

A Unifying Viewpoint of some Clustering Techniques Using Bregman Divergences and Extensions to Mixed Data Sets

Cécile Levasseur, Brandon Burdge, Ken Kreutz-Delgado
University of California, San Diego
Jacobs School of Engineering
La Jolla, California, USA
{clevasseur, bburdge, kreutz}@ucsd.edu

Uwe F. Mayer
University of Utah
Department of Mathematics
Salt Lake City, Utah, USA
mayer@math.utah.edu

Abstract

We present a general viewpoint using Bregman divergences and exponential family properties that contains as special cases the three following algorithms: 1) exponential family Principal Component Analysis (exponential PCA), 2) Semi-Parametric exponential family Principal Component Analysis (SP-PCA) and 3) Bregman soft clustering. This framework is equivalent to a mixed data-type hierarchical Bayes graphical model assumption with latent variables constrained to a low-dimensional parameter subspace. We show that within this framework exponential PCA and SP-PCA are similar to the Bregman soft clustering technique with the addition of a linear constraint in the parameter space. We implement the resulting modifications to SP-PCA and Bregman soft clustering for mixed (continuous and/or discrete) data sets, and add a nonparametric estimation of the point-mass probabilities to exponential PCA. Finally, we compare the relative performances of the three algorithms in a clustering setting for mixed data sets.

1. Introduction

We present a general point of view that relates the exponential family Principal Component Analysis (exponential PCA) technique of [7] to the Semi-Parametric Principal Component Analysis (SP-PCA) technique of [20] and to the Bregman soft clustering method presented in [4]. The proposed standpoint is then illustrated with a clustering problem in mixed data sets.

The three techniques considered here all utilize Bregman divergences and can all be explained within a single hierarchical Bayes graphical model framework shown in Figure 1. They are not separate unrelated algorithms but different manifestations of parameter choices taken within a common

framework. The proposed model is mathematically equivalent to equation (6) and we will demonstrate that various parametric choices symbolically shown in Figure 2 determine the three algorithms. Because of this insight, we will readily extend the algorithms to deal with the important mixed data type case.

There are two ways in which Bregman divergences are important. First, they generalize the squared Euclidean distance to a class of distances that all share similar properties. Second, there exists a bijection between Bregman divergences and exponential family distributions [3, 4]. Recently, researchers have shown that many important algorithms can be generalized from Euclidean metrics to distances defined by a Bregman divergence [7, 20, 4, 6], i.e., from Gaussian distributed data components to data components distributed according to an exponential family, such as binary- or integer-valued.

Data mining techniques seek potentially useful information contained in complex data sets. Complexity of these data generally comes from a high number of components and also from the fact that the components usually are of mixed data types (categorical, count, continuous, etc.). In order to address the latter, techniques should allow for the components to have different parametric forms by using the large range of exponential family distributions and their associated Bregman divergences. While a modification for mixed data sets was presented in [12] for exponential PCA, we implement here a modification for SP-PCA and Bregman soft clustering. The issue of a high number of components of the original data is addressed by exponential PCA and SP-PCA through dimensionality reduction by projecting the data to a low-dimensional parameter subspace.

Finally, we consider synthetic data examples of mixed types. We demonstrate that exponential PCA, with the addition of a nonparametric estimation tool, rivals SP-PCA and Bregman soft clustering in terms of clustering performance.

2. Theoretical background

To motivate theoretical developments, the hierarchical Bayes graphical model for hidden or latent variables shown in Figure 1 is considered.

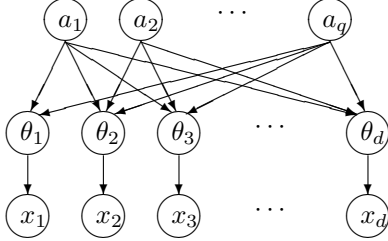


Figure 1. Graphical model for our framework.

The row vector $\mathbf{x} = [x_1, \dots, x_d] \in \mathbb{R}^d$ consists of observed features in a d -dimensional space. It is assumed that training points can be drawn from populations having class-conditional probability density functions

$$p(\mathbf{x}|\boldsymbol{\theta}) = p_1(x_1|\theta_1) \cdot \dots \cdot p_d(x_d|\theta_d), \quad (1)$$

where, when conditioned on the random parameter vector $\boldsymbol{\theta} = [\theta_1, \dots, \theta_d] \in \mathbb{R}^d$, the components of \mathbf{x} are independent. It is further assumed that $\boldsymbol{\theta}$ can be written as

$$\boldsymbol{\theta} = \mathbf{a}\mathbf{V} + \mathbf{b}$$

with the hidden or latent variable $\mathbf{a} = [a_1, \dots, a_q] \in \mathbb{R}^q$ random and unknown with $q < d$ (and ideally $q \ll d$), $\mathbf{V} \in \mathbb{R}^{q \times d}$ and $\mathbf{b} \in \mathbb{R}^d$ deterministic and unknown. The latent variable \mathbf{a} in some way explains part (or all) of the random behavior of the observed variables. The subscript i on $p_i(\cdot|\cdot)$ serves to indicate that the marginal densities can all be different, allowing for the possibility of \mathbf{x} containing categorical, discrete, and continuous valued components. Also, the marginal densities are each assumed to be one-parameter exponential family densities, and θ_i is taken to be the natural parameter (or some simple bijective function of it) of the exponential family density p_i . Hence, each component density $p_i(x_i|\theta_i)$ in (1) for $x_i \in \mathcal{X}_i, i = 1, \dots, d$, is of the form

$$p(x_i|\theta_i) = \exp(\theta_i x_i - G(\theta_i)), \quad (2)$$

where $G(\cdot)$ is the cumulant generating function defined as

$$G(\theta_i) = \log \int_{\mathcal{X}_i} \exp(\theta_i x_i) \nu(dx_i), \quad (3)$$

with $\nu(\cdot)$ a σ -finite measure that generates the exponential family. It can be shown, using Fubini's theorem [11], that $G(\boldsymbol{\theta}) = \sum_{i=1}^d G(\theta_i)$.

The maximum likelihood identification of the blind random effect model

$$p(\mathbf{x}) = \int p(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta} = \int \prod_{i=1}^d p_i(x_i|\theta_i)\pi(\boldsymbol{\theta})d\boldsymbol{\theta}, \quad (4)$$

where $\pi(\boldsymbol{\theta})$ is the probability density function of $\boldsymbol{\theta}$, is quite a difficult problem. It corresponds to identifying $\pi(\boldsymbol{\theta})$, which, under the condition $\boldsymbol{\theta} = \mathbf{a}\mathbf{V} + \mathbf{b}$, corresponds to identifying the matrix \mathbf{V} , the vector \mathbf{b} , and a density function on the random effect \mathbf{a} via a maximization of the likelihood function $p(\mathbf{X})$ with respect to \mathbf{V} , \mathbf{b} , and the random effect density function, where

$$\begin{aligned} p(\mathbf{X}) &= \prod_{k=1}^n p(\mathbf{x}[k]) = \prod_{k=1}^n \int p(\mathbf{x}[k]|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta} \\ &= \prod_{k=1}^n \int \prod_{i=1}^d p_i(x_i[k]|\theta_i)\pi(\boldsymbol{\theta})d\boldsymbol{\theta}, \end{aligned} \quad (5)$$

and \mathbf{X} is the $(n \times d)$ observation matrix

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}[1] \\ \mathbf{x}[2] \\ \vdots \\ \mathbf{x}[n] \end{pmatrix} = \begin{pmatrix} x_1[1] & \dots & x_d[1] \\ x_1[2] & \dots & x_d[2] \\ \vdots & \ddots & \vdots \\ x_1[n] & \dots & x_d[n] \end{pmatrix}.$$

This difficulty can be avoided by NonParametric Maximum Likelihood (NPML) estimation of the random effect distribution, concurrently with the structural model parameters. The NPML estimate is known to be a discrete distribution on a finite number of support points or ‘‘atoms’’ [9, 10, 13]. Finding the NPML estimate is widely regarded as computationally intensive, the particular difficulty being the location of the atoms [1].

As shown in [12], with $\boldsymbol{\theta} = \mathbf{a}\mathbf{V} + \mathbf{b}$, with \mathbf{V} , \mathbf{b} fixed and \mathbf{a} random, the single-sample likelihood (4) is equal to

$$p(\mathbf{x}) = \sum_{l=1}^m p(\mathbf{x}|\boldsymbol{\theta}[l])\pi_l = \sum_{l=1}^m p(\mathbf{x}|\mathbf{a}[l]\mathbf{V} + \mathbf{b})\pi_l$$

and the data likelihood (5) is equal to

$$p(\mathbf{X}) = \prod_{k=1}^n \sum_{l=1}^m p(\mathbf{x}[k]|\boldsymbol{\theta}[l])\pi_l = \prod_{k=1}^n \sum_{l=1}^m p(\mathbf{x}[k]|\mathbf{a}[l]\mathbf{V} + \mathbf{b})\pi_l, \quad (6)$$

with point-mass probability estimates π_l , unknown point-mass support points $\mathbf{a}[l]$, and the linear predictor $\boldsymbol{\theta}[l] = \mathbf{a}[l]\mathbf{V} + \mathbf{b}$ in the l th mixture component, $l = 1, \dots, m$.

The data likelihood is thus approximately the likelihood of a finite mixture of exponential family densities with unknown mixture proportions or point-mass probability estimates π_l and unknown point-mass support points $\mathbf{a}[l]$,

with the linear predictor $\theta[l] = \mathbf{a}[l]\mathbf{V} + \mathbf{b}$ in the l th mixture component [2]. The combined problem of Maximum Likelihood Estimation (MLE) of the parameters \mathbf{V} , \mathbf{b} , the point-mass support points (atoms) $\mathbf{a}[l]$ and the point-mass probability estimates $\pi_l, l = 1, \dots, m$, is known as the Semiparametric Maximum Likelihood mixture density Estimation (SMLE) problem [14, 16]. It can be attacked by using the Expectation-Maximization (EM) algorithm [8, 10, 13, 19, 15, 2, 5, 16, 18], as done in particular in the Semi-Parametric Principal Component Analysis (SP-PCA) technique proposed in [20] and discussed below. Note that, historically, [10] appears to be generally acknowledged as the first paper that proposed the EM algorithm for NPML estimation in the mixture density context; then, [13] improved upon the theoretical foundations of the NPML estimation approach and later [15] further explored some of the fundamental issues raised in [13].

However, the SMLE problem can also be attacked by simply considering the special case of uniform point-mass probabilities, i.e., $\pi_l = 1/m$ for $l = 1, \dots, m$, for which the number of support points equals the number of data samples, i.e., $m = n$. It was demonstrated in [12] that this special uniform case corresponds to exponential PCA.

3. A unifying framework

Within the proposed hierarchical Bayes graphical model framework, exponential PCA, SP-PCA and Bregman soft clustering are not separate uncorrelated algorithms but different manifestations of parameter choices.

Figure 2 considers the number of atoms as a common characteristic for comparison purposes. It symbolizes how various parametric choices determine the three algorithms. Whereas the exponential PCA approach requires a number of atoms equal to the number of data points and hence can be seen as an extreme case of the NPML technique, SP-PCA deals with a smaller number of atoms. Finally, the Bregman soft clustering approach considers an even smaller number of atoms, viewed as cluster centers, since its primary goal is clustering. Furthermore, both exponential PCA and SP-PCA impose a low-dimensional (unknown) latent variable subspace in their structure. However, Bregman soft clustering does not impose this lower dimensional constraint and hence can be seen as a degenerate case.

It becomes clear while looking at Table 1, Table 2 and Table 3 shown below that both SP-PCA and Bregman soft clustering utilize the EM algorithm for estimation purposes whereas exponential PCA does not. This is because exponential PCA assumes uniform point-mass probabilities and does not need to estimate them. Both exponential PCA and SP-PCA impose the low-dimensional parameter subspace

constraint, hence the need for the Newton-Raphson iterative algorithm to find it as discussed below.

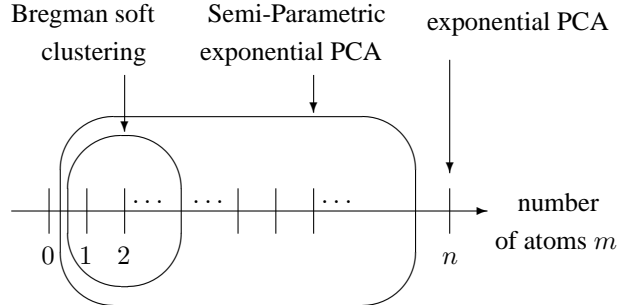


Figure 2. General point of view on SP-PCA, exponential PCA and Bregman soft clustering based on the number of NPML atoms.

4. Extensions to mixed data sets

Now that the relationship between exponential PCA, SP-PCA and Bregman soft clustering has been demonstrated within a unifying mixed data type framework, we can extend SP-PCA and Bregman soft clustering to deal with the important mixed data case. Recalling the modifications on exponential PCA presented in [12], we also introduce a penalty function that improves exponential PCA convergence efficiency. Towards these ends, following the notation in [5], we use the *mixing distribution* $\mathcal{Q} = \{\theta[l], \pi_l\}_{l=1}^m$ which contains the parameters $\theta[l]$ and their associated point-mass probabilities π_l . The mixing distribution needs to be estimated in all three algorithms.

4.1. Semi-parametric exponential PCA

The Semi-Parametric exponential family Principal Component Analysis (SP-PCA) approach presented in [20] attacks the Semiparametric Maximum Likelihood mixture density Estimation (SMLE) problem by using the Expectation-Maximization (EM) algorithm [8]. As done previously in [12] for exponential PCA, we present an SP-PCA modified approach for mixed data types. For simplicity of presentation, we consider that the f first attributes are distributed according to the exponential family distribution $p^{(1)}$ and the $(d - f)$ last attributes are distributed according to the exponential family distribution $p^{(2)}$. The following notation is used:

$$\begin{aligned} \mathbf{x}[k] &= [x_1[k], \dots, x_f[k], x_{f+1}[k], \dots, x_d[k]] \\ &= [\mathbf{x}^{(1)}[k] | \mathbf{x}^{(2)}[k]], \end{aligned}$$

for $k = 1, \dots, n$, and similarly for the observation matrix,

$$\mathbf{X} = (\mathbf{X}^{(1)} | \mathbf{X}^{(2)}).$$

The EM approach introduces a *missing* (unobserved) variable $\mathbf{z}_k = [z_{k1}, \dots, z_{km}]$, for $k = 1, \dots, n$. This variable is an m -dimensional binary vector whose l th component equals 1 if the observed variable $\mathbf{x}[k]$ was drawn from the l th mixture component and 0 otherwise; its value is estimated during the E-step. Using this information, a *complete* log-likelihood function is defined as follows:

$$L^{(c)}(\mathcal{Q}, \{\mathbf{z}_k\}_{k=1}^n) = \log \prod_{k=1}^n \prod_{l=1}^m p(\mathbf{x}[k]|\boldsymbol{\theta}[l])^{z_{kl}} \pi_l^{z_{kl}},$$

with $\boldsymbol{\theta}[l] = \mathbf{a}[l]\mathbf{V} + \mathbf{b}$. Its maximization during the M-step yields parameters \mathbf{A} , \mathbf{V} and \mathbf{b} estimates. Then,

$$p(\mathbf{x}[k]|\boldsymbol{\theta}[l]) = p^{(1)}(\mathbf{x}^{(1)}[k]|\boldsymbol{\theta}^{(1)}[l]) \cdot p^{(2)}(\mathbf{x}^{(2)}[k]|\boldsymbol{\theta}^{(2)}[l]) \quad (7)$$

using the assumption from (1), where, for $l = 1, \dots, m$,

$$\boldsymbol{\theta}[l] = [\theta_1[l], \dots, \theta_f[l], \theta_{f+1}[l], \dots, \theta_d[l]] = [\boldsymbol{\theta}^{(1)}[l] | \boldsymbol{\theta}^{(2)}[l]],$$

$$\boldsymbol{\Theta} = \begin{pmatrix} \boldsymbol{\theta}[1] \\ \boldsymbol{\theta}[2] \\ \vdots \\ \boldsymbol{\theta}[n] \end{pmatrix} = (\boldsymbol{\Theta}^{(1)} | \boldsymbol{\Theta}^{(2)}).$$

The following decompositions arise:

$$\mathbf{V} = (\mathbf{V}^{(1)} | \mathbf{V}^{(2)}),$$

$$\mathbf{B} = (\mathbf{B}^{(1)} | \mathbf{B}^{(2)}),$$

where $\mathbf{B}^{(1)} = [\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(1)}]^T$ and $\mathbf{b}^{(1)} = [b_1, \dots, b_f]$, $\mathbf{B}^{(2)} = [\mathbf{b}^{(2)}, \dots, \mathbf{b}^{(2)}]^T$ and $\mathbf{b}^{(2)} = [b_{f+1}, \dots, b_d]$. Hence,

$$\boldsymbol{\Theta} = (\mathbf{A}\mathbf{V}^{(1)} + \mathbf{B}^{(1)} | \mathbf{A}\mathbf{V}^{(2)} + \mathbf{B}^{(2)}).$$

Note that there is no such split for \mathbf{A} . The complete log-likelihood function becomes:

$$\begin{aligned} L^{(c)}(\mathcal{Q}, \{\mathbf{z}_k\}_{k=1}^n) &= \sum_{k=1}^n \sum_{l=1}^m z_{kl} \log \pi_l \\ &+ \sum_{k=1}^n \sum_{l=1}^m z_{kl} \log p^{(1)}(\mathbf{x}^{(1)}[k]|\boldsymbol{\theta}^{(1)}[l]) \cdot p^{(2)}(\mathbf{x}^{(2)}[k]|\boldsymbol{\theta}^{(2)}[l]). \end{aligned} \quad (8)$$

The E-step yields for $k = 1, \dots, n$ and $l = 1, \dots, m$:

$$\begin{aligned} \hat{z}_{kl} &= \mathbb{E}\{z_{kl}|\mathbf{x}[k], \pi_1, \dots, \pi_m\} \\ &= \frac{p^{(1)}(\mathbf{x}^{(1)}[k]|\boldsymbol{\theta}^{(1)}[l]) \cdot p^{(2)}(\mathbf{x}^{(2)}[k]|\boldsymbol{\theta}^{(2)}[l]) \pi_l}{\sum_{r=1}^m p^{(1)}(\mathbf{x}^{(1)}[k]|\boldsymbol{\theta}^{(1)}[r]) \cdot p^{(2)}(\mathbf{x}^{(2)}[k]|\boldsymbol{\theta}^{(2)}[r]) \pi_r}. \end{aligned}$$

For all l and all k , each data point $\mathbf{x}[k]$ has an estimated probability \hat{z}_{kl} of belonging to the l th mixture component.

The M-step first yields the estimates for the point-mass probabilities:

$$\pi_l^{(new)} = \frac{\sum_{k=1}^n \hat{z}_{kl}}{n}.$$

The second part of the M-step, i.e., the estimation of the parameters \mathbf{V} , \mathbf{b} , and the point-mass support points $\mathbf{A} = [\mathbf{a}[1]^T, \dots, \mathbf{a}[m]^T]^T \in \mathbb{R}^{m,q}$, is affected by the mixed data type assumption. It consists of maximizing the complete log-likelihood function (8) with respect to these parameters:

$$\arg \max_{\mathbf{A}, \mathbf{V}, \mathbf{b}} \mathbb{E} \left\{ L^{(c)}(\{\boldsymbol{\theta}[l], \pi_l^{(new)}\}_{l=1}^m, \{\hat{\mathbf{z}}_k\}_{k=1}^n) \right\}.$$

Following the notation in [20], we set, for $l = 1, \dots, m$:

$$\tilde{\mathbf{x}}[l] = \frac{\sum_{k=1}^n \hat{z}_{kl} \mathbf{x}[k]}{\sum_{k=1}^n \hat{z}_{kl}},$$

the l th mixture component center. It can be shown, using exponential family properties, that the loss function is:

$$\begin{aligned} L(\mathbf{A}, \mathbf{V}, \mathbf{b}) &= \sum_{l=1}^m \left\{ G^{(1)}(\mathbf{a}[l]\mathbf{V}^{(1)} + \mathbf{b}^{(1)}) - (\mathbf{a}[l]\mathbf{V}^{(1)} + \mathbf{b}^{(1)}) \tilde{\mathbf{x}}[l]^T \right\} \\ &+ \sum_{l=1}^m \left\{ G^{(2)}(\mathbf{a}[l]\mathbf{V}^{(2)} + \mathbf{b}^{(2)}) - (\mathbf{a}[l]\mathbf{V}^{(2)} + \mathbf{b}^{(2)}) \tilde{\mathbf{x}}[l]^T \right\}, \end{aligned} \quad (9)$$

where $G^{(1)}(\cdot)$, $G^{(2)}(\cdot)$ respectively, is the cumulant generating function associated with the exponential family distribution $p^{(1)}(\cdot)$, $p^{(2)}(\cdot)$ respectively, as seen in equation (3).

Following the derivations in [12], the Newton-Raphson method is used for the iterative minimization of the loss function (9) and the resulting update equations are as follows. First at iteration t , for $l = 1, \dots, m$,

$$\begin{aligned} \mathbf{a}^{(t+1)}[l]^T &= \mathbf{a}^{(t)}[l]^T - \alpha_{\mathbf{a}}^{(t+1)} \\ &\cdot \left\{ \mathbf{V}^{(1)(t)} G^{(1)''}(\mathbf{a}^{(t)}[k]\mathbf{V}^{(1)(t)} + \mathbf{b}^{(1)(t)}) \mathbf{V}^{(1)(t),T} \right. \\ &+ \left. \mathbf{V}^{(2)(t)} G^{(2)''}(\mathbf{a}^{(t)}[k]\mathbf{V}^{(2)(t)} + \mathbf{b}^{(2)(t)}) \mathbf{V}^{(2)(t),T} \right\}^{-1} \\ &\cdot \left\{ \mathbf{V}^{(1)(t)} (G^{(1)'})'(\mathbf{a}^{(t)}[k]\mathbf{V}^{(1)(t)} + \mathbf{b}^{(1)(t)}) - \tilde{\mathbf{x}}[k]^T \right\} \\ &+ \left. \mathbf{V}^{(2)(t)} (G^{(2)'})'(\mathbf{a}^{(t)}[k]\mathbf{V}^{(2)(t)} + \mathbf{b}^{(2)(t)}) - \tilde{\mathbf{x}}[k]^T \right\}. \end{aligned} \quad (10)$$

For the second step, the two sets of row vectors $\{\mathbf{v}_j^{(1)}\}_{j=1}^q$ and $\{\mathbf{v}_j^{(2)}\}_{j=1}^q$ are updated separately. The update equations can then be used for $\{\mathbf{v}_j^{(1)}\}_{j=1}^q$ and $\{\mathbf{v}_j^{(2)}\}_{j=1}^q$ by changing \mathbf{v}_j to $\mathbf{v}_j^{(1)}$, respectively to $\mathbf{v}_j^{(2)}$, \mathbf{b} to $\mathbf{b}^{(1)}$, respectively to $\mathbf{b}^{(2)}$, $G(\cdot)$, $G'(\cdot)$, and

$G''(\cdot)$ to $G^{(1)}(\cdot)$, $G^{(1)'(\cdot)}$, and $G^{(1)''}(\cdot)$, respectively to $G^{(2)}(\cdot)$, $G^{(2)'(\cdot)}$, and $G^{(2)''}(\cdot)$. For $j = 1, \dots, q$:

$$\begin{aligned} \mathbf{v}_j^{(t+1),T} &= \mathbf{v}_j^{(t),T} \\ &- \alpha_{\mathbf{v}}^{(t+1)} \left(\sum_{l=1}^m a_j^{(t+1)} [l]^2 G''(\mathbf{a}^{(t+1)}[l] \mathbf{V}^{(t)} + \mathbf{b}^{(t)}) \right)^{-1} \\ &\cdot \left(\sum_{k=1}^n a_j^{(t+1)} [l] \{ G'(\mathbf{a}^{(t+1)}[l] \mathbf{V}^{(t)} + \mathbf{b}^{(t)}) - \tilde{\mathbf{x}}[k]^T \} \right). \end{aligned} \quad (11)$$

And finally for the last step, the update equations can then be used for $\mathbf{b}^{(1)}$ and $\mathbf{b}^{(2)}$ by changing \mathbf{b} to $\mathbf{b}^{(1)}$, respectively to $\mathbf{b}^{(2)}$, \mathbf{V} to $\mathbf{V}^{(1)}$, respectively to $\mathbf{V}^{(2)}$, $G(\cdot)$, $G'(\cdot)$, and $G''(\cdot)$ to $G^{(1)}(\cdot)$, $G^{(1)'(\cdot)}$, and $G^{(1)''}(\cdot)$, respectively to $G^{(2)}(\cdot)$, $G^{(2)'(\cdot)}$, and $G^{(2)''}(\cdot)$.

$$\begin{aligned} \mathbf{b}^{(t+1),T} &= \mathbf{b}^{(t),T} \\ &= \mathbf{b}^{(t),T} - \alpha_{\mathbf{b}}^{(t+1)} \left(\sum_{l=1}^m G''(\mathbf{a}^{(t+1)}[l] \mathbf{V}^{(t+1)} + \mathbf{b}) \right)^{-1} \\ &\cdot \left(\sum_{l=1}^m \{ G'(\mathbf{a}^{(t+1)}[l] \mathbf{V}^{(t+1)} + \mathbf{b}) - \tilde{\mathbf{x}}[k]^T \} \right). \end{aligned} \quad (12)$$

Table 1 summarizes the SP-PCA algorithm.

4.2. Exponential family PCA

The work proposed in [12] is based on the generalization of Principal Component Analysis to the exponential family technique presented in [7], often referred to as exponential Principal Component Analysis.

As stated earlier, instead of the fastidious estimation of the point-mass probabilities, exponential PCA considers the special uniform case where $\pi_l = 1/m$ for $l = 1, \dots, m$, for which the number of support points equals the number of data samples ($m = n$). Hence, the point-mass probabilities do not need to be estimated and the EM algorithm is unnecessary. Then, to each vector \mathbf{x} corresponds a vector \mathbf{a} , i.e., a vector $\boldsymbol{\theta}$, and they can share the same index $k = 1, \dots, n$.

The update equations were presented in [12]. It is easily noticed that they are the same as equations (10), (11) and (12), the only difference being that the exponential PCA update equations use \mathbf{x} instead of $\tilde{\mathbf{x}}$, i.e., data points instead of mixture component centers. Each data point \mathbf{x} is its own mixture component center.

As noted in [7], it is possible for the atoms obtained with exponential PCA to diverge since the optimum may be at infinity. To avoid such behavior, we introduce a penalty function that defines and places a set of constraints into the

Algorithm: Semi-Parametric exp. family PCA [20]

Input: a set of observations $\{\mathbf{x}[k]\}_{k=1}^n \subseteq \mathbb{R}^d$, two exponential family distribution $p^{(1)}, p^{(2)}$ defined by their cumulant generating functions $G^{(1)}, G^{(2)}$, a number of atoms $m, q \ll d$ the dimension of the latent variable lower dimensional subspace.

Output: the NPML estimator that maximizes the complete log-likelihood function $L^{(c)}(\mathcal{Q}, \{\mathbf{z}_k\}_{k=1}^n)$: $\hat{\mathcal{Q}} = \{\hat{\boldsymbol{\theta}}[l], \hat{\pi}_l\}_{l=1}^m$ with $\hat{\boldsymbol{\theta}}[l] = \hat{\mathbf{a}}[l] \hat{\mathbf{V}} + \hat{\mathbf{b}}$ for all l , $\{\hat{\mathbf{a}}[l]\}_{l=1}^m \in \mathbb{R}^q$, $\hat{\mathbf{V}} \in \mathbb{R}^{q \times d}$ and $\hat{\mathbf{b}} \in \mathbb{R}^d$.

Method:

Initialize \mathbf{V} , \mathbf{b} and $\{\mathbf{a}[l], \pi_l\}_{l=1}^m$ with $\pi_l \geq 0$ for all l and $\sum_{l=1}^m \pi_l = 1$; $\boldsymbol{\theta}[l] = \mathbf{a}[l] \mathbf{V} + \mathbf{b} \in \Theta$ for all l ; $p(\mathbf{x}[k]|\boldsymbol{\theta}[l])$ as defined in (7) for all k and l ;

repeat

{The Expectation Step}

for $k = 1$ to n do

for $l = 1$ to m do

$\hat{z}_{kl} \leftarrow p(\mathbf{x}[k]|\boldsymbol{\theta}[l])\pi_l / \sum_{r=1}^m p(\mathbf{x}[k]|\boldsymbol{\theta}[r])\pi_r$

end for

end for

{The Maximization Step}

for $l = 1$ to m do

$\pi_l \leftarrow (1/n) \sum_{k=1}^n \hat{z}_{kl}$

end for

{The Newton-Raphson iterative algorithm}

for $l = 1$ to m do

$\mathbf{a}[l] \leftarrow$ update equation (10)

end for

for $j = 1$ to q do

$\mathbf{v}_j \leftarrow$ update equation (11)

end for

$\mathbf{b} \leftarrow$ update equation (12)

until convergence;

return $\hat{\mathcal{Q}} = \{\hat{\boldsymbol{\theta}}[l] = \hat{\mathbf{a}}[l] \hat{\mathbf{V}} + \hat{\mathbf{b}}, \hat{\pi}_l\}_{l=1}^m$.

Table 1. Semi-Parametric PCA algorithm.

loss function via a penalty parameter in a way that penalizes any divergence to infinity. The penalty function

$$\begin{aligned} \psi(\boldsymbol{\theta}) &= \sum_{i=1}^d \left\{ \exp(-\beta_{min}(\theta_i - \theta_{min,i})) \right. \\ &\quad \left. + \exp(\beta_{max}(\theta_i - \theta_{max,i})) \right\}, \end{aligned}$$

is designed so that $\psi(\boldsymbol{\theta})$ is close to zero for $\theta_{min} \leq \boldsymbol{\theta} \leq \theta_{max}$ and reaches infinity otherwise. The penalty function modifications on the update equations (10), (11) and (12), i.e., an additional additive term within each of the large parentheses, are omitted here for sake of brevity.

Table 2 summarizes exponential PCA algorithm.

Algorithm: Exponential PCA [7]

Input: a set of observations $\{\mathbf{x}[k]\}_{k=1}^n \subseteq \mathbb{R}^d$, two exponential family distribution $p^{(1)}, p^{(2)}$ defined by their cumulant generating functions $G^{(1)}, G^{(2)}$, a number of atoms $n, q \ll d$ the dimension of the latent variable lower dimensional subspace.

Output: the NPML estimator that maximizes the log-likelihood function $L(\mathcal{Q}): \hat{\mathcal{Q}} = \{\hat{\boldsymbol{\theta}}[k]\}_{k=1}^n$ with $\hat{\boldsymbol{\theta}}[k] = \hat{\mathbf{a}}[k]\hat{\mathbf{V}} + \hat{\mathbf{b}}$ for all k , $\{\hat{\mathbf{a}}[k]\}_{k=1}^n \in \mathbb{R}^q$, $\hat{\mathbf{V}} \in \mathbb{R}^{q \times d}$ and $\hat{\mathbf{b}} \in \mathbb{R}^d$.

Method:

Initialize \mathbf{V}, \mathbf{b} and $\{\mathbf{a}[k]\}_{k=1}^n; \boldsymbol{\theta}[k] = \mathbf{a}[k]\mathbf{V} + \mathbf{b} \in \Theta$ for all $k; p(\mathbf{x}[k]|\boldsymbol{\theta}[k])$ as defined in (7) for all k ;

repeat

{The Newton-Raphson iterative algorithm}

for $k = 1$ to n **do**

$\mathbf{a}[l] \leftarrow$ penalty-modified update equation (10)
with \mathbf{x} instead of $\tilde{\mathbf{x}}$

end for**for** $j = 1$ to q **do**

$\mathbf{v}_j \leftarrow$ penalty-modified update equation (11)
with \mathbf{x} instead of $\tilde{\mathbf{x}}$

end for

$\mathbf{b} \leftarrow$ penalty-modified update equation (12)
with \mathbf{x} instead of $\tilde{\mathbf{x}}$

until convergence;

return $\hat{\mathcal{Q}} = \{\hat{\boldsymbol{\theta}}[k] = \hat{\mathbf{a}}[k]\hat{\mathbf{V}} + \hat{\mathbf{b}}\}_{k=1}^n$.

Table 2. Exponential PCA algorithm.

4.3. Bregman soft clustering

The Bregman soft clustering approach presented in [4] utilizes an alternative interpretation of the EM algorithm for learning models involving mixtures of exponential family distributions. It is a simple soft clustering algorithm for all Bregman divergences, i.e., for all exponential family distributions. We choose here to present this technique without referring to the Bregman divergence as done in [4] but by using its corresponding exponential family probability distribution for the sake of comparison with SP-PCA and exponential PCA.

Given a data set of observations $\{\mathbf{x}[k]\}_{k=1}^n$, Bregman soft clustering aims at modeling the statistical structure of the data as a mixture of m densities of the same exponential family. The clusters correspond to the components of the mixture model and the soft membership of a data point in each cluster is proportional to the probability of the data point being generated by the corresponding density function. The Bregman soft clustering problem is based on a

Algorithm: Bregman Soft Clustering [4]

Input: a set of observations $\{\mathbf{x}[k]\}_{k=1}^n \subseteq \mathbb{R}^d$, two exponential family distribution $p^{(1)}, p^{(2)}$ defined by their cumulant generating functions $G^{(1)}, G^{(2)}$, a number of atoms m .

Output: the NPML estimator that maximizes the complete log-likelihood function $L^{(c)}(\mathcal{Q}, \{\mathbf{z}_k\}_{k=1}^n): \hat{\mathcal{Q}} = \{\hat{\boldsymbol{\theta}}[l], \hat{\pi}_l\}_{l=1}^m$.

Method:

Initialize $\{\boldsymbol{\theta}[l], \pi_l\}_{l=1}^m$ with $\pi_l \geq 0$ for all l and $\sum_{l=1}^m \pi_l = 1; p(\mathbf{x}[k]|\boldsymbol{\theta}[l])$ as defined in (7) for all k and $l; \boldsymbol{\theta}[l] \in \Theta$ for all l ;

repeat

{The Expectation Step}

for $k = 1$ to n **do****for** $l = 1$ to m **do**

$\hat{z}_{kl} \leftarrow p(\mathbf{x}[k]|\boldsymbol{\theta}[l])\pi_l / \sum_{r=1}^m p(\mathbf{x}[k]|\boldsymbol{\theta}[r])\pi_r$

end for**end for**

{The Maximization Step}

for $l = 1$ to m **do**

$\pi_l \leftarrow (1/n) \sum_{k=1}^n \hat{z}_{kl}$

$\boldsymbol{\theta}[l] \leftarrow$ solve for $\boldsymbol{\theta}[l]$:

$$G'(\boldsymbol{\theta}[l]) = \sum_{k=1}^n \hat{z}_{kl} \mathbf{x}[k] / \sum_{k=1}^n \hat{z}_{kl}$$

end for**until convergence;**

return $\hat{\mathcal{Q}} = \{\hat{\boldsymbol{\theta}}[l], \hat{\pi}_l\}_{l=1}^m$.

Table 3. Bregman soft clustering algorithm.

maximum likelihood estimation of the cluster parameters $\{\boldsymbol{\theta}[l], \pi_l\}_{l=1}^m$ satisfying the following mixture structure:

$$p(\mathbf{x}) = \sum_{l=1}^m p(\mathbf{x}|\boldsymbol{\theta}[l])\pi_l,$$

where $p(\mathbf{x}|\cdot)$ is an exponential family distribution. The data likelihood function takes the following form:

$$p(\mathbf{X}) = \prod_{k=1}^n \sum_{l=1}^m p(\mathbf{x}[k]|\boldsymbol{\theta}[l])\pi_l. \quad (13)$$

The data likelihood function in (13) is similar to the data likelihood function in (6) without the linear constraint $\boldsymbol{\theta}[l] = \mathbf{a}[l]\mathbf{V} + \mathbf{b}$ for $l = 1, \dots, m$. Hence, the Bregman soft clustering problem is similar to the SP-PCA problem without the lower dimensional subspace constraint and a simple EM algorithm is used to estimate the cluster parameters. We consider again the mixed data type case. The E-step and the first part of the M-step yield the same results as for SP-PCA. In the second part of the M-step, the component parameters

$\theta[l], l = 1, \dots, m$, are estimated in the following way:

$$\begin{aligned} \theta[l]^{(new)} &= \arg \max_{\theta[l]} \sum_{k=1}^n \sum_{r=1}^m \hat{z}_{kr} \log p(\mathbf{x}[k]|\theta[r]) \\ &= \arg \max_{\theta[l]} \left\{ \sum_{k=1}^n \sum_{r=1}^m \hat{z}_{kr} \log p^{(1)}(\mathbf{x}^{(1)}[k]|\theta^{(1)}[r]) \right. \\ &\quad \left. + \sum_{k=1}^n \sum_{r=1}^m \hat{z}_{kr} \log p^{(2)}(\mathbf{x}^{(2)}[k]|\theta^{(2)}[r]) \right\}, \end{aligned}$$

with $\log p(\mathbf{x}[k]|\theta[r]) = \theta[r]\mathbf{x}[k]^T - G(\theta[r])$. Using the convexity properties of $G(\cdot)$, it is easily shown that:

$$G'(\theta[l]^{(new)}) = \left(\sum_{k=1}^n \hat{z}_{kl} \mathbf{x}[k] \right) / \left(\sum_{k=1}^n \hat{z}_{kl} \right)$$

can be solved for $\theta[l]^{(new),(1)}$ and $\theta[l]^{(new),(2)}$ by changing \mathbf{x} to $\mathbf{x}^{(1)}$, respectively to $\mathbf{x}^{(2)}$, $G'(\cdot)$ to $G^{(1)'}(\cdot)$, respectively to $G^{(2)'}(\cdot)$.

Table 3 summarizes the Bregman soft clustering algorithm.

5. Experimental results on synthetic data

We compare the relative performances of exponential PCA, SP-PCA and Bregman soft clustering in a mixed data set clustering problem with two data types and demonstrate how exponential PCA with the addition of a nonparametric estimation of the point-mass probabilities exceeds SP-PCA in performance.

We first consider a synthetic $d = 3$ -dimensional data set with a lower dimensional subspace of dimension $q = 1$. The first data feature is Poisson distributed, the second and third features are Gaussian distributed. The data has $n = 500$ points and is composed of two mixture components with parameters $\theta[1]$ and $\theta[2]$ constrained to the lower dimensional subspace.

We first use exponential PCA. However, exponential PCA does not estimate point-mass probabilities. We use a nonparametric density estimation technique based on a kernel smoothing method to estimate the point-mass probabilities using the support points values $\mathbf{a}[k], k = 1, \dots, n$, obtained by exponential PCA. Figure 3 shows that the nonparametric density estimation exhibits a definite two-component shape. The dotted lines represent the correct values $\mathbf{a}[1]$ and $\mathbf{a}[2]$. We can then estimate the values of $\mathbf{a}[1]$ and $\mathbf{a}[2]$ as well as their mixing distributions π_1 and π_2 using a simple kmeans algorithm, with the $\pi_1 + \pi_2 = 1$ assumption.

Figure 4 presents the histogram of the estimated point-mass probabilities obtained with SP-PCA, $m = 2$.

Table 4 shows detailed results for this synthetic data setting (“modified” means the extension to mixed data sets of the algorithm): the mixing distributions or point-mass probabilities π_1 and π_2 , the latent variable or point of support values $\mathbf{a}[1]$ and $\mathbf{a}[2]$, the parameter values $\theta[1]$ and $\theta[2]$ as well as the sine of the angle between the estimated lower dimensional subspace and the correct subspace. Bregman soft clustering does not have the lower dimensional subspace constraint, and hence does not exhibit a sine or the latent variables values in Table 4. The estimation quality of the $\theta[1]$, $\theta[2]$ and π_1, π_2 values defines the clustering performance. For this simple Poisson-Gaussian mixed data setting, both exponential PCA and Bregman soft clustering seem to perform better than SP-PCA: the SP-PCA obtained parameter values for $\theta[2]$ are far from the original values, contrary to exponential PCA and Bregman soft clustering.

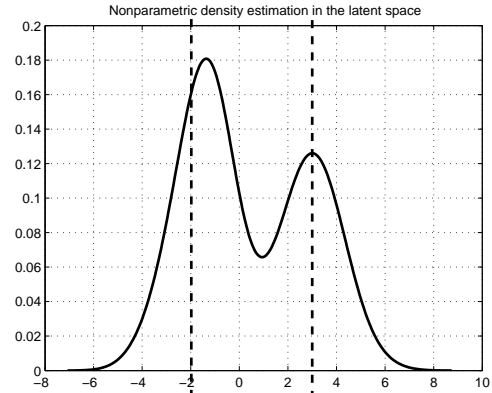


Figure 3. Nonparametric estimation of the point-mass probabilities obtained with exponential PCA (dotted: correct cluster centers).

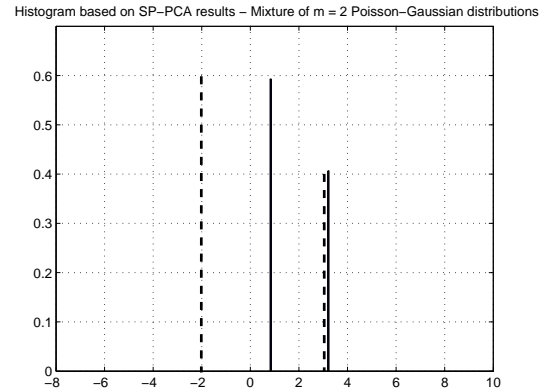


Figure 4. Histogram of the estimated point-mass probabilities obtained with SP-PCA (dotted: correct cluster values).

Results for a second experiment are shown in Table 5 for a Binomial-Gaussian mixed data set created in a similar fashion as the Poisson-Gaussian mixed data set (the parameter N is set to 10 for the Binomial component). Again, exponential PCA exceeds SP-PCA in clustering performance.

	$\pi_1; \pi_2$	$\mathbf{a}[1]; \mathbf{a}[2]$	$\boldsymbol{\theta}[1]; \boldsymbol{\theta}[2]$	sin
correct model values	0.4 0.6	3 -2	[1.9404, 1.6148, 1.6210] [-1.2936, -1.0765, -1.0806]	
modified exponential PCA	0.4107 0.5893	3.0009 -1.3725	[1.6235, 1.8648, 1.7007] [-0.7425, -0.8529, -0.7778]	0.1368
modified SP-PCA	0.3724 0.6276	3.2170 0.8355	[2.1732, 1.5715, 1.7768] [0.5644, 0.4081, 0.4614]	0.058663
modified Bregman soft clustering	0.4069 0.5931		[1.9317, 1.7162, 1.5585] [-1.1061, -1.0802, -1.0304]	

Table 4. Clustering results for a Poisson-Gaussian mixed data set.

	$\pi_1; \pi_2$	$\mathbf{a}[1]; \mathbf{a}[2]$	$\boldsymbol{\theta}[1]; \boldsymbol{\theta}[2]$	sin
correct model values	0.4 0.6	1 -2	[0.8914, 0.1688, 0.4206] [-1.7828, -0.3375, -0.8412]	
modified exponential PCA	0.4475 0.5525	0.8559 -1.9972	[0.7796, 0.1166, 0.3334] [-1.8193, -0.2721, -0.7779]	0.049038
modified SP-PCA	0.3978 0.6022	-0.9548 -3.1821	[-0.9046, -0.0989, -0.2890] [-3.0148, -0.3296, -0.9633]	0.1455
modified Bregman soft clustering	0.3973 0.6027		[0.82252, 0.144, 0.41004] [-1.8072, -0.3089, -0.9816]	

Table 5. Clustering results for a Binomial-Gaussian mixed data set.

6. Conclusion

We presented a mixed data-type hierarchical Bayesian graphical model framework that adds clarity and perspective to our understanding of exponential PCA, SP-PCA and Bregman soft clustering. We demonstrated that these techniques are not separate unrelated algorithms but different manifestations of parameter choices taken within a common framework. Because of this insight, we were able to extend the algorithms to readily derive novel extensions that deal with the important mixed data type case. Our framework has the critical advantage of allowing one to go from high-dimensional mixed-type data components to low-dimensional common-type latent variables that are then used to perform clustering in a much simpler manner using well-known continuous-parameter clustering techniques.

References

- [1] M. Aitkin, A general maximum likelihood analysis of overdispersion in generalized linear models, *Statistics and Comput.*, vol. 6, pp. 251-262, 1996.
- [2] M. Aitkin, A maximum likelihood analysis of variable components in generalized linear models, *Biometrics*, vol. 55, pp. 117-128, 1999.
- [3] K. S. Azoury and M. K. Warmuth, Relative loss bounds for on-line density estimation with the exponential family of distributions, *Mach. Learning*, vol. 43, pp. 211-246, 2001.
- [4] A. Banerjee, S. Merugu, I. Dhillon and J. Ghosh, Clustering with Bregman divergences, *J. Mach. Learning Research*, vol. 6, pp. 1705-1749, 2005.
- [5] D. Boehning, *Computer-Assisted Analysis of Mixtures and Applications: Meta-Analysis, Disease Mapping, and Others*, Chapman and Hall/CRC, 2000.
- [6] L. Cayton, Fast nearest neighbor retrieval for Bregman divergences, *25th Int. Conf. Mach. Learning*, 2008.
- [7] M. Collins, S. Dasgupta and R. Shapire, A generalization of principal component analysis to the exponential family, *NIPS*, 2001.
- [8] A. P. Dempster, N. M. Laird and D. B. Rubin, Maximum likelihood estimation from incomplete data via the EM algorithm, *J. Royal Statistical Soc., B*, vol. 39, pp. 1-38, 1977.
- [9] J. Kiefer and J. Wolfowitz, Consistency of the maximum likelihood estimator in the presence of infinitely many incidental parameters, *The Annals of Mathematical Statistics*, vol. 27, pp. 887-906, 1956.
- [10] N. Laird, Nonparametric Maximum Likelihood Estimation of a Mixing Distribution, *J. American Statistical Soc.*, vol. 73, 1978.
- [11] E. L. Lehmann and G. Castella, *Theory of Point Estimation*, Springer, 2nd edition, 1998.
- [12] C. Levasseur, U. F. Mayer, B. Burdige and K. Kreutz-Delgado, Generalized statistical methods for unsupervised minority class detection in mixed data sets, *IAPR Workshop on Cognitive Inf.*, 2008.
- [13] B. G. Lindsay, The geometry of mixture likelihoods: a general theory, *Annals of Stat.*, vol. 11, no. 1, pp. 86-94, 1983.
- [14] B. G. Lindsay and M. L. Lesperance, A review of semiparametric mixture models, *J. Statistical Planning and Inference*, vol. 47, pp. 29-99, 1995.
- [15] A. Mallet, A maximum likelihood estimation method for random coefficient regression models, *Biometrika*, vol. 73, no. 3, pp. 645-656, 1986.
- [16] G. McLachlan and D. Peel, *Finite Mixture Models*, Wiley-Interscience, 2000.
- [17] F. Nielsen, J.-. Boissonnat and R. Nock, On Visualizing Bregman Voronoi Diagrams, *Proc. of the 23rd ACM Conf. on Comput. Geometry*, 2007.
- [18] R. S. Pilla and B. Lindsay, Alternative EM methods for nonparametric finite mixture models, *Biometrika*, vol. 88, no. 2, pp. 535-550, 2001.
- [19] R. A. Redner and H. F. Walker, Mixture densities, maximum likelihood and the EM algorithm, *SIAM Review*, vol. 26, no. 2, pp. 195-239, 1984.
- [20] Sajama and A. Orlitsky, Semi-parametric exponential family PCA, *NIPS*, 2004.