



**HAL**  
open science

# Learning nonlinear image manifolds by global alignment of local linear models

Jakob Verbeek

► **To cite this version:**

Jakob Verbeek. Learning nonlinear image manifolds by global alignment of local linear models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2006, 28 (8), pp.1236–1250. 10.1109/TPAMI.2006.166 . inria-00321131v2

**HAL Id: inria-00321131**

**<https://inria.hal.science/inria-00321131v2>**

Submitted on 8 Apr 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Learning non-linear image manifolds by global alignment of local linear models

Jakob Verbeek

**Abstract**—Appearance based methods, based on statistical models of the pixels values in an image (region) rather than geometrical object models, are increasingly popular in computer vision. In many applications the number of degrees of freedom (DOF) in the image generating process is much lower than the number of pixels in the image. If there is a smooth function that maps the DOF to the pixel values, then the images are confined to a low dimensional manifold embedded in the image space. We propose a method based on probabilistic mixtures of factor analyzers to (i) model the density of images sampled from such manifolds and (ii) recover global parameterizations of the manifold. A globally non-linear probabilistic two-way mapping between coordinates on the manifold and images is obtained by combining several, locally valid, linear mappings. We propose a parameter estimation scheme that improves upon an existing scheme, and experimentally compare the presented approach to self-organizing maps, generative topographic mapping, and mixtures of factor analyzers. In addition, we show that the approach also applies to find mappings between different embeddings of the same manifold.

**Index Terms**—Feature extraction or construction, Machine learning, Statistical image representation.

## I. INTRODUCTION

OVER the last two decades appearance based methods have become increasingly popular to solve computer vision problems such as object recognition and pose estimation [1]. These methods are based on statistical models of the pixel gray or color values rather than on geometrical models of the depicted objects. One can think of the images as vectors in a high dimensional image space which is spanned by the values of the individual pixels. Note that the dimensionality of the image space is typically large; e.g. when modelling small  $64 \times 64$  pixel images the image space has 4096 dimensions. It turns out that in many applications the images of interest are confined to a low dimensional manifold embedded in the high dimensional image space. Our interest in this paper is to model image collections sampled from such manifolds. We build upon a recent method [2] to (i) model the density of images sampled from such manifolds and (ii) recover global parameterizations of these manifolds in an unsupervised manner. A globally non-linear probabilistic mapping between coordinates on the manifold to images—and vice-versa—is obtained by combining several, locally valid, linear mappings.

In general, reliable estimation of the parameters of densities that implement an appearance based model requires many images due to the high dimensionality of the image space. To alleviate this problem, often the simplifying assumption is made that, although the images are high dimensional, the set of

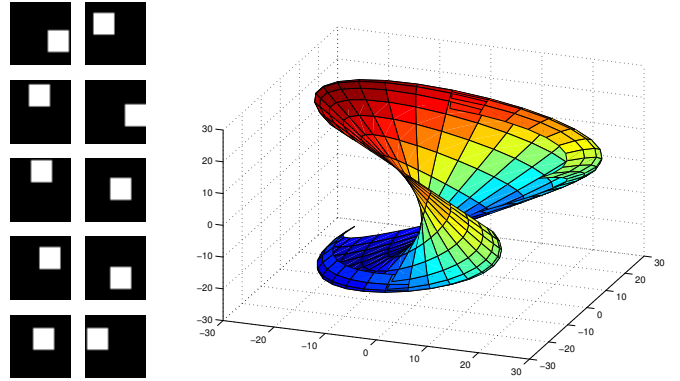


Fig. 1. Ten randomly selected images of the shifted square (left). Images projected on the first three principal components, projections of images of similar shifts are connected (right).

images considered in a given application are confined to a low dimensional subspace of the high dimensional image space. This idea was introduced in [3], where principal component analysis (PCA) [4] was used to identify, for a given set of images, a linear subspace which contains most of the variance in the images. The low dimensional representation obtained by projecting an image on this subspace then gives an informative characterization of the actual image. Later, others proposed to use the projections on a low dimensional linear subspace for tasks such as recognition and pose estimation [5], [6]. The idea that the images are confined to a low, say  $d$ , dimensional subspace of the image space (say of dimension  $D$ ) is supported by the following reasoning. If there are only  $d \ll D$  degrees of freedom (DOF) in the system that generates the images, and small changes in the DOF lead to small changes in the images, then the images will be confined to a  $d$  dimensional manifold smoothly embedded in the  $D$  dimensional image space.<sup>1</sup>

In many applications images have indeed been found to be confined to a low dimensional manifold. However, the manifold is often non-linearly embedded in the image space. Fig. 1 gives a simple illustration. A set of 400 images of  $29 \times 29$  pixels containing a white square of  $10 \times 10$  pixels against a black background is considered. Here, there are two DOF: corresponding to the horizontal and vertical position of the white square. The right panel of Fig. 1 depicts the projections of the images on the first three principal components (note: this projection is linear). Images where the position of the square differs only one pixel are connected to reveal the manifold structure. Clearly, the manifold spanned by the DOF is non-

<sup>1</sup>With the DOF in the system that generates the images we mean both those in the imaged object and those in the camera.

linearly embedded in the image space.

This paper is concerned with modelling the DOF by recovering, or ‘learning’, the underlying manifold structure given just a set of high dimensional points (images) sampled from this manifold. In specific, we want to find a function that maps points in the image space to the corresponding low dimensional coordinates on the manifold and vice versa. In this manner we obtain a compact representation of the images in terms of the DOF rather than the, somewhat arbitrary and redundant, coordinates of the image space. The learning task is ‘unsupervised’ in the sense that for none of the images the coordinates on the manifold are known in advance.

A wide range of techniques has been proposed for unsupervised learning of non-linear manifolds, such as auto-encoder neural networks [7], principal curves and surfaces [8], [9], [10], self-organizing maps [11], [12], and generative topographic mapping [13]. For more extensive overviews see [14], [15], [16]. All these methods suffer from one or more of the following drawbacks: (i) parameter estimation can return suboptimal solution since it is based on minimizing an error function with local minima that are not global minima, and (ii) there is no notion of the uncertainty in the estimated low dimensional coordinates of the images. Recently various techniques, such as Isomap [17], Locally Linear Embedding (LLE) [18], Kernel PCA [19], charting [20], Locality Preserving Projections [21], [22], Laplacian Eigenmaps [23], and a semi-definite programming approach [24], have been proposed that formalize manifold learning as minimizing a convex error function that encodes how well certain inter-point relationships are preserved. Due to the convexity of the error function there is no risk that minimization yields poor local minima rather than the global minimum of the error function, so drawback (i) is resolved for these methods. However, these methods do suffer from drawback (ii) mentioned above. In addition, these methods do not provide a parametric function that maps between the image space and the low dimensional manifold. Non-parametric methods can be used to map between the spaces, but in principle they require storage and access to all training data [25], which is costly for large high dimensional data sets. If the goal is to find the manifold structure in a given training set only, e.g. for visualization, then this is not a problem. However, if we want to generalize the recovered manifold structure to new data (i.e. represent new images in terms of the DOF) then this is problematic because of storage and computational requirements.

In this paper, in contrast to these non-parametric methods, we consider a probabilistic model that combines several, locally valid, linear mappings to define a globally non-linear mapping between the image space and the low dimensional manifold. Once the parameters have been estimated from a set of images, a non-linear mapping is obtained which can be used to map previously unseen images  $\mathbf{x}$  to their low dimensional coordinates on the manifold  $\mathbf{z}$  and vice versa, without the need to store the original training data. In fact, the model yields a conditional density on  $\mathbf{z}$  given  $\mathbf{x}$  and vice-versa, and thus also gives a notion of uncertainty in the mapping. We note here that the required amount of training data to learn the model necessarily increases with the dimensionality of the manifold,

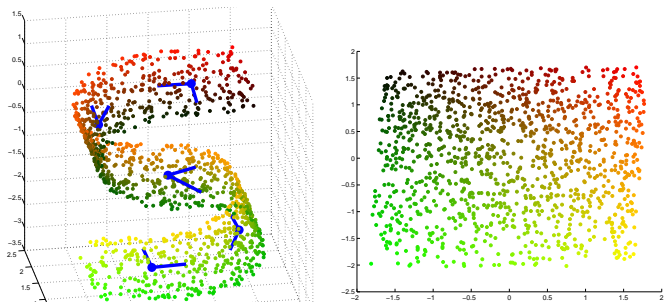


Fig. 2. Data in  $\mathbb{R}^3$  with local two dimensional subspaces indicated by the axes (left). The Desired global two dimensional data representation unfolds the non-linear manifold (right).

for reasons of numerical stability.<sup>2</sup>

Intuitively, the approach taken in this paper can be understood as a combination of clustering and PCA. The idea is illustrated in Fig. 2; although globally the data is spread out in all three dimensions, locally the data is concentrated around a two dimensional linear subspace. Thus, the data can be clustered in such a manner that the data within each cluster can be accurately reconstructed from a two dimensional PCA projection. Several authors, e.g. [26], [27], [28], have reported that such a combination of clustering and PCA allows for significantly better reconstructions of images as compared to a reconstruction using one single linear PCA subspace. Others [29] successfully used such a ‘mixture of PCAs’ (MPCA) to classify images of hand written digits; a MPCA model was learned on images of each digit  $i = 0, \dots, 9$ , yielding density models  $p_i$  ( $i = 0, \dots, 9$ ). The learned models were used to classify new images  $\mathbf{x}$  of handwritten digits by classifying them as digit  $i = \arg \max_j p_j(\mathbf{x})$ .

An important observation is that in all these works [26], [27], [28], [29] each cluster provides a *separate* low dimensional coordinate system that is only used for data in that cluster. Therefore, it is not possible to directly interpolate between images that are assigned to different clusters since a global low dimensional structure is lacking. Finding a global low dimensional representation is also useful as a preprocessing stage for supervised learning problems since the low dimensionality safeguards against the curse of dimensionality [30]. Thus, once a global low dimensional representation is found based on many unsupervised examples, it is possible to learn functions on the manifold using relatively few supervised examples.

In the next section we describe probabilistic mixtures of factor analyzers, and how they can be used to obtain a global low dimensional representation. We present a parameter estimation technique, which improves upon earlier work by Roweis et al. [2]. Then, in Section III, we show how this approach applies to correspondence learning, a recently introduced learning problem [31]. The learning data consists of two sets of high dimensional points related by correspondences, i.e. some points in one set are known to have the same low dimensional coordinate as a point in the other set. For example, each set could contain images of a different object as viewed

<sup>2</sup>However, to our knowledge it is unclear how the amount of data should depend on the manifold dimensionality. In our experiments, with manifold dimensionalities up to three, we did not encounter numerical instabilities.

from different directions. The correspondences are pairs of images where the two objects are viewed from the same direction. Non-linear dimension reduction for the two sets of points simultaneously allows us to predict the correspondences for other points. In Section IV we present the experimental results to evaluate the presented methods. Finally, in Section V we present our conclusions.

## II. MODELLING DATA SAMPLED FROM MANIFOLDS

In this section we describe a probabilistic model for data sampled from non-linear manifolds. In Section II-A we describe linear subspace densities which can be combined in a mixture density to model non-linear manifolds. Then, in Section II-B, we describe a framework to simultaneously estimate the parameters of such mixtures and to obtain a global parameterization of the manifold. In Section II-C we present our improved parameter estimation procedure. Since the optimization procedure is only guaranteed to find local optima of the objective function, careful initialization of the parameters is required. In Section II-D we consider parameter initialization methods.

### A. Mixtures of linear models

We assume that the high dimensional data vectors  $\mathbf{x}_n \in \mathbb{R}^D$  ( $n = 1, \dots, N$ ) are sampled independently and identically from a smooth non-linear manifold and are potentially corrupted by additive uncorrelated Gaussian noise. Thus, the data is distributed on, or near, a low dimensional manifold embedded in a high dimensional space. If the manifold is sufficiently smooth, i.e. locally it is linearly embedded in the high dimensional space, then we can model the data density with a mixture of densities that concentrate their mass in a linear subspace. To this end we can use mixtures of factor analyzers (MFA).

A MFA density over data vectors  $\mathbf{x} \in \mathbb{R}^D$  can be specified by introducing for each observed data vector two ‘hidden’ (or ‘unobserved’) variables:  $c$  and  $\mathbf{z}$ . The first variable,  $c \in \{1, \dots, C\}$ , is an index over the  $C$  components of the mixture. The second variable  $\mathbf{z} \in \mathbb{R}^d$  is a coordinate in the  $d$ -dimensional subspace associated with the mixture component with index given by  $c$ . Thus,  $\mathbf{z}$  may be interpreted as a coordinate on one of several local linear approximations of the manifold. The density on vectors  $\mathbf{x}$  follows from the joint density over  $\mathbf{x}, \mathbf{z}$ , and  $c$ :

$$p(\mathbf{x}, \mathbf{z}, c) = p(\mathbf{x}|\mathbf{z}, c)p(\mathbf{z}|c)p(c), \quad (1)$$

$$p(\mathbf{z}|c) = \mathcal{N}(\mathbf{z}; \boldsymbol{\kappa}_c, \boldsymbol{\Sigma}_c), \quad (2)$$

$$p(\mathbf{x}|c, \mathbf{z}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_c + \boldsymbol{\Lambda}_c(\mathbf{z} - \boldsymbol{\kappa}_c), \boldsymbol{\Psi}_c), \quad (3)$$

where  $\mathcal{N}(\mathbf{z}; \boldsymbol{\kappa}, \boldsymbol{\Sigma})$  denotes the Gaussian density on  $\mathbf{z}$  with mean  $\boldsymbol{\kappa}$  and covariance matrix  $\boldsymbol{\Sigma}$ . The density on  $\mathbf{x}$  given  $c$  is obtained by marginalizing over  $\mathbf{z}$  and given by:

$$p(\mathbf{x}|c) = \int p(\mathbf{x}|\mathbf{z}, c)p(\mathbf{z}|c) d\mathbf{z} = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_c, \boldsymbol{\Gamma}_c), \quad (4)$$

$$\boldsymbol{\Gamma}_c = \boldsymbol{\Lambda}_c \boldsymbol{\Sigma}_c \boldsymbol{\Lambda}_c^\top + \boldsymbol{\Psi}_c. \quad (5)$$

Finally, the marginal density on  $\mathbf{x}$  is then given by the mixture:

$$p(\mathbf{x}) = \sum_{c=1}^C p(c)p(\mathbf{x}|c). \quad (6)$$

The columns of the matrix  $\boldsymbol{\Lambda}_c$  span the subspace of component  $c$  which has an offset  $\boldsymbol{\mu}_c$  from the origin. The uncertainty in the latent coordinate  $\mathbf{z}$  is mapped by  $\boldsymbol{\mu}_c$  and  $\boldsymbol{\Lambda}_c$  into uncertainty in  $\mathbf{x}$  within the subspace. The diagonal matrix  $\boldsymbol{\Psi}_c$ , with positive entries, adds variance in all dimensions, and makes  $\boldsymbol{\Gamma}_c$  a proper, i.e. positive definite, covariance matrix.<sup>3</sup>

The subspace spanned by  $\boldsymbol{\Lambda}_c$  provides a coordinate system which can be used to reconstruct the data from low dimensional coordinates  $\mathbf{z} \in \mathbb{R}^d$ , at least for data that receives high likelihood under  $p(\mathbf{x}|c)$ . The joint model  $p(\mathbf{x}, \mathbf{z}, c)$  induces a Gaussian density on  $\mathbf{z}$  given  $\mathbf{x}$  and  $c$ :

$$p(\mathbf{z}|\mathbf{x}, c) = \mathcal{N}(\mathbf{z}; \mathbf{m}_c, \mathbf{V}_c^{-1}), \quad (7)$$

$$\mathbf{V}_c = \boldsymbol{\Sigma}_c^{-1} + \boldsymbol{\Lambda}_c^\top \boldsymbol{\Psi}_c^{-1} \boldsymbol{\Lambda}_c, \quad (8)$$

$$\mathbf{m}_c = \boldsymbol{\kappa}_c + \mathbf{V}_c^{-1} \boldsymbol{\Lambda}_c^\top \boldsymbol{\Psi}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c), \quad (9)$$

which can be used to infer the coordinates  $\mathbf{z}$  in the local subspace given a data point  $\mathbf{x}$  and a mixture component  $c$ . The EM algorithm [34], [35] can be used to estimate the parameters of the mixture to find a (local) maximum of the data log-likelihood:

$$\mathcal{L} = \sum_{n=1}^N \log p(\mathbf{x}_n). \quad (10)$$

The inverse and determinant of  $\boldsymbol{\Gamma}_c$ , required to evaluate  $p(\mathbf{x}|c)$ , are efficiently computed as:

$$\boldsymbol{\Gamma}_c^{-1} = \boldsymbol{\Psi}_c^{-1} (\boldsymbol{\Psi}_c - \boldsymbol{\Lambda}_c \mathbf{V}_c^{-1} \boldsymbol{\Lambda}_c^\top) \boldsymbol{\Psi}_c^{-1}, \quad (11)$$

$$|\boldsymbol{\Gamma}_c| = |\boldsymbol{\Psi}_c| \times |\boldsymbol{\Sigma}_c| \times |\mathbf{V}_c|. \quad (12)$$

Note that the model is ‘over-parameterized’:  $\boldsymbol{\kappa}_c$  does not appear in the likelihood  $p(\mathbf{x}|c)$  and any positive definite matrix  $\boldsymbol{\Sigma}_c = \mathbf{U}^\top \mathbf{U} \neq \mathbf{I}$  can be absorbed in  $\boldsymbol{\Lambda}_c$  by setting  $\boldsymbol{\Lambda}_c \leftarrow \boldsymbol{\Lambda}_c \mathbf{U}^\top$  ( $\mathbf{U}$  is the Cholesky factor of  $\boldsymbol{\Sigma}_c$ ). Therefore, the likelihood is invariant to linear transformations of the local coordinates  $\mathbf{z}$ ; in principle we can set  $\boldsymbol{\kappa}_c = \mathbf{0}$  and  $\boldsymbol{\Sigma}_c = \mathbf{I}$ .

### B. Aligning local linear models

Above we defined a density model for data sampled from non-linear manifolds. Next we describe how we can use this model in order to recover a global parameterization of the manifold by aligning the local linear models. We define a learning criterion that combines the data log-likelihood and a measure of the alignment. Therefore, the estimates of the local models can be modified to accommodate a better alignment.

In contrast, other approaches for alignment of local models [36], [37], [38] proceed in two steps: first a mixture of local linear models is estimated (typically on the basis of maximum likelihood), and second the local models are aligned on the

<sup>3</sup>If the data contains outliers, robust estimates of the local subspaces can be obtained using a t-distribution noise model. See [32] for a derivation of the EM algorithm for t-distributions, and [33] for an application of mixtures of local linear models based on t-distributions to image manifold learning.

basis of a second criterion. Such approaches thus do not allow modification of the local models to improve alignment. In [39] a mixture of linear models is proposed, but this assumes that the global low dimensional coordinates are known. Thus in the latter approach the local models can be adapted, but the global low dimensional coordinates cannot. In the approach described below both the local models and the estimates of the global low dimensional coordinates of the data are iteratively re-estimated to obtain an alignment of the local models.

We assume that each data point has a unique (but unknown) global coordinate on the manifold and that locally (i.e. in a region that is assigned high likelihood by a single mixture component) there is a linear transformation that maps the local coordinates  $\mathbf{z}$  to the global coordinates on the manifold. Since the likelihood is invariant to linear transformations of the local coordinates  $\mathbf{z}$  we may choose  $\kappa_c$  and  $\Sigma_c$  to implement this linear transformation. Thus, under these assumptions, there are settings of  $\kappa_c$  and  $\Sigma_c$  such that  $\mathbf{z}$  corresponds to the global coordinates on the manifold.

Since the global coordinates are assumed to be unique, if a data point  $\mathbf{x}$  is assigned high likelihood by two or more mixture components, then the global coordinates as predicted with  $p(\mathbf{z}|\mathbf{x}, c)$  using different mixture components should be, approximately, the same. In other words: mixture components with high posterior  $p(c|\mathbf{x})$  should ‘agree’ on the global coordinate  $\mathbf{z}$  of point  $\mathbf{x}$ . Since  $\mathbf{z}$  is not known with certainty, a formal notion of ‘agreement’ between the mixture components needs to take into account the uncertainty in  $\mathbf{z}$  as given by  $p(\mathbf{z}|\mathbf{x}, c)$ . The conditional density on  $\mathbf{z}$  given  $\mathbf{x}$  takes the form of a mixture of Gaussian densities:

$$p(\mathbf{z}|\mathbf{x}) = \sum_{c=1}^C p(c|\mathbf{x})p(\mathbf{z}|\mathbf{x}, c). \quad (13)$$

If several mixture components with non-negligible posterior  $p(c|\mathbf{x})$  for a data point  $\mathbf{x}$  yield quite different corresponding densities  $p(\mathbf{z}|\mathbf{x}, c)$  on the global coordinates, then they do not ‘agree’ and the mixture  $p(\mathbf{z}|\mathbf{x})$  exhibits several modes. On the other hand, if all components with non-negligible posterior  $p(c|\mathbf{x})$  yield exactly the same density  $p(\mathbf{z}|\mathbf{x}, c)$ , then the components are in perfect agreement and the mixture  $p(\mathbf{z}|\mathbf{x})$  is Gaussian and exhibits a single mode.

Our goal is to find (i) a density model for the high dimensional data, and (ii) a global parameterization of the manifold from which the data was sampled. The above observations suggest that to find the parameters  $\kappa_c$  and  $\Sigma_c$  such that the local coordinates on the linear subspaces correspond to the global coordinates on the manifold, we can add a penalty term to the data log-likelihood that penalizes  $p(\mathbf{z}|\mathbf{x})$  that depart from a (unimodal) Gaussian distribution. The penalty term we use measures for each data point  $\mathbf{x}_n$  how much the mixture  $p(\mathbf{z}|\mathbf{x}_n)$  resembles a Gaussian  $q_n(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{z}_n, \Sigma_n)$ , through the Kullback-Leibler (KL) divergence

$$\mathcal{D}(q_n(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}_n)) = \int q_n(\mathbf{z}) \log \frac{q_n(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}_n)} d\mathbf{z}. \quad (14)$$

The penalized log-likelihood objective function reads:

$$\mathcal{L}' = \sum_{n=1}^N \left[ \log p(\mathbf{x}_n) - \mathcal{D}(q_n(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}_n)) \right] \quad (15)$$

$$= \sum_{n=1}^N \left[ \mathcal{H}(q_n(\mathbf{z})) + \int q_n(\mathbf{z}) \log p(\mathbf{x}_n, \mathbf{z}) d\mathbf{z} \right], \quad (16)$$

where  $\mathcal{H}$  denotes the entropy of a distribution, defined as:

$$\mathcal{H}(q_n(\mathbf{z})) = - \int q_n(\mathbf{z}) \log q_n(\mathbf{z}) d\mathbf{z}. \quad (17)$$

Note that the objective  $\mathcal{L}'$  does not only depend on the parameters of the mixture model  $p$ , but also on the distributions  $q_n$  which can be regarded as probabilistic estimates of the global coordinates of the data.<sup>4</sup>

The ‘mixture of experts’ (MoE) approach to classification and regression [41] is related to the method presented here. The MoE approach constructs a complex (e.g. non-linear) classification or regression function by combining several ‘experts’. Each expert implements a simple (e.g. linear) classification or regression function that is suitable for some subset of the data space. To determine which expert should be used to classify a new data point, a gating network is used. The gating network produces weighting factors (that are positive and sum to one) by which the predictions of the individual experts are averaged. The output of both the experts and the gating network depend on the presented data point. Analogously, in the model described above, the non-linear mapping between high dimensional coordinates  $\mathbf{x}$  and latent coordinates  $\mathbf{z}$  defined in (13) is also a weighted average of simple linear Gaussian dependencies  $p(\mathbf{z}|\mathbf{x}, c)$  weighted by factors  $p(c|\mathbf{x})$  that switch between the ‘experts’. The main difference is that in this paper we consider unsupervised learning while MoE is used for supervised learning and parameters are found that maximize the probability of the desired output given the inputs.

### C. Parameter estimation

The second term in the summands of our objective function defined in (16) is an expectation of the logarithm of the mixture likelihood  $p(\mathbf{x}_n, \mathbf{z}) = \sum_{c=1}^C p(c)p(\mathbf{x}_n, \mathbf{z}|c)$ , and as a result  $\mathcal{L}'$  exhibits local maxima that are not global maxima. Below, we consider how we can use an EM-like algorithm to find a (local) maximum of the objective  $\mathcal{L}'$ . For each data point we introduce a distribution  $q_n(c)$  over the mixture components, in order to bound the mixture log-likelihoods:

$$\log p(\mathbf{x}_n, \mathbf{z}) \geq \log p(\mathbf{x}_n, \mathbf{z}) - \mathcal{D}(q_n(c)\|p(c|\mathbf{x}_n, \mathbf{z})) \quad (18)$$

$$= \mathcal{H}(q_n(c)) + \sum_{c=1}^C q_n(c) \log p(\mathbf{x}_n, \mathbf{z}, c). \quad (19)$$

<sup>4</sup>It is possible to weight the penalty term by a factor different from one, yielding an objective function similar to the one used in [40] for accelerated maximum likelihood estimation. It is straightforward to derive an optimization algorithm similar to the one presented here if the penalty term is multiplied by a factor in  $[0, 1]$ .

The bounds on the individual mixture log-likelihoods can be combined to bound the complete objective function:

$$\begin{aligned} \mathcal{L}' \geq \Phi = & \sum_{n=1}^N \mathcal{H}(q_n(c)) + \mathcal{H}(q_n(\mathbf{z})) \\ & + \sum_{n=1}^N \sum_{c=1}^C q_n(c) \int q_n(\mathbf{z}) \log p(\mathbf{x}_n, \mathbf{z}, c) \, d\mathbf{z}. \end{aligned} \quad (20)$$

We can now iteratively increase  $\Phi$ , analogous to the EM algorithm, by maximizing it in turn with respect to the parameters of  $q_n(c)$ ,  $q_n(\mathbf{z})$  and  $p$  respectively. We call this EM-like maximization of  $\Phi$  the Coordinated Factor Analysis (CFA) algorithm.

Using the abbreviations,  $q_{nc} = q_n(c)$ ,  $\mathbf{x}_{nc} = \mathbf{x}_n - \boldsymbol{\mu}_c$ , and  $\mathbf{z}_{nc} = \mathbf{z}_n - \boldsymbol{\kappa}_c$ , we can write  $\Phi$  as:

$$\Phi = \sum_{n=1}^N \sum_{c=1}^C q_{nc} [\mathcal{S}_{nc} - \mathcal{E}_{nc}], \quad (21)$$

$$\mathcal{S}_{nc} = \frac{1}{2} \log |\boldsymbol{\Sigma}_n| - \log q_{nc} + \frac{d}{2} \log(2\pi), \quad (22)$$

$$\begin{aligned} \mathcal{E}_{nc} = & -\log p(c) + \frac{D+d}{2} \log(2\pi) + \frac{1}{2} \log |\boldsymbol{\Sigma}_c| + \frac{1}{2} |\boldsymbol{\Psi}_c| \\ & + \frac{1}{2} \text{Tr}\{\boldsymbol{\Sigma}_c(\boldsymbol{\Sigma}_n + \mathbf{z}_{nc}\mathbf{z}_{nc}^\top)\} + \frac{1}{2} \text{Tr}\{\boldsymbol{\Sigma}_n \boldsymbol{\Lambda}_c^\top \boldsymbol{\Psi}_c^{-1} \boldsymbol{\Lambda}_c\} \\ & + \frac{1}{2} (\mathbf{x}_{nc} - \boldsymbol{\Lambda}_c \mathbf{z}_{nc})^\top \boldsymbol{\Psi}_c^{-1} (\mathbf{x}_{nc} - \boldsymbol{\Lambda}_c \mathbf{z}_{nc}). \end{aligned} \quad (23)$$

To maximize  $\Phi$  with respect to the distributions  $q_n(c)$  we equate the corresponding partial derivatives to zero and find the maximizing assignment as:

$$q_{nc} \leftarrow \left[ \sum_{c'=1}^C e^{-\mathcal{E}_{nc'}} \right]^{-1} e^{-\mathcal{E}_{nc}}. \quad (24)$$

Similarly, to maximize  $\Phi$  with respect to the  $q_n(\mathbf{z})$  we set:

$$\boldsymbol{\Sigma}_n^{-1} \leftarrow \sum_{c=1}^C q_{nc} \mathbf{V}_c, \quad \mathbf{z}_n \leftarrow \boldsymbol{\Sigma}_n \sum_{c=1}^C q_{nc} \mathbf{V}_c \mathbf{m}_{nc}.$$

Given  $q_n(c)$  and  $q_n(\mathbf{z})$ , we maximize  $\Phi$  with respect to the parameters of  $p$ . With  $\tilde{q}_{nc} = q_{nc} / \sum_{n'=1}^N q_{n'c}$ , the maximizing assignments for the mixing weights and means are:

$$p(c) \leftarrow \frac{1}{N} \sum_{n=1}^N q_{nc}, \quad \boldsymbol{\kappa}_c \leftarrow \sum_{n=1}^N \tilde{q}_{nc} \mathbf{z}_n, \quad \boldsymbol{\mu}_c \leftarrow \sum_{n=1}^N \tilde{q}_{nc} \mathbf{x}_n.$$

In [2] the authors suggested to use an fixed-point algorithm to find the stationary points of  $\Phi$  w.r.t. the remaining parameters of the model:  $\boldsymbol{\Sigma}_c$ ,  $\boldsymbol{\Lambda}_c$  and  $\boldsymbol{\Psi}_c$ . However, the parameters can be found in closed-form as shown in [42] for a constrained version of the above model and in [14] for the unconstrained model.<sup>5</sup> In Section IV-A we experimentally show that the closed-form algorithm is computationally more efficient than the fixed-point algorithm, a more detailed comparison can be found in [14]. Using the weighted correlations and variances:

$$\mathbf{S}_c = \sum_{n=1}^N \tilde{q}_{nc} \mathbf{x}_{nc} \mathbf{z}_{nc}^\top, \quad \mathbf{G}_c = \sum_{n=1}^N \tilde{q}_{nc} [\mathbf{z}_{nc} \mathbf{z}_{nc}^\top + \boldsymbol{\Sigma}_n],$$

<sup>5</sup>In [2] the model was parameterized in a slightly different manner, which somewhat obscured the closed-form solution.

the closed-form equations are:

$$\begin{aligned} \boldsymbol{\Sigma}_c & \leftarrow \mathbf{G}_c, \quad \boldsymbol{\Lambda}_c \leftarrow \mathbf{S}_c \boldsymbol{\Sigma}_c^{-1}, \\ [\boldsymbol{\Psi}_c]_{ii} & \leftarrow \sum_{n=1}^N \tilde{q}_{nc} \left( [\mathbf{x}_{nc} - \boldsymbol{\Lambda}_c \mathbf{z}_{nc}]_i^2 + [\boldsymbol{\Lambda}_c \boldsymbol{\Sigma}_n \boldsymbol{\Lambda}_c^\top]_{ii} \right), \end{aligned}$$

where we used  $[\cdot]_{ii}$  and  $[\cdot]_i$  to denote the  $(i, i)^{th}$  entry of a matrix and the  $i^{th}$  entry of a vector respectively. The computational complexity of once updating the distributions and all parameters is  $O(NCDd + Nd^3)$ . Unless there are fewer mixture components than latent dimensions ( $C < d$ ), this is linear in the number of data points  $N$ , the number of mixture components  $C$ , the dimensionality of the data  $D$  and the dimensionality of the manifold  $d$ .

#### D. Parameter initialization

Like most other EM algorithms, the algorithm described above can terminate at poor local optima of the objective function. Careful parameter initialization, especially of the  $\mathbf{z}_n$ , is needed to successfully use the algorithm. In [2] the authors proposed to initialize the  $\mathbf{z}_n$  by using some other unsupervised non-linear dimension reduction technique. We can then apply the EM-like algorithm described above until convergence while keeping the  $\mathbf{z}_n$  fixed at their initialized value. The  $\boldsymbol{\Sigma}_n$  are also kept fixed at a small multiple of identity and the  $q_{nc}$  are initialized at random near uniform. In this manner we learn the model on the ‘complete’ data, i.e. each data vector  $\mathbf{x}_n$  is augmented with its estimated global coordinate  $\mathbf{z}_n$ . After such an initialization phase, the  $\mathbf{z}_n$  and  $\boldsymbol{\Sigma}_n$  can be updated as well, to obtain better estimates of the global coordinates and their uncertainties. Note that the overall procedure is still unsupervised.

Other non-linear dimension reduction method could be used for the initialization. However, since we will use it to avoid poor local optima in the optimization procedure of the last section, the initialization should not suffer from the same problem. Therefore the non-parametric manifold learning techniques based on minimization of convex error functions, mentioned in Section I, are particularly useful for this purpose. Methods that first fit a mixture of local linear models [36], [37], [38] are computationally more efficient, but we found them to give suboptimal results as compared to the non-parametric methods. In our experiments described in Section IV we use the LLE algorithm to initialize the latent coordinates, mainly because of its computational efficiency.

Finally, the number of mixture components  $C$ , and the latent dimensionality  $d$  have to be set. If the number of degrees of freedom in the data is not known<sup>6</sup>, then  $d$  may be estimated using a variety of techniques [43], [44], [45]. Given a latent dimensionality  $d$ , the number of mixture components  $C$  can be set by employing cross-validation or a penalized log-likelihood criterion, e.g. following [46], [47]. Alternatively, a variational Bayesian learning approach [48] can be taken to estimate both the latent dimensionality and the number of components.

<sup>6</sup>If known,  $d$  can be set to the number of degrees of freedom. However, this practice is somewhat controversial due to the existence of space-filling curves, such as Sierpinski curves, and the possibility of periodic degrees of freedom, see also Section IV-D.

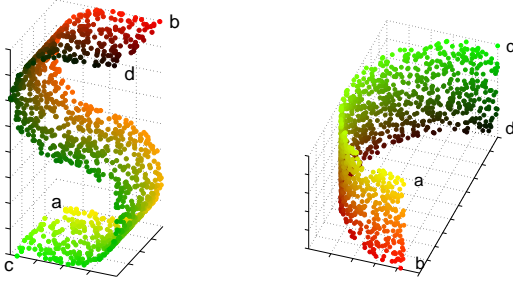


Fig. 3. Two data sets (left and right panel), the four pairs of corresponding points have been labelled ‘a’-‘d’.

### III. MAPPING BETWEEN MANIFOLD EMBEDDINGS

In the previous section we presented the CFA model for non-linear dimension reduction. In this section, we consider how it applies to the case where multiple embeddings of the same underlying low dimensional manifold are observed and we want to learn a mapping between these embeddings [31], [37], [49]. So rather than one set of high dimensional data points, we are now given two sets of high dimensional data points:  $\{\mathbf{x} \in \mathbb{R}^{D_x}\}$  and  $\{\mathbf{y} \in \mathbb{R}^{D_y}\}$ . The two sets are related through a set of correspondences: pairs  $(\mathbf{x}, \mathbf{y})$  which are known to share the same low dimensional coordinate on the manifold. For example, the two sets may be images of objects  $x$  and  $y$  viewed from different directions, or images of the face of persons  $x$  and  $y$  in different poses or with different facial expressions. A mapping between the manifolds makes it possible to map e.g. an image sequence of person  $x$  with a varying facial expression or pose to an image sequence of person  $y$ , in such a manner that only the identity is changed from  $x$  to  $y$  but the sequence of facial expressions or poses is preserved. Such techniques could be of considerable interest for video synthesis applications.

An abstract illustration of this setting is given in Fig. 3: two data sets are plotted respectively in the left and right panel and the correspondences are indicated by the letters ‘a’-‘d’. The goal is to predict missing correspondences, i.e. for a point in one set without a correspondence predict the coordinates of the point in the other set that would correspond to it. In other words, we consider a prediction problem where both the ‘inputs’ and ‘outputs’ are intrinsically low dimensional, but specified as vectors in high dimensional spaces.

One can also think of the data as points in the product space  $\mathbb{R}^{D_x} \times \mathbb{R}^{D_y} = \mathbb{R}^{D_x+D_y}$ . Corresponding pairs  $(\mathbf{x}, \mathbf{y})$  can be regarded as points in the product space, the first  $D_x$  coordinates are given by  $\mathbf{x}$  and the remaining  $D_y$  coordinates are given by  $\mathbf{y}$ . Points  $\mathbf{x}$  for which we do not have a correspondence, can be regarded as points in the product space for which the first  $D_x$  coordinates are given by  $\mathbf{x}$  and the other coordinates are not observed, and similarly for points  $\mathbf{y}$  without a correspondence. The goal is to estimate a density  $p(\mathbf{x}, \mathbf{y})$  on the product space, which induces the conditional densities  $p(\mathbf{x}|\mathbf{y})$  and  $p(\mathbf{y}|\mathbf{x})$  which are used to predict correspondences.

The method described in this section is related to the parameterized self-organizing map (PSOM) [50]. The PSOM

also produces a mapping between two high dimensional spaces through an underlying low dimensional representation. The basic idea is to first find a low dimensional representation<sup>7</sup> of the set of given input-output pairs  $(\mathbf{x}_n, \mathbf{y}_n)$  ( $n = 1, \dots, N$ ) and then to find a vector-valued function  $\mathbf{f}$  that optimally reproduces each input-output pair from its low dimensional representation. Thus if  $\mathbf{z}_n$  is the vector of low-dimensional coordinates of the  $n$ -th pair, then  $\mathbf{f}$  is such that for  $n = 1, \dots, N$  we have  $\mathbf{f}(\mathbf{z}_n) \approx (\mathbf{x}_n, \mathbf{y}_n)$ . The function  $\mathbf{f}$  is taken to be a linear combination of non-linear basis functions. To predict the output for a new input  $\mathbf{x}$ , low dimensional coordinates  $\mathbf{z}^*$  are found such that  $\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{x} - \mathbf{f}_x(\mathbf{z})\|^2$ , where  $\mathbf{f}_x$  is the function  $\mathbf{f}$  restricted to the part corresponding to  $\mathbf{x}$ . The prediction on  $\mathbf{y}$  is then given by  $\mathbf{f}_y(\mathbf{z}^*)$ . Compared to the CFA approach the PSOM has two drawbacks: (i) PSOM can not use observations without a correspondence, and (ii) to map between the two spaces we have to find the low dimensional coordinates  $\mathbf{z}^*$  that optimally reproduce the input. To find  $\mathbf{z}^*$  a non-convex error function has to be minimized, which may produce erroneous results if a local but non-global minimum is found. Another drawback of PSOM is that for a given input it only produces an output where the CFA approach produces a density over possible outputs.

Another approach to the correspondence problem is to fit a MFA to the data in the product space. Note that the data has many missing values: for points without a correspondence the first  $D_x$  or the last  $D_y$  coordinates are not observed. The EM algorithm can be used to estimate the parameters of a MFA in the presence of missing values [35], [51]. The difference with the CFA approach is that the MFA neglects the *global* manifold structure of the data; it does not relate the low dimensional coordinate systems. Note that if a mixture component  $c$  in a MFA does not model any correspondences (i.e. all correspondences have negligible posterior  $p(c|\mathbf{x}, \mathbf{y})$  on component  $c$ ), then it is not possible to (reliably) estimate the dependency between  $\mathbf{x}$  and  $\mathbf{y}$  for this component. Contrary, the CFA approach estimates global low dimension coordinates for all data points. Therefore, even for mixture components which do not receive any posterior probability from any correspondence, the dependency between  $\mathbf{x}$  and  $\mathbf{y}$  can be estimated through the dependence of both  $\mathbf{x}$  and  $\mathbf{y}$  on the global low dimensional coordinates  $\mathbf{z}$ .

In Section III-A we consider how, with some small changes, the CFA model also applies to this new setting. Then, in Section III-B, we consider suitable initialization techniques.

#### A. The probabilistic model and parameter estimation

The probabilistic model and learning algorithm presented in Section II can be adapted to apply to the current problem. Recall that in Section II we assumed the densities:

$$p(\mathbf{z}|c) = \mathcal{N}(\mathbf{z}; \boldsymbol{\kappa}_c, \boldsymbol{\Sigma}_c), \quad (25)$$

and we postulated a density  $p(\mathbf{x}|\mathbf{z}, c)$  where all elements of the vector  $\mathbf{x}$  are distributed independently given  $\mathbf{z}$  and  $c$ . Here

<sup>7</sup>In [50] the authors suggest using self-organizing maps [11].

we postulate a density of the same form on vectors in the product space, which can be written as:

$$p(\mathbf{x}, \mathbf{y} | \mathbf{z}, c) = p(\mathbf{x} | \mathbf{z}, c) p(\mathbf{y} | \mathbf{z}, c), \quad (26)$$

$$p(\mathbf{x} | \mathbf{z}, c) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_c^{\mathbf{x}} + \boldsymbol{\Lambda}_c^{\mathbf{x}}(\mathbf{z} - \boldsymbol{\kappa}_c), \boldsymbol{\Psi}_c^{\mathbf{x}}), \quad (27)$$

$$p(\mathbf{y} | \mathbf{z}, c) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_c^{\mathbf{y}} + \boldsymbol{\Lambda}_c^{\mathbf{y}}(\mathbf{z} - \boldsymbol{\kappa}_c), \boldsymbol{\Psi}_c^{\mathbf{y}}). \quad (28)$$

Note that the data generating model is a MFA model in the product space. The objective function  $\Phi$  defined in (20), is modified to take into account only the observed data.

To facilitate notation, we introduce the binary variables  $o_n^{\mathbf{x}}$  and  $o_n^{\mathbf{y}}$ . If and only if  $o_n^{\mathbf{x}} = 1$  the first  $D_{\mathbf{x}}$  coordinates are observed for the  $n$ -th data point, and similarly for  $o_n^{\mathbf{y}}$ . Thus, only for correspondences we have  $o_n^{\mathbf{x}} = o_n^{\mathbf{y}} = 1$ . We can now write the objective function as:

$$\begin{aligned} \Phi = & \sum_{n=1}^N \sum_{c=1}^C q_n(c) \int q_n(\mathbf{z}) [\log p(\mathbf{z}, c) + o_n^{\mathbf{x}} \log p(\mathbf{x}_n | \mathbf{z}, c) \\ & + o_n^{\mathbf{y}} \log p(\mathbf{y}_n | \mathbf{z}, c)] d\mathbf{z} + \sum_{n=1}^N \mathcal{H}(q_n(c)) + \mathcal{H}(q_n(\mathbf{z})). \end{aligned}$$

In terms of the parameters of  $p$ ,  $q_n(c)$ , and  $q_n(\mathbf{z})$  the objective reads (ignoring constants):

$$\Phi = \sum_{n=1}^N \sum_{c=1}^C q_{nc} [\mathcal{S}_{nc} - \mathcal{E}_{nc}], \quad (29)$$

$$\mathcal{S}_{nc} = \frac{1}{2} \log |\boldsymbol{\Sigma}_n| - \log q_{nc}, \quad (30)$$

$$\begin{aligned} \mathcal{E}_{nc} = & -\log p(c) + \frac{1}{2} \log |\boldsymbol{\Sigma}_c| + \frac{1}{2} \text{Tr}\{\boldsymbol{\Sigma}_c^{-1}(\boldsymbol{\Sigma}_n + \mathbf{z}_{nc} \mathbf{z}_{nc}^{\top})\} \\ & + \frac{o_n^{\mathbf{x}}}{2} (\mathbf{x}_{nc} - \boldsymbol{\Lambda}_c^{\mathbf{x}} \mathbf{z}_{nc})^{\top} \boldsymbol{\Psi}_c^{\mathbf{x}-1} (\mathbf{x}_{nc} - \boldsymbol{\Lambda}_c^{\mathbf{x}} \mathbf{z}_{nc}) \\ & + \frac{o_n^{\mathbf{y}}}{2} (\mathbf{y}_{nc} - \boldsymbol{\Lambda}_c^{\mathbf{y}} \mathbf{z}_{nc})^{\top} \boldsymbol{\Psi}_c^{\mathbf{y}-1} (\mathbf{y}_{nc} - \boldsymbol{\Lambda}_c^{\mathbf{y}} \mathbf{z}_{nc}) \\ & + \frac{o_n^{\mathbf{x}}}{2} \left[ \text{Tr}\{\boldsymbol{\Sigma}_n \boldsymbol{\Lambda}_c^{\mathbf{x}\top} \boldsymbol{\Psi}_c^{\mathbf{x}-1} \boldsymbol{\Lambda}_c^{\mathbf{x}}\} + |\boldsymbol{\Psi}_c^{\mathbf{x}}| \right] \\ & + \frac{o_n^{\mathbf{y}}}{2} \left[ \text{Tr}\{\boldsymbol{\Sigma}_n \boldsymbol{\Lambda}_c^{\mathbf{y}\top} \boldsymbol{\Psi}_c^{\mathbf{y}-1} \boldsymbol{\Lambda}_c^{\mathbf{y}}\} + |\boldsymbol{\Psi}_c^{\mathbf{y}}| \right]. \end{aligned} \quad (31)$$

For  $\boldsymbol{\Sigma}_n$  and  $\mathbf{z}_n$  the updates in the M-step become:

$$\begin{aligned} \boldsymbol{\Sigma}_n^{-1} \leftarrow & \sum_{c=1}^C q_{nc} \left[ \boldsymbol{\Sigma}_c^{-1} + o_n^{\mathbf{x}} \boldsymbol{\Lambda}_c^{\mathbf{x}\top} \boldsymbol{\Psi}_c^{\mathbf{x}-1} \boldsymbol{\Lambda}_c^{\mathbf{x}} + o_n^{\mathbf{y}} \boldsymbol{\Lambda}_c^{\mathbf{y}\top} \boldsymbol{\Psi}_c^{\mathbf{y}-1} \boldsymbol{\Lambda}_c^{\mathbf{y}} \right], \\ \mathbf{z}_n \leftarrow & \boldsymbol{\Sigma}_n \sum_{c=1}^C q_{nc} \left[ \boldsymbol{\Sigma}_c^{-1} \boldsymbol{\kappa}_c + o_n^{\mathbf{x}} \boldsymbol{\Lambda}_c^{\mathbf{x}\top} \boldsymbol{\Psi}_c^{\mathbf{x}-1} (\mathbf{x}_{nc} + \boldsymbol{\Lambda}_c^{\mathbf{x}} \boldsymbol{\kappa}_c) \right. \\ & \left. + o_n^{\mathbf{y}} \boldsymbol{\Lambda}_c^{\mathbf{y}\top} \boldsymbol{\Psi}_c^{\mathbf{y}-1} (\mathbf{y}_{nc} + \boldsymbol{\Lambda}_c^{\mathbf{y}} \boldsymbol{\kappa}_c) \right]. \end{aligned}$$

The updates for  $q_{nc}$ ,  $p(c)$ ,  $\boldsymbol{\kappa}_c$ , and  $\boldsymbol{\Sigma}_c$  remain as before, using  $\tilde{q}_{nc} = q_{nc} / \sum_{n'=1}^N q_{n'c}$ , they are:

$$\begin{aligned} q_{nc} \leftarrow & \frac{e^{-\mathcal{E}_{nc}}}{\sum_{c'=1}^C e^{-\mathcal{E}_{nc'}}}, \quad p(c) \leftarrow \frac{1}{N} \sum_{n=1}^N q_{nc}, \\ \boldsymbol{\kappa}_c \leftarrow & \sum_{n=1}^N \tilde{q}_{nc} \mathbf{z}_n, \quad \boldsymbol{\Sigma}_c \leftarrow \sum_n \tilde{q}_{nc} [\boldsymbol{\Sigma}_n + \mathbf{z}_{nc} \mathbf{z}_{nc}^{\top}]. \end{aligned}$$

With  $q_{nc}^{\mathbf{x}} = o_n^{\mathbf{x}} q_{nc} / \sum_{n'=1}^N o_{n'}^{\mathbf{x}} q_{n'c}$ ,  $\bar{\mathbf{x}}_c^{\mathbf{x}} = \sum_{n=1}^N q_{nc}^{\mathbf{x}} \mathbf{x}_{nc}$ , and  $\bar{\mathbf{z}}_c^{\mathbf{x}} = \sum_{n=1}^N q_{nc}^{\mathbf{x}} \mathbf{z}_{nc}$ , the remaining updates are:

$$\begin{aligned} \boldsymbol{\Lambda}_c^{\mathbf{x}} \leftarrow & \left( \sum_{n=1}^N q_{nc}^{\mathbf{x}} \mathbf{x}_{nc} \mathbf{z}_{nc}^{\top} - \bar{\mathbf{x}}_c^{\mathbf{x}} \bar{\mathbf{z}}_c^{\mathbf{x}\top} \right) \\ & \times \left( \sum_{n=1}^N q_{nc}^{\mathbf{x}} [\boldsymbol{\Sigma}_n + \mathbf{z}_{nc} \mathbf{z}_{nc}^{\top}] - \bar{\mathbf{z}}_c^{\mathbf{x}} \bar{\mathbf{z}}_c^{\mathbf{x}\top} \right)^{-1}, \\ \boldsymbol{\mu}_c^{\mathbf{x}} \leftarrow & \sum_{n=1}^N q_{nc}^{\mathbf{x}} \mathbf{x}_n - \boldsymbol{\Lambda}_c^{\mathbf{x}} \bar{\mathbf{z}}_c^{\mathbf{x}}, \\ [\boldsymbol{\Psi}_c^{\mathbf{x}}]_{ii} \leftarrow & \sum_{n=1}^N q_{nc}^{\mathbf{x}} \left( [\mathbf{x}_{nc} - \boldsymbol{\Lambda}_c^{\mathbf{x}} \mathbf{z}_{nc}]_i^2 + [\boldsymbol{\Lambda}_c^{\mathbf{x}} \boldsymbol{\Sigma}_n \boldsymbol{\Lambda}_c^{\mathbf{x}\top}]_{ii} \right). \end{aligned}$$

The updates for  $\boldsymbol{\mu}_c^{\mathbf{y}}$ ,  $\boldsymbol{\Lambda}_c^{\mathbf{y}}$  and  $\boldsymbol{\Psi}_c^{\mathbf{y}}$  are analogous. Note that the above is straightforwardly generalized to settings with more than two sets of points.

## B. Parameter initialization

As before, proper initialization of the  $\mathbf{z}_n$  is crucial in order to find good parameter estimates with our EM-like algorithm. In the current setting we cannot directly apply the LLE based parameter initialization used before because the missing values do not allow us to compute nearest neighbors and reconstruction weights in the product space.

In [31] an approach was proposed to extend the LLE algorithm to the current setting. The idea is to perform LLE on both sets of points simultaneously, and to constrain the latent coordinates of corresponding points to be identical. The error function is defined as the sum of the LLE error function for both sets of points. The error function is quadratic and can be minimized, subject to the correspondence constraints, by computing the  $(d+1)$  eigenvectors with smallest eigenvalues of a sparse matrix.

## IV. EXPERIMENTAL RESULTS

In this section we present experimental results obtained with the approaches described in the previous sections. In Section IV-A we present qualitative results of applying the CFA algorithm to data sampled from a single manifold. We compare CFA performance with that of several alternative approaches in Section IV-B, and in Section IV-C we compare our closed-form CFA algorithm with the fixed-point algorithm of Roweis et al. In Section IV-D we use the CFA model to predict correspondences between two manifolds and compare its performance with that obtained using an MFA model.

### A. Manifold modelling with CFA: qualitative results

The CFA model can be used in two directions. First, the model can be used for dimension reduction, i.e. mapping a point  $\mathbf{x}$  in the data space to a point  $\mathbf{z}$  in the latent space. Second, the model can be used for data reconstruction, i.e. mapping a point  $\mathbf{z}$  in the latent space to a point  $\mathbf{x}$  in the



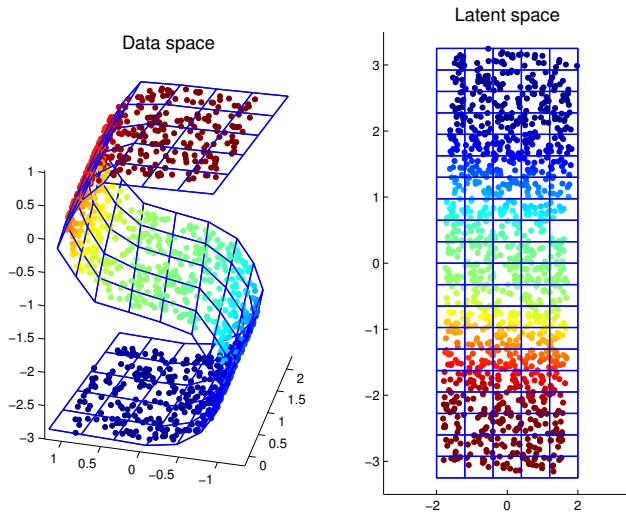


Fig. 4. Original data in  $\mathbb{R}^3$  and a grid mapped from the latent space to the data space with the model (left). Original grid and data mapped to latent space (right).

data space. The mappings can be performed by evaluating  $E_{p(\mathbf{z}|\mathbf{x})}[\mathbf{z}]$  and  $E_{p(\mathbf{x}|\mathbf{z})}[\mathbf{x}]$ .<sup>8</sup>

In Fig. 4 we illustrate the two directions in which we can use the model using the ‘S’ data set of Fig. 2. We initialized the latent coordinates with the LLE algorithm using 10 nearest neighbors and we used a mixture of  $C = 10$  components (below we consider the effect of  $C$  on the results). In the left panel of Fig. 4 the original three dimensional data is depicted on which the model was trained. The right panel depicts the two-dimensional latent representation of the original data and a rectangular grid in the same space. The left panel also depicts the grid when mapped to the data space with the trained model. The illustration shows that the non-linear subspace containing the data was indeed recovered from the training data and that accurate non-linear mapping between the data space and latent space is possible using  $C = 10$  components.

Since for the ‘S’ data set of Fig. 2 the true latent coordinates are available, we can measure how well they are recovered by CFA and LLE respectively. We applied the LLE algorithm, and subsequently the CFA algorithm, on a random subset that contains 80% of the available 1240 data points. We then used the CFA model to map the remaining 20% to the latent space. For LLE, the remaining points were mapped to the latent space by the non-parametric method described in [39]. Since both CFA and LLE find the latent representation only upto a linear transformation, we first optimally linearly transform the recovered latent coordinates before we measure their squared distance to the known true latent coordinates. We call the sum of these squared distances the ‘embedding error’. For LLE we used  $k = 6, \dots, 20$  neighbors, and we used the LLE results to learn CFA models with  $C = 6, \dots, 20$  mixture components. For all combinations of  $k$  and  $C$  we found that the average (over ten random selections of test and train data) difference

<sup>8</sup>It is an interesting direction for further research to consider alternatives to average over the local projections. For example, it might be possible to use notions like the Karger mean [52], or methods used in [27] for interpolation on non-linear manifolds.

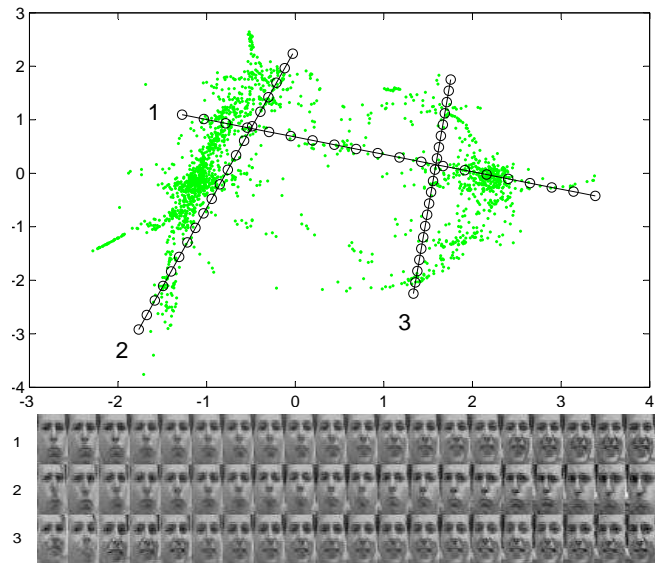


Fig. 5. Recovered latent representation of the images (top). Images generated with model along the line segments shown in the top panel (bottom).

in embedding error for LLE and CFA is very small: about two orders of magnitude smaller than the standard deviations in the embedding errors.

We repeated this experiment using the data shown in Fig. 1, here the two dimensional manifold (the true latent coordinates are given by the horizontal and vertical position of the square, which both range 1 to 20) is non-linearly embedded in a 841 dimensional space. Here we used for LLE  $k = 5, 10, \dots, 45, 50$  neighbors and  $C = 10, 20, \dots, 60, 70$  mixture components for CFA. Again, except for a few combinations, the average difference in the embedding error is smaller than the standard deviation in the embedding errors (although the differences and standard deviations are of the same order of magnitude for this data set). Both LLE and CFA make, on average, an error of about one pixel in mapping the image to the position of the square.

Next, we show results obtained when applying the model to a data set of facial images, which vary in the pose and expression of the face. This data set consists of 1965 images of  $20 \times 28 = 560$  pixels each. The data set was also used in [2] where the authors showed similar results as shown here (although the results here are obtained with the faster closed-form solution algorithm).<sup>9</sup> For initialization we used the LLE algorithm with 14 nearest neighbors.

The recovered latent representation of the data is shown in Fig. 5, where each dot represents the (expected) coordinates  $\mathbf{z}$  of an image (we use the model to map the high dimensional images to the two dimensional latent space). We used the trained model also to generate images  $\mathbf{x}$  from latent coordinates  $\mathbf{z}$  along three straight trajectories in the latent space (the model is used to map two dimensional coordinates along the trajectories to 560 dimensional images). The trajectories are plotted in the figure and the corresponding generated images are shown in the bottom panel of Fig. 5. The latent coordinates form two clusters; the first trajectory passes through both

<sup>9</sup>The data is available at <http://www.cs.toronto.edu/~roweis>.

clusters and the other two trajectories stay within a single cluster. From the generated images along the trajectories we can see that the clusters roughly correspond to images of smiling and non-smiling faces. The reconstructed images along the second and third line segment show that the variation within the clusters corresponds to pose variation of the faces.

### B. Quantitative comparisons

We compared CFA with a probabilistic variant of self-organizing maps called self-organizing mixture models (SOMM) [12], and generative topographic mapping (GTM) [13]. As CFA, these are probabilistic methods that allow for a mapping between data space and latent space in both directions. Therefore it is possible to make a direct comparison between the generalization performance in terms of data log-likelihood and reconstruction errors. To measure the generalization performance of these methods we use a separate set of test data not included when estimating the model parameters. The reconstruction  $\hat{\mathbf{x}}_n$  of a point  $\mathbf{x}_n$  using the model is obtained by mapping  $\mathbf{x}_n$  to a single point  $\mathbf{z}_n$  in the latent space and then mapping  $\mathbf{z}_n$  back to the data space<sup>10</sup>. The reconstruction error is defined as:

$$E_{rec} = \sum_n \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2.$$

All models used a two dimensional latent space; for GTM and SOMM the nodes were placed at a square rectangular grid in the latent space such that the node locations had zero mean and identity covariance. For GTM and SOMM we use Gaussian component densities with isotropic covariance matrix. For GTM we used as many basis functions as mixture components. The basis functions were of the form  $\phi_s(\mathbf{z}) = \exp(-\|\mathbf{z} - \boldsymbol{\kappa}_s\|^2/2\sigma^2)$ , we used  $\sigma^2 = 1/10$  which yielded the best results of several values that were tried. For CFA we constrained the noise variance matrices  $\boldsymbol{\Psi}_c$  to be of the form  $\boldsymbol{\Psi}_c = \sigma_c^2 \mathbf{I}$  so that also here the variance outside the subspaces is isotropic. Note that the CFA model has  $(d+1)$  times more parameters per mixture component than the SOMM and GTM models to encode  $\Lambda_c$ . To obtain a fair comparison we compare models with equal number of parameters. For SOMM and GTM we used  $C \in \{2^2, 4^2, 6^2, \dots, 16^2\}$ , for CFA the number of components was taken to be  $C' = \lfloor C/(d+1) \rfloor$ .

The comparison is based on two data sets: the synthetic ‘S’ data set used before (using a training and test set of 600 points) and a data set of 2000 images of  $40 \times 40$  pixels each of a face looking in different directions (using 1500 training images and 500 test images).<sup>11</sup> The reported results are averages and standard deviations over 10 randomly drawn train and test sets. Results obtained using the synthetic data set are summarized in Fig. 6, and those obtained using the image data set in Fig. 7. In Fig. 8 we also plotted some of the reconstructions obtained for the second data set using the different models (using  $C = 256$  components for SOMM and GTM and  $C' = 85$  components for CFA). As expected

<sup>10</sup>For GTM we map latent coordinates to the data space using the generalized linear model. For SOMM we assign the high dimensional coordinates of the mixture component with nearest mean in the latent space.

<sup>11</sup>The data is available at <http://lear.inrialpes.fr/~verbeek>.

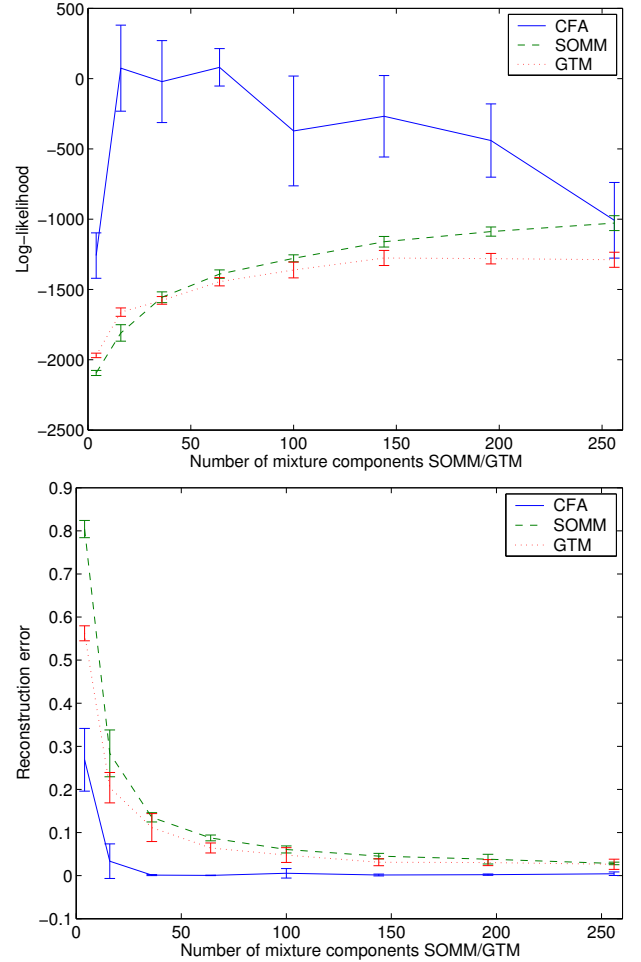


Fig. 6. Comparison of CFA, GTM, and SOMM using synthetic data.

from the results in Fig. 7, the reconstructions of the different methods appear quite similar.

The results lead to the following conclusions. The CFA model is able to achieve significantly higher log-likelihood on the test data than the SOMM and GTM models, especially when a moderate model capacity is used (small number of mixture components/parameters). Compared to log-likelihood, the differences in reconstruction errors are relatively small. The discrepancy between the reconstruction and log-likelihood scores can be understood as follows. SOMM and GTM are based on mixtures where the components are Gaussians with isotropic covariance. The isotropic noise makes that SOMM and GTM have to ‘leak’ a considerable amount of mass to parts of the data space that are not filled by the low dimensional manifold, even though they use about three times more mixture components. The CFA model on the other hand can locally concentrate its mass in the linear subspace approximating the manifold. The reconstruction error measures the squared Euclidean distance between a data point and its reconstruction. Note that the reconstruction is the mean of the distribution on  $\mathbf{x}$  given its low dimensional representation  $\mathbf{z}$ ; the reconstruction error does not take into account the uncertainty in this distribution.

The standard deviation in the results obtained with the CFA

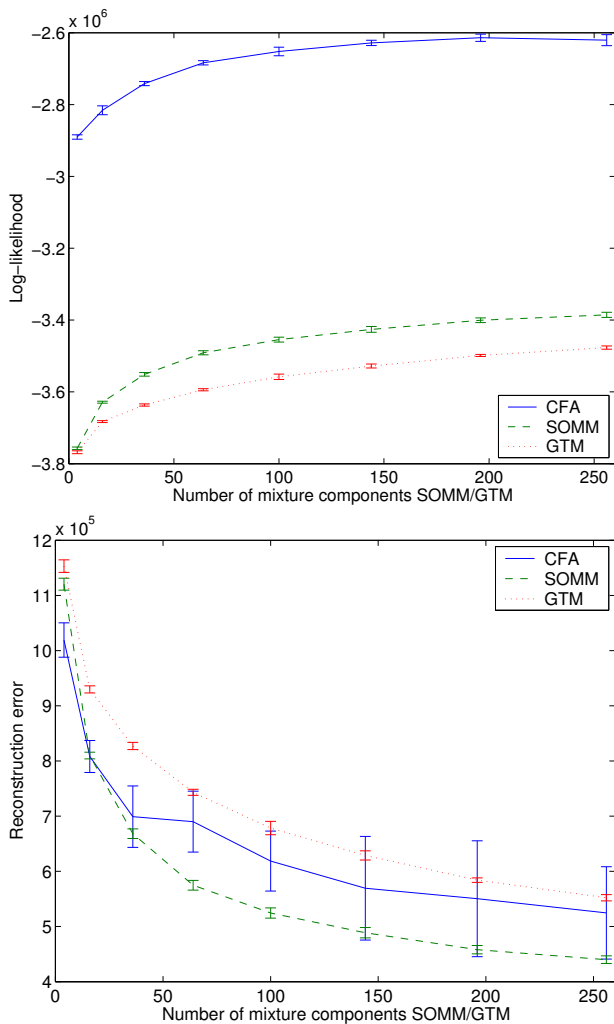


Fig. 7. Comparison of CFA, GTM, and SOMM, results for image data.

model are generally larger or comparable to those obtained for SOMM and GTM. This is caused by the sensitivity to the initialization of the CFA model. This sensitivity might be resolved by using more robust initialization techniques, e.g. it is possible to minimize the sum of LLE error function using several numbers of neighbors or to use bootstrap-like procedures as in [31].

Using many mixture components, CFA shows overfitting on both data sets in terms of log-likelihood. This effect is well known [53] and is caused by the fact that using many components more parameters need to be estimated. However, since the amount of available data is limited, the parameters cannot be accurately estimated and as a result the model is erroneous. The overfitting effect is not observed for SOMM, and for GTM only on the synthetic data set. This is due to the fact that the latter two models use an isotropic noise model.

### C. Comparison of closed-form and fixed-point algorithm

Next, we compare the EM algorithm using our closed-form M-step and using the fixed-point M-step proposed in [2]. The complexity of both M-steps is  $O(Dd^2)$ . Thus, the factor of speed-up depends on the number of iterations needed by the

TABLE I

COMPARISON BETWEEN CLOSED-FORM AND FIXED-POINT ALGORITHMS.

Results for the 'S' data set, $d = 2, C = 10$ .			
	# EM steps	time	$\Phi$
closed-form	$123.5 \pm 49.5$	$6.6 \pm 3.9$	$0.9 \pm 0.6$
fixed-point	$125.5 \pm 33.8$	$77.8 \pm 29.3$	$0.2 \pm 0.8$
Results for the isomap face data set, $d = 3, C = 50$ .			
	# EM steps	time	$\Phi$
closed-form	$78.4 \pm 7.0$	$30.0 \pm 2.8$	$9.7 \pm 0.7$
fixed-point	$76.7 \pm 9.1$	$198.0 \pm 32.7$	$8.1 \pm 1.1$

iterative approach to reach a fixed-point. We measured the performance when applying the algorithms to two data sets. The first one is the 'S' data set of Fig. 2. The second data set is a set of 698 computer generated images of  $64 \times 64$  pixel each of a face seen from different directions and under different lighting conditions.<sup>12</sup> To speed-up experimentation, the images were first projected on the 30-dimensional principal components subspace which contains over 90% of the total data variance.

To start the fixed-point iterations we initialized the parameters at their value found in the previous M-step. We considered the fixed-point algorithm to have converged if the maximum relative change in the parameters dropped below  $10^{-4}$ . If the fixed-point algorithm did not converge in 1000 iterations, it was terminated and the parameter values of the previous M-step were retained. To quantify the speed-up of the closed-form equations in practice, and to see whether they lead to better final parameter estimates, we measured: (i) runtime of the algorithm<sup>13</sup>, (ii) value of  $\mathcal{L}'$  attained after convergence and (iii) number of EM steps required to converge.

The results are summarized in Tab. I by averages and standard deviations from ten runs. The results show that the closed-form algorithm is significantly faster and yields higher values for the objective function in a comparable number of EM iterations. We conclude that the closed-form update equations are to be preferred over the fixed-point algorithm.

### D. Mapping between multiple embeddings of a manifold

In this section we experimentally compare our CFA approach against a MFA to predict high dimensional correspondences. We show how the performance of the approaches depends on the amount of correspondences and the number of mixture components. First, we compare prediction accuracy on two synthetic data sets to gain insight in the differences between the two approaches. Then we compare them using a data set consisting of gray scale images of two objects to determine if the differences observed with the synthetic data sets are also apparent in more realistic data sets.

In the first experiment we used the data set shown in Fig. 3. Both sets contain 1240 points in  $\mathbb{R}^3$ , in each experiment we used a random subset of 600 points from each set to fit the models and a second set of 600 corresponding points to assess the quality of the fitted models. The results reported below are averages over 20 experiments. We trained CFA and MFA

<sup>12</sup>This data set is available online at <http://isomap.stanford.edu>.

<sup>13</sup>To measure convergence of the EM algorithm we checked whether the relative change in  $\mathcal{L}'$  was smaller than  $10^{-4}$ .



Fig. 8. From top to bottom: original images, and reconstructions using CFA, SOMM, and GTM.

models using latent dimensionality of  $d = 2$ , while varying the percentage of training data for which correspondences are available (ranging from 5% up to 55% with 10% intervals, i.e. using from 30 up to 330 correspondences) and the number  $C$  of mixture components (ranging from  $C = 2$  to  $C = 8$ ). To assess the quality of the fitted models we used the test data, for which all 6 coordinates are known. For the test data we predicted the last three coordinates given the first three coordinates with the trained models, and vice versa.

When using a MFA, each mixture component  $c$  is a Gaussian density on the product space which induces the conditional Gaussian  $p(\mathbf{y}|\mathbf{x}, c)$ . Combining the different components we obtain the conditional (mixture) density  $p(\mathbf{y}|\mathbf{x}) = \sum_{c=1}^C p(\mathbf{y}|\mathbf{x}, c)p(c|\mathbf{x})$ . Hence, the expected value of  $\mathbf{y}$  under this distribution, which we denote by  $\hat{\mathbf{y}}$ , is the average of the means of  $p(\mathbf{y}|\mathbf{x}, c)$  where each mean is weighted by the corresponding posterior probability  $p(c|\mathbf{x})$ . The same holds for the CFA models and the predictions on  $\mathbf{x}$  given  $\mathbf{y}$ .

As an error measure of the models we used the squared difference between the predicted coordinates  $\hat{\mathbf{x}}$  (or  $\hat{\mathbf{y}}$ ) and the true value of the coordinates of  $\mathbf{x}$  (or  $\mathbf{y}$ ), averaged over all test data points and over the coordinates. Thus for the  $N = 600$  test points we measure:

$$E_{rec} = \frac{1}{ND_{\mathbf{x}}} \sum_{n=1}^N \|\hat{\mathbf{x}}_n - \mathbf{x}_n\|^2 + \frac{1}{ND_{\mathbf{y}}} \sum_{n=1}^N \|\hat{\mathbf{y}}_n - \mathbf{y}_n\|^2.$$

In Tab. II we tabulated, for different amounts of correspondences and mixture components, the error  $E_{rec}$  obtained with CFA and MFA models. Results where the CFA error is significantly<sup>14</sup> smaller are printed bold.

The results show that if many correspondences are available CFA and MFA give comparable results and if only a few correspondences are given CFA performs significantly better than MFA. This difference is more pronounced as the number of mixture components becomes larger. The explanation of the effect is that if the average number of correspondences per mixture component becomes small, then in a MFA model there can be mixture components which have (almost) no responsibility for correspondences (i.e. observations without

TABLE II

RECONSTRUCTION ERRORS FOR SYNTHETIC DATA SET. FOR EACH PERCENTAGE CFA RESULTS ARE PRINTED ABOVE MFA RESULTS.

$C$	2	3	4	5	6	7	8
5%	<b>0.27</b> <b>0.74</b>	<b>0.34</b> <b>0.98</b>	<b>0.34</b> <b>1.26</b>	<b>0.22</b> <b>1.13</b>	<b>0.25</b> <b>1.45</b>	<b>0.25</b> <b>1.40</b>	<b>0.23</b> <b>1.28</b>
15%	<b>0.18</b> <b>0.34</b>	0.21 0.33	<b>0.14</b> <b>0.57</b>	<b>0.15</b> <b>0.51</b>	<b>0.13</b> <b>0.71</b>	<b>0.09</b> <b>0.67</b>	<b>0.09</b> <b>0.67</b>
25%	0.16 0.18	<b>0.13</b> <b>0.19</b>	<b>0.12</b> <b>0.25</b>	<b>0.08</b> <b>0.30</b>	<b>0.08</b> <b>0.37</b>	<b>0.09</b> <b>0.40</b>	<b>0.06</b> <b>0.42</b>
35%	<b>0.16</b> <b>0.18</b>	0.13 0.18	<b>0.12</b> <b>0.16</b>	<b>0.09</b> <b>0.20</b>	<b>0.08</b> <b>0.23</b>	<b>0.09</b> <b>0.23</b>	<b>0.06</b> <b>0.26</b>
45%	0.16 0.17	<b>0.13</b> <b>0.15</b>	<b>0.11</b> <b>0.17</b>	<b>0.10</b> <b>0.18</b>	<b>0.07</b> <b>0.18</b>	<b>0.06</b> <b>0.20</b>	<b>0.05</b> <b>0.19</b>
55%	0.16 0.16	<b>0.13</b> <b>0.15</b>	<b>0.10</b> <b>0.14</b>	<b>0.08</b> <b>0.16</b>	<b>0.06</b> <b>0.16</b>	<b>0.05</b> <b>0.14</b>	<b>0.05</b> <b>0.13</b>

missing values). In such cases it is not possible to determine the dependencies between coordinates in the two different spaces, since for each data point with non-negligible responsibility for the component either only the first  $D_1$  coordinates are observed or only the last  $D_2$  coordinates. In comparison, the CFA model can exploit the incomplete observations to determine these dependencies; this is possible through the dependencies between the incomplete observations and the estimated global low dimensional coordinates on the manifold.

In the second experiment we also compare the reconstruction errors using MFA and CFA models, but we use a more realistic data set of much higher dimensionality. The data consists of 2500 gray-scale images of  $64 \times 64$  pixels of each of two toy puppets as viewed from different directions. In Fig. 9 some corresponding views of the two puppets are depicted. The images, originally used in [54], were provided by G. Peters who recorded them at the Institute for Neural Computation of the Ruhr-Universität-Bochum, Germany. Images of the objects were recorded while moving the camera over the *viewing hemisphere*, see the left panel of Fig. 10. The viewing hemisphere was sampled at  $3.6^\circ$  intervals in longitude (yielding 100 steps to complete the circle) and at  $3.6^\circ$  intervals in latitude (yielding 25 steps to go from the equator of the hemisphere to the pole). Note that the images recorded near the top of the hemisphere differ considerable, since these are

<sup>14</sup>To determine significance we used a  $t$ -test with 19 DOF and  $p = 0.05$ .



Fig. 9. Examples of the images of the two toy puppets, the top row contains images of the cat figure and the bottom row images of the dwarf figure. Corresponding views are displayed above each other.

different rotations of the top view of the object.

There are only two degrees of freedom in each set of images since the images are determined by position of the camera, which is in turn determined by its longitude and latitude on the hemisphere. Since the longitude is a periodic degree of freedom, the images can be embedded on the surface of a cylinder in a Euclidean space. In principle the images can also be embedded, while preserving nearest neighbor relations, in a two dimensional space by embedding images with equal latitude on concentric circles with a radius that is monotonically increasing with the latitude. However, such a two dimensional embedding is not returned by the LLE based algorithm since it is not directly aiming at the preservation of nearest neighbor relations. Three dimensions are required to recover a cylinder-like embedding of the images in which only similar images are embedded nearby. Therefore, we used a three dimensional latent space for the CFA models, and to obtain a fair comparison we also use  $d = 3$  for the MFA models. In Fig. 10 we illustrated the three dimensional latent representation recovered by the LLE based initialization method when using all 2500 images of each set and 125 correspondences. The recovered coordinates indeed trace-out a cylinder-like shape.

In order to speed-up the experimentation, the images were first projected to the 100 dimensional linear PCA subspace of the  $64 \times 64 = 4096$  dimensional space spanned by the pixel values. Over 90% of the original variance in the data was contained in this subspace with approximately 40 times fewer dimensions. Since the discarded dimensions contain only a small fraction of the data variance, projecting the original data to the PCA subspace is expected to have little effect on the obtained results, below we also show results obtained when using 1000 dimensional PCA projection. We trained CFA and MFA models with a number of mixture components ranging from  $C = 10$  to  $C = 60$  with steps of 10. The amount of images with a correspondence was set to 2%, 5%, 10%, 20% and 50%. Of each object 2000 images were used for training and 500 to assess the reconstruction error  $E_{rec}$  from (32).

The obtained errors, averaged over six random selections of training and test data, are tabulated in Tab. III. Statistically significant differences are printed bold (using a  $t$ -test with 5 degrees of freedom and  $p = 0.05$ ). In Fig. 11 we plotted averages of the error of CFA models and MFA models using  $C = 40$  components against the percentage of correspondences that was used. From this figure it can be clearly observed that already with relatively few correspondences CFA

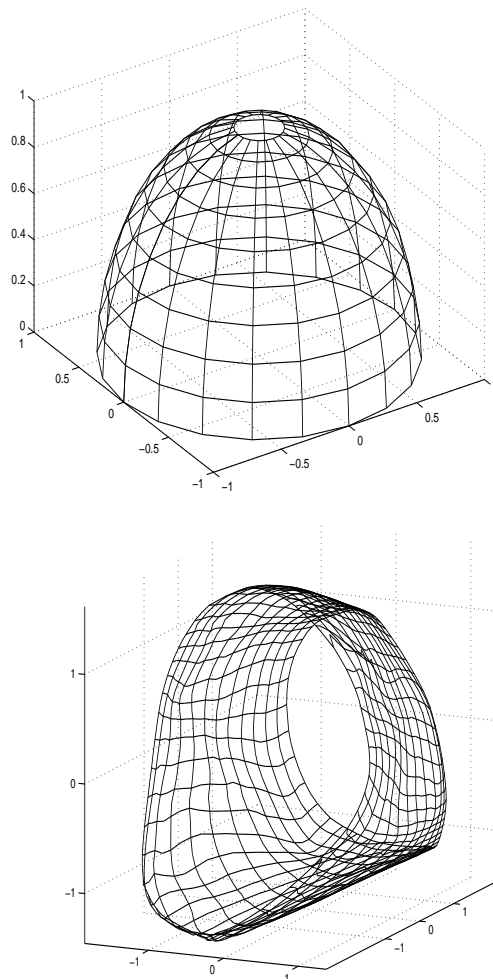


Fig. 10. The viewing hemisphere: the object was fixed in the center of the sphere and images were recorded when the camera was placed a different locations on the hemisphere and directed toward the object (left panel). Three dimensional embedding of 625 of the 2500 images (right panel).

models make accurate predictions, MFA models require many more correspondences to obtain similar errors.

To compare MFA and CFA qualitatively, we plotted the true and predicted correspondence for some of the test examples in Fig. 12. For reference, we also included reconstructions of the true correspondences obtained with a three dimensional linear

TABLE III

RECONSTRUCTION ERRORS FOR THE IMAGE DATA SET. FOR EACH PERCENTAGE, CFA RESULTS ARE PRINTED ABOVE MFA RESULTS.

$C$	1	10	20	30	40	50	60
2%	<b>4.78</b>	<b>2.31</b>	<b>1.44</b>	<b>1.15</b>	<b>1.51</b>	<b>2.32</b>	<b>1.49</b>
	<b>6.47</b>	<b>8.16</b>	<b>7.21</b>	<b>7.14</b>	<b>7.06</b>	<b>7.94</b>	<b>7.36</b>
5%	<b>4.64</b>	<b>1.99</b>	<b>1.18</b>	<b>0.96</b>	<b>0.74</b>	<b>0.73</b>	<b>0.55</b>
	<b>6.47</b>	<b>5.25</b>	<b>4.89</b>	<b>4.71</b>	<b>4.20</b>	<b>4.30</b>	<b>4.55</b>
10%	<b>4.21</b>	<b>1.78</b>	<b>1.08</b>	<b>0.84</b>	<b>0.68</b>	<b>0.54</b>	<b>0.59</b>
	<b>6.42</b>	<b>3.74</b>	<b>4.29</b>	<b>3.35</b>	<b>2.71</b>	<b>3.27</b>	<b>2.95</b>
20%	<b>4.04</b>	<b>1.54</b>	<b>0.93</b>	<b>0.73</b>	<b>0.66</b>	<b>0.50</b>	<b>0.46</b>
	<b>6.42</b>	<b>2.51</b>	<b>2.55</b>	<b>2.40</b>	<b>2.26</b>	<b>2.34</b>	<b>1.78</b>
50%	<b>3.89</b>	<b>1.47</b>	<b>0.84</b>	<b>0.64</b>	<b>0.64</b>	<b>0.44</b>	<b>0.36</b>
	<b>6.32</b>	<b>1.74</b>	<b>1.02</b>	<b>1.07</b>	<b>0.86</b>	<b>0.72</b>	<b>0.89</b>

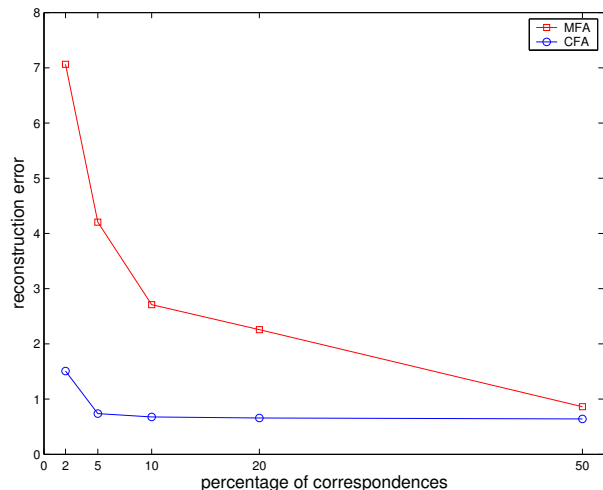


Fig. 11. Average reconstruction errors obtained with CFA and MFA models using 40 components.

PCA projection.<sup>15</sup> The depicted results were obtained by training models with  $C = 65$  mixture components on 2000 images of each object with 100 images in correspondence. Before fitting the models, the images in each set were projected on the first 1000 principal components derived from the images in that set, preserving over 99.6% of the variance in each set. The average errors (per image and per dimension) were 0.127 for CFA and 0.732 for MFA. Clearly, the CFA model—which exploits the global manifold structure of the data—yields superior predictions as compared to those obtained with the MFA model.

## V. CONCLUSIONS

In this paper we considered the Coordinated Factor Analysis (CFA) approach to combine linear models to form a non-linear model. Compared to GTM and SOMM, the advantage of this method is that the low dimensional representation is not restricted to a discrete set of points (the nodes in SOMM and GTM) but offers a continuous low dimensional representation. Furthermore, the conditional density  $p(\mathbf{z}|\mathbf{x})$

<sup>15</sup>The PCA results were obtained by performing PCA on both sets of images separately and reconstructing the *true* correspondence from its three dimensional PCA representation. This is the best possible result that could be obtained using a linear reconstruction of the images from a three dimensional representation.

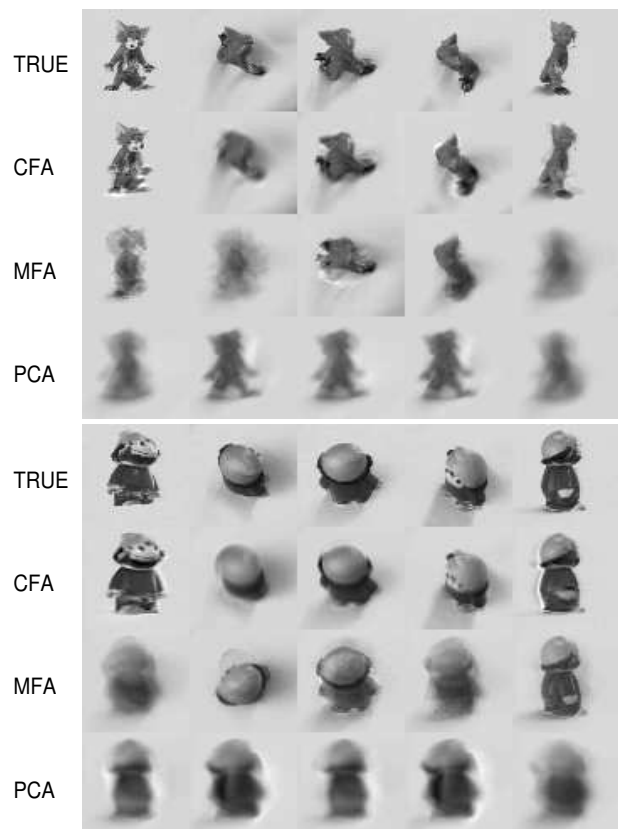


Fig. 12. Top panel shows (predicted) corresponding view of the cat given a view of the dwarf and bottom panel shows (predicted) corresponding view of the dwarf given a view of the cat. In each panel, from top to bottom: true correspondence, reconstruction with CFA, MFA and 3D PCA projection.

on latent coordinates given a data vector (and vice versa  $p(\mathbf{x}|\mathbf{z})$ ) is a Gaussian mixture density whose parameters and expectation are readily computed. We presented an improved parameter estimation scheme, as compared to that proposed in [2], which replaces a fixed-point iteration within the M-step of the EM-like algorithm by a closed-form solution. The experimental results presented in Section IV indicated that the CFA approach compares favorably to SOMM and GTM in the sense that fewer mixture components are needed to obtain similar reconstruction errors (errors in predicting the high dimensional points from their low dimensional representation).

In Section III we generalized the CFA approach to a setting where the goal is to find a model that enables us to predict between two high-dimensional variables. In this setting, the training data consists of two sets of high dimensional points. For some points the corresponding point in the other set is given. We compared the CFA approach with a mixture of factor analyzers approach, and experimentally found that when only a few correspondences are given the CFA approach leads to significantly better prediction accuracy, both for synthetic and natural data. The difference in performance is explained by the fact that CFA successfully uses the manifold structure of the data to determine the dependencies between the high dimensional variables, even if few correspondences are given.

A line of research we want to pursue in the future is the use of non-parametric semi-supervised learning techniques, in

combination with the models presented here for regression problems. In semi-supervised learning for some images the response variables are known and are used to infer a density on the response variables of the other images. The response variables can be either the DOF that span the image manifold or some smooth function on the manifold (e.g. the absolute angular deviation from the frontal pose of a face). In this manner the unsupervised data can be used to bias the learning of the regression function. This is similar to the way in which unsupervised data was used for correspondence learning here.

#### ACKNOWLEDGMENT

The author would like to thank the reviewers for their detailed comments and suggestions. The author is also indebted to Nikos Vlassis and Ben Kröse for many useful discussions.

#### REFERENCES

- [1] A. Leonardis and H. Bischof, "Kernel and subspace methods for computer vision," *Pattern Recognition*, vol. 36, no. 9, pp. 1925–1927, 2003.
- [2] S. T. Roweis, L. K. Saul, and G. E. Hinton, "Global coordination of local linear models," in *Advances in Neural Information Processing Systems*, T. G. Dietterich and S. Becker and Z. Ghahramani, Ed. Cambridge, MA, USA: MIT Press, 2002, vol. 14, pp. 889–896.
- [3] L. Sirovich and M. Kirby, "Low-dimensional procedure for the characterization of human faces," *Journal of Optical Society of America A*, vol. 4, no. 3, pp. 519–524, 1987.
- [4] I. T. Jolliffe, *Principal Component Analysis*. New-York, NY, USA: Spinger-Verlag, 1986.
- [5] M. Turk and A. Pentland, "Eigenfaces for recognition," *Journal of Cognitive Neuroscience*, vol. 3, no. 1, pp. 71–86, 1991.
- [6] H. Murase and S. K. Nayar, "Visual learning and recognition of 3-d objects from appearance," *International Journal of Computer Vision*, vol. 14, pp. 5–24, 1995.
- [7] E. Oja, "Data compression, feature extraction, and autoassociation in feedforward neural networks," in *Proceedings of the International Conference on Artificial Neural Networks*, T. Kohonen, M. Mäkisara, O. Simula, and J. Kangas, Eds., vol. 1. Amsterdam, The Netherlands: North-Holland, 1991, pp. 737–745.
- [8] T. Hastie and W. Stuetzle, "Principal curves," *Journal of the American Statistical Association*, vol. 84, no. 406, pp. 502–516, 1989.
- [9] B. Kégl, A. Krzyzak, T. Linder, and K. Zeger, "Learning and design of principal curves," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 3, pp. 281–297, 2000.
- [10] K. Chang and J. Ghosh, "A unified model for probabilistic principal surfaces," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 1, pp. 22–41, 2001.
- [11] T. Kohonen, *Self-Organizing Maps*. New-York, NY, USA: Spinger-Verlag, 1995.
- [12] J. J. Verbeek, N. Vlassis, and B. J. A. Kröse, "Self-organizing mixture models," *Neurocomputing*, vol. 63, pp. 99–123, 2005.
- [13] C. M. Bishop, M. Svensén, and C. K. I. Williams, "GTM: the generative topographic mapping," *Neural Computation*, vol. 10, pp. 215–234, 1998.
- [14] J. J. Verbeek, "Mixture models for clustering and dimension reduction," Ph.D. dissertation, University of Amsterdam, 2004.
- [15] I. K. Fodor, "A survey of dimension reduction techniques," Lawrence Livermore National Laboratory, Center for Applied Scientific Computing, Tech. Rep. UCRL-ID-148494, 2002.
- [16] M. Á. Carreira-Perpiñán, "A review of dimension reduction techniques," Department of Computer Science, University of Sheffield, Tech. Rep. CS-96-09, 1997.
- [17] J. B. Tenenbaum, V. de Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000.
- [18] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, 2000.
- [19] B. Schölkopf, A. J. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation*, vol. 10, pp. 1299–1319, 1998.
- [20] M. Brand, "Charting a manifold," in *Advances in Neural Information Processing Systems*, S. Becker and S. Thrun and K. Obermayer, Ed. Cambridge, MA, USA: MIT Press, 2003, vol. 15, pp. 961–968.
- [21] X. He and P. Niyogi, "Locality preserving projections," in *Advances in Neural Information Processing Systems*, S. Thrun and L. K. Saul and B. Schölkopf, Ed. Cambridge, MA, USA: MIT Press, 2004, vol. 16, pp. 153–160.
- [22] X. He, S. Yan, Y. Hu, P. Niyogi, and H.-J. Zhang, "Face recognition using Laplacianfaces," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 3, pp. 328–340, 2005.
- [23] M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering," in *Advances in Neural Information Processing Systems*, T. G. Dietterich and S. Becker and Z. Ghahramani, Ed. Cambridge, MA, USA: MIT Press, 2002, vol. 14, pp. 585–591.
- [24] K. Q. Weinberger and L. K. Saul, "Unsupervised learning of image manifolds by semidefinite programming," in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, 2004, pp. 988–995.
- [25] Y. Bengio, J.-F. Paiement, P. Vincent, O. Delalleau, N. Le Roux, and M. Ouimet, "Out-of-sample extensions for LLE, isomap, MDS, eigenmaps, and spectral clustering," in *Advances in Neural Information Processing Systems*, S. Thrun and L. K. Saul and B. Schölkopf, Ed. Cambridge, MA, USA: MIT Press, 2004, vol. 16, pp. 177–184.
- [26] N. Kambhatla and T. K. Leen, "Fast non-linear dimension reduction," in *Advances in Neural Information Processing Systems*, J. D. Cowan and G. Tesauro and J. Alspector, Ed. San Mateo, CA, USA: Morgan Kaufmann, 1994, vol. 6, pp. 152–159.
- [27] C. Bregler and S. M. Omohundro, "Nonlinear image interpolation using manifold learning," in *Advances in Neural Information Processing Systems*, G. Tesauro and D. S. Touretzky and T. K. Leen, Ed. Cambridge, MA, USA: MIT Press, 1995, vol. 7, pp. 973–980.
- [28] M. E. Tipping and C. M. Bishop, "Mixtures of probabilistic principal component analysers," *Neural Computation*, vol. 11, no. 2, pp. 443–482, 1999.
- [29] G. E. Hinton, P. Dayan, and M. Revow, "Modelling the manifolds of images of handwritten digits," *IEEE Transactions on Neural Networks*, vol. 8, no. 1, pp. 65–74, 1997.
- [30] R. Bellman, *Adaptive control processes: a guided tour*. Princeton University Press, 1961.
- [31] J. H. Ham, D. D. Lee, and L. K. Saul, "Learning high dimensional correspondences from low dimensional manifolds," in *ICML 2003, workshop on the continuum from labeled to unlabeled data in machine learning and data mining*, 2003.
- [32] X. Meng and D. van Dyk, "The EM algorithm – an old folk song sung to a fast new tune," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 59, no. 1, pp. 511–567, 1997.
- [33] D. de Ridder and V. Franc, "Robust subspace mixture models using t-distributions," in *Proceedings of the British Machine Vision Conference*, R. Harvey and J. Bangham, Eds., 2003, pp. 319–328.
- [34] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 39, no. 1, pp. 1–38, 1977.
- [35] Z. Ghahramani and G. E. Hinton, "The EM algorithm for mixtures of factor analyzers," University of Toronto, Tech. Rep. CRG-TR-96-1, 1996.
- [36] Y. W. Teh and S. T. Roweis, "Automatic alignment of local representations," in *Advances in Neural Information Processing Systems*, S. Becker and S. Thrun and K. Obermayer, Ed. Cambridge, MA, USA: MIT Press, 2003, vol. 15, pp. 841–848.
- [37] J. J. Verbeek, S. T. Roweis, and N. Vlassis, "Non-linear CCA and PCA by alignment of local models," in *Advances in Neural Information Processing Systems*, S. Thrun and L. K. Saul and B. Schölkopf, Ed. Cambridge, MA, USA: MIT Press, 2004, vol. 16, pp. 297–304.
- [38] J. Wieghardt, "Learning the topology of views: from images to objects," Ph.D. dissertation, Ruhr-Universität-Bochum, Bochum, Germany, 2001.
- [39] L. K. Saul and S. T. Roweis, "Think globally, fit locally: unsupervised learning of low dimensional manifolds," *Journal of Machine Learning Research*, vol. 4, pp. 119–155, 2003.
- [40] S. Chretien and A. O. Hero, "Kullback proximal algorithms for maximum likelihood estimation," *IEEE Transactions on Information Theory*, vol. 46, no. 5, pp. 1800–1810, 2000.
- [41] R. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton, "Adaptive mixtures of local experts," *Neural Computation*, vol. 3, pp. 79–87, 1991.
- [42] J. J. Verbeek, N. Vlassis, and B. J. A. Kröse, "Coordinating principal component analyzers," in *Proceedings of the International Conference*

- on *Artificial Neural Networks*, J. R. Dorronsoro, Ed. New-York, NY, USA: Spinger-Verlag, 2002, vol. 12, pp. 914–919.
- [43] B. Kégl, “Intrinsic dimension estimation using packing numbers,” in *Advances in Neural Information Processing Systems*, S. Becker and S. Thrun and K. Obermayer, Ed. Cambridge, MA, USA: MIT Press, 2003, vol. 15, pp. 681–688.
- [44] J. A. Costa and A. O. Hero, “Geodesic entropic graphs for dimension and entropy estimation in manifold learning,” *IEEE Transactions on Signal Processing*, vol. 52, no. 8, pp. 2210–2221, 2004.
- [45] E. Levina and P. J. Bickel, “Maximum likelihood estimation of intrinsic dimension,” in *Advances in Neural Information Processing Systems*, L. K. Saul and Y. Weiss and L. Bottou, Ed. Cambridge, MA, USA: MIT Press, 2005, vol. 17, pp. 777–784.
- [46] J. J. Oliver, R. A. Baxter, and C. S. Wallace, “Unsupervised learning using MML,” in *Proceedings of the International Conference on Machine Learning*, L. Saitta, Ed. San Mateo, CA, USA: Morgan Kaufmann, 1996, vol. 13, pp. 364–374.
- [47] M. A. T. Figueiredo and A. K. Jain, “Unsupervised learning of finite mixture models,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 3, pp. 381–396, 2002.
- [48] Z. Ghahramani and M. J. Beal, “Variational inference for Bayesian mixtures of factor analysers,” in *Advances in Neural Information Processing Systems*, S. A. Solla and T. K. Leen and K.-R. Müller, Ed. Cambridge, MA, USA: MIT Press, 2000, vol. 12, pp. 449–455.
- [49] J. H. Ham, D. D. Lee, and L. K. Saul, “Semisupervised alignment of manifolds,” in *Proceedings of the Annual Conference on Uncertainty in Artificial Intelligence*, Z. Ghahramani and R. Cowell, Eds., vol. 10, 2005, pp. 120–127.
- [50] H. Ritter, “Parametrized self-organizing maps,” in *Proceedings of the International Conference on Artificial Neural Networks*, S. Gielen and B. Kappen, Eds. New-York, NY, USA: Spinger-Verlag, 1993, vol. 3, pp. 568–577.
- [51] Z. Ghahramani and M. I. Jordan, “Supervised learning from incomplete data via an EM approach,” in *Advances in Neural Information Processing Systems*, J. D. Cowan and G. Tesauro and J. Alspector, Ed. San Mateo, CA, USA: Morgan Kaufmann, 1994, vol. 6, pp. 120–127.
- [52] H. Karger, “Riemann center of mass and mollifier smoothing,” *Communications on Pure and Applied Mathematics*, vol. 3, pp. 509–541, 1977.
- [53] A. R. Webb, *Statistical pattern recognition*. New-York, NY, USA: Wiley, 2002.
- [54] G. Peters, B. Zitova, and C. von der Malsburg, “How to measure the pose robustness of object views,” *Image and Vision Computing*, vol. 20, no. 4, pp. 249–256, 2002.



**Jakob Verbeek** received a cum laude MSc in Artificial Intelligence (1998) and a cum laude MSc in Logic (2000), both at the University of Amsterdam, The Netherlands. In 2004 he received a PhD degree from the same university for his research on clustering and dimension reduction. After a postdoctoral researcher position at the University of Amsterdam he is currently working as a postdoctoral researcher at INRIA Rhone-Alpes. His research interests include latent variable models, their applications in computer vision, and machine learning in general.