

ABSTRACT

Subspace analysis (which includes PCA) seeks for feature subspace (which corresponds to the eigenspace), given multivariate input data and has been widely used in computer vision and pattern recognition. Typically data space belongs to very high dimension, but only a few principal components need to be extracted. In this paper I present a fast sequential algorithm for subspace analysis or tracking. Useful behavior of the algorithm is confirmed by numerical experiments.

1. INTRODUCTION

The task of pattern recognition is to classify each data in an unlabelled set of data (test set), given a set of data labelled with its class (training set). Statistical pattern recognition relies on information coding (data representation). For example, linear data representation decomposes a given set of data into a linear sum of basis vectors. Feature vectors are obtained by projecting the data onto the subspace spanned by the basis vectors. Popular methods are factor analysis and principal component analysis (PCA). The eigenface method [8, 11] might be one exemplary subspace analysis technique in pattern recognition.

In computer vision and pattern recognition, one often encounters into a set of huge dimensional data and wish to extract a small number of features which is able to represent the data as well as possible. The singular value decomposition (SVD) is a numerically robust method which calculates the eigenvectors of the covariance matrix, however, it is computationally expensive, especially for the case of data with high dimension. For adaptive computation of eigenvectors, a variety of PCA neural networks have been developed [1], most of which are gradient-based learning algorithms, so their convergence is very slow.

Recently probabilistic model-based methods for subspace analysis have been proposed. These include probabilistic PCA (PPCA) [10], EM-PCA [5], mixtures of factor analyzers [2], and mixtures of probabilistic principal component analyzers [9]. All these algorithms employ the EM learning which is an iterative maximum likelihood estimation method in the presence of hidden variables. PPCA and EM-PCA are batch algorithms, thus, when a new data arrives, whole calculation should be carried out again. In order to overcome this drawback, I present a sequential EM learning algorithm.

2. LINEAR GENERATIVE MODEL

The linear generative model assumes that the set of m dimensional observed vectors $\{\mathbf{x}_t\}$ is generated from a corresponding set of latent variables $\{\mathbf{s}_t\}$ by

$$\mathbf{x}_t = \mathbf{A}\mathbf{s}_t + \mathbf{v}_t, \quad (1)$$

where $\mathbf{s}_t \in \mathbb{R}^n$ ($n \leq m$) and \mathbf{v}_t is Gaussian noise vector that is assumed to be statistically independent of \mathbf{s}_t .

In standard factor analysis, latent variables \mathbf{s} are assumed to have a unit isotropic Gaussian distribution, i.e., $\mathbf{s} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The noise model is Gaussian, i.e., $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$ with $\mathbf{\Sigma}$ being a diagonal matrix. Given this formulation, the model for \mathbf{x} is also Gaussian, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ where the covariance matrix $\mathbf{C} = \mathbf{\Sigma} + \mathbf{A}\mathbf{A}^T$. Due to the diagonality of $\mathbf{\Sigma}$, the observed variables \mathbf{x} are conditionally independent given the values of the latent variables \mathbf{s} . Thus the reduced-dimensional distribution \mathbf{s} is intended to model the dependencies between the observed variables. This is in contrast to PCA which treats the inter-variable dependencies and the independent noise being identical. Factor analysis seeks for a factor loading matrix which best model the covariance structure of the observation data. In general, the columns of the factor loading matrix do not correspond to the principal subspace of the data. Maximum likelihood solution to factor analysis can be found in [7].

3. PROBABILISTIC PCA

This section reviews PPCA [10] for the case of isotropic Gaussian noise model and zero noise limit. In the limit of zero noise, PPCA is known as EM-PCA [5].

3.1. Isotropic Gaussian Noise

In general factor loadings \mathbf{A} differ from the principal axes due to the diagonal noise model $\mathbf{\Sigma}$. Principal components emerges when independent noise variables have common variance σ^2 , i.e., noise is isotropic Gaussian. Recently Tipping and Bishop [10] showed that under an isotropic noise structure, the maximum likelihood estimator \mathbf{A}_{ML} spans a principal subspace (which consists of scaled and rotated principal eigenvectors of the sample covariance matrix \mathbf{R}) even when the covariance model is approximate.

We assume an isotropic noise model, $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, which implies a probability distribution over data space for a given \mathbf{s} given by

$$p(\mathbf{x}|\mathbf{s}) = \frac{1}{(2\pi\sigma^2)^{\frac{m}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{A}\mathbf{s}\|^2 \right\}. \quad (2)$$

A Gaussian prior over the latent variables is used, i.e.,

$$p(\mathbf{s}) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} \|\mathbf{s}\|^2 \right\}. \quad (3)$$

The marginal distribution of the data has the form of

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}|\mathbf{s})p(\mathbf{s})d\mathbf{s} \\ &= \frac{1}{(2\pi)^{\frac{m}{2}} |\mathbf{C}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}\mathbf{s}^T \mathbf{C}^{-1} \mathbf{s}\right\}, \end{aligned} \quad (4)$$

where the model covariance is

$$\mathbf{C} = \mathbf{A}\mathbf{A}^T + \sigma^2 \mathbf{I}. \quad (5)$$

Then the log-likelihood of observing the data under this model is

$$\begin{aligned} \mathcal{L} &= \sum_{t=1}^N \log p(\mathbf{x}_t|\theta) \\ &= -\frac{Nm}{2} \log(2\pi) - \frac{N}{2} \log |\mathbf{C}| - \frac{N}{2} \text{tr}\{\mathbf{C}^{-1} \mathbf{R}\}, \end{aligned} \quad (6)$$

where \mathbf{R} is the sample covariance matrix given by

$$\mathbf{R} = \frac{1}{N} \sum_{t=1}^N \mathbf{x}_t \mathbf{x}_t^T. \quad (7)$$

It was shown in [10] that with \mathbf{C} given by (5), the only non-zero stationary points of $\frac{\partial \mathcal{L}}{\partial \mathbf{A}}$ occur for

$$\mathbf{A} = \mathbf{U}_n (\mathbf{\Lambda}_n - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{Q}, \quad (8)$$

where the n column vectors in \mathbf{U}_n are eigenvectors of the sample covariance matrix \mathbf{R} , with corresponding eigenvalues in the diagonal matrix $\mathbf{\Lambda}_n$, and \mathbf{Q} is an arbitrary $n \times n$ orthogonal rotation matrix.

Following Rubin and Thayer [7], Tipping and Bishop derived an EM algorithm for maximizing the log-likelihood (6). Here I briefly review the EM algorithm for probabilistic PCA. See [10] for more details.

In the framework of EM, the latent variables $\{\mathbf{s}_t\}$ are treated as *missing data*. The complete-data log-likelihood \mathcal{L}_c is given by

$$\begin{aligned} \mathcal{L}_c &= \sum_{t=1}^N \log p(\mathbf{x}_t, \mathbf{s}_t|\theta) \\ &= \sum_{t=1}^N \left[\text{const} - \frac{m}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{x}_t - \mathbf{A}\mathbf{s}_t\|^2 \right. \\ &\quad \left. - \frac{1}{2} \mathbf{s}_t^T \mathbf{s}_t \right]. \end{aligned} \quad (9)$$

Algorithm Outline: PPCA

E-step Compute sufficient statistics

$$\langle \mathbf{s}_t \rangle = \mathbf{M}^{-1} \mathbf{A}^T \mathbf{x}_t, \quad (10)$$

$$\langle \mathbf{s}_t \mathbf{s}_t^T \rangle = \sigma^2 \mathbf{M}^{-1} + \langle \mathbf{s}_t \rangle \langle \mathbf{s}_t^T \rangle, \quad (11)$$

where $\mathbf{M} = \sigma^2 \mathbf{I} + \mathbf{A}^T \mathbf{A}$.

M-step Re-estimate the parameters \mathbf{A} and σ^2 by

$$\hat{\mathbf{A}} = \left(\sum_{t=1}^N \mathbf{x}_t \langle \mathbf{s}_t^T \rangle \right) \left(\sum_{t=1}^N \langle \mathbf{s}_t \mathbf{s}_t^T \rangle \right)^{-1}, \quad (12)$$

$$\hat{\sigma}^2 = \frac{1}{m} \left\{ \text{tr} \left[\mathbf{R} - \hat{\mathbf{A}} \mathbf{M}^{-1} \hat{\mathbf{A}}^T \right] \right\}. \quad (13)$$

3.2. Zero Noise Limit

Probabilistic PCA algorithm described in previous section is able to find scaled and rotated principal eigenvectors of the sample covariance matrix of the observed variables. It was derived in the framework of factor analysis with isotropic Gaussian noise model. PCA is a limiting case of the linear Gaussian model [6] (which factor analysis is based on) as the covariance of the noise \mathbf{v} becomes infinitesimally small and equal in all directions. Hence a simple EM algorithm for PCA can be obtained by taking the zero noise limit into account. In the zero noise limit ($\sigma^2 \rightarrow 0$), the likelihood of a data point \mathbf{x} is dominated solely by the squared distance between it and its reconstruction $\mathbf{A}\mathbf{s}$. In such a case, the posterior collapses to a single point and the covariance becomes zero, i.e.,

$$\begin{aligned} p(\mathbf{s}|\mathbf{x}) &= \mathcal{N}((\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{x}, 0) \\ &= \delta(\mathbf{s} - (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{x}). \end{aligned} \quad (14)$$

Now inference reduces to simple least squares projection, which leads to a simple EM algorithm [5] that is summarized below.

Algorithm Outline: PPCA (zero noise limit)

E-step Inference is carried out by LS projection,

$$\mathbf{S} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{X}, \quad (15)$$

where

$$\begin{aligned} \mathbf{S} &= [\mathbf{s}_1, \dots, \mathbf{s}_N], \\ \mathbf{X} &= [\mathbf{x}_1, \dots, \mathbf{x}_N]. \end{aligned} \quad (16)$$

M-step Re-estimate the matrix \mathbf{A} by

$$\hat{\mathbf{A}} = \mathbf{X} \mathbf{S} (\mathbf{S} \mathbf{S}^T)^{-1}. \quad (17)$$

4. OJA'S SUBSPACE RULE

Let us consider a linear feedforward neural networks whose output $\mathbf{y}_t \in \mathbb{R}^n$ is described by

$$\mathbf{y}_t = \mathbf{W} \mathbf{x}_t, \quad (18)$$

where $\mathbf{W} \in \mathbb{R}^{n \times m}$ is a synaptic weight matrix.

The Oja's subspace learning rule [4] is given by

$$\begin{aligned} \mathbf{W}_{t+1} &= \mathbf{W}_t + \eta_t \mathbf{y}_t \mathbf{x}_t^T \left\{ \mathbf{I} - \mathbf{W}_t^T \mathbf{W}_t \right\} \\ &= \mathbf{W}_t + \eta_t \left\{ \mathbf{y}_t \mathbf{x}_t^T - \mathbf{y}_t \mathbf{y}_t^T \mathbf{W}_t \right\}, \end{aligned} \quad (19)$$

where $\eta_t > 0$ is the learning rate. For the case of $n = 1$, the subspace rule (19) becomes the well known Oja's algorithm [3] which is known to extract the eigenvector associated the largest eigenvalue of the input data covariance matrix. When the convergence of (19) is achieved, the synaptic weight matrix \mathbf{W} corresponds to the eigenspace which is spanned by the principal eigenvectors of the input data covariance matrix.

Besides the Oja's subspace rule, a variety of PCA neural networks have been developed (see [1] and references therein).

5. SEQUENTIAL EM FOR PSA

5.1. Separable Least Squares

The PPCA algorithm for the case of zero noise limit can be also derived in the framework of separable LS method. Moreover we can employ the sequential LS (also known as recursive LS) in order to develop an on-line algorithm which learns principal subspace of the observed variables.

As pointed out in [5], in the zero noise limit, the likelihood of a data point \mathbf{x} is dominated solely by the squared distance between it and its reconstruction $\mathbf{A}\mathbf{s}$. In such a case, ML estimation of both \mathbf{A} and \mathbf{s} becomes a separable LS minimization problem. The LS estimates, $\hat{\mathbf{A}}$ and $\hat{\mathbf{S}}$ are computed by

$$\hat{\mathbf{A}}, \hat{\mathbf{S}} = \min_{\mathbf{A}, \mathbf{S}} \|\mathbf{X} - \mathbf{A}\mathbf{S}\|_F^2. \quad (20)$$

The separable LS minimization is carried out in two steps. First we minimize (20) with respect to \mathbf{A} with \mathbf{S} being fixed. It leads to

$$\hat{\mathbf{A}} = \mathbf{X}\mathbf{S} \left(\mathbf{S}\mathbf{S}^T \right)^{-1} \quad (21)$$

which corresponds to the M-step in PPCA (for the case of zero noise limit).

The estimate $\hat{\mathbf{A}}$ is substituted back into (20), then we obtain a new criterion which is a function of \mathbf{S} only

$$\min_{\mathbf{S}} \left\| \mathbf{X}\mathbf{P}_{\mathbf{S}}^{\perp} \right\|_F^2, \quad (22)$$

where $\mathbf{P}_{\mathbf{S}}^{\perp}$ is the orthogonal projection matrix given by

$$\mathbf{P}_{\mathbf{S}}^{\perp} = \mathbf{I} - \mathbf{S}^T \left(\mathbf{S}\mathbf{S}^T \right)^{-1} \mathbf{S}. \quad (23)$$

The minimization of (22) is achieved when

$$\mathbf{S} = \left(\hat{\mathbf{A}}^T \hat{\mathbf{A}} \right)^{-1} \hat{\mathbf{A}}^T \mathbf{X}, \quad (24)$$

which correspond to the E-step.

5.2. Sequential LS

For sequential estimation of \mathbf{A} and \mathbf{s} , we consider the weighted LS minimization problem where the objective function is given by

$$\mathcal{E} = \sum_{k=1}^t \beta^{t-k} \|\mathbf{x}_k - \mathbf{A}\mathbf{s}_k\|^2, \quad (25)$$

where $0 < \beta \leq 1$ is the forgetting factor.

Our objective is to compute \mathbf{A}_t and \mathbf{s}_t , assuming a good estimate of \mathbf{s}_{t-1} (or equivalently \mathbf{A}_{t-1}) is available. The exponential weighting is used to de-emphasize old data in a time-varying environment. Setting the derivative of \mathcal{E} with respect to \mathbf{A} to be zero, then we have

$$\mathbf{A}_t = \mathbf{R}_{\mathbf{x}\mathbf{s},t} \left[\mathbf{R}_{\mathbf{s}\mathbf{s},t} \right]^{-1}, \quad (26)$$

where

$$\begin{aligned} \mathbf{R}_{\mathbf{x}\mathbf{s},t} &= \sum_{k=1}^t \beta^{t-k} \mathbf{x}_k \mathbf{s}_k^T, \\ \mathbf{R}_{\mathbf{s}\mathbf{s},t} &= \sum_{k=1}^t \beta^{t-k} \mathbf{s}_k \mathbf{s}_k^T. \end{aligned} \quad (27)$$

Define $\mathbf{P}_t = \mathbf{R}_{\mathbf{s}\mathbf{s},t}^{-1}$ and apply the matrix inversion lemma. Then we have a recursion equation for updating \mathbf{P}_t

$$\mathbf{P}_t = \frac{1}{\beta} \left\{ \mathbf{P}_{t-1} - \frac{\mathbf{P}_{t-1} \mathbf{s}_t \mathbf{s}_t^T \mathbf{P}_{t-1}}{\beta + \mathbf{s}_t^T \mathbf{P}_{t-1} \mathbf{s}_t} \right\}. \quad (28)$$

Using the recursion (28), the adaptation for \mathbf{A} is given by

$$\begin{aligned} \mathbf{A}_t &= \mathbf{R}_{\mathbf{x}\mathbf{s},t} \mathbf{P}_t \\ &= \mathbf{A}_{t-1} + [\mathbf{x}_t - \mathbf{A}_{t-1} \mathbf{s}_t] \frac{\mathbf{s}_t^T \mathbf{P}_{t-1}}{\beta + \mathbf{s}_t^T \mathbf{P}_{t-1} \mathbf{s}_t}. \end{aligned} \quad (29)$$

Algorithm Outline: Sequential EM

E-step Estimate \mathbf{s}_t by the LS projection

$$\mathbf{s}_t = \left(\mathbf{A}_{t-1}^T \mathbf{A}_{t-1} \right)^{-1} \mathbf{A}_{t-1}^T \mathbf{x}_t. \quad (30)$$

M-step Estimate \mathbf{A}_t by

$$\mathbf{A}_t = \mathbf{A}_{t-1} + \epsilon_t \frac{\mathbf{s}_t^T \mathbf{P}_{t-1}}{\beta + \mathbf{s}_t^T \mathbf{P}_{t-1} \mathbf{s}_t}, \quad (31)$$

where

$$\epsilon_t = \mathbf{x}_t - \mathbf{A}_{t-1} \mathbf{s}_t, \quad (32)$$

6. A NUMERICAL EXAMPLE

A simple numerical example is given here to confirm fast convergence and high performance of the proposed algorithm. The algorithm is compared with the Oja's subspace rule [4]. The 3-dimensional observation vector \mathbf{x} was generated with its covariance matrix given by

$$\begin{bmatrix} 1.391 & 0.173 & -0.536 \\ 0.173 & 0.032 & -0.078 \\ -0.536 & -0.078 & 2.584 \end{bmatrix}. \quad (33)$$

As a performance measure, I use the subspace error (SE) which is defined by

$$\text{SE} = \frac{1}{\sqrt{n}} \left\| \mathbf{P}_{\hat{\mathbf{A}}}^{\perp} \mathbf{U}_n \mathbf{U}_n^T \right\|, \quad (34)$$

where $\mathbf{P}_{\hat{\mathbf{A}}}^{\perp}$ is the projection matrix onto the orthogonal subspace, i.e.,

$$\mathbf{P}_{\hat{\mathbf{A}}}^{\perp} = \mathbf{I} - \mathbf{A} \left(\mathbf{A}^T \mathbf{A} \right)^{-1} \mathbf{A}^T, \quad (35)$$

which is decided by the estimated values of \mathbf{A} and $\mathbf{U}_n \mathbf{U}_n^T$ is the projection matrix onto the signal subspace that is computed from the SVD of the covariance matrix of the input data. When the estimated \mathbf{A} spans the true signal subspace, $\mathbf{P}_{\hat{\mathbf{A}}}^{\perp}$ should be orthogonal to $\mathbf{U}_n \mathbf{U}_n^T$. Hence the SE defined above is able to serve as a performance measure.

The matrices \mathbf{A}_0 (for the sequential EM) and \mathbf{W}_0 (for Oja's subspace rule) were initialized as a random matrix whose elements are drawn from uniformly distributed random variables over [0,1]. The learning rate in the Oja's subspace rule was $\eta_t = .01$. In Fig. 1, one can observe that the sequential EM algorithm converges to a solution much faster than the Oja's subspace rule and even after convergence, the former shows slight better performance than the latter.

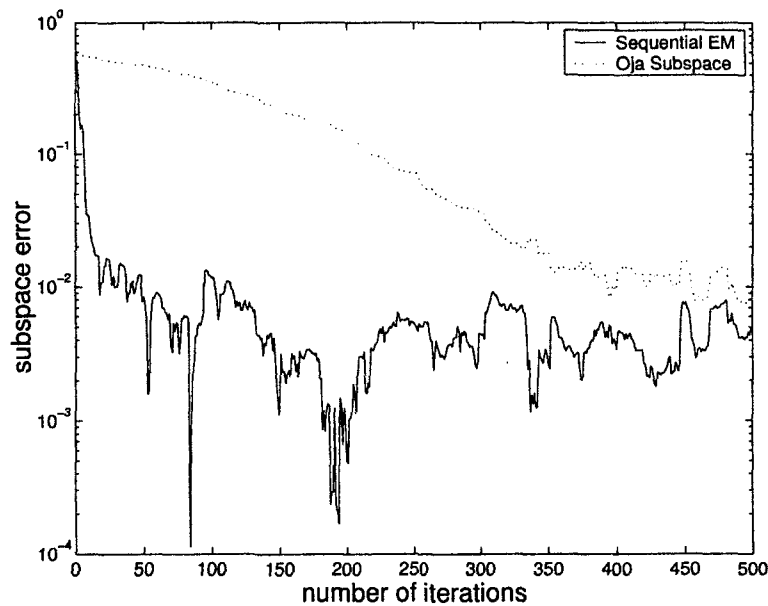


Figure 1: Convergence comparison for sequential EM algorithm and Oja's subspace rule in terms of subspace error.

7. CONCLUSION

I have described some relations between the EM algorithm and separable LS fitting in PPCA with zero noise limit. Based on this observation, I have developed a sequential EM learning for subspace analysis which is nothing but a sequential LS algorithm. The sequential EM algorithm can be easily extended to mixtures of PCA, which is currently under investigation.

8. ACKNOWLEDGMENT

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

9. REFERENCES

- [1] K. I. Diamantaras and S. Y. Kung. *Principal Component Neural Networks: Theory and Applications*. John Wiley & Sons, INC, 1996.
- [2] Z. Ghahramani and G. E. Hilton. The EM algorithm for mixtures of factor analyzers. Technical Report CRG-TR-96-1, Department of Computer Science, University of Toronto, 1997.
- [3] E. Oja. A simplified neuron model as a principal component analyzer. *J. Mathematical Biogology*, 15:267-273, 1982.
- [4] E. Oja. Neural networks, principal component analysis, and subspaces. *International Journal of Neural Systems*, 1:61-68, 1989.
- [5] S. Roweis. EM algorithms for PCA and SPCA. In *Advances in Neural Information Processing Systems*, volume 10, pages 626-632. MIT press, 1998.
- [6] S. Roweis and Z. Ghahramani. A unifying review of linear Gaussian models. *Neural Computation*, 11(2):305-345, 1999.
- [7] D. Rubin and D. Thayer. EM algorithms for ML factor analysis. *Psychometrika*, 47:69-76, 1982.
- [8] L. Sirovich and M. Kirby. Low-dimensional procedure for the characterization of human faces. *Journal of the Optical Society of America A*, 4(3):519-524, 1987.
- [9] M. E. Tipping and C. M. Bishop. Mixtures of probabilistic principal component analyzers. *Neural Computation*, 11(2):443-482, 1999.
- [10] M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society B*, 61(3):611-622, 1999.
- [11] M. Turk and A. Pentland. Eigenfaces for recognition. *Journal of Cognitive Neuroscience*, 3(1):71-86, 1991.