ON THE CONVERGENCE OF ICA ALGORITHMS WITH SYMMETRIC ORTHOGONALIZATION

Alper T. Erdogan

Koc University, EE Department Istanbul, Turkey

ABSTRACT

We study the convergence behavior of Independent Component Analysis (ICA) algorithms that are based on the contrast function maximization and that employ symmetric orthogonalization method to guarantee the orthogonality property of the search matrix. In particular, the characterization of the critical points of the corresponding optimization problem and the stationary points of the conventional gradient ascent and fixed point algorithms are obtained. As an interesting and a useful feature of the symmetrical orthogonalization method, we show that the use of symmetric orthogonalization enables the monotonic convergence for the fixed point ICA algorithms that are based on the convex contrast functions.

Index Terms— Independent Component Analysis, Blind Source Separation, Symmetric Orthogonalization, Fixed Point Algorithms, Convergence

1. INTRODUCTION

In the area of Independent Component Analysis (ICA), and Blind Source Separation (BSS), the contrast function maximization based algorithms have attracted special attention. The fast fixed point (FastICA) algorithm introduced in [1] has played an important role in the growth of such interest.

The convergence behavior of these algorithms is of both theoretical and practical concern. Recently, Regalia & Kofidis [2] shown the monotonic convergence property of the fixed point ICA algorithms for extracting a single source from a mixture. In addition to the deflationary case, where sources (components) are sequentially obtained, the convergence behavior in the case of simultaneous extraction through a algorithm using symmetrical orthogonalization is of special interest. Oja [3] has shown that a particular set of matrices, where nonzero entries at each column has a constant magnitude, are the fixed points of the FastICA algorithm for the Kurtosis cost function, and among these fixed points, only the matrices corresponding to perfect separation condition are stable.

In this article, our goal is to generalize the characterization of fixed points in [3] for more general contrast functions and for both gradient ascent and fixed point algorithms using symmetric orthogonalization. Furthermore, as an important result, we provide the extension of the work in [2] and show that symmetric orthogonalization enables monotonic convergence to a fixed point in the simultaneous extraction of sources.

2. ICA (BSS) SETUP

We consider the following ICA setup: the p independent components (sources) s_1, \ldots, s_p , which are assumed to have unity variance and zero mean, are mixed through a linear memoryless mapping

$$\underbrace{\begin{bmatrix} y_1(k) \\ \vdots \\ y_q(k) \end{bmatrix}}_{\mathbf{y}(k)} = \mathbf{H} \underbrace{\begin{bmatrix} s_1(k) \\ \vdots \\ s_p(k) \end{bmatrix}}_{\mathbf{s}(k)} \quad k = 1, \dots, \Omega, \quad (1)$$

where

- y_k, k = 1,..., q are the mixtures and q is the number of mixtures (we assume overdetermined case, i.e., q ≥ p),
- $\mathbf{H} \in \Re^{q \times p}$ is the full rank (tall or square) mixing matrix,
- Ω is the number of available samples.

We assume that the mixture sequence $\mathbf{y}(k)$ is prewhitened through the matrix \mathbf{W}_{pre} such that the whitened observations

$$\mathbf{x}(k) = \mathbf{W}_{pre}\mathbf{y}(k) = \underbrace{\mathbf{W}_{pre}\mathbf{H}}_{\mathbf{C}}\mathbf{s}$$
 (2)

has the covariance matrix $E(\mathbf{x}\mathbf{x}^T) = \mathbf{C}\mathbf{C}^T = \mathbf{I}$.

The goal is to obtain an orthogonal separator matrix $\Theta \in \Re^{p \times p}$ such that the overall mapping from the sources to the separator outputs, which is given by

$$\mathbf{G} = \boldsymbol{\Theta} \mathbf{C} \tag{3}$$

is equal to the product of a permutation matrix and a diagonal matrix with unity magnitude entries, i.e.,

$$\mathbf{G} = \mathbf{P}\boldsymbol{\Lambda} \tag{4}$$

This work is supported in part by TUBITAK Career Award, Contract No:104E073.

where the diagonal matrix Λ represents the sign ambiguity and the permutation matrix \mathbf{P} represents the permutation ambiguity inherent in ICA problem.

The algorithms to achieve this goal can be derived from an optimization setting of the form

optimize
$$\mathcal{J}(\Theta)$$

subject to $\Theta\Theta^T = \mathbf{I},$ (5)

which is typically posed as a maximization problem.

In the deflation approach [4], each row is obtained separately where the above optimization problem is divided into a sequence of sub-optimization problems for each row:

optimize
$$\mathcal{J}_k(\Theta_{k,:})$$

subject to $\|\Theta_{k,:}\|_2 = 1,$ (6)

During this process, once a row of Θ is obtained, the corresponding source signal is subtracted from the mixture. The deflation approach has the advantage that the individual suboptimization problems have a unity-2-norm constraint which is much easier to handle compared to the original orthogonal matrix constraint in (5). However, a clear disadvantage of the deflation process is the error accumulation problem, where estimation inaccuracies in the earlier optimization stages cause growing errors for the later stages.

Among the alternative approaches to solve the optimization problem in (5), a class of algorithms simultaneously train all the rows of Θ using a two step procedure:

$$\underline{\Theta}^{(k+1)} = f(\Theta^{(k)}) \qquad (i)
\Theta^{(k)} = \mathcal{M}_O(\underline{\Theta}^{(k+1)}) \qquad (ii)$$
(7)

where

- the Step-(i) is the gradient update, where the two typical choices are
 - The Conventional Gradient Ascent:

$$f(\mathbf{\Theta}^{(k)}) = \mathbf{\Theta}^{(k)} + \mu^{(k)} \mathbf{\Delta}_{\mathbf{\Theta}^{(k)}}^{(k)}$$
(8)

with $\Delta_{\Theta^{(k)}}^{(k)}$ is the gradient of the cost function with respect to $\Theta^{(k)}$.

- Fixed Point Update:

$$f(\mathbf{\Theta}^{(k)}) = \mathbf{\Delta}_{\mathbf{\Theta}^{(k)}}^{(k)} \tag{9}$$

- the Step-(ii) is the mapping to the set of orthogonal matrices, which is typically implemented using
 - Gram-Schmidt Orthogonalization
 - Projection to the closest orthogonal Matrix (according to Frobenius norm sense): In this case, if

the singular value decomposition of $\underline{\Theta}^{(k)}$ is given by,

$$\underline{\Theta}^{(k)} = \underline{\mathbf{U}}^{(k)} \underline{\boldsymbol{\Sigma}}^{(k)} \underline{\mathbf{V}}^{(k)^{T}}$$
(10)

then the minimum distance orthogonal mapping is given by

$$\mathcal{M}_O(\underline{\Theta}^{(k)}) = \underline{\mathbf{U}}^{(k)} \underline{\mathbf{V}}^{(k)T}.$$
 (11)

Note that, alternatively and equivalently, \mathcal{M} can be written as

$$\mathcal{M}(\underline{\Theta}^{(k)})_O = (\underline{\Theta}^{(k)}\underline{\Theta}^{(k)}^T)^{-1/2}\underline{\Theta}^{(k)}$$
(12)

where the matrix-square root above is the symmetric square root. This approach is also referred as symmetric orthogonalization.

In the next section, we investigate some convergence related properties of ICA (BSS) algorithms employing symmetric (minimum distance) orthogonalization.

3. CONVERGENCE PROPERTIES OF SYMMETRICAL ICA ALGORITHMS

In analyzing the convergence behavior of symmetrical ICA algorithms, we first look at the characterization of the potential convergence (or stationary) points. Following that, we'll investigate the monotonic convergence property of the fixed point ICA algorithm with symmetric orthogonalization. Throughout the article, we assume that the cost function \mathcal{J} in (5) is differentiable with full-rank gradient over the set of orthogonal matrices. The case of rank deficient gradient can be similarly approached with a more careful treatment.

3.1. Critical Points of ICA Algorithms with Symmetric Orthogonalization

We start with statement of a fact related to the critical points of the optimization problem in (5):

Theorem 1. Let Θ_* be a local maximum of the problem (5). Then there exists an $\mathbf{S} \in C^{p \times p}$ such that

$$\Delta_{\Theta_*} = \Theta_* \mathbf{S} \text{ and } \mathbf{S} = \mathbf{S}^T.$$
(13)

Proof: According to the Proposition 4.7.3 in [5], if Θ_* is a local maxima of (5), then the gradient Δ_{Θ_*} should be a member of the polar cone of the tangent cone of the constraint set at Θ_* , which is the vector space given by $\{\Theta_* \mathbf{S} \mid \mathbf{S} = \mathbf{S}^T\}$.

Using the linear relation in (3), we can rewrite the condition in (13), in terms of the overall mapping \mathbf{G} as

$$\Delta_{\mathbf{G}_*} = \mathbf{G}_* \mathbf{S}' \text{ and } \mathbf{S}' = {\mathbf{S}'}^T, \tag{14}$$

where $\mathbf{S}' = \mathbf{C}^T \mathbf{S} \mathbf{C}$.

A similar characterization can be provided for the stationary points of the conventional gradient ascent (i.e., (8)) based ICA algorithm with symmetric orthogonalization:

Theorem 2. An orthogonal matrix Θ_* is the stationary point of the ICA algorithm with update rule given in (8) using symmetric orthogonalization if and only if there exists an $\mathbf{S} \in \Re^{p \times p}$ for which the condition in (13) holds and $\lambda_{min}(\mathbf{S})\mu^{(k)} \geq -1 \quad \forall k$.

Proof: Θ_* is the stationary point of the algorithm if and only if the Θ_* is mapped back to itself after the two step update in (7), which leads to

$$\boldsymbol{\Theta}_{*} = \mathcal{M}_{O}(\boldsymbol{\Theta}_{*} + \boldsymbol{\mu}^{(k)} \boldsymbol{\Delta}_{\boldsymbol{\Theta}_{*}})$$
(15)

$$= \mathcal{M}_O(\mathbf{\Theta}_*(\mathbf{I} + \boldsymbol{\mu}^{(k)} \mathbf{\Theta}_*^T \mathbf{\Delta}_{\mathbf{\Theta}_*})) \qquad (16)$$

$$= \boldsymbol{\Theta}_* \mathcal{M}_O (\mathbf{I} + \boldsymbol{\mu}^{(k)} \boldsymbol{\Theta}_*^T \boldsymbol{\Delta}_{\boldsymbol{\Theta}_*})$$
(17)

where the last equality follows from the definition of the mapping \mathcal{M}_O defined by equations (10)-(11) and from the fact that Θ_* is an orthogonal matrix.

As a result, from (17) we obtain,

$$\mathcal{M}_O(\mathbf{I} + \mu^{(k)} \boldsymbol{\Theta}_*^T \boldsymbol{\Delta}_{\boldsymbol{\Theta}_*}) = \mathbf{I},$$
(18)

which implies that $\mathbf{I} + \mu^{(k)} \Theta_*^T \Delta_{\Theta_*}$ has a singular value decomposition given by

$$\mathbf{I} + \mu^{(k)} \boldsymbol{\Theta}_*^{T} \boldsymbol{\Delta}_{\boldsymbol{\Theta}_*} = \mathbf{U}_* \boldsymbol{\Sigma}_* \mathbf{U}_*^{T}, \qquad (19)$$

where \mathbf{U}_* is a orthogonal matrix and $\boldsymbol{\Sigma}_*$ is a nonnegative diagonal matrix. The condition in (19) also implies that $\mathbf{I} + \mu^{(k)} \boldsymbol{\Theta}_*^T \boldsymbol{\Delta}_{\boldsymbol{\Theta}_*}$ is a Hermitian positive-definite matrix, which further implies that $\boldsymbol{\Theta}_*^T \boldsymbol{\Delta}_{\boldsymbol{\Theta}_*}$ is equal to a Hermitian matrix \mathbf{S} for which $\lambda_{min}(\mathbf{S})\mu^{(k)} + 1 \ge 0$.

As a result, the condition for the stationary point for the symmetric ICA algorithm using conventional gradient ascent with appropriate step size is equivalent to the necessary condition for the local maxima of the corresponding problem.

For the fixed point ICA algorithm employing symmetric orthogonalization, the stationary point condition is more restrictive, as stated by the following theorem:

Theorem 3. An orthogonal matrix Θ_* is a stationary point of the fixed point ICA algorithm with symmetric orthogonalization if and only if there exists a <u>positive-definite</u> $\mathbf{S} \in \Re^{p \times p}$ such that the condition in (13) holds.

Proof: If Θ_* is the stationary point then the following condition must hold:

$$\boldsymbol{\Theta}_* = \mathcal{M}_O(\boldsymbol{\Delta}_{\boldsymbol{\Theta}_*}). \tag{20}$$

Assuming that Δ_{Θ_*} has a singular value decomposition

$$\Delta_{\Theta_*} = \mathbf{U}_{\Delta} \boldsymbol{\Sigma}_{\Delta} \mathbf{V}_{\Delta}^T, \qquad (21)$$

then

$$\Theta_* = \mathbf{U}_{\Delta} \mathbf{V}_{\Delta}^T. \tag{22}$$

Therefore, we have

$$\Theta_*^T \Delta_{\Theta_*} = \mathbf{V}_{\Delta} \mathbf{U}_{\Delta}^T \mathbf{U}_{\Delta} \boldsymbol{\Sigma}_{\Delta} \mathbf{V}_{\Delta}^T$$
(23)

$$= \mathbf{V}_{\Delta} \boldsymbol{\Sigma}_{\Delta} \mathbf{V}_{\Delta}^{I} = \mathbf{S}$$
(24)

where \mathbf{S} is a positive-definite matrix.

3.2. Special Case of Kurtosis Maximization

In order to illustrate the results of the previous section, we concentrate on the popular Kurtosis maximization for which the cost function in (5) can be written as

$$\mathcal{J}(\mathbf{\Theta}) = \sum_{k=1}^{p} kurt(z_k)$$
(25)

where z_k 's are the separator outputs and kurt(.) is the kurtosis of its argument. We assume that all components (sources) have positive kurtosis values κ_k . The gradient of this cost function with respect to **G** is given by

$$\Delta_{\mathbf{G}} = 4\mathbf{G}^{\odot 3}\mathbf{K} \tag{26}$$

where $\mathbf{K} = diag(\kappa_1, \kappa_2, \dots, \kappa_p)$ and $\mathbf{G}^{\odot n}$ stands for n^{th} Hadamard power of \mathbf{G} .

We can now outline some specific examples of critical points based on (14) and (26):

Example 1: If **G** satisfies the perfect separation condition in (4), then

$$\Delta_{\mathbf{G}} = 4\mathbf{G}^{\odot 3}\mathbf{K} = 4\mathbf{P}\mathbf{\Lambda}^{3}\mathbf{K} = 4\mathbf{P}\mathbf{\Lambda}\mathbf{K}\mathbf{\Lambda}^{2} \quad (27)$$

$$4\mathbf{G}\mathbf{K}\mathbf{\Lambda}^2$$
 (28)

Therefore, $\mathbf{G}^T \boldsymbol{\Delta}_{\mathbf{G}} = 4\mathbf{K}\boldsymbol{\Lambda}^2$ is a Hermitian matrix, which is also positive definite. Therefore not surprisingly, perfect separation matrices are the critical points of the optimization problem, the stationary points for the symmetrical gradient ascent algorithm for any positive $\mu^{(k)}$ and the stationary points for the symmetrical fixed point algorithm.

Example 2: If G has columns with nonzero entries having a constant magnitude β_j (the example provided in [3]), then

$$\Delta_{\mathbf{G}} = 4\mathbf{G}^{\odot 3}\mathbf{K} = 4\mathbf{G}\mathbf{B}^{2}\mathbf{K}$$
(29)

where $\mathbf{B} = diag(\beta_1, \beta_2, \dots, \beta_p)$. Therefore, $\mathbf{G}^T \boldsymbol{\Delta}_{\mathbf{G}} = 4\mathbf{B}^2 \mathbf{K}$ is a Hermitian positive-definite matrix, therefore, we can make the same comments as in the previous case.

Example 3: As a more interesting example, consider the case where

$$\mathbf{G} = \frac{1}{7} \begin{bmatrix} -2 & 2 & -2 & 2 & 2 & -5 & -2 \\ -2 & 2 & -2 & 2 & 2 & 2 & 5 \\ 2 & -2 & 2 & 5 & -2 & -2 & 2 \\ 2 & 5 & 2 & -2 & -2 & -2 & 2 \\ -2 & 2 & -2 & 2 & -5 & 2 & -2 \\ -2 & 2 & 5 & 2 & 2 & 2 & -2 \\ 5 & 2 & -2 & 2 & 2 & 2 & -2 \end{bmatrix}$$
(30)

and $\mathbf{K} = \mathbf{I}$. If we look at the product $\mathbf{G}^T \boldsymbol{\Delta}_{\mathbf{G}}$, it is equal to

	Γ 103	30	-30	30	30	30	-30 J	
$\frac{4}{7^3}$	30	103	30	-30	-30	-30	30	
	-30	30	103	30	30	30	-30	
	30	-30	30	103	-30	-30	30	(31)
	30	-30	30	-30	103	-30	30	
	30	-30	30	-30	-30	103	30	
		30	-30	30	30	30	103	

which is clearly a symmetric matrix, and therefore, **G** in (30) is a critical point. Furthermore, since the smallest eigenvalue of the matrix in (31) is equal to $\frac{-44}{49}$, it would be a stationary point for the conventional gradient ascent algorithm if $\mu^{(k)} < \frac{49}{44}$. Due to the fact that the symmetric matrix in (31) isn't positive definite, **G** in (30) is not a stationary point of the fixed point algorithm. In fact, if we apply one iteration of fixed point update (with symmetric orthogonalization) we obtain

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix},$$
(32)

as new G, which is a perfect separation point.

3.3. Monotonic Convergence of Fixed Point ICA Algorithms with Symmetric Orthogonalization

In order to show the monotonic convergence property of the Fixed Point ICA Algorithm employing symmetrical orthogonalization, we follow a treatment similar to the one in [2]. The basic result is summarized by the following theorem:

Theorem 4. Fixed Point ICA algorithm with symmetrical orthogonalization corresponding to the optimization setting in (5) where \mathcal{J} is a convex cost function (bounded on the set of orthogonal matrices), is monotonically convergent to one of the stationary points defined by Theorem 3.

Proof: Given \mathcal{J} is a differentiable convex function, for any $\mathbf{G}^{(k)}, \mathbf{G}^{(k+1)}$ pair, we have

$$\mathcal{J}(\mathbf{G}^{(k+1)}) \geq \mathcal{J}(\mathbf{G}^{(k)}) + \operatorname{Tr}(\mathbf{\Delta}_{\mathbf{G}^{(k)}}^T(\mathbf{G}^{(k+1)} - \mathbf{G}^{(k)})). (33)$$

If we now look at the optimization problem

maximize
$$\operatorname{Tr}(\boldsymbol{\Delta}_{\mathbf{G}}^{T})$$

subject to $\mathbf{G}\mathbf{G}^{T} = \mathbf{I},$ (34)

according to Proposition 4.7.3 in [5], a local optimum G_* of the problem in (34) should satisfy

$$\mathbf{G}_*\mathbf{S} = \boldsymbol{\Delta}_{\mathbf{G}^{(k)}} \quad \text{for some } \mathbf{S} = \mathbf{S}^T.$$
(35)

From (35), we can write

$$\mathbf{S}^T \mathbf{S} = \mathbf{\Delta}_{\mathbf{G}^{(k)}}^T \mathbf{\Delta}_{\mathbf{G}^{(k)}}$$
(36)

$$= \mathbf{V}_{\Delta} \boldsymbol{\Sigma}_{\Delta} \mathbf{U}_{\Delta}^{T} \mathbf{U}_{\Delta} \boldsymbol{\Sigma}_{\Delta} \mathbf{V}_{\Delta}^{T} \qquad (37)$$
$$= \mathbf{V}_{\Delta} \boldsymbol{\Sigma}_{\Delta}^{2} \mathbf{V}_{\Delta}^{T}. \qquad (38)$$

From which we conclude that ${f S}$ is a Hermitian matrix with

$$\mathbf{S} = \mathbf{V}_{\Delta} \mathbf{\Lambda} \mathbf{V}_{\Delta}^T \tag{39}$$

where $\Lambda = \Sigma_{\Delta} J$ and J is a diagonal matrix with 1's,and/or -1s on the diagonal. Note that the cost function value at the critical point G_* is equal to

$$Tr(\mathbf{\Delta}_{\mathbf{G}^{(k)}}^T \mathbf{G}_*) = Tr(\mathbf{S})$$
(40)

$$= \operatorname{Tr}(\mathbf{V}_{\Delta}\Sigma_{\Delta}\mathbf{J}\mathbf{V}_{\Delta}^{T}) \qquad (41)$$

$$= \operatorname{Tr}(\Sigma_{\Delta} \mathbf{J}). \tag{42}$$

Among all critical points specified by (35), the global maximum value is achieved for J = I. Note that this case corresponds to

=

$$\mathbf{G}_* = \mathbf{U}_{\boldsymbol{\Delta}} \mathbf{V}_{\boldsymbol{\Delta}}^T \tag{43}$$

$$= \mathcal{M}_O(\mathbf{\Delta}_{\mathbf{G}^{(\mathbf{k})}}), \tag{44}$$

i.e., the global maximum point for the problem in (34) is obtained by projecting the gradient $\Delta_{\mathbf{G}^{(k)}}$ to the set of orthogonal matrices using mapping \mathcal{M}_O (symmetric orthogonalization). Furthermore, the value of $\operatorname{Tr}(\Delta_{\mathbf{G}^{(k)}}^T \mathbf{G}_*)$ at the global maximum point is strictly greater than its values at other critical points obtained for the choices of $\mathbf{J} \neq \mathbf{I}$. Therefore, given $\mathbf{G}^{(k)}$ is not a stationary point of the algorithm, i.e., $\mathbf{G}^{(k)} \neq \mathcal{M}_O(\Delta_{\mathbf{G}^{(k)}})$, the choice

$$\mathbf{G}^{(k+1)} = \mathcal{M}_O(\mathbf{\Delta}_{\mathbf{G}^{(\mathbf{k})}}), \tag{45}$$

combined with (33), guarantees that $\mathcal{J}(\mathbf{G}^{(k+1)}) > \mathcal{J}(\mathbf{G}^{(k)})$. This fact, together with the boundedness of \mathcal{J} , implies the convergence.

4. REFERENCES

- Aapo Hyvärinen and Erkki Oja, "A fast fixed-point algorithm for independent component analysis," *Neural Computation*, vol. 9, no. 7, pp. 1483–1492, October 1997.
- [2] P. A. Regalia and E. Kofidis, "Monotonic convergence of fixed-point algorithms for ICA," *IEEE Transactions on Neural Networks*, vol. 14, no. 4, pp. 943–949, July 2003.
- [3] Erikki Oja, "Convergence of the symmetrical FastICA algorithm," Proc. of the 9th Int. Conf. on Neural Information Processing, vol. 3, pp. 1368–1372, November 2002.
- [4] P. Loubaton and Phillip Regalia, "Blind deconvolution of multivariate signals: a deflation approach," *IEEE International Conference on Communications*, vol. 2, pp. 1160–1164, May 1993.
- [5] D. P. Bertsekas, A. Nedic, and A. Ozdaglar, *Convex Analysis and Optimization*, Athena Scientific, 2003.