

Transformed Component Analysis: Joint Estimation of Spatial Transformations and Image Components

Brendan J. Frey

Department of Computer Science
University of Waterloo

Nebojsa Jojic

Electrical and Computer Engineering
University of Illinois at Urbana-Champaign

Abstract

A simple, effective way to model images is to represent each input pattern by a linear combination of “component” vectors, where the amplitudes of the vectors are modulated to match the input. This approach includes principal component analysis, independent component analysis and factor analysis. In practice, images are subjected to randomly selected transformations of a known nature, such as translation, rotation and scale. Direct application of the above methods will lead to severely blurred components and even to components that only account for the transformations and ignore the more interesting and useful structure. We propose a method called *transformed component analysis*, which incorporates a discrete, hidden variable that accounts for transformations and uses a linear-time expectation maximization algorithm to jointly extract components and normalize for transformations. We illustrate the algorithm using a shading problem, facial expression modeling and handwritten digit recognition.

1 Introduction

Many popular ways of modeling images use a linear combination of vectors to represent each input. Principal components analysis (PCA, a.k.a. “*eigen-whatever*”) is a way of representing a class of images using a small set of component vectors in the vector space of image pixel intensities [1]. An image can be described by a linear combination of these components plus some distortion. The first principal component is the direction in which the projected training data has greatest variance, the second principal component is the direction of greatest variance *after* the first component is subtracted off, and so on.

A different technique that is gaining in popularity is independent component analysis (ICA), which tries to find components such that when the training data is projected on these components, the component activities are independent (not just uncorrelated) [2]. This

method has been used for blind separation and blind deconvolution.

These methods are “feedforward” in the sense that they are trained to map an input image to a hidden (unobserved) representation. In contrast, *generative models* are trained to generate patterns that look similar to the training data. The model specifies a distribution over some hidden variables and a conditional distribution over the input image given the values of the hidden variables. The representation for an input is the posterior distribution over the hidden variables.

Factor analysis (FA) [3] is a generative model that is similar in spirit to PCA. The distribution over a small set of real-valued hidden variables is a zero-mean unit-covariance Gaussian and the distribution over the inputs given the hidden variables is also Gaussian, with a diagonal covariance matrix and a mean given by a linear combination (matrix of components) of the hidden variables. A factor analyzer can be fit to training data using the EM algorithm [4].

The above methods are useful in some applications. However, real data is often subjected to randomly selected transformations of a known nature, such as translation, rotation and scale in images. In these cases, direct application of the above methods will lead to severely blurred components and even to components that only account for the transformations and ignore the more interesting and useful structure.

In this paper, we propose a method called *transformed component analysis*, which incorporates a discrete, hidden variable that accounts for transformations and uses a linear-time expectation maximization algorithm to jointly extract components and normalize for transformations. After illustrating the algorithm using a toy problem and facial expression modeling, we give results on classifying images of handwritten digits and compare the components found by transformed component analysis with the components found by PCA and FA.

2 The Transformed Component Analyzer (TCA)

A transformed component analyzer (TCA) is a probability model that specifies how the linear combination of a set of component vectors is transformed in different ways to model input patterns. The component activities (amplitudes) \mathbf{y} , which form a subspace representation of an image, are assumed to be independent and Gaussian:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{I}), \quad (1)$$

where \mathbf{I} is the identity covariance matrix. A latent image \mathbf{z} is produced by combining these components linearly using a “factor loading matrix” $\mathbf{\Lambda}$, offsetting the resulting image by an image mean $\boldsymbol{\mu}$ and then adding independent Gaussian noise to each pixel:

$$p(\mathbf{z}|\mathbf{y}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu} + \mathbf{\Lambda}\mathbf{y}, \boldsymbol{\Phi}), \quad (2)$$

where $\boldsymbol{\Phi}$ is a diagonal covariance matrix.

We assume the image produced by the subspace model is further transformed to obtain the observed image. Although real-valued transformations (*e.g.*, rotation) are common, real-valued latent variables introduce complicated integrals (instead of summations) into the EM learning algorithm. So, we will assume there is a fixed set of possible transformations and that this set is specified beforehand. Also, the algorithm is simplified by assuming that the vector of pixel values for the transformed image is obtained by multiplying the vector of pixel values for the latent image by a *sparse* matrix. This permits a broad class of transformations, including translation, scale, in-plane rotation and out-of-plane rotation.

Let $\ell \in \{1, \dots, L\}$ index the set of L transformations represented by matrices $\mathbf{\Gamma}_1, \dots, \mathbf{\Gamma}_L$. The probability density of the vector of pixel values \mathbf{x} for the image corresponding to transformation ℓ and latent image \mathbf{z} is

$$p(\mathbf{x}|\ell, \mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{\Gamma}_\ell \mathbf{z}, \boldsymbol{\Psi}), \quad (3)$$

where $\boldsymbol{\Psi}$ is a diagonal covariance matrix that specifies the noise on the observed pixels.

The variances $\boldsymbol{\Phi}$ of the pixels in the latent image are quite different from the variances $\boldsymbol{\Psi}$ of the pixels in the observed image. The noise modeled by $\boldsymbol{\Phi}$ coheres to the latent image during transformations. In contrast, $\boldsymbol{\Psi}$ models noise in the observed pixels and this noise does not depend on the transformation.

The translation, scale, rotation, *etc.* corresponding to each ℓ are selected ahead of time so that the set of transformations effectively covers the range of possible

transformations in the data and at the same time is small enough to keep the learning time reasonable. Transformation ℓ has a prior probability $P(\ell) = p_\ell$, which may be fixed to be uniform or learned from the training data.

The joint distribution over the observed image, the transformation index, the latent image and the subspace representation is

$$\begin{aligned} p(\mathbf{x}, \ell, \mathbf{z}, \mathbf{y}) &= p(\mathbf{y})P(\ell)p(\mathbf{z}|\mathbf{y})p(\mathbf{x}|\ell, \mathbf{z}) \\ &= p_\ell \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{I}) \mathcal{N}(\mathbf{z}; \boldsymbol{\mu} + \mathbf{\Lambda}\mathbf{y}, \boldsymbol{\Phi}) \mathcal{N}(\mathbf{x}; \mathbf{\Gamma}_\ell \mathbf{z}, \boldsymbol{\Psi}). \end{aligned} \quad (4)$$

In this generative model, the subspace representation, transformation index and latent image are unobserved variables. Generating an image from a TCA consists of drawing a subspace representation from $p(\mathbf{y})$, drawing a latent image from $p(\mathbf{z}|\mathbf{y})$, drawing a transformation index from $P(\ell)$ and then drawing an image from $p(\mathbf{x}|\ell, \mathbf{z})$.

The number of parameters in a TCA is roughly equal to the number of parameters in a standard factor analyzer and the number of “parameters” used in PCA. In Sec. 3, we show how the ML estimate of this model can be obtained using the EM algorithm.

2.1 Inferring the likelihood for an image and the responsibilities of the transformations

For a given image and set of model parameters, it is useful to know how appropriate each transformation is for matching the image to the transformed components (the “responsibility” of the transformation) and overall how well the components match the image (the likelihood).

The responsibilities are the posterior probabilities of the transformation indices,

$$P(\ell|\mathbf{x}) = p(\mathbf{x}, \ell)/p(\mathbf{x}), \quad (5)$$

where $p(\mathbf{x})$ is the likelihood:

$$p(\mathbf{x}) = \sum_{\ell=1}^L p(\mathbf{x}, \ell). \quad (6)$$

Both the responsibilities and the likelihood can be easily computed from the joint distribution over the image and transformation index, $p(\mathbf{x}, \ell)$.

To obtain $p(\mathbf{x}, \ell)$, it turns out that the component activities \mathbf{y} and the latent image \mathbf{z} can be integrated out in closed form, giving

$$\begin{aligned} p(\mathbf{x}, \ell) &= \int_{\mathbf{z}} \int_{\mathbf{y}} d\mathbf{z} d\mathbf{y} p(\mathbf{x}, \ell, \mathbf{z}, \mathbf{y}) \\ &= p_\ell \mathcal{N}(\mathbf{x}; \mathbf{\Gamma}_\ell \boldsymbol{\mu}, \mathbf{\Gamma}_\ell (\mathbf{\Lambda}\mathbf{\Lambda}' + \boldsymbol{\Phi}) \mathbf{\Gamma}_\ell' + \boldsymbol{\Psi}), \end{aligned} \quad (7)$$

where “ ’ ” indicates transpose. Each transformation ℓ has a corresponding mean image $\Gamma_\ell \boldsymbol{\mu}$ and covariance matrix $\Gamma_\ell(\boldsymbol{\Lambda}\boldsymbol{\Lambda}' + \boldsymbol{\Phi})\Gamma_\ell' + \boldsymbol{\Psi}$.

For an N -pixel image, K components and L transformations, it turns out the responsibilities and the likelihood for an image can be computed in $\mathcal{O}(LKN)$ time, if we assume $|\Gamma_\ell(\boldsymbol{\Lambda}\boldsymbol{\Lambda}' + \boldsymbol{\Phi})\Gamma_\ell' + \boldsymbol{\Psi}| \approx |\Gamma_\ell(\boldsymbol{\Lambda}\boldsymbol{\Lambda}' + \boldsymbol{\Phi})\Gamma_\ell'|$. This corresponds to assuming the image noise is significantly lesser than the variability introduced by the transformed components (*i.e.*, the image structure).

2.2 Inferring the subspace representation

For a given image \mathbf{x} , the posterior distribution over the component activities is given by

$$p(\mathbf{y}|\mathbf{x}) = \sum_{\ell=1}^L p(\mathbf{y}|\mathbf{x}, \ell)P(\ell|\mathbf{x}), \quad (8)$$

where $P(\ell|\mathbf{x})$ is the responsibility of transformation ℓ from above. Given ℓ , the TCA is a multinomial Gaussian model, so $p(\mathbf{y}|\mathbf{x}, \ell)$ is Gaussian and the above posterior is a mixture of Gaussians, where each Gaussian corresponds to a different transformation.

In some situations, the posterior distribution over \mathbf{y} may be highly multimodal, since significantly different transformations may be appropriate. For example, if the subspace represents the 3-D position of a point light source for an object that looks the same when an in-plane rotation of 180° is applied, then there will be two quite different possible light source positions corresponding to two transformations that differ by a 180° in-plane rotation.

In contrast with the single point representation obtained from PCA and ICA, TCA gives a set of weighted points. The weight for point ℓ is $P(\ell|\mathbf{x})$ and the corresponding estimate is

$$\mathbf{E}[\mathbf{y}|\mathbf{x}, \ell] = \boldsymbol{\beta}_\ell \boldsymbol{\Lambda}' \boldsymbol{\Phi}^{-1} [\boldsymbol{\Omega}_\ell \Gamma_\ell' \boldsymbol{\Psi}^{-1} \mathbf{x} - (\mathbf{I} - \boldsymbol{\Omega}_\ell \boldsymbol{\Phi}^{-1}) \boldsymbol{\mu}], \quad (9)$$

where

$$\begin{aligned} \boldsymbol{\Omega}_\ell &= \text{COV}(\mathbf{z}|\mathbf{x}, \mathbf{y}, \ell) = (\boldsymbol{\Phi}_c^{-1} + \Gamma_\ell' \boldsymbol{\Psi}^{-1} \Gamma_\ell)^{-1}, \\ \boldsymbol{\beta}_\ell &= (\mathbf{I} + \boldsymbol{\Lambda}' \boldsymbol{\Phi}^{-1} \boldsymbol{\Lambda} - \boldsymbol{\Lambda}' \boldsymbol{\Phi}^{-1} \boldsymbol{\Omega}_\ell \boldsymbol{\Phi}^{-1} \boldsymbol{\Lambda})^{-1}. \end{aligned} \quad (10)$$

To compare the subspace representations of two images, the two corresponding sets of weighted points can be compared. Alternatively, the MAP point (the point with the highest weight) can be selected for each image and a Euclidean distortion can be measured.

3 Transformed Component Analysis using the EM Algorithm

In this section, we describe an iterative expectation maximization (EM) algorithm for estimating the maximum likelihood parameters of a transformed component analyzer. We refer to this procedure as transformed component analysis (TCA). The EM algorithm consists of an E-step, in which statistics that are sufficient for estimating the model parameters are accumulated and an M-step, in which the model parameters are updated using the sufficient statistics. Each iteration increases the likelihood of the training data.

Notation: $\text{diag}(\mathbf{A})$ is the vector containing the diagonal elements of matrix \mathbf{A} ; $\text{diag}(\mathbf{a})$ is the diagonal matrix whose diagonal equals the vector \mathbf{a} ; $\mathbf{a} \circ \mathbf{b}$ is the element-wise product of vectors \mathbf{a} and \mathbf{b} .

3.1 Expectation step

The sufficient statistics for the M-Step are computed in the E-Step during a single pass through the training set. Before making this pass, the following matrices are computed:

$$\begin{aligned} \boldsymbol{\Omega}_\ell &= \text{COV}(\mathbf{z}|\mathbf{x}, \mathbf{y}, \ell) = (\boldsymbol{\Phi}_c^{-1} + \Gamma_\ell' \boldsymbol{\Psi}^{-1} \Gamma_\ell)^{-1}, \\ \boldsymbol{\beta}_\ell &= (\mathbf{I} + \boldsymbol{\Lambda}' \boldsymbol{\Phi}^{-1} \boldsymbol{\Lambda} - \boldsymbol{\Lambda}' \boldsymbol{\Phi}^{-1} \boldsymbol{\Omega}_\ell \boldsymbol{\Phi}^{-1} \boldsymbol{\Lambda})^{-1}. \end{aligned} \quad (11)$$

For each case during the pass through the training set, $P(\ell|\mathbf{x}_t)$, $\mathbf{E}[\mathbf{y}|\mathbf{x}_t, \ell]$ and $\mathbf{E}[\mathbf{y}|\mathbf{x}_t] = \sum_{\ell} P(\ell|\mathbf{x}_t) \mathbf{E}[\mathbf{y}|\mathbf{x}_t, \ell]$ are first computed as described in Sec. 2.1.

For each ℓ , we then compute

$$\begin{aligned} \mathbf{E}[\mathbf{z}|\mathbf{x}_t, \ell] &= \boldsymbol{\mu} + \boldsymbol{\Omega}_\ell \Gamma_\ell' \boldsymbol{\Psi}^{-1} (\mathbf{x}_t - \Gamma_\ell \boldsymbol{\mu}) \\ &+ \boldsymbol{\Omega}_\ell \boldsymbol{\Phi}^{-1} \boldsymbol{\Lambda} \boldsymbol{\beta}_\ell \boldsymbol{\Lambda}' \boldsymbol{\Phi}^{-1} \boldsymbol{\Omega}_\ell \Gamma_\ell' \boldsymbol{\Psi}^{-1} (\mathbf{x}_t - \Gamma_\ell \boldsymbol{\mu}) \end{aligned} \quad (12)$$

and obtain the following expectation:

$$\mathbf{E}[\mathbf{z} - \boldsymbol{\Lambda} \mathbf{y} | \mathbf{x}_t] = \sum_{\ell=1}^L P(\ell|\mathbf{x}_t) (\mathbf{E}[\mathbf{z}|\mathbf{x}_t, \ell] - \boldsymbol{\Lambda} \mathbf{E}[\mathbf{y}|\mathbf{x}_t, \ell]). \quad (13)$$

Two other expectations,

$$\begin{aligned} &\mathbf{E}[(\mathbf{z} - \boldsymbol{\mu}) \circ (\mathbf{z} - \boldsymbol{\mu}) | \mathbf{x}_t, \ell] \\ &= (\mathbf{E}[\mathbf{z}|\mathbf{x}_t, \ell] - \boldsymbol{\mu}) \circ (\mathbf{E}[\mathbf{z}|\mathbf{x}_t, \ell] - \boldsymbol{\mu}) \\ &\quad + \text{diag}(\boldsymbol{\Omega}_\ell) + \text{diag}(\boldsymbol{\Omega}_\ell \boldsymbol{\Phi}^{-1} \boldsymbol{\Lambda} \boldsymbol{\beta}_\ell \boldsymbol{\Lambda}' \boldsymbol{\Phi}^{-1} \boldsymbol{\Omega}_\ell), \end{aligned} \quad (14)$$

and

$$\begin{aligned} &\mathbf{E}[(\mathbf{z} - \boldsymbol{\mu}) \mathbf{y}' | \mathbf{x}_t, \ell] \\ &= (\mathbf{E}[\mathbf{z}|\mathbf{x}_t, \ell] - \boldsymbol{\mu}) \mathbf{E}[\mathbf{y}|\mathbf{x}_t, \ell]' + \boldsymbol{\Omega}_\ell \boldsymbol{\Phi}^{-1} \boldsymbol{\Lambda} \boldsymbol{\beta}_\ell, \end{aligned} \quad (15)$$

are used to compute the expectation,

$$\begin{aligned} \mathbb{E}[(\mathbf{z} - \boldsymbol{\mu} - \boldsymbol{\Lambda}\mathbf{y}) \circ (\mathbf{z} - \boldsymbol{\mu} - \boldsymbol{\Lambda}\mathbf{y}) | \mathbf{x}_t] &= \sum_{\ell=1}^L P(\ell | \mathbf{x}_t) \\ &\cdot \{ \mathbb{E}[(\mathbf{z} - \boldsymbol{\mu}) \circ (\mathbf{z} - \boldsymbol{\mu}) | \mathbf{x}_t, \ell] + \text{diag}(\boldsymbol{\Lambda}\boldsymbol{\beta}_\ell\boldsymbol{\Lambda}') \\ &\quad - 2\text{diag}(\boldsymbol{\Lambda}\mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})\mathbf{y}' | \mathbf{x}_t, \ell]) \\ &\quad + (\boldsymbol{\Lambda}\mathbb{E}[\mathbf{y} | \mathbf{x}_t, \ell]) \circ (\boldsymbol{\Lambda}\mathbb{E}[\mathbf{y} | \mathbf{x}_t, \ell]) \}. \end{aligned} \quad (16)$$

The following expectations are also computed:

$$\begin{aligned} \mathbb{E}[(\mathbf{x}_t - \boldsymbol{\Gamma}_\ell\mathbf{z}) \circ (\mathbf{x}_t - \boldsymbol{\Gamma}_\ell\mathbf{z}) | \mathbf{x}_t] &= \sum_{\ell=1}^L P(\ell | \mathbf{x}_t) \\ &\cdot \{ (\mathbf{x}_t - \boldsymbol{\Gamma}_\ell\mathbb{E}[\mathbf{z} | \mathbf{x}_t, \ell]) \circ (\mathbf{x}_t - \boldsymbol{\Gamma}_\ell\mathbb{E}[\mathbf{z} | \mathbf{x}_t, \ell]) \\ &\quad + \text{diag}(\boldsymbol{\Gamma}_\ell\boldsymbol{\Omega}_\ell\boldsymbol{\Gamma}_\ell') \\ &\quad + \text{diag}(\boldsymbol{\Gamma}_\ell\boldsymbol{\Omega}_\ell\boldsymbol{\Phi}^{-1}\boldsymbol{\Lambda}\boldsymbol{\beta}_\ell\boldsymbol{\Lambda}'\boldsymbol{\Phi}^{-1}\boldsymbol{\Omega}_\ell\boldsymbol{\Gamma}_\ell') \}, \end{aligned} \quad (17)$$

$$\mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})\mathbf{y}' | \mathbf{x}_t] = \sum_{\ell=1}^L P(\ell | \mathbf{x}_t) \mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})\mathbf{y}' | \mathbf{x}_t, \ell], \quad (18)$$

$$\mathbb{E}[\mathbf{y}\mathbf{y}' | \mathbf{x}_t] = \sum_{\ell=1}^L P(\ell | \mathbf{x}_t) (\boldsymbol{\beta}_\ell + \mathbb{E}[\mathbf{y} | \mathbf{x}_t, \ell] \mathbb{E}[\mathbf{y} | \mathbf{x}_t, \ell]'). \quad (19)$$

3.2 Maximization step

We use $\langle \cdot \rangle = \frac{1}{T} \sum_{t=1}^T (\cdot)$ to indicate a statistic computed by averaging the above expectations over the training set and “ $\tilde{\cdot}$ ” to denote the updated parameters:

$$\tilde{\boldsymbol{\mu}} = \langle \mathbb{E}[\mathbf{z} - \boldsymbol{\Lambda}\mathbf{y} | \mathbf{x}_t] \rangle, \quad (20)$$

$$\tilde{\boldsymbol{\Phi}} = \text{diag}(\langle \mathbb{E}[(\mathbf{z} - \boldsymbol{\mu} - \boldsymbol{\Lambda}\mathbf{y}) \circ (\mathbf{z} - \boldsymbol{\mu} - \boldsymbol{\Lambda}\mathbf{y}) | \mathbf{x}_t] \rangle), \quad (21)$$

$$\tilde{\boldsymbol{\Psi}} = \text{diag}(\langle \mathbb{E}[(\mathbf{x}_t - \boldsymbol{\Gamma}_\ell\mathbf{z}) \circ (\mathbf{x}_t - \boldsymbol{\Gamma}_\ell\mathbf{z}) | \mathbf{x}_t] \rangle), \quad (22)$$

$$\tilde{\boldsymbol{\Lambda}} = \langle \mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})\mathbf{y}' | \mathbf{x}_t] \rangle \langle \mathbb{E}[\mathbf{y}\mathbf{y}' | \mathbf{x}_t] \rangle^{-1}, \quad (23)$$

$$\tilde{p}_l = \langle P(\ell = l | \mathbf{x}_t) \rangle. \quad (24)$$

In order to avoid overfitting the noise variances, it is sometimes useful to set the diagonal elements of $\boldsymbol{\Phi}$ and $\boldsymbol{\Psi}$ that are below some ϵ equal to ϵ .

4 Learning shape and lighting representations from noisy unaligned images of an object

Fig. 1 shows a training set of 144×9 noisy images of a uniformly colored pyramid (gray) at randomly selected positions and illuminated by parallel light rays with randomly selected angle and intensity. A cluttered background was simulated by randomly selecting pixel values from a uniform distribution.

The first 8 principal components of the training data, scaled by the standard deviation of the projected data, are shown in Fig. 2a. It appears the components implement a multiresolution approximation to model shifts of the object.

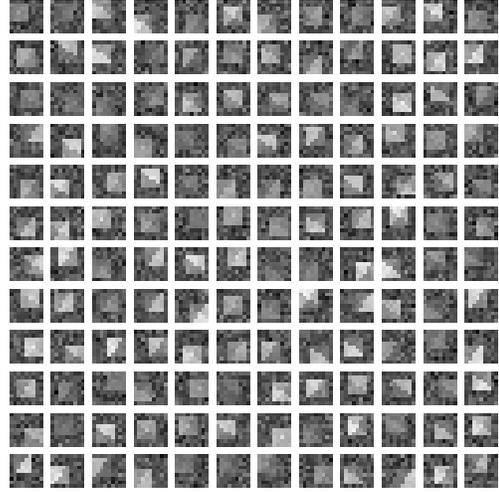


Figure 1: Noisy images of a pyramid at different locations and under different lighting conditions.

We trained a TCA with 3 components and 81 transformations implementing 9 horizontal and 9 vertical shifts using 10 iterations of the EM algorithm. To initialize the parameters, the mean and variance of each pixel was first computed from the training data. The parameters were then initialized to random values, using the mean and variance as a 1st order guide. The transformation probabilities p_ℓ were set equal.

Fig. 2b shows the mean latent image $\boldsymbol{\mu}$, the 3 columns of $\boldsymbol{\Lambda}$ (shown as 3 images), the latent image noise $\boldsymbol{\Phi}$ (shown as an image where the pixel intensity is equal to 4 times the standard deviation) and the observed image noise $\boldsymbol{\Psi}$. The mean clearly shows that the outline of the object has been determined and that the uniform coloring has been determined (except at the point of the pyramid). Linear combinations of the 3 components produce different lighting conditions (see the following paragraph) which implies that the 3-element *rows* of $\boldsymbol{\Lambda}$ are proportional to the object surface normals, up to some rotation in 3-dimensional space. The variance map for the latent image shows that the model predicts low variance for pixels belonging to the object, but high variance for other pixels (the background clutter). Finally, the variance map for the observed image accounts for the small amount of noise that is present in the images.

The TCA can be simulated in a noise-free transformation-free fashion, by drawing a subspace representation from $p(\mathbf{y})$, computing $\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{y}$ and then computing $\mathbf{x} = \boldsymbol{\Gamma}_1\mathbf{z}$. Fig. 2c shows 144 examples simulated in this way. These “fantasies” show that the TCA can simulate the different lighting conditions.

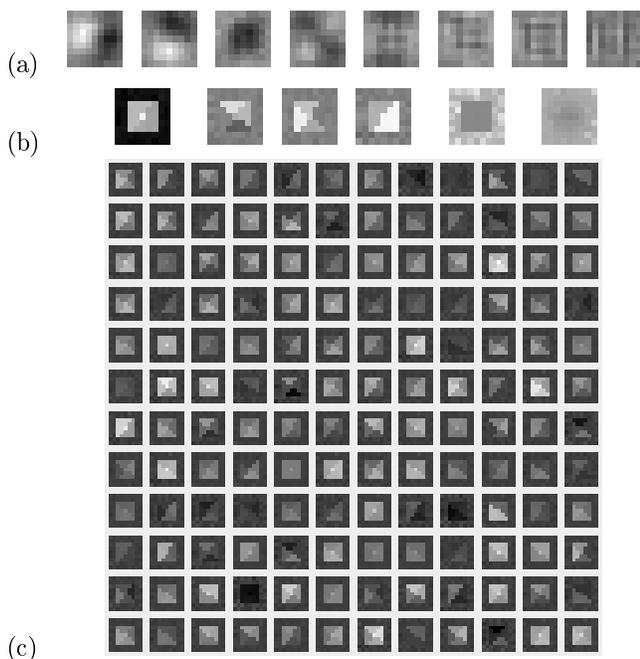


Figure 2: (a) The first 8 scaled principal components. A pixel colored gray – halfway between black and white – corresponds to a component value of 0. (b) The mean, components, and noise deviation of a TCA with 3 components, after 10 iterations of EM. Pixels for the mean and the noise deviations are colored using the same scale as the training images. (c) Examples simulated from the TCA, without noise and without transformations.

5 Learning a subspace representation of facial expressions from imperfectly aligned images

Fig. 3 shows a training set of 100 16×24 images of automatically aligned faces with different expressions. The accuracy of the face detection algorithm used to align the images is ± 2 pixels in each direction.

Fig. 4a shows the mean of the training data and the first 10 principal components, scaled by the standard deviation of the projected data. The first 5 components obviously account for vertical, horizontal and diagonal shifts in the data and the remaining components are very blurred.

To see how well the PCA subspace represents the data, we can draw a subspace point from an axis-aligned Gaussian with variances determined from the projected training data, and then use the principal components to map the point to image space. 100 examples simulated in this manner are shown in Fig. 4b. Although the “faces” do appear to be shifted around



Figure 3: Imperfectly aligned images of faces with different expressions.

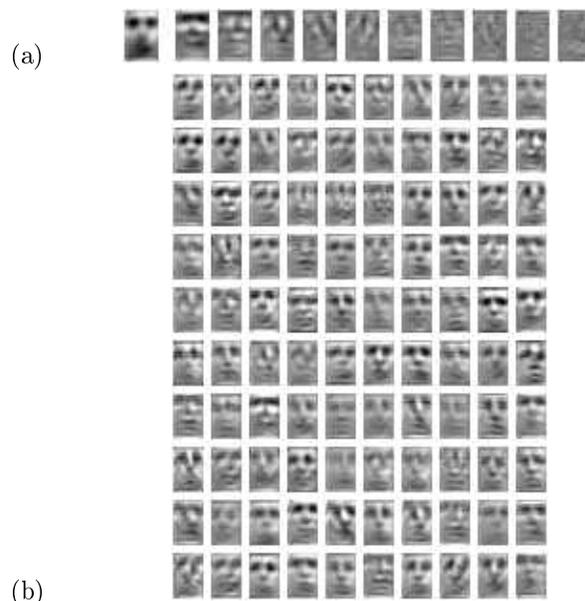


Figure 4: (a) The mean and first 10 scaled principal components. The mean is colored using the same scale as the training images, whereas the components are colored using intermediate gray to represent 0. (b) Examples simulated using the PCA subspace.

the field of vision, they are also severely blurred.

Fig. 5a and b show the parameters and simulated examples (without adding the sensor noise) for a FA model (a TCA with only the identity transformation) trained using 70 iterations of EM. The parameters were initialized using the mean and variance of each

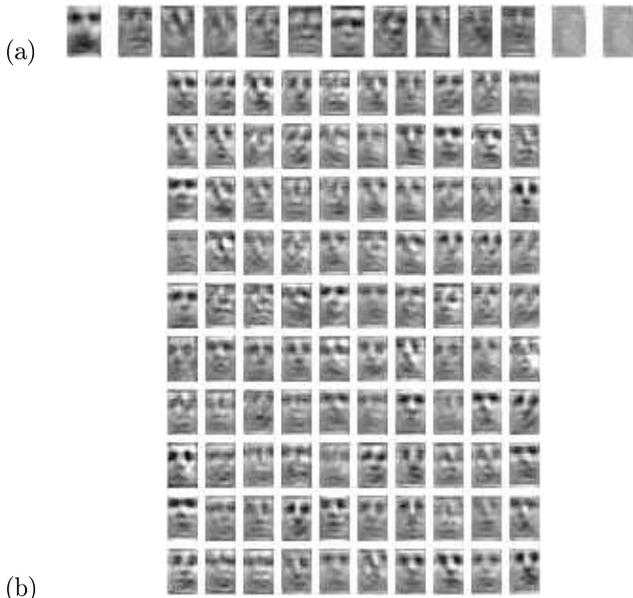


Figure 5: (a) The mean, 10 components and noise deviations found by FA (TCA with only the identity transformation). Pixels for the mean and noise deviations are colored using the same scale as the training images. (b) Examples simulated using the factor analyzer.

pixel in the training data. The sum of the two images on the far right of Fig. 5a show the variance map for FA. In contrast to PCA, different components represent similar amounts of energy (variance). This is because FA does not find a preferred set of basis vectors (factors) for the subspace. In fact, FA is not identifiable in this regard.

We trained a TCA with 10 components and 25 transformations implementing 5 horizontal and 5 vertical shifts using 70 iterations of the EM algorithm. The parameters were initialized using the mean and variance of each pixel in the training data. The transformation probabilities p_ℓ were set equal.

Fig. 6a shows the mean, components and variances maps. Unlike PCA and FA, TCA extracts clear components. The first component appears to expose some teeth, the second component appears to raise the eyebrows, raise the upper lip and expose a “tongue”, and so on. The components found by TCA are unique up to a unitary transformation, so each component often includes more than one feature. A further processing step can be applied to find a unitary transformation that produces components with spatially localized energy.

Fig. 6b shows 100 examples of images simulated us-

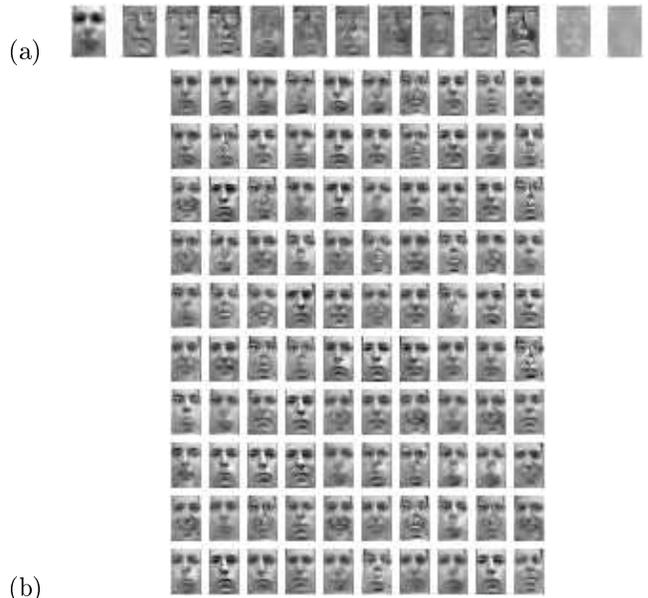


Figure 6: (a) The mean, 10 components and noise deviations found by TCA. Pixels for the mean and noise deviations are colored using the same scale as the training images. (b) Examples simulated using the TCA.

ing the TCA, without the latent image and observed image noise and without randomly selected transformations. The images are much clearer than those simulated using the PCA subspace and the factor analyzer. The expressions in the training set are reproduced and the model also generates novel realistic expressions that are not present in the training set, such as the one in the 5th column of the 1st row.

6 Using shear transformations to improve handwritten digit recognition

Fig. 7 shows 230 images of handwritten digits taken from envelopes mailed by people in Buffalo [5]. A wide variety of handwriting styles are evident. The original images were affinely transformed to fit snugly in an 8×8 gray level image. Although this preprocessing helps to normalize for uniform scaling in the horizontal and vertical directions, it does not normalize for local scaling (*e.g.*, some 8s have smaller loops than others) or for shearing (*e.g.*, the vertical stroke of the 7 is at different angles).

For each class of digit, we trained a 10 component factor analyzer on 200 images using 30 iterations of EM. The parameters of each factor analyzer were initialized using the mean and variance of each pixel in

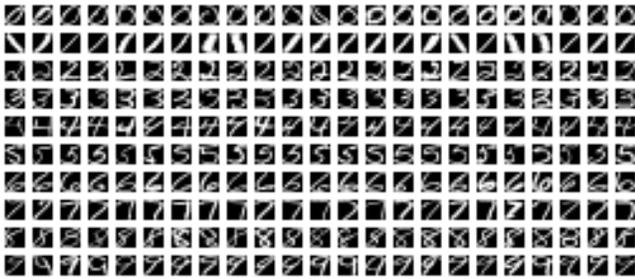


Figure 7: Example images from the training set of handwritten digits show a wide variety in writing style.

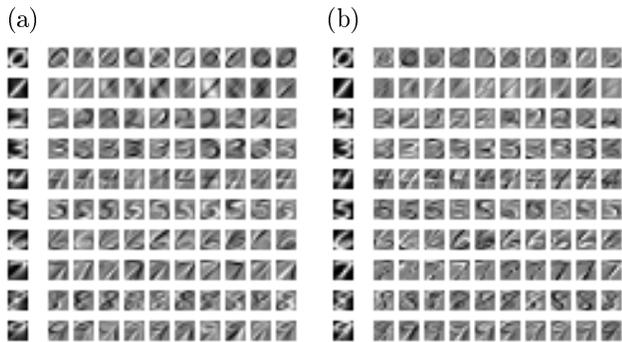


Figure 8: The means and components learned by (a) FA and (b) TCA applied to 200 images of each handwritten digit. Pixels for the means are colored using the same scale as the training images; pixels for the components are colored using intermediate gray to represent 0.

the respective class. Fig. 8a shows the means and components found by FA and Fig. 9 shows images simulated from the trained models, without adding noise.

Some components account for interesting variability in writing style, but many clearly account for simple spatial transformations. For example, most of the components for 7 simply negate the vertical stroke in the image mean and produce a new vertical stroke at a different angle.

Shearing is a simple spatial transformation to model using a sparse matrix. Using 30 iterations of EM, we trained one TCA on each class of 200 digits using 29 transformations that included horizontal shearing plus horizontal translations. The first column of images in Fig. 11 shows how the 29 transformations modify a contrived pattern. Some transformations are quite severe and erase a large number of pixels, so in this experiment the transformation probabilities, p_t , were learned. The parameters were initialized using the mean and variance of each pixel in the training

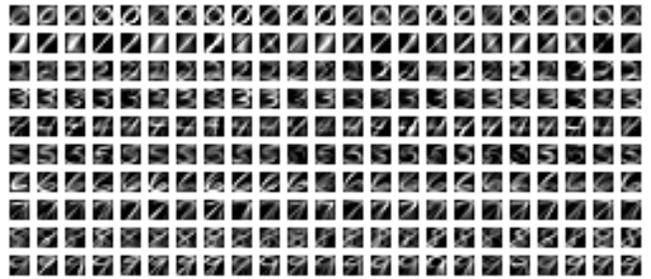


Figure 9: Images simulated from the FA models.

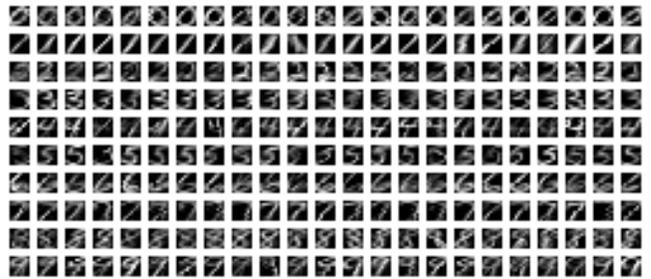


Figure 10: Images simulated from the TCA models.

set.

Fig. 8b shows the means and components found by TCA. The 2nd component for 7 models the accent on the nose of the 7; the 3rd component for 7 rounds out the angle in the upper right. These features are not clearly represented in the components found by FA. Fig. 10 shows digits simulated from the TCA models, without adding noise onto the latent image and the observed image. These images show a greater variation in style and fewer inappropriate artifacts than the images simulated from the FA model.

The 10 columns of images on the right in Fig. 11 show the means and corresponding sheared images from the TCA model. Those transformed images whose average expected responsibilities are less than 0.0001 in the training set are faded (shown with low contrast). A wide variation in writing style is captured by the transformations. For example, whereas FA uses components to model the different angles of vertical strokes in images of the digit 7, TCA models the different angles using the transformations. This allows TCA to use its components to model more subtle variations in writing style, such as the accent on the nose of the 7, described above.

6.1 Improvement in recognition rate

We performed FA (20 components) and TCA (20 components) on 200 training cases from each class of handwritten digit using 50 iterations of EM. The noise



Figure 11: The mean image learned by TCA for each class of digit, along with the corresponding sheared images. The first column of images uses a contrived pattern to illustrate the 29 shear and translation transformations that were used by TCA. For each class of digit, the mean learned by TCA is shown followed by the 29 transformed versions of the mean. Those versions whose average expected responsibilities are less than 0.0001 in the training set are faded.

variances were not allowed to drop below 10^{-2} to prevent overfitting a pixel that happens to always be off in the training data. Bayes rule was then used to classify 1000 test patterns. The results are summarized in Table 1 and compared with a standard feedforward method, k -nearest neighbors, where k was chosen using leave-one-out cross validation. TCA has a lower error rate than the other two methods.

The probability of each transformation p_ℓ in TCA was learned and we believe there was some overfitting. For example, some of the sheared image means in Fig. 11a that are faded (have average responsibilities less than 0.0001) are good generalizations. We are currently running experiments that regularize p_ℓ and we are also running experiments to find the performance of support vector machines on this data.

Table 1: Handwritten digit recognition rates, for a training set of 2000 8×8 images and a test set of 1000 images.

Method	Error rate
k -nearest neighbors	7.6%
Factor analysis	3.2%
Transformed component analysis	2.7%

7 Summary

We introduced a technique called “transformed component analysis” (TCA), that can jointly learn components from data and normalize for transformations. We explored transformations for translation and shearing, but the method can be applied to other transformations, such as rotation, scale and warping. When applied to noisy images of a uniformly colored pyramid at randomly selected positions and illuminated by parallel light rays with randomly selected angle and intensity, TCA is able to extract representations of shape and lighting, despite a significant level of background clutter. TCA also found a clear representation of facial expression from imperfectly aligned images, whereas PCA and factor analysis (FA) found severely blurred components. By including shearing transformations in the TCA model, we showed that handwritten digit recognition performance was improved over FA and that the images simulated from the TCA model contained greater variation in style and fewer inappropriate artifacts than the images simulated from FA.

References

- [1] I. T. Jolliffe, *Principal Component Analysis*, Springer-Verlag, New York NY., 1986.
- [2] A. J. Bell and T. J. Sejnowski, “An information maximization approach to blind separation and blind deconvolution,” *Neural Computation*, vol. 7, pp. 1129–1159, 1995.
- [3] B. S. Everitt, *An Introduction to Latent Variable Models*, Chapman and Hall, New York NY., 1984.
- [4] D. Rubin and D. Thayer, “EM algorithms for ML factor analysis,” *Psychometrika*, vol. 47, no. 1, pp. 69–76, 1982.
- [5] J. J. Hull, “A database for handwritten text recognition research,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 16, no. 5, pp. 550–554, 1994.