apply EVD decomposition.

Suppose $d_1^{(i)}, \dots, d_k^{(i)}$ are the largest k eigenvalues and $\mathbf{v}_1^{(i)}, \dots, \mathbf{v}_k^{(i)}$ are the corresponding eigenvectors, we get $\mathbf{E}_i = [\mathbf{v}_1^{(i)}, \dots, \mathbf{v}_k^{(i)}], i = 1, \dots, K$.

3.2) Update the column vector $\hat{\mathbf{b}}_j$ of $\hat{\mathbf{B}}$. Without losing generality, suppose $\hat{\mathbf{b}}_j$ is the column vector of the $m \times k$ submatrice $\hat{\mathbf{B}}_{j1}, \dots, \hat{\mathbf{B}}_{jC_{n-1}^{k-1}}$ of $\hat{\mathbf{B}}$, respectively. We can update the column vector $\hat{\mathbf{b}}_j$ by solving the following optimization problem:

$$J(\hat{\mathbf{b}}_j) = \min_{\hat{\mathbf{b}}_j} \sum_{r=1}^{C_{n-1}^{k-1}} d^2(\hat{\mathbf{b}}_j, \mathbf{E}_{jr}), \text{ subject to } \hat{\mathbf{b}}_j^H \hat{\mathbf{b}}_j = 1. \quad (10)$$

The solution of problem (10) is the eigenvector corresponding to the smallest eigenvalue of matrix $\sum_{r=1}^{C_{n-1}^{k-1}} \mathbf{E}_{jr} (\mathbf{E}_{jr}^H \mathbf{E}_{jr})^{-1} \mathbf{E}_{jr}^H$. From formula (9), problem (10) can be changed into

$$\max_{\hat{\mathbf{b}}_j} \hat{\mathbf{b}}_j^H \left(\sum_{r=1}^{C_{n-1}^{k-1}} \mathbf{E}_{jr} (\mathbf{E}_{jr}^H \mathbf{E}_{jr})^{-1} \mathbf{E}_{jr}^H \right) \hat{\mathbf{b}}_j, \text{ subject to } \hat{\mathbf{b}}_j^H \hat{\mathbf{b}}_j = 1. \quad (11)$$

(11) can be changed into the following optimization problem without constraint:

$$\max_{\hat{\mathbf{b}}_j} J(\hat{\mathbf{b}}_j) = \max_{\hat{\mathbf{b}}_j} \frac{\hat{\mathbf{b}}_j^H \left(\sum_{r=1}^{C_{n-1}^{k-1}} \mathbf{E}_{jr} (\mathbf{E}_{jr}^H \mathbf{E}_{jr})^{-1} \mathbf{E}_{jr}^H \right) \hat{\mathbf{b}}_j}{\hat{\mathbf{b}}_j^H \hat{\mathbf{b}}_j}. \quad (12)$$

From equation (12), let $\frac{\partial J}{\partial \hat{\mathbf{b}}_j^*} = 0$ and we have

$$\begin{aligned} \frac{\partial J}{\partial \hat{\mathbf{b}}_j^*} \cdot \hat{\mathbf{b}}_j^H \hat{\mathbf{b}}_j + J \cdot \frac{\partial (\hat{\mathbf{b}}_j^H \hat{\mathbf{b}}_j)}{\partial \hat{\mathbf{b}}_j^*} &= \left(\sum_{r=1}^{C_{n-1}^{k-1}} \mathbf{E}_{jr} (\mathbf{E}_{jr}^H \mathbf{E}_{jr})^{-1} \mathbf{E}_{jr}^H \right) \hat{\mathbf{b}}_j \Rightarrow \\ \frac{\partial J}{\partial \hat{\mathbf{b}}_j^*} \cdot \hat{\mathbf{b}}_j^H \hat{\mathbf{b}}_j &= \left[\left(\sum_{r=1}^{C_{n-1}^{k-1}} \mathbf{E}_{jr} (\mathbf{E}_{jr}^H \mathbf{E}_{jr})^{-1} \mathbf{E}_{jr}^H \right) - J \right] \hat{\mathbf{b}}_j = 0. \end{aligned} \quad (13)$$

From equation (13), obviously $\hat{\mathbf{b}}_j$ is the eigenvector of $\sum_{r=1}^{C_{n-1}^{k-1}} \mathbf{E}_{jr} (\mathbf{E}_{jr}^H \mathbf{E}_{jr})^{-1} \mathbf{E}_{jr}^H$.

In this way $\hat{\mathbf{b}}_1, \hat{\mathbf{b}}_2, \dots, \hat{\mathbf{b}}_n$ respectively are updated, so $\hat{\mathbf{B}} = [\hat{\mathbf{b}}_1, \hat{\mathbf{b}}_2, \dots, \hat{\mathbf{b}}_n]$ is updated.

(4) Return to step (2) and repeat until $\hat{\mathbf{B}}$ converges.

(5) Output the estimation of basis matrix $\mathbf{B}^* = [\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_n]$.

Remark In analogy to K-means, we call this algorithm ‘‘K-subspace clustering’’.

5 Sparse component analysis

Sparse Component Analysis (SCA) can be expressed as follows:

$$\mathbf{X} = \mathbf{A} \mathbf{S} + \mathbf{V}, \quad (14)$$

where $\mathbf{X} = [\mathbf{x}(1), \dots, \mathbf{x}(T)] \in C^{m \times T}$ ($T \gg m$) is the given data (observation) matrix,

$\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m] \in C^{m \times m}$ is unknown mixing (basis) matrix (not necessary sparse) and

$\mathbf{S} \in C^{m \times T}$ is also unknown matrix representing sparse sources or hidden components,

T is the number of available sample, m the number of observations and n the number of sources. Our main objective is to find a reasonable basis matrix \mathbf{A} such that the coefficients in matrix \mathbf{S} are as sparse as possible. Here we mainly discuss the $k(>1)$ dominant components SCA. We use the two-stage methods to solve the SCA problem: estimate at the first step the basis matrix \mathbf{A} using K-subspace clustering approach and in the next step to estimate the coefficient matrix \mathbf{S} using the following Minimum Distance Decomposition (MDD):

- (1) Denote $K(K = C_n^k)$ $m \times k$ submatrice of \mathbf{A} as $\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_K$.
- (2) Calculate each distance $d(\mathbf{x}(t), \mathbf{A}_i), i = 1, \dots, K$ from the observation sample point $\mathbf{x}(t)$ to linear space $\text{span}\{\mathbf{A}_i\}$ using the distance formula (9). The observation sample point $\mathbf{x}(t) \in \theta(\mathbf{A}_i)$ if and only if $d(\mathbf{x}(t), \mathbf{A}_i) = \min\{d(\mathbf{x}(t), \mathbf{A}_j), j = 1, \dots, K\}$, where $\theta(\mathbf{A}_i)$ is a data vector set. By this means, the T observation sample points $\mathbf{x}(1), \dots, \mathbf{x}(T)$ are assigned into K different clusters $\theta(\mathbf{A}_i), i = 1, \dots, K$.
- (3) Suppose $\mathbf{A}_i = [\mathbf{a}_{i1}, \dots, \mathbf{a}_{ik}]$. If $\mathbf{x}(t) \in \theta(\mathbf{A}_i)$, the estimation of the column vector $s(t), t = 1, \dots, T$ of coefficient matrix \mathbf{S} is as follows:

$$\begin{cases} \begin{pmatrix} s(i1, t) \\ \vdots \\ s(ik, t) \end{pmatrix} = (\mathbf{A}_i^H \mathbf{A}_i)^{-1} \mathbf{A}_i^H \mathbf{x}(t), \\ s(j, t) = 0, \text{ for } j \neq i1, \dots, ik. \end{cases} \quad (15)$$

6 Numerical experiments and result analysis

We give a SCA experiment. To check how well the mixing matrix is estimated, we introduce the following the Biased Angles Sum (BAS), the sum of angles between the column vectors (of mixing matrix) and their corresponding estimations:

$$BAS(\mathbf{A}, \hat{\mathbf{A}}) = \sum_{i=1}^n \text{acos}(\langle \mathbf{a}_i, \hat{\mathbf{a}}_i \rangle), \quad (16)$$

where $\text{acos}(\bullet)$ is the inverse cosine function, $\langle \bullet, \bullet \rangle$ the inner product and $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$.

In addition, the Signal to Interference Ratio (SIR) is employed to measure the accuracy of the estimations of coefficient matrix. It is shown as follows:

$$SIR(s, \hat{s}) = 10 \log \left[\frac{\|s\|_2^2}{\|\hat{s}\|_2^2} \right] (\text{dB}). \quad (17)$$

Consider that \hat{s} is permitted to scale up to a nonzero constant factor $c (c \neq 0)$ with the corresponding s , we rescale \hat{s} to be the same energy level as s before computing its SIR. Usually, when $SIR > 18\text{dB}$, the estimation is acceptable.

The basis matrix is taken by randomly as follows:

$$\mathbf{A} = \begin{pmatrix} 0.7930 & -0.7428 & 0.1404 & 0.9021 \\ 0.1480 & -0.5901 & 0.7010 & -0.3691 \\ -0.5910 & -0.3161 & -0.6992 & -0.2235 \end{pmatrix}, \quad (18)$$

The two dominant components sparse source matrix $S \in R^{4 \times 10000}$ was generated artificially and added Gaussian noise. After 7 iterations, the algorithm converged. It took less than 2 seconds. The initial value $A^{(0)}$ and estimation \hat{A} of A are respectively as follows:

$$A^{(0)} = \begin{pmatrix} 0.3137 & -0.7592 & 1.4124 & 0.9980 \\ 1.3431 & -0.6031 & -0.5058 & 0.1446 \\ -1.3657 & -0.3231 & -0.4132 & -0.7089 \end{pmatrix} \quad \hat{A} = \begin{pmatrix} -0.1405 & -0.7428 & 0.9021 & 0.7931 \\ -0.7010 & -0.5901 & -0.3691 & 0.1481 \\ 0.6992 & -0.3163 & -0.2236 & -0.5908 \end{pmatrix}$$

$BAS(A, \hat{A}) = 5.3933e-004$, and the estimations of 4 components of coefficient matrix are 36.5774dB, 108.0649dB, 56.4985dB and 47.0995. The result is very good.

7 Conclusions

Conventional SCA methods only can deal with single dominant component SCA problem. To overcome the drawback of conventional SCA methods, in this paper the K-subspace clustering algorithm is proposed to perform multiple dominant components SCA. In fact the K-subspace algorithm can be used to sparsely represent the complex valued data. Here we are sorry that we can not give the examples of complex sparse representation for the limited space.

References

- [1] P. G. Georgiev, F. Theis, and A. Cichocki, Sparse component analysis and blind source separation of underdetermined mixtures, *IEEE Trans. on Neural Networks*, July 2005, Vol. 16(4), pp. 992-996.
- [2] L. Parra, C. Spence. Convolutional blind separation of non-stationary sources. *IEEE Trans. Speech and audio processing*, vol. 8(3), pp. 320-327, May 2000.
- [3] M. Aharon, M. Elad, A.M. Bruckstein. K-SVD and its non-negative variant for dictionary design. *Proceedings of the SPIE conference wavelets*, Vol. 5914, July 2005.
- [4] Y.Q. Li, S. Amari, A. Cichocki, and D. W. C. Ho: "Underdetermined Blind Source Separation Based on Sparse Representation", *IEEE Trans. on Signal Processing* (in print).
- [5] Y. Q. Li, A. Cichocki, and S. Amari, "Blind estimation of channel parameters and source components for EEG signals: A sparse factorization approach, *IEEE Transactions on Neural Networks*, 2006, (accepted for publication)
- [6] D. Brandwood, A complex gradient operator and its application in adaptive array theory. *Proc. Inst. Elect. Eng.*, vol 130, pp.11-16, Feb. 1983.
- [7] K. Janich, *Einführung in Die Funktionentheorie*. Berlin, Germany: Springer-Verlag, 1997, ch. 2.
- [8] P. Bofill, M. Zibulevsky. Undetermined blind source separation using sparse representations. *Signal processing*, vol.81, 2353-2362, 2001.
- [9] Y.Q. Li, A. Cichocki, S. Amari. Analysis of sparse representation and blind source separation. *Neural computation*, vol.16, 1193-1234, 2004.
- [10] I. Takigawa, M. Kudo, J. Toyama, Performance analysis of Minimum l_1 -norm solutions for underdetermined source separation, *IEEE Trans. Signal processing*, vol.52(3):582-591, March 2004.
- [11] O. Yilmaz, S. Rickard. Blind separation of speech mixtures via time-frequency masking. *IEEE Trans. Signal Processing*, vol.52(7), pp.1830-1847, 2004.