

DIFFERENTIAL LEARNING AND ICA

Seungjin Choi

Department of Computer Science, POSTECH, Korea
seungjin@postech.ac.kr

ABSTRACT

Differential learning relies on the differentiated values of nodes, whereas the conventional learning depends on the values themselves of nodes. In this paper, I elucidate the differential learning in the framework of maximum likelihood learning of linear generative model with latent variables obeying random walk. I apply the idea of differential learning to the problem independent component analysis (ICA), which leads to *differential ICA*. Algorithm derivation using the natural gradient and local stability analysis are provided. Usefulness of the algorithm is emphasized in the case of blind separation of temporally correlated sources and is demonstrated through a simple numerical example.

1. INTRODUCTION

Independent component analysis (ICA) is a statistical method, the goal of which is to learn non-orthogonal basis vectors from a set of observation data with basis coefficients being statistically independent. In the framework of linear transform, ICA finds a representation of the form

$$\begin{aligned} \mathbf{x}(t) &= \sum_{i=1}^n \mathbf{a}_i s_i(t) \\ &= \mathbf{A} \mathbf{s}(t), \end{aligned} \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the observation data (which is given) and $\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_n] \in \mathbb{R}^{n \times n}$ (which is known as a *mixing matrix* in source separation) consists of basis vectors $\{\mathbf{a}_i\}$ and $\mathbf{s} = [s_1 \cdots s_n]$ is an n -dimensional vector containing basis coefficients $\{s_i\}$ (which are called *independent components* and are also known as *sources*).

It is known that ICA performs source separation, the goal of which is to restore unknowns sources without resorting to any prior knowledge, given only a set of observation data. Source separation is achieved by estimating the mixing matrix \mathbf{A} or its inverse $\mathbf{W} = \mathbf{A}^{-1}$ (which is known as *demixing matrix*).

Let $\mathbf{y}(t)$ be the output of demixing transform, i.e.,

$$\mathbf{y}(t) = \mathbf{W} \mathbf{x}(t). \quad (2)$$

Either maximum likelihood estimation or the minimization of mutual information leads to the well-known natural gradient ICA algorithm [1] whose updating rule has the form

$$\mathbf{W}(t+1) = \mathbf{W}(t) + \eta \left\{ \mathbf{I} - \varphi(\mathbf{y}(t)) \mathbf{y}^T(t) \right\} \mathbf{W}(t), \quad (3)$$

where η is a learning rate and $\varphi(\mathbf{y}) = [\varphi_1(y_1) \cdots \varphi_n(y_n)]^T$ is an n -dimensional vector, each element of which corresponds

to the negative score function, i.e., $\varphi_i(y_i) = -\frac{d \log p_i(y_i)}{dy_i}$ where $p_i(\cdot)$ is the hypothesized probability density function for s_i . More details on ICA or source separation can be found in [2, 3] (and references therein).

In a wide sense, most of ICA algorithms based on unsupervised learning belong to Hebb-type rule or its generalization with adopting nonlinear functions. Motivated from differential Hebb's rule [4] and differential decorrelation [5], we develop an ICA algorithm which employs differential learning where learning resorts to differentiated values (or difference of values in discrete-time counterpart).

Differential Hebb's rule was studied as an alternative to the Hebb's rule. The motivation of the differential Hebb's rule is that concurrent change, rather than just concurrent activation, more accurately captures the *concomitant variation*. The differential learning was introduced in the framework of ICA [6] and decorrelation [5] recently. In this paper we derive a differential ICA algorithm in the framework of maximum likelihood estimation and random walk model. In fact, our differential ICA algorithm can be viewed as a simpler form of ICA algorithms which exploit the temporal structure of sources [7, 8].

2. RANDOM WALK MODEL FOR LATENT VARIABLES

Given a set of observation data, $\{\mathbf{x}(t)\}$, the task of learning the linear generative model (1) under a constraint that latent variables being statistically independent, is a semi-parametric estimation problem. The maximum likelihood estimation of basis vectors $\{\mathbf{a}_i\}$ is involved with a probabilistic model for latent variables which are treated as nuisance parameters.

In order to show a link between the differential learning and maximum likelihood estimation, we consider a random walk model for latent variables which is a simple Markov chain, i.e.,

$$s_i(t) = s_i(t-1) + \epsilon_i(t), \quad (4)$$

where the innovation $\epsilon_i(t)$ is assumed to have zero mean with a density function $q_i(\epsilon_i(t))$. In addition, innovation sequences $\{\epsilon_i(t)\}$ are assumed to be mutually independent.

Let us consider the latent variables $s_i(t)$ over N -point time block. We define the vector \underline{s}_i as

$$\underline{s}_i = [s_i(0), \dots, s_i(N-1)]^T. \quad (5)$$

Then the joint probability density function of \underline{s}_i can be

written as

$$\begin{aligned} p_i(\underline{s}_i) &= p_i(s_i(0), \dots, s_i(N-1)) \\ &= \prod_{t=0}^{N-1} p_i(s_i(t)|s_i(t-1)), \end{aligned} \quad (6)$$

where $s_i(t) = 0$ for $t < 0$ and the statistical independence of innovation sequences was taken into account.

It follows from the random walk model (4) that the conditional probability density of $s_i(t)$ given its past samples can be written as

$$p_i(s_i(t)|s_i(t-1)) = q_i(\epsilon_i(t)). \quad (7)$$

Combining (6) and (7) leads to

$$\begin{aligned} p_i(\underline{s}_i) &= \prod_{t=0}^{N-1} q_i(\epsilon_i(t)) \\ &= \prod_{t=0}^{N-1} q_i(s'_i(t)), \end{aligned} \quad (8)$$

where $s'_i(t) = s_i(t) - s_i(t-1)$ which is the first-order approximation of differentiation.

Take the statistical independence of latent variables and (8) into account, then we can write the joint density $p(\underline{s}_1, \dots, \underline{s}_n)$ as

$$\begin{aligned} p(\underline{s}_1, \dots, \underline{s}_n) &= \prod_{i=1}^n p_i(\underline{s}_i) \\ &= \prod_{t=0}^{N-1} \prod_{i=1}^n q_i(s'_i(t)). \end{aligned} \quad (9)$$

The factorial model given in (9) will be used as a optimization criterion to derive the proposed algorithm.

3. DIFFERENTIAL ICA ALGORITHM

Denote a set of observation data by

$$\mathcal{X} = \{\underline{x}_1, \dots, \underline{x}_n\}, \quad (10)$$

where

$$\underline{x}_i = \{x_i(0), \dots, x_i(N-1)\}. \quad (11)$$

Then the normalized log-likelihood is given by

$$\begin{aligned} \frac{1}{N} \log p(\mathcal{X}|\mathbf{A}) &= -\log |\det \mathbf{A}| + \frac{1}{N} \log p(\underline{s}_1, \dots, \underline{s}_n) \\ &= -\log |\det \mathbf{A}| + \frac{1}{N} \sum_{t=0}^{N-1} \sum_{i=1}^n \log q_i(s'_i(t)). \end{aligned} \quad (12)$$

Let us denote the inverse of \mathbf{A} by $\mathbf{W} = \mathbf{A}^{-1}$. The estimate of latent variables is denoted by $\mathbf{y}(t) = \mathbf{W}\mathbf{x}(t)$. With these defined variables, the objective function (that is the negative normalized log-likelihood) is given by

$$\begin{aligned} \mathcal{J}_2 &= -\frac{1}{N} \log p(\mathcal{X}|\mathbf{A}) \\ &= -\log |\det \mathbf{W}| - \frac{1}{N} \sum_{t=0}^{N-1} \sum_{i=1}^n \log q_i(y'_i(t)), \end{aligned} \quad (13)$$

where s_i is replaced by its estimate y_i .

For on-line learning, the sample average is replaced by instantaneous value. Hence the objective function (13) becomes

$$\mathcal{J}_3 = -\log |\det \mathbf{W}| - \sum_{i=1}^n \log q_i(y'_i(t)), \quad (14)$$

Note that objective function (14) is slightly different from the one used in the conventional ICA based on the minimization of mutual information or the maximum likelihood estimation.

We derive a natural gradient learning algorithm which finds a minimum of (14). To this end, we follow the way that was discussed in [1, 9, 10]. We calculate the total differential $d\mathcal{J}_3(\mathbf{W})$ due to the change $d\mathbf{W}$

$$\begin{aligned} d\mathcal{J}_3 &= \mathcal{J}_3(\mathbf{W} + d\mathbf{W}) - \mathcal{J}_3(\mathbf{W}) \\ &= d\{-\log |\det \mathbf{W}|\} + d\left\{-\sum_{i=1}^n \log q_i(y'_i(t))\right\} \end{aligned} \quad (15)$$

Define

$$\varphi_i(y'_i) = -\frac{d \log q_i(y'_i)}{dy'_i}. \quad (16)$$

and construct a vector $\varphi(\mathbf{y}') = [\varphi_1(y'_1) \dots \varphi_n(y'_n)]^T$.

With this definition, we have

$$\begin{aligned} d\left\{-\sum_{i=1}^n \log q_i(y'_i(t))\right\} &= \sum_{i=1}^n \varphi_i(y'_i(t)) dy'_i(t) \\ &= \varphi^T(\mathbf{y}'(t)) d\mathbf{y}'(t). \end{aligned} \quad (17)$$

One can easily see that

$$d\{-\log |\det \mathbf{W}|\} = \text{tr}\{d\mathbf{W}\mathbf{W}^{-1}\}. \quad (18)$$

Define a modified differential matrix $d\mathbf{V}$ by

$$d\mathbf{V} = d\mathbf{W}\mathbf{W}^{-1}. \quad (19)$$

Then, with this modified differential matrix, the total differential $d\mathcal{J}_3(\mathbf{W})$ is computed as

$$d\mathcal{J}_3 = -\text{tr}\{d\mathbf{V}\} + \varphi^T(\mathbf{y}'(t)) d\mathbf{V}\mathbf{y}'(t). \quad (20)$$

A gradient descent learning algorithm for updating \mathbf{V} is given by

$$\begin{aligned} \mathbf{V}(t+1) &= \mathbf{V}(t) - \eta_t \frac{d\mathcal{J}_3}{d\mathbf{V}} \\ &= \eta_t \left\{ \mathbf{I} - \varphi(\mathbf{y}'(t))\mathbf{y}'^T(t) \right\}. \end{aligned} \quad (21)$$

Hence, it follows from the relation (19) that the updating rule for \mathbf{W} has the form

$$\mathbf{W}(t+1) = \mathbf{W}(t) + \eta_t \left\{ \mathbf{I} - \varphi(\mathbf{y}'(t))\mathbf{y}'^T(t) \right\} \mathbf{W}(t). \quad (22)$$

Remarks

- The algorithm (22) was originally derived in an *ad hoc* manner in [6]. Here we show that the algorithm (22) can be derived in the framework of maximum likelihood estimation and a random walk model.

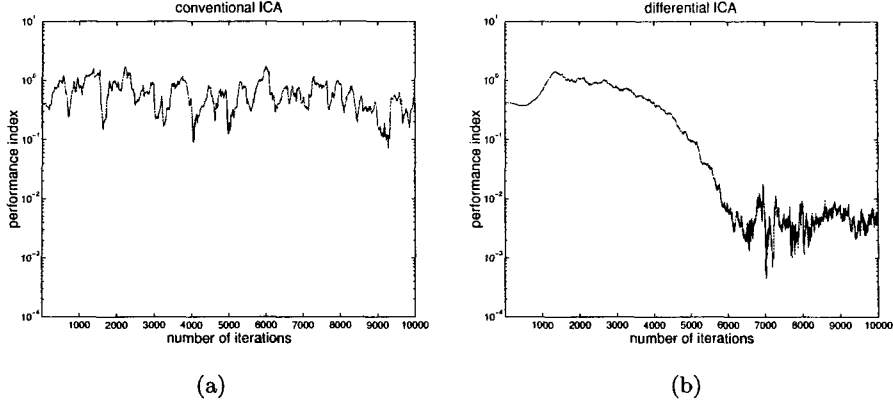


Figure 1: Evolution of performance index: (a) conventional ICA; (b) differential ICA.

- The algorithm (22) can be viewed as a special case of temporal ICA algorithm [7] where the spatiotemporal generative model was employed.
- In the conventional ICA algorithm, the nonlinear function $\varphi_i(\cdot)$ depends on the probability distribution of source. However, in the differential ICA algorithm, the nonlinear function is chosen, depending on the probability distribution of $\epsilon_i(t) = s_i(t) - s_i(t-1)$, i.e., the difference of adjacent latent variables in time domain. In general, the innovation is more non-Gaussian, compared to the signal itself. In this sense, the differential ICA algorithm works better than the conventional ICA algorithm when source was generated by a linear combination of innovation and its time-delayed replica (e.g., moving average). This is confirmed by a simple numerical example.
- As in the flexible ICA [10], we can adopt a flexible nonlinear function based on the generalized Gaussian distribution.

4. LOCAL STABILITY ANALYSIS

The differential ICA algorithm (22) can be obtained by replacing $\mathbf{y}(t)$ by $\mathbf{y}'(t)$ in the conventional ICA algorithm (3). Thus the local stability analysis of the algorithm (22) can be done similarly, following the result in [1]. As in [1], we calculate the expected Hessian $E\{d^2\mathcal{J}_3\}$ (in which the expectation is taken at $\mathbf{W} = \mathbf{A}^{-1}$) in terms of the modified differential matrix $d\mathbf{V}$. For shorthand notation, we omit the time index t in the following analysis.

The expected Hessian $E\{d^2\mathcal{J}_3\}$ is given by

$$\begin{aligned}
 E\{d^2\mathcal{J}_3\} &= E\left\{\mathbf{y}'d\mathbf{V}^T\Phi d\mathbf{y}' + \varphi^T(\mathbf{y}')d\mathbf{V}d\mathbf{y}'\right\} \\
 &= E\left\{\mathbf{y}'d\mathbf{V}^T\Phi d\mathbf{V}\mathbf{y}' + \varphi^T(\mathbf{y}')d\mathbf{V}d\mathbf{V}\mathbf{y}'\right\} \\
 &= \sum_{j \neq i} [\sigma_i^2 \kappa_j (dv_{ji})^2 + dv_{ij} dv_{ji}] \\
 &\quad + \sum_i (\zeta_i + 1) (dv_{ii})^2, \tag{23}
 \end{aligned}$$

where the statistical expectation is taken at the solution so

that $\{y_i\}$ are mutually independent and

$$\Phi = \begin{bmatrix} \varphi_1(y'_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \varphi_n(y'_n) \end{bmatrix} \tag{24}$$

$$\dot{\varphi}_i(y'_i) = \frac{d\varphi_i(y'_i)}{dy'_i} \tag{25}$$

$$\sigma_i^2 = E\{y_i'^2\} \tag{26}$$

$$\kappa_i = E\{\dot{\varphi}_i(y'_i)\} \tag{27}$$

$$\zeta_i = E\{y_i'^2 \dot{\varphi}_i(y'_i)\}. \tag{28}$$

It follows from (23) that $E\{d^2\mathcal{J}_3\}$ is positive if and only if

$$\kappa_i > 0 \tag{29}$$

$$\zeta_i + 1 > 0 \tag{30}$$

$$\sigma_i^2 \sigma_j^2 \kappa_i \kappa_j > 1. \tag{31}$$

5. NUMERICAL EXAMPLE

We present a simple numerical example to show the usefulness of our differential ICA algorithm which is described in (22). Three independent innovation sequences were drawn from Laplacian distribution. Each innovation sequence was convolved with a moving average filter (with exponentially decreasing impulse response) in order to generate colored sources. These sources were linearly mixed via 3×3 mixing matrix \mathbf{A}

We compare the performance of our differential ICA algorithm with that of the conventional natural gradient ICA algorithm in terms of the performance index (PI) which is defined as

$$\begin{aligned}
 \text{PI} = \frac{1}{2(n-1)} \sum_{i=1}^n \left\{ \left(\sum_{k=1}^n \frac{|g_{ik}|^2}{\max_j |g_{ij}|^2} - 1 \right) \right. \\
 \left. + \left(\sum_{k=1}^n \frac{|g_{ki}|^2}{\max_j |g_{ji}|^2} - 1 \right) \right\}, \tag{32}
 \end{aligned}$$

where g_{ij} is the (i, j) -element of the global system matrix $\mathbf{G} = \mathbf{W}\mathbf{A}$ and $\max_j g_{ij}$ represents the maximum value

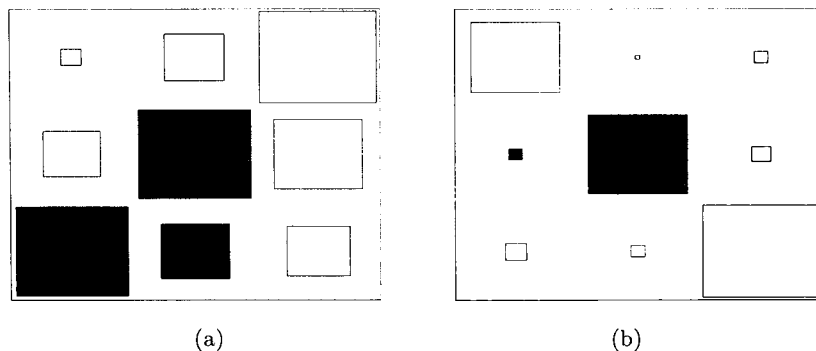


Figure 2: Hinton's diagram for the global matrix \mathbf{G} : (a) conventional ICA; (b) differential ICA. Each square's area represents the magnitude of the element of the matrix \mathbf{G} . White square is for positive sign and black square is for negative sign.

among the elements in the i th row vector of \mathbf{G} , $\max_j g_{ji}$ does the maximum value among the elements in the i th column vector of \mathbf{G} . The performance index defined in (32) tells us how far the global system matrix \mathbf{G} is from a generalized permutation matrix.

It is expected that the conventional ICA algorithm would have difficulty in separating these sources because they are close to Gaussian. The differential ICA algorithm inherently resort to the innovation sequence rather than the source itself (since it is motivated by a simple Markov model). The result of a numerical example is shown in Fig. 1.

6. DISCUSSION

In this paper we have presented a natural gradient learning algorithm for differential decorrelation, the goal of which is to minimize the correlation between differentiated random variables. We showed that the differential decorrelation algorithm could be derived from learning a linear generative model by the maximum likelihood estimation under a random walk model. We also discussed a differential version of the natural gradient ICA algorithm and showed that it could also be derived under the random walk model. The differential correlation algorithm (22) or the differential ICA algorithm (22) could be generalized by adopting higher-order differentiation. This generalization is currently under investigation.

7. ACKNOWLEDGMENTS

This work was supported by Korea Ministry of Science and Technology under Brain Science and Engineering Research Program and an International Cooperative Research Project, by KOSEF 2000-2-20500-009-5, and by Ministry of Education of Korea for its financial support toward the Electrical and Computer Engineering Division at POSTECH through its BK21 program.

8. REFERENCES

- [1] S. Amari, T. P. Chen, and A. Cichocki, "Stability analysis of learning algorithms for blind source separation," *Neural Networks*, vol. 10, no. 8, pp. 1345–1351, 1997.
- [2] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*. John Wiley & Sons, Inc., 2001.
- [3] A. Cichocki and S. Amari, *Adaptive Blind Signal and Image Processing: Learning Algorithms and Applications*. John Wiley & Sons, Inc., 2002.
- [4] B. Kosko, "Differential Hebbian learning," in *Proc. American Institute of Physics: Neural Networks for Computing*, 1986, pp. 277–282.
- [5] S. Choi, "Adaptive differential decorrelation: A natural gradient algorithm," in *Proc. ICANN*, Madrid, Spain, Aug. 2002, pp. 1168–1173.
- [6] —, "Differential Hebbian-type learning algorithms for decorrelation and independent component analysis," *Electronics Letters*, vol. 34, no. 9, pp. 900–901, 1998.
- [7] H. Attias and C. E. Schreiner, "Blind source separation and deconvolution: The dynamic component analysis algorithms," *Neural Computation*, vol. 10, pp. 1373–1424, 1998.
- [8] S. Amari, "Estimating functions of independent component analysis for temporally correlated signals," *Neural Computation*, vol. 12, no. 9, pp. 2083–2107, 2000.
- [9] —, "Natural gradient works efficiently in learning," *Neural Computation*, vol. 10, no. 2, pp. 251–276, 1998.
- [10] S. Choi, A. Cichocki, and S. Amari, "Flexible independent component analysis," *Journal of VLSI Signal Processing*, vol. 26, no. 1/2, pp. 25–38, Aug. 2000.