Tensor Analyzers

Anonymous Author(s) Affiliation Address email

Abstract

Factor Analysis is a statistical method that seeks to explain linear variations in data by using unobserved latent variables. Due to its *additive* nature, it is not suitable for modeling data that is generated by multiple groups of latent factors which interact *multiplicatively*. In this paper, we introduce Tensor Analyzers which are a multilinear generalization of Factor Analyzers. We describe a fairly efficient way of sampling from the posterior distribution over factor values and we demonstrate that these samples can be used in the EM algorithm for learning interesting mixture models of natural image patches and of images containing a variety of simple shapes that vary in size and color. Tensor Analyzers can also accurately recognize a face under significant pose and illumination variations when given only one previous image of that face. We also show that mixtures of Tensor Analyzers outperform mixtures of Factor Analyzers at modeling natural image patches and artificial data produced using multiplicative interactions.

1 Introduction

000

002

008

009

010

011 012 013

014 015

016

017

018

019

021

025

026

027 028 029

031

Exploratory Factor Analysis is widely used in statistics to identify underlying linear factors. Mixtures of Factor Analyzers have been used successfully for unsupervised learning [20, 18]. Factor Analyzers (FAs) model each observation vector as a weighted linear combination of the unobserved factor values plus additive uncorrelated noise. For many types of data, this additive generative process is less suitable than a generative process that also contains multiplicative interactions between latent factors.

In this paper, we introduce Tensor Analyzers (TAs) which are one way of generalizing FAs to include multiplicative interactions between latent variables from different groups. A TA generates data in the following way: First we sample all of the factor values in all of the groups independently from a zero-040 mean unit-variance Gaussian. Then, for each possible way of choosing one factor from each latent 041 group, we take the product of the factor values and contribute this product times a learned parameter 042 to each component of the generated data-vector. Finally, we add independent Gaussian noise to every 043 component of the generated vector with the variance of the noise being learned separately for each 044 component. The number of "tensor loading" parameters is the product of the sizes of all the factor groups and the dimensionality of the data-vector. We may additionally use lower order products that do not include a factor from every group, but in this paper the only lower-order products we use are 046 the individual factors of every group. FAs are special cases of TAs that contain only one group of 047 factors. 048

Conditioned on all but one group of factors, a TA reduces to an ordinary FA in which the factor
loadings are a function of the factor values in the groups we are conditioning on. The posterior
distribution of the factor values in a FA can be computed analytically, so by cycling through each
group of factors, efficient alternating Gibbs sampling is therefore possible in the TA. A TA is a proper
density model so the extension to a mixture of TAs (MTA) is straightforward. When performing
inference or learning in a TA it is easy to make use of a supervisory signal that specifies the fact

that 2 or more different observations are generated from the same factor values in some of the factor groups.

The related tensor decomposition has been widely applied to signal processing and computer vi-057 sion [16, 2, 9, 19]. The SVD algorithm was used to learn a bilinear model to separate style and content [15] (henceforth referred to as the S&C model). Tucker decomposition was applied to a 5-mode array of face images in [17], finding multilinear bases called TensorFaces. However, tensor 060 decomposition based methods lack a probabilistic formulation and are not density models. Inference 061 given a *single* new test case is ill-posed and can be ad hoc¹. Moreover, training data must be arranged 062 into a tensor [11]. More recently, bilinear models with priors on the latent variables have been pro-063 posed. In [6], sparsity is induced in the codes of a bilinear model to learn translational invariant 064 representation from video. However, the model does not try to maximize the $\log p(\mathbf{x})$ but instead finds the MAP estimate of the code activations à la sparse coding. [3] describe an outer-product 065 factorization of the bilinear model that is trained as a density model using EM, but expensive Hamil-066 tonian dynamics is required for sampling from the posterior. In addition, their model only admits an 067 approximate M-step, and does not make it easy to incorporate label information. 068

069 In contrast, TAs do not have any of the above deficiencies. A TA is a density model that can be learned directly from data vectors in an entirely unsupervised manner, but it can also make use of 071 supervision in the form of equality constraints that specify that one group of factors should have the same vector of values for a subset of the training cases (Sec. 3.4). As an extension to FA, TA 072 inherits an efficient inference algorithm that is used in each step of alternating Gibbs sampling and 073 a closed-form M-step during learning. Unlike bilinear models, it can handle multilinear cases with 074 3 or more groups of latent factors (Sec. 4.3). It can also be easily extended to a mixture model, 075 provided we are willing to compute approximate densities as described in (Sec. 3.3). In addition, 076 posterior inference for a *single* test case is simple and accurate, as demonstrated by one-shot face 077 recognition experiments of Sec. 4.4.

079

2 Preliminaries

Following [8], we refer to the number of dimensions of the tensor as its *order* (also known as modes). We will use bold lowercase letters to denote vectors (tensors of order one), e.g. x; bold uppercase letters for matrices (tensors of order two), e.g. W. We use the notation $\mathbf{w}_{(i,:)}$ to denote the *i*-th row of matrix W. Higher order tensors are denoted by Euler script letters, e.g. a third-order tensor with dimensions of *I*, *J*, and $K: \mathfrak{T} \in \mathbb{R}^{I \times J \times K}$.

Fibers: Fibers are higher-order generalization of row/column vectors. Elements of a tensor fiber is found by fixing all but one index. Specifically, $t_{(:,j,k)}$ is the mode-1 fiber of the tensor \mathcal{T} . Row and column vectors are the mode-2 and mode-1 fiber of a 2nd-order tensor, respectively.

Matricization: Matricization is the process of "flattening" a tensor into a matrix, by reordering the elements of the tensor. It is denoted by $\mathbf{T}_{(n)}$, where the mode-*n* fibers of \mathfrak{T} are placed in the columns of the resulting matrix $\mathbf{T}_{(n)}$. For example, given $\mathfrak{T} \in \mathbb{R}^{I \times J \times K}$, $\mathbf{T}_{(1)} \in \mathbb{R}^{I \times J K}$.

n-mode vector product: By multiplying a vector $\mathbf{y} \in \mathbb{R}^{D_n}$ with a tensor $\boldsymbol{\mathcal{T}} \in \mathbb{R}^{D_1 \times D_2 \times \cdots \times D_N}$ along the mode-*n*, the *n-mode (vector) product* is denoted by $\boldsymbol{\mathcal{T}} \times_n \mathbf{y}$. The resulting tensor is of size $D_1 \times \cdots \times D_{n-1} \times D_{n+1} \times \cdots \times D_N$.

096 2.1 Factor Analyzers

Let $\mathbf{x} \in \mathbb{R}^D$ denote the *D*-dimensional data, let $\{\mathbf{z} \in \mathbb{R}^d : d \leq D\}$ denote *d*-dimensional latent factors. FA is a directed model, defined as by a prior and likelihood:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}), \quad p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{\Lambda} \mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}), \tag{1}$$

I is the $d \times d$ identity matrix; $\mathbf{\Lambda} \in \mathbb{R}^{D \times d}$ is the factor loading matrix, $\boldsymbol{\mu}$ is the mean; and a diagonal $\boldsymbol{\Psi} \in \mathbb{R}^{D \times D}$ represents the variance of the observation noise. By integrating out the latent variable \mathbf{z} , a FA model becomes a Gaussian with constrained covariance:

$$p(\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Gamma), \quad \Gamma = \mathbf{\Lambda} \mathbf{\Lambda}^{\mathsf{T}} + \boldsymbol{\Psi}$$
(2)

095

097 098

099 100 101

102

103

¹E.g., the asymmetric model in (S&C) requires EM learning of a separate model during test time.

108 109 110

111

112

113

114

115

For inference, we are interested in the posterior, which is also a multivariate Gaussian:

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mathbf{m}, \mathbf{V}^{-1}), \tag{3}$$

where $\mathbf{V} = \mathbf{I} + \mathbf{\Lambda}^{\mathsf{T}} \Psi^{-1} \mathbf{\Lambda}$, and $\mathbf{m} = \mathbf{V}^{-1} \mathbf{\Lambda}^{\mathsf{T}} \Psi^{-1} (\mathbf{x} - \boldsymbol{\mu})$. Maximum likelihood estimation of the parameters is straightforward using the EM algorithm [14]. During the E-step, Eqs. 3 are used to compute the posterior sufficient statistics given the current setting of the model parameters. During the M-step, the expected complete-data log-likelihood $\mathbb{E}_{p(\mathbf{z}|\mathbf{x};\theta_{old})}[\log p(\mathbf{x},\mathbf{z};\theta)]$ is maximized with respect to the model parameters $\theta = \{\Lambda, \mu, \Psi\}$.

116 117 118

121

125 126

127 128 129

130

131 132

133

134

135 136

137

138

139

140

141

142

143

149

150

151

152

156

3 **Tensor Analyzers**

119 FA generates data by linear combining factors, as there are no multiplicative interactions involving 120 terms such as $z_i z_j, i \neq j$. By using a (J+1)-order "loading" tensor $\mathfrak{T} \in \mathbb{R}^{D \times d_1 \times \cdots \times d_J}$, a TA can model multiplicative interactions among its J groups of factors: $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_J\}$. For $j = 1, \dots, J$, 122 $\mathbf{z}_i \in \mathbb{R}^{d_j}$. We will use the notation TA $\{D, d_1, d_2, \dots, d_J\}$ to denote the aforementioned TA with 123 latent factors $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_J\}$ and interaction tensor \mathcal{T} . We give each group of factors a standard 124 Normal prior:

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|0, \mathbf{I}), \quad j = 1, 2, \dots, J \tag{4}$$

For clarity and WLOG, we assume J = 3 for the following equations. The likelihood is defined as:

$$p(\mathbf{x}|\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3) = \mathcal{N}(\mathbf{x}|\mathbf{m} + \mathbf{W}_1\mathbf{z}_1 + \mathbf{W}_2\mathbf{z}_2 + \mathbf{W}_3\mathbf{z}_3 + \mathbf{T}_{(1)}(\mathbf{z}_3 \otimes \mathbf{z}_2 \otimes \mathbf{z}_1), \Psi),$$
(5)



Figure 1: Diagram of TA's (J 2) generative process. Vector-tensor multiplications determines the mean of $p(\mathbf{x}|\mathbf{z}_1,\mathbf{z}_2).$

where $\mathbf{x} \in \mathbb{R}^{D}$, \mathbf{m} , and Ψ are same as in FA. $\mathbf{W}_{j} \in \mathbb{R}^{D \times d_{j}}$ are the "biases" factor loadings, $\mathbf{T}_{(1)} \in \mathbb{R}^{D \times (d_1 d_2 d_3)}$ is the matricization of the tensor T as discussed in Sec. 2, and " \otimes " is the Kronecker product operator. Multiplicative interactions are due to the term: $\mathbf{z}_3 \otimes \mathbf{z}_2 \otimes \mathbf{z}_1$, which is a vector with dimensionality of $d_1 \times d_2 \times d_3$.

For clarity, we can concatenate the factors and loading matrices: let $\mathbf{y} \in \mathbb{R}^{d_1+d_2+d_3+1} \triangleq [\mathbf{z}_1; \mathbf{z}_2; \mathbf{z}_3; 1]; \mathbf{W} \in$ $\mathbb{R}^{D \times (d_1 + d_2 + d_3 + 1)} \triangleq [\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3, \mathbf{m}]; \text{ and } \mathbf{u} \in \mathbb{R}^{d_1 d_2 d_3} =$ $\mathbf{z}_3 \otimes \mathbf{z}_2 \otimes \mathbf{z}_1$. We note that $\mathbf{T}_{(1)}(\mathbf{z}_3 \otimes \mathbf{z}_2 \otimes \mathbf{z}_1)$ is mathematically equivalent to $\sum_{i,j,k} \mathbf{t}_{(:,i,j,k)} \mathbf{z}_1(i) \mathbf{z}_2(j) \mathbf{z}_3(k)$, where $\mathbf{z}_1(i)$ is the *i*-th element of vector \mathbf{z}_1 , and \mathbf{t} is the mode-1 fiber of T.

The joint/complete log-likelihood of the TA is:

$$\log p(\mathbf{x}, \mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3) = \sum_{j=1}^3 \left(-\frac{d_j}{2} \log(2\pi) - \frac{1}{2} \mathbf{z}_j^{\mathsf{T}} \mathbf{z}_j \right) - \frac{D}{2} \log(2\pi) - \frac{1}{2} \log|\Psi| - \frac{1}{2} (\mathbf{x} - \mathbf{e})^{\mathsf{T}} \Psi^{-1} (\mathbf{x} - \mathbf{e})$$
(6)

where $\mathbf{e} = \mathbf{W}\mathbf{y} + \mathbf{T}_{(1)}\mathbf{u}$. See Fig. 1 for a visual diagram of the TA's generative process. The last term in Eq. 6 indicates that the TA models contain higher-order interactions (squared of the outer product of all factors). In comparison, FAs have only 2nd-order interactions among its latent factors.

Conditioned on any two of the three groups of factors, e.g. z_2 and z_3 , the log-likelihood of x and z_1 153 becomes: 154

$$\log p(\mathbf{x}, \mathbf{z}_1 | \mathbf{z}_2, \mathbf{z}_3) = -\frac{d_1}{2} \log(2\pi) - \frac{1}{2} \mathbf{z}_1^{\mathsf{T}} \mathbf{z}_1 - \frac{D}{2} \log(2\pi) - \frac{1}{2} \log|\mathbf{\Psi}| - \frac{1}{2} (\mathbf{x} - \mathbf{e})^{\mathsf{T}} \mathbf{\Psi}^{-1} (\mathbf{x} - \mathbf{e})$$
(7)

157 Here, e can be re-written as $(m + W_2 z_2 + W_3 z_3) + (W_1 + \Im \overline{x}_3 z_3 \overline{x}_2 z_2) z_1$. We can see that 158 conditioned on z_2 and z_3 , we have a FA with parameters (c.f. Eq. 1): 159

161

$$oldsymbol{\mu} = \mathbf{m} + \mathbf{W}_2 \mathbf{z}_2 + \mathbf{W}_3 \mathbf{z}_3, ~~ oldsymbol{\Lambda} = \mathbf{W}_1 + oldsymbol{ au} ildsymbol{ar{ extsf{x}}}_3 ~ \mathbf{z}_3 ~oldsymbol{ar{ extsf{x}}}_2 ~ \mathbf{z}_2$$

(8)

The marginal probability density function is a Gaussian: $p(\mathbf{x}|\mathbf{z}_2, \mathbf{z}_3) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}\boldsymbol{\Lambda}^{\mathsf{T}} + \boldsymbol{\Psi}).$

¹⁶² 3.1 Inference

167 168 169

171 172

178 179

180 181

182

183

185

186

187

188

189

190

191

192

193

194

205 206 207

211

Higher order interaction in the TA means that inference is more complicated as the joint posterior $p(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_3 | \mathbf{x})$ has no closed-form solution. We resort to alternating Gibbs sampling: Step 1: $p(\mathbf{z}_1 | \mathbf{x}, \mathbf{z}_2, \mathbf{z}_3)$; Step 2: $p(\mathbf{z}_2 | \mathbf{x}, \mathbf{z}_1, \mathbf{z}_3)$; Step 3: $p(\mathbf{z}_3 | \mathbf{x}, \mathbf{z}_1, \mathbf{z}_2)$.

Conditioned on two groups of factors, the posterior of third is simple as the model reduces to a FA:

$$p(\mathbf{z}_1|\mathbf{x},\mathbf{z}_2,\mathbf{z}_3) = \mathcal{N}(\mathbf{z}_1|\mathbf{V}^{-1}\mathbf{\Lambda}^{\mathsf{T}}\boldsymbol{\Psi}^{-1}(\mathbf{x}-\boldsymbol{\mu}),\mathbf{V}^{-1}),$$
(9)

where $\mathbf{V} = \mathbf{I} + \mathbf{\Lambda}^{\mathsf{T}} \boldsymbol{\Psi}^{-1} \mathbf{\Lambda}$. $\boldsymbol{\mu}$, and $\boldsymbol{\Lambda}$ are defined by Eq. 8.

3.2 Learning

173 Maximum likelihood learning of a TA is similar to FA and is straightforward using a variant of 174 the EM algorithm. During the E-step, MCMC samples are drawn from the posterior distribu-175 tion using alternating Gibbs sampling. In the M-step, the samples are used to approximate the 176 sufficient statistics involving u and y, followed by closed-form updates of the model parameters, 177 $\theta = \{\mathbf{W}, \mathbf{T}_{(1)}, \Psi\}$. The expected joint log-likelihood function is:

$$Q = E \left[\log \prod_{i}^{N} (2\pi)^{-D/2} |\Psi|^{-1/2} \exp\{-\frac{1}{2} (\mathbf{x}_{i} - \mathbf{e}_{i})^{\mathsf{T}} \Psi^{-1} (\mathbf{x}_{i} - \mathbf{e}_{i})\}\right]$$
(10)

Setting $\frac{\partial Q}{\partial \theta} = 0$, we have update equations:

$$\mathbf{W} = \left(\sum_{i}^{N} \mathbf{x}_{i} E[\mathbf{y}_{i}^{\mathsf{T}}] - \mathbf{T}_{(1)} \sum_{i}^{N} E[\mathbf{u}_{i} \mathbf{y}_{i}^{\mathsf{T}}]\right) \left(\sum_{i}^{N} E[\mathbf{y}_{i} \mathbf{y}_{i}^{\mathsf{T}}]\right)^{-1} (11)$$

$$\mathbf{T}_{(1)} = \left(\sum_{i}^{N} \mathbf{x}_{i} E[\mathbf{u}_{i}^{\mathsf{T}}] - \mathbf{W} \sum_{i}^{N} E[\mathbf{y}_{i} \mathbf{u}_{i}^{\mathsf{T}}]\right) \left(\sum_{i}^{N} E[\mathbf{u}_{i} \mathbf{u}_{i}^{\mathsf{T}}]\right)^{-1} (12)$$

$$\Psi = \frac{1}{N} diag \left\{ \sum_{i}^{N} \left(\mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} - 2\mathbf{T}_{(1)} \left(E[\mathbf{u}_{i}] \mathbf{x}_{i}^{\mathsf{T}} - E[\mathbf{u}_{i} \mathbf{y}_{i}^{\mathsf{T}}] \mathbf{W}^{\mathsf{T}} - \frac{1}{2} E[\mathbf{u} \mathbf{u}^{\mathsf{T}}] \mathbf{T}_{(1)}^{\mathsf{T}} \right) - 2\mathbf{W} \left(E[\mathbf{y}_{i}] \mathbf{x}_{i}^{\mathsf{T}} - \frac{1}{2} E[\mathbf{y}_{i} \mathbf{y}_{i}^{\mathsf{T}}] \mathbf{W}^{\mathsf{T}} \right) \right\}$$
(13)

(See Supplementary Materials for the derivation.)

3.3 Likelihood Computation

196 For model comparison, we are interested in evaluat-197 ing the data log-likelihood $\log p(\mathbf{x}|\theta)$. As noted in Sec. 3, a TA with J groups of factors reduces to a 199 FA when conditioned on J-1 factor groups. Uti-200 lizing the fact that $\log p(\mathbf{x}|\theta)$ can be easily computed 201 (Eq. 2), a Monte Carlo estimation of data log-likelihood 202 in TA can be performed by sampling from the prior of 203 the J-1 groups of factors. For example, in a model $TA\{D, d_1, d_2\}, J = 2:$ 204

Algorithm 1 EM Learning for TA

- 1: Given training data with N samples: $\mathbf{X} \in \mathbb{R}^{D \times N}$
- 2: Initialize θ : {**W**, **T**₍₁₎} ~ $\mathcal{N}(0, .01^2)$, $\Psi \leftarrow 10 * \text{stddev}(\mathbf{X})$.

repeat

//Approximate E-step:

for n = 1 to N do 3: Sample $\{\mathbf{z}_1^{(n)}, \mathbf{z}_2^{(n)}, \mathbf{z}_3^{(n)}\}$ from $p(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3 | \mathbf{x}^{(n)})$ using Eq. 9, and alternating between $\mathbf{z}_1, \mathbf{z}_2, \& \mathbf{z}_3$. end for

//*M*-step:

using samples:
$$E[\mathbf{y}_n] \simeq \mathbf{y}^n$$
,
 $E[\mathbf{u}_n \mathbf{y}_n^{\mathsf{T}}] \simeq \mathbf{u}^n \mathbf{y}^{n\mathsf{T}}$, etc.

6: Update $\{\mathbf{W}, \mathbf{T}_{(1)}, \Psi\}$ using Eqs. 11, 12, and 13. **until** convergence

$$\log p(\mathbf{x}) = \log \int_{\mathbf{z}_2} p(\mathbf{x}|\mathbf{z}_2) p(\mathbf{z}_2) d\mathbf{z}_2 \simeq \log \frac{1}{K} \sum_{k=1}^K p(\mathbf{x}|\mathbf{z}_2^{(k)}), \quad \mathbf{z}_2^{(k)} \sim \mathcal{N}(0, I)$$
(14)

This simple estimator is asymptotically unbiased but has high variance unless the dimensionality of z_2 , or d_2 , is very small. Since z_1 can be analytically integrated out, the Monte Carlo technique can be accurate when only one factor group has large dimensionality.

212 AIS estimation of likelihood

For large d_j , however, simple Monte Carlo estimation is very inefficient as random samples from the prior will mostly yield close to zero log-likelihoods, therefore giving an estimator with large variance. In this situation, Annealed Importance Sampling [13] is a much better alternative. We can treat the problem of estimating $\log p(\mathbf{x})$ as calculating the partition function of unnormalized



Figure 2: TA vs. FA on 2D synthetic datasets. Training data log-likelihood of each model is in parenthesis. TA better models more complex densities.

posterior distribution $p^*(\mathbf{z}|\mathbf{x}) \triangleq p(\mathbf{x}, \mathbf{z})$, where $p^*(\cdot)$ denotes an unnormalized distribution. The basic Importance Sampling gives:

$$p(\mathbf{x}) = \int_{\mathbf{z}} d\mathbf{z} \; \frac{p^*(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} q(\mathbf{z}) \simeq \frac{1}{M} \sum_{i}^{M} w^{(i)}, \text{ where } w^{(i)} = \frac{p^*(\mathbf{z}^{(i)}|\mathbf{x})}{q(\mathbf{z}^{(i)})}, \quad z^{(i)} \sim q(\mathbf{z})$$
(15)

AIS provides a better estimate by first sampling from a tractable base distribution q(z). Subsequent MCMC steps are taken in a set of intermediate distribution, annealing to the distribution of interest: $p(\mathbf{z}|\mathbf{x})$. Annealing allows for a much better estimate of $w^{(i)}$. For TAs, we assume the base distribution is the prior over the factors: $q(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3) = \prod_{i=1}^{3} p(\mathbf{z}_i)$. An intermediate distribution is:

$$p_{\beta}(\{\mathbf{z}_j\}) \propto q(\{\mathbf{z}_j\})^{1-\beta} p^*(\{\mathbf{z}_j\} | \mathbf{x})^{\beta} = p(\{\mathbf{z}_j\}) p^{\beta}(\mathbf{x} | \{\mathbf{z}_j\}),$$
(16)

where β is a scalar which varies from 0.0 to 1.0, as we anneal from the prior to the posterior. Derivations and experiments with the AIS estimator is provided in Supplementary Materials.

3.4 Equality Constraints

An equality constraint indicates that a set of observed vectors $\{\mathbf{x}^{(k)}\}_{k=1}^{K}$ have the same factor values for the *j*-th factor group \mathbf{z}_j . During learning, the availability of equility constraints will only change the inference step. Assuming we have constraints for the factor group j = 1, the posterior for z_i (Eq. 9) will be modified as follows:

$$p(\mathbf{z}_1|\{\mathbf{x}^{(k)}\},\{\mathbf{z}_2^{(k)}\},\{\mathbf{z}_3^{(k)}\}) = \mathcal{N}(\mathbf{z}_1|\tilde{\mathbf{V}}^{-1}\sum_{k=1}^{K}\{\mathbf{\Lambda}^{(k)\mathsf{T}}\mathbf{\Psi}^{-1}(\mathbf{x}^{(k)}-\boldsymbol{\mu}^{(k)})\},\tilde{\mathbf{V}}^{-1})$$
(17)

$$\tilde{\mathbf{V}} = \mathbf{I} + \sum_{k=1}^{K} \mathbf{\Lambda}^{(k)\mathsf{T}} \boldsymbol{\Psi}^{-1} \mathbf{\Lambda}^{(k)}, \qquad \mathbf{\Lambda}^{(k)} = \mathbf{W}_1 + \boldsymbol{\mathfrak{T}} \,\bar{\mathbf{x}}_3 \, \mathbf{z}_3^{(k)} \,\bar{\mathbf{x}}_2 \, \mathbf{z}_2^{(k)}, \qquad \boldsymbol{\mu}^{(k)} = \mathbf{m} + \mathbf{W}_2 \mathbf{z}_2^{(k)} + \mathbf{W}_3 \mathbf{z}_3^{(k)}$$

The M-step is not affected by the presence of equality constraints so TAs can learn when equality constraints are provided for arbitrary subsets of the data.

3.5 Mixture of Tensor Analyzers

Extending TAs to Mixture of Tensor Analyzers (MTAs) is straightforward. Each component c will have its own parameters $\theta_c = \{\mathbf{W}_c, \mathbf{T}_{(1),c}, \Psi_c\}$. The data likelihood is marginalized over the C components: $p(\mathbf{x}) = \sum_{c=1}^{C} p(\mathbf{x}|c)p(c)$. Posterior distribution over the factors and components can be decomposed as:

$$p(\{\mathbf{z}_j\}, c | \mathbf{x}) = p(\{\mathbf{z}_j\} | \mathbf{x}, c) p(c | \mathbf{x})$$
(18)

where $p({\mathbf{z}_i} | \mathbf{x}, c)$ can be sampled using Eq. 9 with θ_c and $p(c|\mathbf{x}) \propto p(\mathbf{x}|c)p(c)$. Sec. 3.3 showed how $p(\mathbf{x}|c)$ can be efficiently approximated, thereby making it feasible to train MTAs.

Experiments

We demonstrate the usefulness of TA on 2 synthetic and 3 real-life datasets. We plan to release source code with implementation details of the experiments in the future.

4.1 Synthetic Data

As a proof of concept, we compared TA to FA on two synthetic datasets (Fig. 2 A & B). Data A is highly structured and is generated using a TA with random parameters. Data B has high kurtosis, with density concentrated at the origin. For both datasets, we learned using a TA{ $D = 2, d_1 =$ $2, d_2 = 2$ and a FA with the same number of parameters as the TA until convergence. The TAs



Figure 3: Learning MTA on natural image patches. (a) Training data. (b) Samples from MTA look realistic. (c) Average test log-likelihood comparisons. ICA: Independent Component Analysis [1]. GRBM: Gaussian Restricted Boltzmann Machine [7]. (d) Each row contains filters from a different MTA component.

performed model recovery nicely. The left panel of Fig. 2(a) plots training points. The data loglikelihood of the true model is -2.58 nats. The middle panel plots the samples of a TA which achieved -2.62 nats on the training data. The right panel plots samples drawn from a FA. Likewise in Fig. 2(b), TA is a significantly better model than the FA: -2.04 ± 0.05 to -2.48 ± 0.09 nats. We also tested mixtures of TAs vs mixtures of FA (MFA) on data generated by randomly initialized MFA models. The performance of MTA and MFA were very similar, demonstrating that (M)TAs can also efficiently emulate (M)FAs when necessary.

288 289 4.2 Natural Images

279

280

Learning a good density model of natural images is useful for image denoising and inpainting. 290 We compared MTAs to MFAs and other models on image patches. 100,000 8×8 patches were 291 extracted from the training set of the Berkeley Natural Images database [12] for training, while 292 20,000 patches from the test set were extracted for testing. PCA is performed to preprocess the 293 data to 30 dimensions, preserving 98% of the variance. For both MTA and MFA², we used 10 294 mixture components and selected the number of factors per component using a validation set. For 295 MFA, 40 factors per component were used, resulting in an overcomplete representation. For MTA, 296 each component is a TA $\{30,8,2\}$. Convergence is achieved when the objective does not improve 297 by more than 0.01 percent. The average test log-likelihood of MTA is better than that of the MFA by **0.8 nats**³. Fig. 3(b) shows that samples of a MTA matches closely to the training patches in 298 Fig. 3(a). Each row of Fig. 3(d) shows filters (fibers of the tensor) of one of the MTA components. 299 Components of MTA specialize to model patches of different spatial frequencies and orientations. 300

4.3 Concept Learning with equality constraints

302 Human intelligence is characterized by the ability to form abstract concepts which allows for gen-303 eralization across widely differing percepts. The color concept of "red" is perceived from both the 304 images of an small red apple and a red car. Concepts can be learned using TAs. We use images of 305 synthetic shapes for training. There are 4 different shapes: circle, triangle, square, and pentagon; 306 5 different colors: red, green, blue, yellow, and purple; 3 different sizes: small, medium, and big. 307 In total, there are 60 images of resolution 24×24 with 3 color channels per image. We learned a $TA\{3 \times 24^2, 5, 3, 4\}$ using 30 iterations of EM. 200 alternating Gibbs steps were taken to estimate 308 the posterior during each E-step. Equality constraints were used during training, implicitly assigning 309 meanings to the latent factor groups, i.e. $\mathbf{z}_1 \triangleq \mathbf{z}_{shape}, \mathbf{z}_2 \triangleq \mathbf{z}_{color}, \text{ and } \mathbf{z}_3 \triangleq \mathbf{z}_{size}.$ 310

311 After learning, to probe into the representation learned by the TA, we first infer the latent factors 312 $\{\mathbf{z}_{shape}, \mathbf{z}_{color}, \mathbf{z}_{size}\}\$ for the training shapes by sampling from $p(\mathbf{z}_{shape}, \mathbf{z}_{color}, \mathbf{z}_{size}|\mathbf{x})$. To see 313 if \mathbf{z}_{color} truly represents the color concept, we fix the inferred factors $\{\mathbf{z}_{shape}, \mathbf{z}_{size}\}$ of a small If \mathbf{z}_{color} thuy represents the control concept, we have the interver factors (*z*_{shape}, *s*_{color}) blue circle, while drawing a new \mathbf{z}_{color}^{new} from its standard Normal prior. The top row of Fig. 4(b) shows the generated images from the factors { $\mathbf{z}_{shape}, \mathbf{z}_{color}^{new}, \mathbf{z}_{size}$ }. Old shapes with the same size 314 315 but novel (and more importantly) homogeneous colors are generated. This type of generalization 316 is possible if and only if \mathbf{z}_{color} models color. We further sample from the priors of both the color 317 and shape factors. Generations mix and match novel colors and novel shapes, shown in the middle 318 row of Fig. 4(b). Due to the fact that TA is a multilinear model, it can only linearly combine the 4 319 training shapes to synthesize new shapes. The bottom row of Fig. 4(b) are filters/detectors for the 320 concept red - tuned to red shapes of differing size and shape. They are the fibers of the tensor of the 321 TA conditioned on $\mathbf{z}_{color} = red$.

322 323

²MFA code [18] is downloaded from http://lear.inrialpes.fr/~verbeek/software.php

³The gain is found to be statistically significant using the paired t-test at p = 0.05.

324 325

326

327 328

330

331 332

334

337

333

(a)

Training data

Small Me (c) Test object

Shape Concept

0.8

0.6 0.4

0.2 0 R

> 0. 0.

0.2

0

Color Concept

Pur Blue Grn Yel

Size Concept

Figure 4: Learning color, shape, and size concepts. (a) Data contains 4 shapes, 5 colors, and 3 sizes. (b) Each row demonstrates a particular generalization achieved by the TA. (c) Posterior distribution over the concept 335 mixture of Gaussians. the model is able to correctly activate previously learned concepts for a novel test object. 336 See text for details, best viewed in color and on-screen.

(b) Generalization

The aggregated posterior $\sum_{n=1}^{N} p(\mathbf{z}_{shape}^{(n)}, \mathbf{z}_{color}^{(n)}, \mathbf{z}_{size}^{(n)} | \mathbf{x}^{(n)})$ is tightly clustered in distinct modes. 338 339 Since training data has 5 colors, 5 clusters will form in the space of \mathbf{z}_{color} . Likewise, there are 3 and 340 4 clusters for the size and shape factors. We train three second layer mixture of Gaussians (MoG) 341 models using the samples from the TA posterior as training data, e.g. $\mathbf{z}_{color}^{(n)} \sim p(\mathbf{z}_{color} | \mathbf{x}^{(n)})$. Using 342 prior knowledge, we set the number of components of the "color" MoG to be 5, "size" MoG to be 3, 343 and the "shape" MoG to be 4. Learning using EM and initializing the means of MoG components 344 to small values, the components of the MoGs become semantically meaningful, each representing a separate concept. Using this hierarchical model, we demonstrate strong generalization by using a 345 24×24 image of a toy duck as the observed test image Fig. 4(c). A two stage inference step first 346 infers the factors using the parameters of the TA. We then compute the posterior probabilities over 347 the components of all 3 MoGs, conditioned on the inferred factors from the first stage. Interestingly, 348 our model perceive a big, yellow object which is mostly square but slightly circular.

349 350

4.4 Face Recognition with equality constraints

351 In a one-shot learning setting, classifiers must be trained from only one example per class. For face 352 recognition, only one example per test subject is used for training. We use the Yale B database and 353 its Extended version in this experiment [10]. The database contains 38 subjects under 45 different 354 lighting conditions. We use 28 subjects for training and test on the 10 subjects from the original 355 Yale B database. The images are first downsampled to 24×24 and we used a TA{576, 80, 4}, which 356 contains 2 groups of factors. Equality constraints specifying which images have the same identity 357 or lighting type are used during training.

358 The learned TA allows for strong generalization to new people under new lighting conditions. It 359 achieved an average log-likelihood of 836 ± 7 nats on the images of the 10 held-out subjects. As 360 a comparison, the best MFA model achieved only 791 ± 10 nats. The number of components and 361 factors of the MFA are optimally selected using grid search. The gain of 45 nats demonstrates a 362 significant win for the TA. Qualitatively, to see how well the TA is able to factor out identity from 363 lighting, we first sample from the posterior distribution conditioned on a single test image, exactly as in Sec. 4.3. One factor group's activation is fixed, while we sample the other group of factors 364 from its standard Normal prior. Results are in Fig. 5. A row in panel (b) shows the same person 365 under different lighting. A row in panel (c) shows a different person but under the same lighting. 366 We emphasize that only a *single* test image from *novel* subjects is used for inference (panel (a)). 367

368 For the recognition task, we first use the 28 training subjects to learn the parameters for the 369 TA{576,80,4} in the training phase. During the testing phase, a single image (under frontal lighting) for each of the 10 test subjects is used to compute the TA's posterior mean using 200 alternating 370 Gibbs steps. For each one of the 10 $z_{identity}$ factors, we use it to collapse the TA into a FA (See 371 Eqs. 7, 8). Each FA has 4 factors and is equivalent to a rank 4 Gaussian in pixel space. During 372 testing, an image is classified to be the same class as the FA which assigns it the highest likelihood. 373

374 We compared TA with Nearest Neighbor (NN), normalized Cross Correlation (CC), Support Vector 375 Machines (SVM), S&C bilinear model, Factor Analyzers, and Human. The CC method first subtract the image mean from all pixels of the image. The image is then normalized to have unit norm. 376 The cosine similarity between test and training images is used for classification. Multiclass linear 377 SVM from the LIBLINEAR package [4] is used, where the hyperparameter C was chosen using



Figure 5: Tensor Analyzer is able to simultaneously decompose a test image into separate identity and lighting factors. (a) 4 test images with frontal, left, right and top light sources. (b) Random samples with identity factor fixed to the inferred values from the test faces in (a). (c) Random samples with lighting factor fixed to the inferred values from (a). (d) Factor loading transferring in a FA model creates artifacts. (e) Recognition error on the one-shot face recognition task.

validation. For the S&C bilinear method, we implemented the exact algorithm as stated in Sec. 394 3.2 of [15]. During one-shot learning, adapting the style and content codes is performed using the 395 algorithm described in Sec. 6.2 of [15]. Recognition using the FA method requires the transfer of 396 the factor loadings from the training phase. We first learned 28 FAs, one for each training subject. 397 During the testing phase, each one of the 10 training images of the test subject is matched with 398 the FA which gives it the largest likelihood. A new FA is created, centered at the training image. 399 The rest of the parameters are transfered from the matched FA. In essence, the transfered loadings 400 models the lighting variations of the training subjects (1 of 28) which is most similar to the test phase 401 training image. After the creation of these 10 new FAs, classification is same as in the TA method. 402 Human error is the average of testing several human subjects on the same one-shot recognition task. 403 Recognition errors are listed in Fig. 5(e). TA outperforms the rest of the field by a significant margin 404 and approaches human performance.

FA do not work as well as TA because the transferred factor loadings are not accurate. Fig. 5(d) shows random samples from a transfered FA of a test subject. This is due to the fact that lighting variations of faces must be a *function* of identity, and can not be simply transfered additively.

409 4.5 Learning with incomplete equality constraints

410 We demonstrate the advantage of the TA in a semisupervised setting on the UMIST face database [5]. It 411 contains 20 subjects with 20 to 40 training images per 412 subject. The variation consists of in-depth head rota-413 tions. We compared the TA to S&C model and NN 414 on one-shot face recognition. Out of the 20 subjects, 415 15 were used for training and 5 for testing. The split 416 was randomized over 10 different trials. The images are 417 downsampled to the resolution of 24×24 . We exper-418 imented with using equality constraints and 3, 4, or 5 419 images per training subject. During the training phase, 420 a TA{576,3,5} model was trained using 30 EM itera-



Figure 6: UMIST Recognition errors.

tions. The algorithms for classification are exactly the same as in Sec. 4.4. Fig. 6 plots recognition
errors as a function of the number of images per subject used during the *training phase*. For the case
where 5 images per training subjects are available, TA achieves significantly lower errors at 12.4%
compared to 24.9% for S&C. If we add an equal number of images without equality constraints
during training, the error is further reduced to 10.9%.

426 427

389

390

391

392

393

5 Conclusions

We have introduced a new density model which extends Factor Analysis to modeling multilinear in teractions. Using efficient alternating sampling and the EM algorithm, we have shown that (M)TAs
 can learn more complex densities, better model natural image patches, and separate factors of vari ation, leading to the learning of simple concepts. Moreover, at the important task of one-shot face
 recognition, TAs outperform a variety of other models.

432	Ref	erences
433	[1]	A Ball and T.I. Sainowski. An information maximization approach to blind separation and blind decon
434	[1]	volution. <i>Neural Computation</i> , 7:1129–1159, 1995.
435	[2]	I Carroll and Jib I Chang. Analysis of individual differences in multidimensional scaling via an n-way.
436	[2]	generalization of 'Eckart-Young' decomposition. <i>Psychometrika</i> , 35(3):283–319, September 1970.
438	[3]	Benjamin J. Culpepper, Jascha Sohl-Dickstein, and Bruno A. Olshausen. Building a better probabilistic model of images by factorization. In <i>ICCV</i> , pages 2011–2017, 2011.
439	[4]	Rong-En Fan Kai-Wei Chang Cho-Jui Hsieh Xiang-Rui Wang and Chih-Jen Lin LIBI INFAR: A
440 441	נדי	library for large linear classification, August 2008.
442 443	[5]	Daniel B Graham and Nigel M Allinson. Characterizing virtual eigensignatures for general purpose face recognition. In <i>Face Recognition: From Theory to Applications</i> , pages 446–456. 1998.
444	[6]	David B. Grimes and Rajesh P. Rao. Bilinear sparse coding for invariant vision. <i>Neural computation</i> , 17(1):47–73, January 2005.
445	[7]	G. E. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. <i>Science</i> , 313:504–507, 2006.
447 448	[8]	Tamara G. Kolda and Brett W. Bader. Tensor decompositions and applications. <i>SIAM Review</i> , 51(3):455–500, 2009
449	101	Lieven De Lathauwer and loos Vandewalle. Dimensionality reduction in higher order signal processing
450 451	[2]	and rank- $(r_1, r_2,, r_n)$ reduction in multilinear algebra, 2004.
452	[10]	K. C. Lee, Jeffrey Ho, and David Kriegman. Acquiring linear subspaces for face recognition under variable lighting. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 27:684–698, 2005.
453		
454	[11]	Haiping Lu, K. N. Plataniotis, and A. N. Venetsanopoulos. Multilinear principal component analysis of
455		tensor objects for recognition. In in Proc. Int. Conf. on Pattern Recognition, pages 776-779, 2006.
456	[12]	D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its
457		application to evaluating segmentation algorithms and measuring ecological statistics. In <i>Proc. 8th Int'l</i>
458		Conf. Computer Vision, volume 2, pages 416–423, July 2001.
459	[13]	R. M. Neal. Annealed importance sampling. <i>Statistics and Computing</i> , 11:125–139, 2001.
460	[14]	Donald B. Rubin and Dorothy T. Thayer. EM algorithms for the ML factor analysis. <i>Psychometrika</i> , 47(1):69–76, 1982.
461	[15]	Joshua B. Tenenbaum and William T. Freeman. Separating style and content with bilinear models. <i>Neural</i>
462	[]	<i>Computation</i> , pages 1247–1283, 2000.
464 465	[16]	L. R. Tucker. Implications of factor analysis of three-way matrices for measurement of change. In C. W. Harris, editor, <i>Problems in measuring change.</i> , pages 122–137. University of Wisconsin Press, Madison WI, 1963.
466 467	[17]	M. A. O. Vasilescu and D. Terzopoulos. Multilinear analysis of image ensembles: Tensorfaces. In <i>ECCV</i> , pages 447–460, 2002
468	F101	Jakeh Varheelt. Learning nonlineer image manifolds hy clobal clianment of least linear models.
469	[10]	Trans. Pattern Analysis and Machine Intelligence, 28:14, 2006.
470	[10]	Hongshang Wang and Narandra Abuja Easial expression decomposition. In ICCV pages 058, 065, 2003
471	[19]	Ministry wang and twatehold a Anuja. Factal expression decomposition. In <i>PCCV</i> , pages 936–905, 2005.
472	[20]	Ming-Hsuan Yang, Narendra Anuja, and David Kriegman. Face detection using a mixture of factor analyzers. In <i>IEEE International Conference on Image Processing (ICIP)</i> Kobe, Japan 1999
473		analyzers. In the I meriational Conjerence on Image Processing (1011), Robe, Japan, 1999.
474		
475		
476		
477		
478		
479		
480		
481		
482		
483		
484		