# A Canonical Coordinate Decomposition Network

Ali Pezeshki, Mahmood R. Azimi-Sadjadi, and Louis L. Scharf
Department of Electrical and Computer Engineering
Colorado State University
Fort Collins, Colorado 80523–1373
Email: ali@engr.colostate.edu

*Abstract*—A network structure for canonical coordinate decomposition is presented. The network consists of two single-layer linear subnetworks that together extract the canonical coordinates of two data channels. The connection weights of the networks are trained by a stochastic gradient descent learning algorithm. Each subnetwork features a hierarchical set of lateral connections among its outputs. The lateral connections perform a deflation process that subtracts the contributions of the already extracted coordinates from the input data subspace. This structure allows for adding new nodes for extracting additional canonical coordinates without the need for retraining the previous nodes. The performance of the network is evaluated on a synthesized data set.

## I. INTRODUCTION

Canonical correlation analysis [1]-[4] provides a minimal description of the correlation between two data channels by concentrating the linear dependence of the channels into a small set of canonical variables. Canonical correlations are maximal invariants to uncoupled linear transformations of two-channel data [3],[4]. The corresponding canonical coordinates resolve the channels into coordinates that are only pairwise correlated [3],[4]. Canonical coordinates have been used to decompose Wiener filters and Gaussian communication channels into their canonical modes, where each mode corresponds to a scalar Gaussian channel or Wiener filter [3],[4]. They provide an elegant framework for analyzing linear dependence and mutual information between two data channels. In this coordinate system, the linear dependence and mutual information between the original channels are decomposed into those of canonical coordinates of the channels, which are determined by the corresponding canonical correlations. The canonical correlation associated with each pair of canonical coordinates determines the contribution of that pair to the linear dependence and mutual information between the channels [3],[4].

The conventional method of finding canonical coordinates [3],[4] involves computation of square-root-inverses of covariance matrices followed by an SVD of a coherence matrix. These operations become computationally intractable and inefficient especially for large dimensional data. In addition all the singular values and singular vectors of the coherence matrix have to be evaluated even though only the most significant singular values and their associated singular vectors are used in most applications. These deficiencies make the conventional scheme inefficient for real-time applications. Consequently, to perform the canonical coordinate decomposition efficiently, a method is required to extract the most significant canonical

coordinate pairs and the corresponding canonical correlations recursively, and without any matrix inversion, matrix square root computation or direct SVD operations.

Neural networks have been proven to be powerful tools for performing complex transformations. Several neural network-based approaches haven been reported for extracting principal components of a stochastic vector process directly from the input data set [5]-[10]. In [11], a neural network structure has been proposed for computing the reduced-rank Wiener filter [4],[12]. A neural network-based approach has been reported in [13] for performing canonical correlation analysis. However, this network only finds the most significant canonical coordinate pair and the corresponding canonical correlation.

In this paper a network and a set of updating rules for performing canonical coordinate decomposition is presented. First, the problem of finding the first canonical coordinate pair is formulated as a constrained minimization problem. Then, given the first $r$ canonical coordinate pairs, the problem of finding the $(r + 1)$th pair is formulated as one of finding the first canonical coordinate pair after the contributions of the first $r$ pairs are deflated from the input data subspace. This formulation is used to propose a network structure that consists of two single-layer linear subnetworks. The weights of the subnetworks are trained using a stochastic gradient descent learning algorithm. Each subnetwork consists of a set of lateral connections that whiten the output. The idea of using lateral connections among the outputs was first exploited in [7] for recursive extraction of principal components. In fact, the structure of each subnetwork is similar to the structure of the network reported in [7]. The lateral connections are trained to deflate the contributions of the already extracted canonical coordinates from the input data subspace. This structure allows for adding new nodes for extracting a new canonical coordinate without the need for retraining the previous nodes. This is very useful since in most cases the number of canonical coordinates or canonical correlations required is not known *a priori*. It is also useful where the statistics of the data channels are slowly varying with time. A simulation example is given to demonstrate the validity of the proposed network and the learning rules.

## II. CANONICAL COORDINATE DECOMPOSITION: A REVIEW

Consider the two random vectors, $\mathbf{x} \in R^{m \times 1}$ and $\mathbf{y} \in R^{n \times 1}$ with $m$ being the smaller dimension ($m \leq n$). Assume that $\mathbf{x}$ and $\mathbf{y}$ have zero means and share the composite covariance

matrix

$$E\left[\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}\begin{pmatrix} \mathbf{x}^T & \mathbf{y}^T \end{pmatrix}\right] = \begin{bmatrix} \mathbf{R}_{xx} & \mathbf{R}_{xy} \\ \mathbf{R}_{yx} & \mathbf{R}_{yy} \end{bmatrix} \quad (1)$$

This composite covariance matrix has the following block tridiagonal decomposition [3],[4].

$$\begin{bmatrix} \mathbf{F}^T & 0 \\ 0 & \mathbf{G}^T \end{bmatrix}\begin{bmatrix} \mathbf{R}_{xx}^{-1/2} & 0 \\ 0 & \mathbf{R}_{yy}^{-1/2} \end{bmatrix}\begin{bmatrix} \mathbf{R}_{xx} & \mathbf{R}_{xy} \\ \mathbf{R}_{yx} & \mathbf{R}_{yy} \end{bmatrix},$$
$$\cdot \begin{bmatrix} \mathbf{R}_{xx}^{-T/2} & 0 \\ 0 & \mathbf{R}_{yy}^{-T/2} \end{bmatrix}\begin{bmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{K} \\ \mathbf{K} & \mathbf{I} \end{bmatrix}$$
$$(2)$$

where $\mathbf{R}_{xx}^{-1/2}\mathbf{R}_{xx}\mathbf{R}_{xx}^{-T/2} = \mathbf{I}$, $\mathbf{R}_{xx}^{1/2}\mathbf{R}_{xx}^{T/2} = \mathbf{R}_{xx}$, and F, G and K are chosen to be the economical SVD of the coherence matrix $\mathbf{C} = \mathbf{R}_{xx}^{-1/2}\mathbf{R}_{xy}\mathbf{R}_{yy}^{-T/2}$. That is,

$$\mathbf{C} = \mathbf{R}_{xx}^{-1/2}\mathbf{R}_{xy}\mathbf{R}_{yy}^{-T/2} = \mathbf{FKG}^T \quad \text{and} \quad \mathbf{F}^T\mathbf{CG} = \mathbf{K},$$
$$\mathbf{F}^T\mathbf{F} = \mathbf{I}, \quad \mathbf{G}^T\mathbf{G} = \mathbf{I}, \quad \mathbf{K} = \mathrm{diag}[k_1, k_2, \ldots, k_m]; \quad (3)$$

The diagonal matrix K is the canonical correlation matrix of canonical correlations $k_i$. The canonical correlations are arranged in descending order $(1 \geq k_1 \geq \ldots \geq k_m > 0)$.

The canonical coordinates of x and y are defined as

$$\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{F}^T & 0 \\ 0 & \mathbf{G}^T \end{bmatrix}\begin{bmatrix} \mathbf{R}_{xx}^{-1/2} & 0 \\ 0 & \mathbf{R}_{yy}^{-1/2} \end{bmatrix}\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \quad (4)$$

where the elements of u are the canonical coordinates of x and the elements of v are the canonical coordinates of y. Correspondingly, the matrices

$$\mathbf{W}^T = \mathbf{F}^T\mathbf{R}_{xx}^{-1/2} \quad \text{and} \quad \mathbf{D}^T = \mathbf{G}^T\mathbf{R}_{yy}^{-1/2} \quad (5)$$

map x and y to their corresponding canonical coordinates u and v. Thus we may rewrite the canonical coordinate map of (4) as

$$\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{W}^T & 0 \\ 0 & \mathbf{D}^T \end{bmatrix}\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}. \quad (6)$$

The composite vector of canonical coordinates, $[\mathbf{u}^T\mathbf{v}^T]^T$ has covariance matrix

$$E\left[\begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix}\begin{pmatrix} \mathbf{u}^T & \mathbf{v}^T \end{pmatrix}\right] = \begin{bmatrix} \mathbf{R}_{uu} & \mathbf{R}_{uv} \\ \mathbf{R}_{vu} & \mathbf{R}_{vv} \end{bmatrix}$$
$$= \begin{bmatrix} \mathbf{W}^T\mathbf{R}_{xx}\mathbf{W} & \mathbf{W}^T\mathbf{R}_{xy}\mathbf{D} \\ \mathbf{D}^T\mathbf{R}_{yx}\mathbf{W} & \mathbf{D}^T\mathbf{R}_{yy}\mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{K} \\ \mathbf{K} & \mathbf{I} \end{bmatrix} \quad (7)$$

which indicates that the canonical coordinates u and v are individually white but diagonally cross-correlated. The canonical correlation matrix is the cross covariance matrix of u and v and is given by

$$E[\mathbf{uv}^T] = \mathbf{K} = \mathbf{W}^T\mathbf{R}_{xy}\mathbf{D} = \mathbf{F}^T\mathbf{R}_{xx}^{-1/2}\mathbf{R}_{xy}\mathbf{R}_{yy}^{-T/2}\mathbf{G} \quad (8)$$

The diagonal elements of this matrix are arranged in descending order. Thus the first canonical coordinate pair $\{u_1, v_1\}$ has the largest correlation.

The conventional method of canonical coordinate decomposition, i.e. (4), requires the computation of the SVD of the coherence matrix $\mathbf{C} = \mathbf{R}_{xx}^{-1/2}\mathbf{R}_{xy}\mathbf{R}_{yy}^{-T/2}$ and the products $\mathbf{F}^T\mathbf{R}_{xx}^{-1/2}$ and $\mathbf{G}^T\mathbf{R}_{yy}^{-1/2}$. However, the major computational

burden is not just the SVD. The procedure before and after the SVD involves computation of square-root-inverses that requires more than $O(\min(m^2n, n^2m))$ flops. These operations become computationally intractable and inefficient for large dimension data. This motivates our next discussion.

## III. CANONICAL COORDINATE DECOMPOSITION NETWORK

This section presents a network and a set of updating rules to recursively extract the canonical coordinates of two data channels. The updating rules are derived so that no matrix inversion or square root computation is required. The network may be trained in either batch or sequential mode and thus may be used for online applications as well.

Let $\{u_i, v_i\}$ denote the $i$th pair of canonical coordinates of x and y. The canonical correlation associated with this pair is $k_i$, which is the $i$th diagonal element of the canonical correlation matrix K. Correspondingly the $i$th columns of matrices W and D, denoted by $\mathbf{w}_i \in R^{m \times 1}$ and $\mathbf{d}_i \in R^{n \times 1}$, are respectively, the mapping vectors that map x and y to their $i$th canonical coordinates $u_i$ and $v_i$. The canonical correlation $k_i$ is then $k_i = E[\mathbf{w}_i^T\mathbf{xy}^T\mathbf{d}_i] = \mathbf{w}_i^T\mathbf{R}_{xy}\mathbf{d}_i$. We further refer to $\mathbf{w}_i$ and $\mathbf{d}_i$ as the $i$th canonical coordinate mappings.

Noting that $k_1 = \mathbf{w}_1^T\mathbf{R}_{xy}\mathbf{d}_1$ is the largest canonical correlation, the problem of finding the first canonical coordinate mappings, $\mathbf{w}_1$ and $\mathbf{d}_1$, may be formulated as the maximization problem

$$\max_{\mathbf{w}_1, \mathbf{d}_1} \mathbf{w}_1^T\mathbf{R}_{xy}\mathbf{d}_1 \quad (9)$$

subject to the constraints

$$\mathbf{w}_1^T\mathbf{R}_{xx}\mathbf{w}_1 = 1 \quad \text{and} \quad \mathbf{d}_1^T\mathbf{R}_{yy}\mathbf{d}_1 = 1. \quad (10)$$

Using the method of Lagrange multipliers we may rewrite the constrained optimization problem defined by (9) and (10) as minimizing the objective function $J_1$ of the form

$$J_1 = -\mathbf{w}_1^T\mathbf{R}_{xy}\mathbf{d}_1 + (\mathbf{w}_1^T\mathbf{R}_{xx}\mathbf{w}_1 - 1)\frac{\lambda_{1,1}}{2} \\ +(\mathbf{d}_1^T\mathbf{R}_{yy}\mathbf{d}_1 - 1)\frac{\lambda_{1,2}}{2}, \quad (11)$$

where $\lambda_{1,1}$ and $\lambda_{1,2}$ are Lagrange multipliers that enforce the constraints in (10).

Now, assume that the first $r < m$ columns of W and D have already been found. Let $\mathbf{W}_r \in R^{m \times r}$ and $\mathbf{D}_r \in R^{n \times r}$ be the matrices that, respectively, contain the first $r$ columns of W and D. That is

$$\mathbf{W}_r = [\mathbf{w}_1, \ldots, \mathbf{w}_r] \quad \text{and} \quad \mathbf{D}_r = [\mathbf{d}_1, \ldots, \mathbf{d}_r]. \quad (12)$$

The first $r$ canonical coordinates of x and y are then given by

$$\mathbf{u}_r = [u_1, \ldots, u_r]^T = \mathbf{W}_r^T\mathbf{x} \\ \mathbf{v}_r = [v_1, \ldots, v_r]^T = \mathbf{D}_r^T\mathbf{y} \quad (13)$$

By introducing deflation we may minimize the deflated version of $J_1$ to find the next canonical coordinate mappings. It can be shown [14] that the $(r+1)$th canonical coordinate pair of x and y, $\{u_{r+1}, v_{r+1}\}$, is the first canonical coordinate pair of $(\mathbf{I} - \mathbf{R}_{xx}\mathbf{W}_r\mathbf{W}_r^T)\mathbf{x}$ and $(\mathbf{I} - \mathbf{R}_{yy}\mathbf{D}_r\mathbf{D}_r^T)\mathbf{y}$. Thus, the problem of finding the $(r+1)$th canonical coordinate mappings $\mathbf{w}_{r+1}$ and

$\mathbf{d}_{r+1}$, may be formulated in the context of (11) with $\mathbf{x}$ and $\mathbf{y}$ being replaced by their deflated versions $(\mathbf{I}-\mathbf{R}_{xx}\mathbf{W}_r\mathbf{W}_r^T)\mathbf{x}$ and $(\mathbf{I}-\mathbf{R}_{yy}\mathbf{D}_r\mathbf{D}_r^T)\mathbf{y}$. We may now write the problem of finding the canonical coordinate mappings $\mathbf{w}_{r+1}$ and $\mathbf{d}_{r+1}$ as minimizing the objective function

$$J_{r+1} = -\mathbf{w}_{r+1}^T(\mathbf{I} - \mathbf{R}_{xx}\mathbf{W}_r\mathbf{W}_r^T)\mathbf{R}_{xy}(\mathbf{I} - \mathbf{R}_{yy}\mathbf{D}_r\mathbf{D}_r^T)\mathbf{d}_{r+1}$$
$$+(\mathbf{w}_{r+1}^T\mathbf{R}_{xx}\mathbf{w}_{r+1} - 1)\frac{\lambda_{r+1,1}}{2} + (\mathbf{d}_{r+1}^T\mathbf{R}_{yy}\mathbf{d}_{r+1} - 1)\frac{\lambda_{r+1,2}}{2}$$
$$(14)$$

where $\lambda_{r+1,1}$ and $\lambda_{r+1,2}$ are Lagrange multipliers that guarantee the unit variance property of the new pair of coordinates;

$$\mathbf{w}_{r+1}^T\mathbf{R}_{xx}\mathbf{w}_{r+1} = 1 \quad \text{and} \quad \mathbf{d}_{r+1}^T\mathbf{R}_{yy}\mathbf{d}_{r+1} = 1. \quad (15)$$

Taking the partial derivatives of $J_c$ with respect to $\mathbf{w}_{r+1}$ and $\mathbf{d}_{r+1}$ yield

$$\frac{\partial J_{r+1}}{\partial \mathbf{w}_{r+1}} = -(\mathbf{I} - \mathbf{R}_{xx}\mathbf{W}_r\mathbf{W}_r^T)\mathbf{R}_{xy}(\mathbf{I} - \mathbf{R}_{yy}\mathbf{D}_r\mathbf{D}_r^T)^T\mathbf{d}_{r+1}$$
$$+\mathbf{R}_{xx}\mathbf{w}_{r+1}\lambda_{r+1,1}$$
$$\frac{\partial J_{r+1}}{\partial \mathbf{d}_{r+1}} = -(\mathbf{I} - \mathbf{R}_{yy}\mathbf{D}_r\mathbf{D}_r^T)\mathbf{R}_{yx}(\mathbf{I} - \mathbf{R}_{xx}\mathbf{W}_r\mathbf{W}_r^T)^T\mathbf{w}_{r+1}$$
$$+\mathbf{R}_{yy}\mathbf{d}_{r+1}\lambda_{r+1,2}$$
$$(16)$$

At the solution the constraints in (15) are satisfied. Moreover, due to deflation, $\mathbf{w}_{r+1}$ and $\mathbf{d}_{r+1}$ are respectively, orthogonal to $\mathbf{R}_{xx}\mathbf{W}_r$ and $\mathbf{R}_{yy}\mathbf{D}_r$. That is

$$\mathbf{w}_{r+1}^T\mathbf{R}_{xx}\mathbf{W}_r = \mathbf{0} \quad \text{and} \quad \mathbf{d}_{r+1}^T\mathbf{R}_{yy}\mathbf{D}_r = \mathbf{0}. \quad (17)$$

Using (15) and (17) the optimal values of Lagrange multipliers in (14) are found to be

$$\lambda_{r+1} = \lambda_{r+1,1} = \lambda_{r+1,2} = \mathbf{w}_{r+1}^T\mathbf{R}_{xy}\mathbf{d}_{r+1} \quad (18)$$

From (6), the $(r+1)$th canonical coordinate pair of $\mathbf{x}$ and $\mathbf{y}$ is given by

$$u_{r+1} = \mathbf{w}_{r+1}^T\mathbf{x}$$
$$v_{r+1} = \mathbf{d}_{r+1}^T\mathbf{y} \quad (19)$$

Using (13) and (17) we may rewrite (19) as

$$u_{r+1} = \mathbf{w}_{r+1}^T(\mathbf{I} - \mathbf{R}_{xx}\mathbf{W}_r\mathbf{W}_r^T)\mathbf{x} = \mathbf{w}_{r+1}^T\mathbf{x} - \mathbf{q}_r^T\mathbf{u}_r$$
$$v_{r+1} = \mathbf{d}_{r+1}^T(\mathbf{I} - \mathbf{R}_{yy}\mathbf{D}_r\mathbf{D}_r^T)\mathbf{y} = \mathbf{d}_{r+1}^T\mathbf{x} - \mathbf{p}_r^T\mathbf{v}_r \quad (20)$$

where

$$\mathbf{q}_r^T = \mathbf{w}_{r+1}^T\mathbf{R}_{xx}\mathbf{W}_r$$
$$\mathbf{p}_r^T = \mathbf{d}_{r+1}^T\mathbf{R}_{yy}\mathbf{D}_r \quad (21)$$

The pair of equations in (20) may be used to define a network structure for extracting the $(r+1)$th pair of canonical coordinates, given the first $r$ pairs. Each equation in (20) defines a single layer subnetwork that features a feedforward set of weights from the input to the output and a set of lateral connections that connects the first $r$ nodes to the $(r+1)$th node. Figure 1 shows the structure of this network. In this structure, $\mathbf{W}_r$ and $\mathbf{D}_r$ are the weight matrices that map $\mathbf{x}$ and $\mathbf{y}$ to their first $r$ canonical coordinates $\mathbf{u}_r$ and $\mathbf{v}_r$. Given these weights, the network may be trained, by minimizing $J_{r+1}$ in (14), to extract the $(r+1)$th canonical coordinate pair and the corresponding mappings. The weight vectors $\mathbf{w}_{r+1}$ and $\mathbf{d}_{r+1}$ are trained to maximize the correlation between the outputs $u_{r+1}$ and $v_{r+1}$ and make them unity variance. The lateral
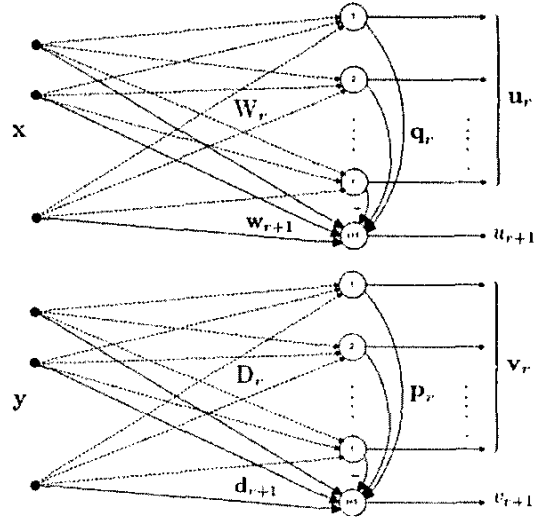


Fig. 1. The structure of the network for recursive extraction of canonical coordinates of x and y.

weight vector $\mathbf{q}_r$ is trained to orthogonalize $\mathbf{u}_r$ (the first $r$ canonical coordinates of $\mathbf{x}$) to $u_{r+1}$ (the $(r+1)$th canonical coordinate of $\mathbf{x}$.) Similarly, the lateral weight vector $\mathbf{p}_r$ is trained to orthogonalize $\mathbf{v}_r$ (the first $r$ canonical coordinates of $\mathbf{y}$) to $v_{r+1}$ (the $(r+1)$th canonical coordinate of $\mathbf{y}$.) The lateral connections perform a deflation process that subtracts the contributions of the already extracted coordinates from the linear subspaces of $\mathbf{x}$ and $\mathbf{y}$. This structures allows for adding new nodes for extracting additional canonical coordinates without the need for retraining the previous nodes.

Using the stochastic gradient descent learning algorithm, with instantaneous values of covariance matrices inserted into (16), we may derive the following updating rules for $\mathbf{w}_{r+1}$, and $\mathbf{d}_{r+1}$:

$$\mathbf{w}_{r+1}(j + 1) = \mathbf{w}_{r+1}(j) + [(\mathbf{x}(j + 1) - \mathbf{S}_r(j + 1)\mathbf{u}_r(j + 1))\cdot$$
$$\cdot v_r(j + 1) - \mathbf{x}(j + 1)\mathbf{x}(j + 1)^T \cdot$$
$$\cdot \mathbf{w}_{r+1}(j)\lambda_{r+1}(j + 1)]\beta(j + 1)$$
$$\mathbf{d}_{r+1}(j + 1) = \mathbf{d}_{r+1}(j) + [(\mathbf{y}(j + 1) - \mathbf{T}_r(j + 1)\mathbf{v}_r(j + 1))\cdot$$
$$\cdot u_r(j + 1) - \mathbf{y}(j + 1)\mathbf{y}(j + 1)^T \cdot$$
$$\cdot \mathbf{d}_{r+1}(j)\lambda_{r+1}(j + 1)]\beta(j + 1)$$
$$(22)$$

where $j$ is the index of iteration. Matrices $\mathbf{S}_r$ and $\mathbf{T}_r$ are updated to asymptotically approximate $\mathbf{R}_{xx}\mathbf{W}_r$ and $\mathbf{R}_{yy}\mathbf{D}_r$ respectively. From (18), the Lagrange multiplier $\lambda_{r+1} = \lambda_{r+1,1} = \lambda_{r+1,2}$ shall be updated to asymptotically approximate $\mathbf{w}_{r+1}^T\mathbf{R}_{xy}\mathbf{d}_{r+1} = k_{r+1}$. Thus the updating rules for $\mathbf{S}_r$, $\mathbf{T}_r$ and $\lambda_{r+1}$ are

$$\mathbf{S}_r(j + 1) = \frac{j}{j+1}\mathbf{S}_r(j) + \frac{1}{j+1}\mathbf{x}(j + 1)\mathbf{u}_r^T(j + 1)$$
$$\mathbf{T}_r(j + 1) = \frac{j}{j+1}\mathbf{T}_r(j) + \frac{1}{j+1}\mathbf{y}(j + 1)\mathbf{v}_r^T(j + 1)$$
$$\lambda_{r+1}(j + 1) = \frac{j}{j+1}\lambda_{r+1}(j) +$$
$$\frac{1}{j+1}\mathbf{w}_{r+1}^T(j)\mathbf{x}(j + 1)\mathbf{y}^T(j + 1)\mathbf{d}_{r+1}(j)$$
$$(23)$$

Finally using (21) the learning rules for the lateral weight vectors $\mathbf{q}_r$, and $\mathbf{p}_r$ may be written as

$$\mathbf{q}_r(j+1) = \mathbf{S}_r^T(j+1)\mathbf{w}_{r+1}(j+1)$$
$$\mathbf{p}_r(j+1) = \mathbf{T}_r^T(j+1)\mathbf{d}_{r+1}(j+1) \tag{24}$$

Thus, we may summarize the training algorithm for extracting the $(r+1)$th canonical coordinate pair for $r = 0, 1, \ldots, m-1$ and the corresponding mapping vectors as

$$\mathbf{u}_r(j+1) = \mathbf{w}_r^T \mathbf{x}(j+1)$$
$$\mathbf{v}_r(j+1) = \mathbf{D}_r^T \mathbf{y}(j+1)$$
$$u_{r+1}(j+1) = \mathbf{w}_{r+1}^T(j+1)\mathbf{x}(j+1) - \mathbf{q}_r^T(j)\mathbf{u}_r(j+1)$$
$$v_{r+1}(j+1) = \mathbf{d}_{r+1}^T(j+1)\mathbf{y}(j+1) - \mathbf{p}_r^T(j)\mathbf{v}_r(j+1)$$
$$\lambda_{r+1}(j+1) = \tfrac{j}{j+1}\lambda_{r+1}(j) + \tfrac{1}{j+1}\mathbf{w}_{r+1}^T(j)\mathbf{x}(j+1)\mathbf{y}^T(j+1)\mathbf{d}_{r+1}(j)$$
$$\mathbf{S}_r(j+1) = \tfrac{j}{j+1}\mathbf{S}_r(j) + \tfrac{1}{j+1}\mathbf{x}(j+1)\mathbf{u}_r^T(j+1)$$
$$\mathbf{T}_r(j+1) = \tfrac{j}{j+1}\mathbf{T}_r(j) + \tfrac{1}{j+1}\mathbf{y}(j+1)\mathbf{v}_r^T(j+1)$$
$$\mathbf{w}_{r+1}(j+1) = \mathbf{w}_{r+1}(j) + [(\mathbf{x}(j+1) - \mathbf{S}_r(j+1)\mathbf{u}_r(j+1))v_r(j+1)$$
$$-\mathbf{x}(j+1)\mathbf{x}(j+1)^T\mathbf{w}_{r+1}(j)\lambda_{r+1}(j+1)]\beta(j+1)$$
$$\mathbf{d}_{r+1}(j+1) = \mathbf{d}_{r+1}(j) + [(\mathbf{y}(j+1) - \mathbf{T}_r(j+1)\mathbf{v}_r(j+1))u_r(j+1)$$
$$-\mathbf{y}(j+1)\mathbf{y}(j+1)^T\mathbf{d}_{r+1}(j)\lambda_{r+1}(j+1)]\beta(j+1)$$
$$\mathbf{q}_r(j+1) = \mathbf{S}_r^T(j+1)\mathbf{w}_{r+1}(j+1)$$
$$\mathbf{p}_r(j+1) = \mathbf{T}_r^T(j+1)\mathbf{d}_{r+1}(j+1) \tag{25}$$

The initial values $\mathbf{w}_{r+1}(0) \in R^{m\times 1}$, $\mathbf{d}_{r+1}(0) \in R^{n\times 1}$, $\mathbf{q}_r(0) \in R^{r\times 1}$, $\mathbf{p}_r(0) \in R^{r\times 1}$, $\mathbf{S}_r(0) \in R^{m\times r}$, $\mathbf{T}_r(0) \in R^{n\times r}$, and $\lambda_{r+1}$ may be chosen randomly. The learning rate $\beta$ may be varied or kept fixed [15]. It is important to note that, owing to the deflation performed by the lateral connections, the outputs within each subnetwork are decoupled. Thus during extraction of the $(r+1)$th pair of canonical coordinates there is no need to retrain the previous nodes and the weight matrices $\mathbf{W}_r$ and $\mathbf{D}_r$ are not changed.

## IV. SIMULATION RESULTS

In this section, the proposed network is used to recursively extract the canonical coordinate mappings for a synthesized data set. The performance of the network is demonstrated by presenting the plots of squared error between the actual canonical coordinate mappings, computed using the direct method in (5), and the ones estimated by the network, along with the plots of squared error for canonical correlations. Let $\hat{\mathbf{w}}_i$ and $\hat{\mathbf{d}}_i$, respectively, denote the estimate of the $i$th pair of the actual canonical coordinate mappings $\mathbf{w}_i$ and $\mathbf{d}_i$. We define $e_{\mathbf{w}_i}^2$ and $e_{\mathbf{d}_i}^2$ as the squared estimation error of the $i$th canonical coordinate mappings $\mathbf{w}_i$ and $\mathbf{d}_i$. That is,

$$e_{\mathbf{w}_i}^2 = \|\mathbf{w}_i - \hat{\mathbf{w}}_i\|^2 \quad \text{and} \quad e_{\mathbf{d}_i}^2 = \|\mathbf{d}_i - \hat{\mathbf{d}}_i\|^2$$

Also, $e_{\mathbf{k}_i}^2 = (k_i - \hat{k}_i)^2$, is defined as the squared estimation error of the $i$th canonical correlation $k_i$. The actual canonical correlation $k_i$ is found from the SVD in (3). From (18), it is seen that the $i$th canonical correlation $k_i$ is estimated by the Lagrange multiplier $\lambda_i$. The data set is formed from 500 samples of two data channels governed by the linear model

$$\mathbf{x} = H_x \eta_x$$
$$\mathbf{y} = H_y \eta_y + H_{yx}\mathbf{x}$$

where $\mathbf{x} \in R^{4\times 1}$, $\mathbf{y} \in R^{5\times 1}$. The matrices $H_x \in R^{4\times 4}$, $H_y \in R^{5\times 5}$ and $H_{yx} \in R^{5\times 4}$ are known matrices, and $\eta_x \in R^{4\times 1}$

and $\eta_y \in R^{5\times 1}$ are two independent white Gaussian vectors. The network is trained for 2500 epochs using the training algorithm in (25). The learning rate is varied linearly from $\beta = 5 \times 10^{-3}$ to $\beta = 5 \times 10^{-6}$ in 2500 steps. All the initial values in (25) are randomly selected.

Figure 2 shows the squared estimation errors $e_{\mathbf{w}_i}^2$, $i \in [1,4]$ vs. the epoch index for 10 independent initializations of the network. It is seen that in all the cases the squared error approaches zero within a misadjustment error and thus the weights of the upper subnetwork (Fig. 1) converge to the actual canonical coordinate mappings that map the first data channel $\mathbf{x}$ into its canonical coordinates $\mathbf{u}$.

The plots of the squared estimation errors $e_{\mathbf{d}_i}^2$, $i \in [1,4]$ vs. epoch index for the 10 initializations are shown in Fig. 3. The convergence behaviors are very similar to those in Fig. 2. It is seen that in all the cases the squared error approaches zero within a misadjustment error and thus the weights of the lower subnetwork (Fig. 1) converge to the actual canonical coordinate mappings that map the second data channel $\mathbf{y}$ into its canonical coordinates $\mathbf{v}$.

Figure 4 shows the squared estimation errors $e_{\mathbf{k}_i}^2$, $i \in [1,4]$ vs. the epoch index for the 10 initializations. The plots show that the squared error decays to zero in all the cases. The estimate of the $i$th canonical correlation is given by the Lagrange multiplier $\lambda_i$. These plots indicate that $\lambda_i$'s converge to the actual canonical correlations $k_i$'s in all the cases.

## V. CONCLUSION

A new network for recursive extraction of canonical coordinates/correlations of two data channels is introduced. The network is based on a constrained minimization problem that exploits a deflation process. The deflation process is performed by incorporating lateral connections into the subnetworks. The learning rules are derived using a stochastic gradient descent algorithm. The structure of the network along with the learning rules allow for adding a new node to the network in order to extract a new pair of canonical coordinates without the need to retrain the previous nodes. Unlike conventional methods, no matrix inversion, matrix square root computation, direct SVD is required during the training. A simulation example demonstrates the validity of the proposed network and learning rules. The results confirm that the extracted canonical coordinate mappings approximate the true ones.

## REFERENCES

[1] H. Hotelling, "Relation between two sets of variates," *Biometrica*, vol. 28, pp. 321-377, 1936.
[2] T. W. Anderson, *An Introduction to Multivariate Statistical Analysis*, New York: Wiley, 1958.
[3] L. L. Scharf and C.T. Mullis, "Canonical coordinates and the geometry of inference, rate and capacity," *IEEE Trans. on Signal Processing*, vol. 48, pp. 824-831, March 2000.
[4] L. L. Scharf and J. T. Thomas, "Wiener filters in canonical coordinates for transform coding, filtering, and quantizing," *IEEE Trans. on Signal Processing*, vol. 46, pp. 647-654, March 1998.
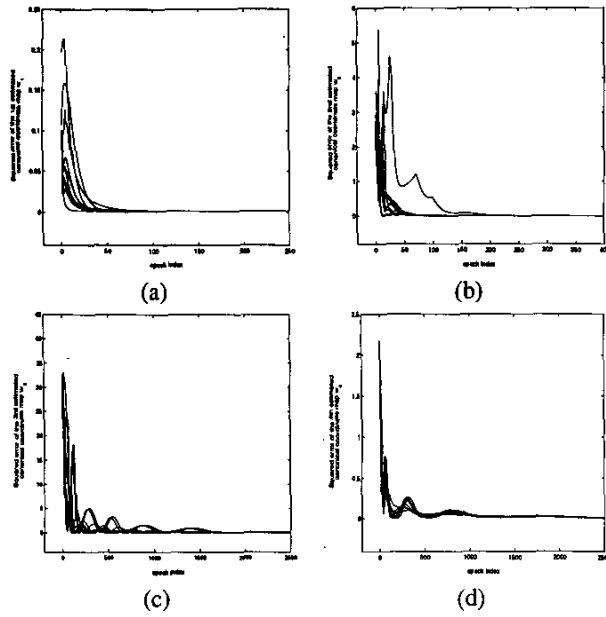
(a)            (b)



(c)            (d)

Fig. 2. The squared error for $w_i$'s, $i \in [1,4]$ vs. the epoch index for 10 independent initializations of the network (a) $i = 1$; $e^2_{\hat{w}_1} = \|w_1 - \hat{w}_1\|^2$. (b) $i = 2$; $e^2_{\hat{w}_2} = \|w_2 - \hat{w}_2\|^2$. (c) $i = 3$; $e^2_{\hat{w}_3} = \|w_3 - \hat{w}_3\|^2$. (d) $i = 4$; $e^2_{\hat{w}_4} = \|w_4 - \hat{w}_4\|^2$. In all the cases the squared error approaches zero and the weights of the upper subnetwork (Fig. 1) converge to the actual canonical coordinate mappings that map the first data channel x into its canonical coordinates u.
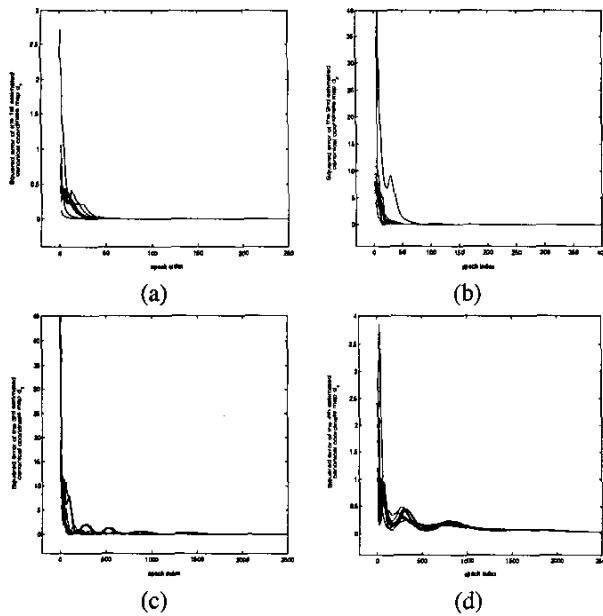


(a)            (b)



(c)            (d)

Fig. 3. The squared error for $d_i$'s, $i \in [1,4]$ vs. the epoch index for 10 independent initializations of the network (a) $i = 1$; $e^2_{d_1} = \|d_1 - \hat{d}_1\|^2$. (b) $i = 2$; $e^2_{d_2} = \|d_2 - \hat{d}_2\|^2$. (c) $i = 3$; $e^2_{d_3} = \|d_3 - \hat{d}_3\|^2$. (d) $i = 4$; $e^2_{d_4} = \|d_4 - \hat{d}_4\|^2$. In all the cases the squared error approaches zero and the weights of the lower subnetwork (Fig. 1) converge to the actual canonical coordinate mappings that map the second data channel y into its canonical coordinates v.
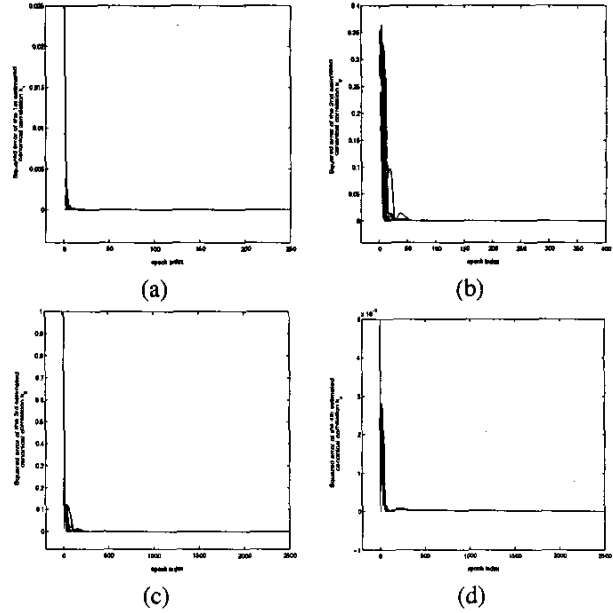


(a)            (b)



(c)            (d)

Fig. 4. The squared error for $k_i$'s, $i \in [1,4]$ vs. the epoch index for 10 independent initializations of the network (a) $i = 1$; $e^2_{k_1} = (k_1 - \hat{k}_1)^2$. (b) $i = 2$; $e^2_{k_2} = (k_2 - \hat{k}_2)^2$. (c) $i = 3$; $e^2_{k_3} = (k_3 - \hat{k}_3)^2$. (d) $i = 4$; $e^2_{k_4} = (k_4 - \hat{k}_4)^2$. The estimate of $k_i$ is given by the Lagrange multiplier $\lambda_i$. The plots show that $\lambda_i$ converges to the actual canonical correlation $k_i$ in all the cases.

[5] K. I. Diamantaras and S. Y. Kung, *Principal Component Neural Networks: Theory and Applications*, 1st ed. John Wiely & Sons Inc. 1996.

[6] S. Bannour and M. R. Azimi-Sadjadi, "Principal component extraction using recursive least squares learning," *IEEE Trans. on Neural Networks*, vol. 6, pp. 457-469, March 1995.

[7] S. Y. Kung and K. I. Diamantaras, "Adaptive principal component extraction (APEX) and applications," *IEEE Trans. on Signal Processing*, vol. 42, pp. 1202-1217, May 1994.

[8] P. Baldi, and K. Hornik, "Neural networks and principal component analysis: Learning from examples without local minima," *Neural Networks*, vol. 2, pp. 53-58, 1989.

[9] T. D. Sanger, "Optimal unsupervised learning in a single-layer linear feedforward neural network," *Neural Networks*, vol. 2, pp. 459-473, 1989.

[10] E. Oja, "A simplified neuron model as principal component analyzer," *J. Math. Biology*, vol. 15, pp. 267-273, 1982.

[11] K. I. Diamantaras and S. Y. Kung, "Multi-layer neural networks for reduced-rank approximation," *IEEE Trans. Neural Networks*, vol. 5, pp. 684-697, Sept. 1994.

[12] L. L. Scharf, *Statistical Signal Processing*. Reading, MA: Addison-Wesley, 1991, pp. 330-331.

[13] P. L. Lai and C. Fyfe, "A neural network implementation of canonical correlation analysis," *Neural Networks*, vol. 12, pp. 1391-1397, 1999.

[14] A. Pezeshki, M. R. Azimi-Sadjadi, and L. L. Scharf, "A network for recusive extraction of canonical coordinates," *To appear in the Special Issue of Journal Neural Networks*, 2003.

[15] S. Haykin, *Neural Networks: A Comprehensive Foundation*. Upper Saddle River, NJ: Printice Hall, second edition, 1991.