
Multiplicative Updating Rule for Blind Separation Derived from the Method of Scoring

Howard Hua Yang
Department of Computer Science
Oregon Graduate Institute
PO Box 91000, Portland, OR 97291, USA
hyang@cse.ogi.edu

Abstract

For blind source separation, when the Fisher information matrix is used as the Riemannian metric tensor for the parameter space, the steepest descent algorithm to maximize the likelihood function in this Riemannian parameter space becomes the serial updating rule with equivariant property. This algorithm can be further simplified by using the asymptotic form of the Fisher information matrix around the equilibrium.

1 Introduction

The relative gradient was introduced by (Cardoso and Laheld, 1996) to design multiplicative updating algorithms with equivariant property for blind separation problems. The idea is to calculate differentials by using a relative increment instead of an absolute increment in the parameter space. This idea has been extended to compute the relative Hessian by (Pham, 1996).

For a matrix function $f = f(\mathbf{W})$, the relative gradient is defined by

$$\hat{\nabla} f = \frac{\partial f}{\partial \mathbf{W}} \mathbf{W}^T. \quad (1)$$

From the differential of $f(\mathbf{W})$ based on the relative gradient, the following learning rule is given by (Cardoso and Laheld, 1996) to maximize the function f :

$$\frac{d\mathbf{W}}{dt} = \eta \hat{\nabla} f \mathbf{W} = \eta \frac{\partial f}{\partial \mathbf{W}} \mathbf{W}^T \mathbf{W} \quad (2)$$

Also motivated by designing blind separation algorithms with equivariant property,

the natural gradient defined by

$$\tilde{\nabla} f = \frac{\partial f}{\partial \mathbf{W}} \mathbf{W}^T \mathbf{W} \quad (3)$$

was introduced in (Amari et al, 1996) which yields the same learning rule (2). The geometrical meaning of the natural gradient is given by (Amari, 1996). More details about the natural gradient can be found in (Yang and Amari, 1997) and (Amari, 1997).

The framework of the natural gradient learning was proposed by (Amari, 1997). In this framework, the ordinary gradient descent learning algorithm in the Euclidean space is not optimal in minimizing a function defined in a Riemannian space. The ordinary gradient should be replaced by the natural gradient which is defined by operating the inverse of the metric tensor in the Riemannian space on the ordinary gradient. Let \mathbf{w} denote a parameter vector. It is proved by (Amari, 1997) that if $C(\mathbf{w})$ is a loss function defined on a Riemannian space $\{\mathbf{w}\}$ with a metric tensor \mathbf{G} , the negative natural gradient of $C(\mathbf{w})$, namely, $-\mathbf{G}^{-1} \frac{\partial C}{\partial \mathbf{w}}$ is the steepest descent direction to decrease this function in the Riemannian space. Therefore, the steepest descent algorithm in this Riemannian space has the following form:

$$\frac{d\mathbf{w}}{dt} = -\eta \mathbf{G}^{-1} \frac{\partial C}{\partial \mathbf{w}}.$$

If the Fisher information matrix is used as the metric tensor for the Riemannian space and $C(\mathbf{w})$ is replaced by the negative log-likelihood function, the above learning rule becomes the method of scoring (Kay, 1993) which is the focus of this paper.

Both the relative gradient $\hat{\nabla}$ and the natural gradient $\tilde{\nabla}$ were proposed in order to design the multiplicative updating algorithms with the equivariant property. The former is due to a multiplicative increment in calculating differential while the latter is due to an increment based on a nonholonomic basis (Amari, 1997). Neither $\hat{\nabla}$ nor $\tilde{\nabla}$ depends on the data model. The Fisher information matrix is a special and important choice for the Riemannian metric tensor for statistical estimation problems. It depends on the data model. Operating the inverse of the Fisher information matrix on the ordinary gradient, we have another gradient operator. It is called a natural gradient induced by the Fisher information matrix.

In this paper, we show how to derive a multiplicative updating algorithm from the method of scoring. This approach is different from those based on the relative gradient and the natural gradient defined by (3).

2 Fisher Information Matrix For Blind Separation

Consider a linear mixing system:

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

where $\mathbf{A} \in \mathfrak{R}^{n \times n}$, $\mathbf{x} = (x_1, \dots, x_n)^T$ and $\mathbf{s} = (s_1, \dots, s_n)^T$. Assume that sources are independent with a factorized joint pdf:

$$r(\mathbf{s}) = \prod_{i=1}^n r(s_i).$$

The likelihood function is

$$p(\mathbf{x}; \mathbf{A}) = \frac{r(\mathbf{A}^{-1} \mathbf{x})}{|\mathbf{A}|}$$

where $|\mathbf{A}| = |\det(\mathbf{A})|$. Let $\mathbf{W} = \mathbf{A}^{-1}$ and $\mathbf{y} = \mathbf{W}\mathbf{x}$ (a demixing system), then we have the log-likelihood function

$$L(\mathbf{W}) = \sum_{i=1}^n \log r_i(\mathbf{y}_i) + \log |\mathbf{W}|.$$

It is easy to obtain

$$\frac{\partial L}{\partial w_{ij}} = \frac{r'_i(\mathbf{y}_i)}{r_i(\mathbf{y}_i)} x_j + \mathbf{W}_{ij}^{-T} \quad (4)$$

where \mathbf{W}_{ij}^{-T} is the (i, j) entry in $\mathbf{W}^{-T} = (\mathbf{W}^{-1})^T$. Writing (4) in a matrix form, we have

$$\frac{\partial L}{\partial \mathbf{W}} = \mathbf{W}^{-T} - \Phi(\mathbf{y})\mathbf{x}^T = (\mathbf{I} - \Phi(\mathbf{y})\mathbf{y}^T)\mathbf{W}^{-T} = \mathbf{F}(\mathbf{y})\mathbf{W}^{-T} \quad (5)$$

where $\Phi(\mathbf{y}) = (\phi_1(\mathbf{y}_1), \dots, \phi_n(\mathbf{y}_n))^T$, $\phi_i(\mathbf{y}_i) = -\frac{r'_i(\mathbf{y}_i)}{r_i(\mathbf{y}_i)}$ and $\mathbf{F}(\mathbf{y}) = \mathbf{I} - \Phi(\mathbf{y})\mathbf{y}^T$.

The maximum likelihood algorithm based on the ordinary gradient $\frac{\partial L}{\partial \mathbf{W}}$ is

$$\frac{d\mathbf{W}}{dt} = \eta(\mathbf{I} - \Phi(\mathbf{y})\mathbf{y}^T)\mathbf{W}^{-T} = \eta\mathbf{F}(\mathbf{y})\mathbf{W}^{-T}$$

which has the high computational complexity due to the matrix inverse \mathbf{W}^{-1} . The maximum likelihood algorithm based on the natural gradient of matrix functions is

$$\frac{d\mathbf{W}}{dt} = \eta\tilde{\nabla}L = \eta(\mathbf{I} - \Phi(\mathbf{y})\mathbf{y}^T)\mathbf{W}. \quad (6)$$

The same algorithm is obtained from $\frac{d\mathbf{W}}{dt} = \eta\hat{\nabla}L\mathbf{W}$ by using the relative gradient. An apparent reason for using this algorithm is to avoid the matrix inverse \mathbf{W}^{-1} . Another good reason for using it is due to the fact that the matrix \mathbf{W} driven by (6) never becomes singular if the initial matrix \mathbf{W} is not singular. This is proved by (Yang and Amari, 1997). In fact, this property holds for any learning rule of the following type:

$$\frac{d\mathbf{W}}{dt} = \mathbf{H}(\mathbf{y})\mathbf{W}. \quad (7)$$

Let $\langle \mathbf{U}, \mathbf{V} \rangle = \text{Tr}(\mathbf{U}^T\mathbf{V})$ denote the inner product of \mathbf{U} and $\mathbf{V} \in \Re^{n \times n}$. When $\mathbf{W}(t)$ is driven by the equation (7), we have

$$\begin{aligned} \frac{d|\mathbf{W}|}{dt} &= \langle \frac{\partial |\mathbf{W}|}{\partial \mathbf{W}}, \frac{d\mathbf{W}}{dt} \rangle = \langle |\mathbf{W}|(\mathbf{W}^{-1})^T, \frac{d\mathbf{W}}{dt} \rangle \\ &= \text{Tr}(|\mathbf{W}|\mathbf{W}^{-1}\mathbf{H}(\mathbf{y})\mathbf{W}) = \text{Tr}(\mathbf{H}(\mathbf{y}))|\mathbf{W}|. \end{aligned}$$

Therefore,

$$|\mathbf{W}(t)| = |\mathbf{W}(0)| \exp\left\{ \int_0^t \text{Tr}(\mathbf{H}(\mathbf{y}(\tau)))d\tau \right\} \quad (8)$$

which is non-singular when the initial matrix $\mathbf{W}(0)$ is non-singular.

The matrix function $\mathbf{F}(\mathbf{y})$ is also called an estimating function. At the equilibrium of the system (6), it satisfies the zero condition $E[\mathbf{F}(\mathbf{y})] = \mathbf{0}$, i.e.,

$$E[\phi_i(\mathbf{y}_i)\mathbf{y}_j] = \delta_{ij} \quad (9)$$

where $\delta_{ij} = 1$ if $i = j$ and 0 otherwise.

To calculate the Fisher information matrix, we need a vector form of the equation (5). Let $\text{Vec}(\cdot)$ denote an operator on a matrix which cascades the columns of the

matrix from the left to the right and forms a column vector. This operator has the following property:

$$\text{Vec}(\mathbf{ABC}) = (\mathbf{C}^T \otimes \mathbf{A})\text{Vec}(\mathbf{B}) \quad (10)$$

where \otimes denotes the Kronecker product. Applying this property, we first rewrite (5) as

$$\frac{\partial L}{\partial \text{Vec}(\mathbf{W})} = \text{Vec}\left(\frac{\partial L}{\partial \mathbf{W}}\right) = (\mathbf{W}^{-1} \otimes \mathbf{I})\text{Vec}(\mathbf{F}(\mathbf{y})), \quad (11)$$

and then obtain the Fisher information matrix

$$\begin{aligned} \mathbf{G} &= E\left[\frac{\partial L}{\partial \text{Vec}(\mathbf{W})}\left(\frac{\partial L}{\partial \text{Vec}(\mathbf{W})}\right)^T\right] \\ &= (\mathbf{W}^{-1} \otimes \mathbf{I})E[\text{Vec}(\mathbf{F}(\mathbf{y}))\text{Vec}^T(\mathbf{F}(\mathbf{y}))](\mathbf{W}^{-T} \otimes \mathbf{I}). \end{aligned} \quad (12)$$

The inverse of \mathbf{G} is

$$\mathbf{G}^{-1} = (\mathbf{W}^T \otimes \mathbf{I})\mathbf{D}^{-1}(\mathbf{W} \otimes \mathbf{I}) \quad (13)$$

where $\mathbf{D} = E[\text{Vec}(\mathbf{F}(\mathbf{y}))\text{Vec}^T(\mathbf{F}(\mathbf{y}))]$.

3 Natural Gradient Induced By Fisher Information Matrix

Define a Riemannian space

$$\mathcal{V} = \{\text{Vec}(\mathbf{W}); \mathbf{W} \in \text{Gl}(n)\}$$

in which the Fisher information matrix \mathbf{G} is used as its metric. Here, $\text{Gl}(n)$ is the space of all the $n \times n$ invertible matrices.

Let $C(\mathbf{W})$ be a matrix function to be minimized. It is shown by (Amari, 1997) that the steepest descent direction in the Riemannian space \mathcal{V} is $-\mathbf{G}^{-1} \frac{\partial C}{\partial \text{Vec}(\mathbf{W})}$.

Let us define the natural gradient in \mathcal{V} by

$$\bar{\nabla} C(\mathbf{W}) = (\mathbf{W}^T \otimes \mathbf{I})\mathbf{D}^{-1}(\mathbf{W} \otimes \mathbf{I}) \frac{\partial C}{\partial \text{Vec}(\mathbf{W})} \quad (14)$$

which is called the natural gradient induced by the Fisher information matrix. The time complexity of computing the natural gradient in the space \mathcal{V} is high since inverting the matrix \mathbf{D} of $n^2 \times n^2$ is needed.

Using the natural gradient in \mathcal{V} to maximize the likelihood function $L(\mathbf{W})$ or the method of scoring, from (11) and (14) we have the following learning rule

$$\text{Vec}\left(\frac{d\mathbf{W}}{dt}\right) = \eta(\mathbf{W}^T \otimes \mathbf{I})\mathbf{D}^{-1}\text{Vec}(\mathbf{F}(\mathbf{y})) \quad (15)$$

We shall prove that the above learning rule has the equivariant property.

Denote Vec^{-1} the inverse of the operator Vec . Let matrices \mathbf{B} and \mathbf{A} be of $n^2 \times n^2$ and $n \times n$, respectively. Denote $\mathbf{B}(i, \cdot)$ the i -th row of \mathbf{B} and $\mathbf{B}_i = \text{Vec}^{-1}(\mathbf{B}(i, \cdot))$, $i = 1, \dots, n^2$. Define an operator $\mathbf{B} \star$ as a mapping from $\mathfrak{R}^{n \times n}$ to $\mathfrak{R}^{n \times n}$:

$$\mathbf{B} \star \mathbf{A} = \begin{bmatrix} \langle \mathbf{B}_1, \mathbf{A} \rangle & \cdots & \langle \mathbf{B}_{n^2-n+1}, \mathbf{A} \rangle \\ \cdots & \cdots & \cdots \\ \langle \mathbf{B}_n, \mathbf{A} \rangle & \cdots & \langle \mathbf{B}_{n^2}, \mathbf{A} \rangle \end{bmatrix}$$

where $\langle \cdot, \cdot \rangle$ is the inner product in $\mathfrak{R}^{n \times n}$. With the operation \star , we have

$$\mathbf{B}\text{Vec}(\mathbf{A}) = \begin{bmatrix} \langle \mathbf{B}_1, \mathbf{A} \rangle \\ \vdots \\ \langle \mathbf{B}_{n^2}, \mathbf{A} \rangle \end{bmatrix} = \text{Vec}\left(\text{Vec}^{-1}\left(\begin{bmatrix} \langle \mathbf{B}_1, \mathbf{A} \rangle \\ \vdots \\ \langle \mathbf{B}_{n^2}, \mathbf{A} \rangle \end{bmatrix}\right)\right) = \text{Vec}(\mathbf{B} \star \mathbf{A}),$$

i.e.,

$$B\text{Vec}(\mathbf{A}) = \text{Vec}(\mathbf{B} \star \mathbf{A}).$$

Applying the above relation, we first rewrite the equation (15) as

$$\text{Vec}\left(\frac{d\mathbf{W}}{dt}\right) = \eta(\mathbf{W}^T \otimes \mathbf{I})\text{Vec}(\mathbf{D}^{-1} \star \mathbf{F}(\mathbf{y})),$$

then applying (10) to the above equation we obtain

$$\frac{d\mathbf{W}}{dt} = \eta(\mathbf{D}^{-1} \star \mathbf{F}(\mathbf{y}))\mathbf{W}. \quad (16)$$

Theorem 1 *For the blind separation problem, the maximum likelihood algorithm based on the natural gradient induced by the Fisher information matrix or the method of scoring has the form (16) which is a multiplicative updating rule with the equivariant property.*

To implement the algorithm (16), we estimate \mathbf{D} by sample average. Let $f_{ij}(\mathbf{y})$ be the (i, j) entry in $\mathbf{F}(\mathbf{y})$. A general form for the entries in \mathbf{D} is

$$d_{ij,kl} = E[f_{ij}(\mathbf{y})f_{kl}(\mathbf{y})]$$

which depends on the source pdfs $r_i(s_i)$. When the source pdfs are unknown, in practice we choose $r_i(s_i)$ as our prior assumptions about the source pdfs. To simplify the algorithm (16), we replace \mathbf{D} by its asymptotic form at the solution points $\mathbf{a} = (c_1 s_{\sigma(1)}, \dots, c_n s_{\sigma(n)})^T$ where $(\sigma(1), \dots, \sigma(n))$ is a permutation of $(1, \dots, n)$.

Regarding the structure of the asymptotic \mathbf{D} , we have the following theorem:

Theorem 2 *Assume that the pdfs of the sources s_i are even functions.*

Then at the solution point $\mathbf{a} = (c_1 s_{\sigma(1)}, \dots, c_n s_{\sigma(n)})^T$, \mathbf{D} is a diagonal matrix and its n^2 diagonal entries have two forms, namely,

$$\begin{aligned} E[f_{ij}(\mathbf{a})f_{ij}(\mathbf{a})] &= \mu_i \lambda_j, \quad \text{for } i \neq j \text{ and} \\ E[(f_{ii}(\mathbf{a}))^2] &= \nu_i \end{aligned}$$

where $\mu_i = E[\phi_i^2(a_i)]$, $\lambda_i = E[a_i^2]$ and $\nu_i = E[\phi_i^2(a_i)a_i^2] - 1$. More concisely, we have

$$\mathbf{D} = \text{diag}(\text{Vec}(\mathbf{H})) \quad (17)$$

where

$$\mathbf{H} = (\mu_i \lambda_j)_{n \times n} - \text{diag}(\mu_1 \lambda_1, \dots, \mu_n \lambda_n) + \text{diag}(\nu_1, \dots, \nu_n)$$

The proof of Theorem 2 is given in Appendix 1.

Let $\mathbf{H} = (h_{ij})_{n \times n}$. Since all μ_i, λ_i , and ν_i are positive, and so are all h_{ij} . We define

$$\frac{1}{\mathbf{H}} = \left(\frac{1}{h_{ij}}\right)_{n \times n}.$$

Then from (17), we have

$$\mathbf{D}^{-1} = \text{diag}\left(\text{Vec}\left(\frac{1}{\mathbf{H}}\right)\right).$$

The results in Theorem 2 enable us to simplify the algorithm (16) to obtain a low complexity learning rule. Since \mathbf{D}^{-1} is a diagonal matrix, for any $n \times n$ matrix \mathbf{A} we have

$$\mathbf{D}^{-1}\text{Vec}(\mathbf{A}) = \text{Vec}\left(\frac{1}{\mathbf{H}} \odot \mathbf{A}\right) \quad (18)$$

where \odot denotes the componentwise multiplication of two matrices of the same dimension. Applying (18) to the learning rule (15), we obtain the following learning rule

$$\text{Vec}\left(\frac{d\mathbf{W}}{dt}\right) = \eta(\mathbf{W}^T \otimes \mathbf{I})\text{Vec}\left(\frac{1}{\mathbf{H}} \odot \mathbf{F}(\mathbf{y})\right).$$

Again, applying (10) to the above equation we have the following learning rule

$$\frac{d\mathbf{W}}{dt} = \eta\left(\frac{1}{\mathbf{H}} \odot \mathbf{F}(\mathbf{y})\right)\mathbf{W}. \quad (19)$$

Like the learning rule (16), the algorithm (19) is also multiplicative; but unlike (16), there is no need to inverse the $n^2 \times n^2$ matrix in (19). The computation of $\frac{1}{\mathbf{H}}$ is straightforward by computing the reciprocals of the entries in \mathbf{H} .

$(\mu_i, \lambda_i, \nu_i)$ are $3n$ unknowns in \mathbf{G} . Let us impose the following constraint

$$\nu_i = \mu_i \lambda_i. \quad (20)$$

Under this constraint, the number of unknowns in \mathbf{G} is $2n$, and \mathbf{D} can be written as

$$\mathbf{D} = \mathbf{D}_\lambda \otimes \mathbf{D}_\mu \quad (21)$$

where $\mathbf{D}_\lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ and $\mathbf{D}_\mu = \text{diag}(\mu_1, \dots, \mu_n)$.

From (14), using (21) we have the natural gradient descent rule in the Riemannian space \mathcal{V}

$$\frac{d\text{Vec}(\mathbf{W})}{dt} = -\eta(\mathbf{W}^T \mathbf{D}_\lambda^{-1} \mathbf{W} \otimes \mathbf{D}_\mu^{-1}) \frac{\partial C}{\partial \text{Vec}(\mathbf{W})}. \quad (22)$$

Applying the property (10), we rewrite the above equation in a matrix form

$$\frac{d\mathbf{W}}{dt} = -\eta \mathbf{D}_\mu^{-1} \frac{\partial C}{\partial \mathbf{W}} \mathbf{W}^T \mathbf{D}_\lambda^{-1} \mathbf{W}. \quad (23)$$

Since μ_i and λ_i are unknown, \mathbf{D}_μ and \mathbf{D}_λ are replaced by the identity matrix in practice. Therefore, the algorithm (2) is an approximation of the algorithm (23).

Taking $C = -L(\mathbf{W})$ as the negative likelihood function and applying the expression (5), we have the following maximum likelihood algorithm based on the natural gradient in \mathcal{V} :

$$\frac{d\mathbf{W}}{dt} = \eta \mathbf{D}_\mu^{-1} (\mathbf{I} - \Phi(\mathbf{y})\mathbf{y}^T) \mathbf{D}_\lambda^{-1} \mathbf{W}. \quad (24)$$

Again, replacing \mathbf{D}_μ and \mathbf{D}_λ by the identity matrix we obtain the maximum likelihood algorithm (6) based on the relative gradient or natural gradient of matrix functions.

In the context of the blind separation, the source pdfs are unknown. The prior assumption $r_i(s_i)$ used to define the functions $\phi_i(\mathbf{y}_i)$ may not match the true pdfs of the sources. However, the algorithm (24) is generally robust to the mismatch between the true pdfs and the pdfs employed by the algorithm if the mismatch is not too large. See (Cardoso, 1997) and (Pham, 1996) for example.

4 Conclusion

In the context of blind separation, when the Fisher information matrix is used as the Riemannian metric tensor for the parameter space, maximizing the likelihood function in this Riemannian space based on the steepest descent method is the method of scoring. This method yields a multiplicative updating rule with the equivariant property. It is further simplified by using the asymptotic form of the Fisher information matrix around the equilibrium.

5 Appendix

Appendix 1 Proof of Theorem 2:

By definition $f_{ij}(\mathbf{y}) = \delta_{ij} - \phi_i(\mathbf{y}_i)y_j$. At the equilibrium $\mathbf{a} = (c_1 s_{\sigma(1)}, \dots, c_n s_{\sigma(n)})^T$, we have $E[\phi_i(\mathbf{a}_i)a_j] = 0$ for $i \neq j$ and $E[\phi_i(\mathbf{a}_i)a_i] = 1$. So $E[f_{ij}(\mathbf{a})] = 0$. Since the source pdfs are even functions, we have $E[a_i] = 0$ and $E[\phi_i(\mathbf{a}_i)] = 0$. Applying these equalities, it is not difficult to verify that

$$E[f_{ij}(\mathbf{a})f_{kl}(\mathbf{a})] = 0, \text{ for } (i, j) \neq (k, l). \quad (25)$$

So, \mathbf{D} is a diagonal matrix and

$$E[f_{ii}(\mathbf{a})f_{ii}(\mathbf{a})] = E[(1 - \phi_i(\mathbf{a}_i)a_i)^2] = E[\phi_i^2(\mathbf{a}_i)a_i^2] - 1,$$

$$E[f_{ij}(\mathbf{a})f_{ij}(\mathbf{a})] = E[\phi_i^2(\mathbf{a}_i)a_j^2] = \mu_i \lambda_j$$

for $i \neq j$.

Q.E.D.

References

- [1] S. Amari. Natural gradient works efficiently in learning. *Accepted by Neural Computation*, 1997.
- [2] S. Amari. Neural learning in structured parameter spaces – natural Riemannian gradient. In *Advances in Neural Information Processing Systems, 9*, ed. M. C. Mozer, M. I. Jordan and T. Petsche, The MIT Press: Cambridge, MA., pages 127–133, 1997.
- [3] S. Amari, A. Cichocki, and H. H. Yang. A new learning algorithm for blind signal separation. In *Advances in Neural Information Processing Systems, 8*, eds. David S. Touretzky, Michael C. Mozer and Michael E. Hasselmo, MIT Press: Cambridge, MA., pages 757–763, 1996.
- [4] J.-F. Cardoso. Infomax and maximum likelihood for blind source separation. *IEEE Signal Processing Letters*, April 1997.
- [5] J.-F. Cardoso and B. Laheld. Equivariant adaptive source separation. *IEEE Trans. on Signal Processing*, 44(12):3017–3030, December 1996.
- [6] S. M. Kay. *Fundamentals of Statistical Signal Processing: Estimation Theory*. PTR Prentice Hall, Englewood Cliffs, 1993.
- [7] D. T. Pham. Blind separation of instantaneous mixture of sources via an ica. *IEEE Trans. on Signal Processing*, 44(11):2768–2779, November 1996.
- [8] H. H. Yang and S. Amari. Adaptive on-line learning algorithms for blind separation: Maximum entropy and minimum mutual information. *Neural Computation*, 9(7):1457–1482, 1997.