# Hierarchical mixtures of Gaussians for combined dimensionality reduction and clustering

Sacha Sokoloski

University of Tübingen

sacha.sokoloski@mailbox.org

Philipp Berens

University of Tübingen

philipp.berens@uni-tuebingen.de

June 13, 2022

**Abstract**

To avoid the curse of dimensionality, a common approach to clustering high-dimensional data is to first project the data into a space of reduced dimension, and then cluster the projected data. Although effective, this two-stage approach prevents joint optimization of the dimensionality-reduction and clustering models, and obscures how well the complete model describes the data. Here, we show how a family of such two-stage models can be combined into a single, hierarchical model that we call a hierarchical mixture of Gaussians (HMoG). An HMoG simultaneously captures both dimensionality-reduction and clustering, and its performance is quantified in closed-form by the likelihood function. By formulating and extending existing models with exponential family theory, we show how to maximize the likelihood of HMoGs with expectation-maximization. We apply HMoGs to synthetic data and RNA sequencing data, and demonstrate how they exceed the limitations of two-stage models. Ultimately, HMoGs are a rigorous generalization of a common statistical framework, and provide researchers with a method to improve model performance when clustering high-dimensional data.

## 1 Introduction

Clustering — the grouping of individuals or items based on their similarity — is an essential part of knowledge discovery. The growing complexity of modern datasets thus poses a challenge, because many clustering algorithms — such as mixtures of Gaussians (MoG) — suffer from the so-called "curse of dimensionality", and exhibit limited performance when classifying high-dimensional data. This arises not only because model complexity scales with dimension, but also because

similarity metrics used to evaluate model performance tend to break down in high-dimensional spaces [1, 2].

A common approach for addressing this challenge is to first apply dimensionality reduction techniques — such as principle component analysis (PCA) or factor analysis (FA) — to project the data into a lower dimensional space, and then cluster the projected data. Such two-stage algorithms are widely applied in problem domains such as image processing [3, 4] and time-series analysis [5], and fields such as neuroscience [6, 7] and bioinformatics [8, 9].

Nevertheless, two-stage approaches struggle with two primary limitations: Firstly, the dimensionality-reduction is typically implemented as a preprocessing step, and is not further optimized based on the results of the subsequent clustering. This can lead to suboptimal projections — if we consider PCA, for example, the directions in the data that best separate the clusters might not be the directions of maximum variance that define the PCA projection [10, 11].

Secondly, two-stage algorithms rely on optimizing two distinct objectives, and it is therefore non-trivial to quantify how well the complete model describes the data. As such, even if we consider algorithms that update projections based on clustering results [12, 13], it is difficult to to be sure that these updates improve overall model performance, because we cannot evaluate the likelihood of the complete model given the data.

To address these limitations, we consider a family of two-stage models that implement dimensionality reduction with a linear Gaussian model (e.g. PCA or FA), and clustering with a mixture of Gaussians (MoG). By applying a novel theory of conjugate priors in latent variable exponential family models, we show how these two stages can be combined and generalized into a single, hierarchical probabilistic model that we call a hierarchical mixture of Gaussians (HMoG).

An HMoG captures both dimensionality reduction and clustering in a single model, and has a single objective given by the likelihood function. Critically, the tractability of HMoGs allows us to derive closed-form expressions for the HMoG density function, and a novel expectation-maximization (EM) algorithm for maximizing the HMoG likelihood. Although the derivation of HMoG theory requires a number of technical innovations, the upshot of the theory is straightforward: a closed-form expression for the HMoG log-likelihood allows us to rigorously compare various algorithms for dimensionality reduction and clustering, and the HMoG EM algorithm allows HMoGs to strictly exceed the performance of the two-stage models that they generalize.

To validate our theory we apply HMoGs to several datasets. In synthetic data, we show how HMoGs can overcome fundamental limitations of standard two-stage algorithms. Finally, in a RNA sequencing (RNA-Seq) data set we show how our methods can significantly exceed the predictive performance of standard two-stage approaches. Our work is related to previous methods on mixtures of PCA/FA [3, 4, 14, 15]. In contrast with these methods, however, our model has a hierarchical structure that affords a more compact parameterization, and ensures that a shared, low-dimensional feature space is extracted during training.

## 2   Theory

Our focus will be on HMoGs where the dimensionality reduction is based on either PCA or FA, and the clustering is based on MoGs. Historically, each of these techniques has fairly different theoretical foundations — PCA finds directions of maximum variance in the data, FA explains the data as a linear combination of hidden factors, and MoGs model the data distribution as a weighted sum of multivariate normal distributions. Nevertheless, each technique may be formulated as a probabilistic latent variable model, that may be fit to data with EM. This unified, probabilistic foundation is essential for combining PCA/FA and MoGs into an HMoG, and ultimately deriving a single EM algorithm that jointly optimizes dimensionality-reduction and clustering.

### 2.1   Linear Gaussian models and mixture models are probabilistic latent variable models

Suppose we make $d_S$ observations $x^{(1)}, \ldots, x^{(d_S)}$ of the random variable $X$, and we hypothesize that $X$ is influenced by some latent random variable $Y$. In the maximum likelihood framework, a statistical model $p(x; \theta)$ with parameters $\theta$ aims to capture the distribution of $X$ by maximizing the log-likelihood objective $\frac{1}{d_S} \sum_{i=1}^{d_S} \log p(x^{(i)}; \theta)$ with respect to $\theta$. We may extend this framework to a latent variable model $p(x, y; \theta)$ that captures both the distribution of $X$ and the effect of the latent variable $Y$, by maximizing the log-likelihood of the marginal distribution $p(x; \theta)$ of $p(x, y; \theta)$. As we will see, both linear Gaussian models and mixture models afford a probabilistic latent variable formulation, and linear Gaussian models include FA and PCA as a special case [16, 17].

On one hand, where $X$ and $Y$ are continuous variables of dimensions $n$ and $m$, respectively, a linear Gaussian model $p(\mathbf{x}, \mathbf{y})$ is simply an $n + m$ dimensional multivariate normal distribution. As a consequence, the observable distribution $p(\mathbf{x})$ and the prior $p(\mathbf{y})$, as well as the likelihood[1] $p(\mathbf{x} \mid \mathbf{y})$ and the posterior $p(\mathbf{y} \mid \mathbf{x})$ are also multivariate normal distributions [see 17]. Where $p(\mathbf{x} \mid \mathbf{y}; \mu, \Sigma) = N(\mu + \mathbf{W} \cdot \mathbf{y}, \Sigma)$ and $p(\mathbf{y}) = N(\mathbf{0}, \mathbf{I})$, we may define probabilistic PCA as the case where the columns of $\mathbf{W}$ are proportional to the ranked eigenvectors of the sample covariance matrix and $\Sigma = \sigma \mathbf{I}$ is isotropic. FA is then the case where the columns of $\mathbf{W}$ are the so-called factor loadings and $\Sigma$ is a diagonal matrix. Finally, in both cases, the mean of the posterior distribution $p(\mathbf{y} \mid \mathbf{x})$ is given by $\mathbb{E}[Y \mid X = \mathbf{x}] = \mathbf{W^T} \cdot (\mathbf{W} \cdot \mathbf{W^T} + \Sigma)^{-1} \cdot \mathbf{x}$, and probabilistically formalizes projecting a datapoint $\mathbf{x}$ into the reduced space ("classic" PCA is recovered in the limit as $\sigma \to 0$).

A mixture model, on the other hand, is a latent variable model where the latent variable represents an index that ranges from 1 to $k$. Since our goal for HMoGs is ultimately to cluster low-dimensional features $Y$, let us introduce a mixture as the model $p(\mathbf{y}, z)$ over $Y$ and an index-valued latent

---

[1] Unfortunately there is no naming convention for distinguishing the likelihood $p(\mathbf{x}; \theta)$ of non-Bayesian model parameters and the likelihood $p(\mathbf{x} \mid \mathbf{y})$ of random variables. We will refer to $\log p(\mathbf{x}; \theta)$ as the log-likelihood *objective* where necessary to help distinguish them.
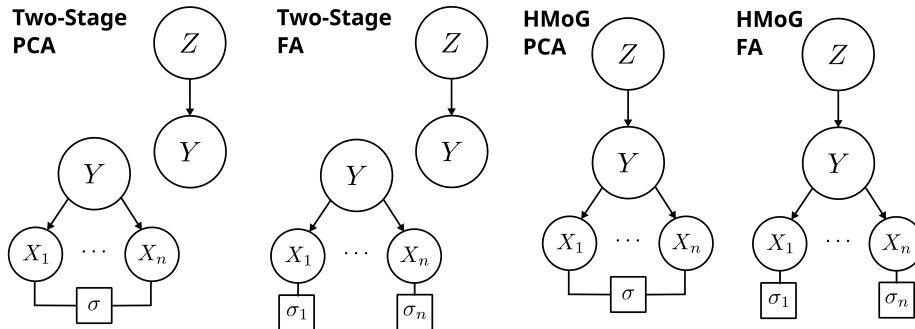
Figure 1: *Dimensionality reduction and clustering as a hierarchical graphical model.* A graphical representation of four methods, where $X = (X_1, \ldots, X_n)$ is the observation, $Y$ is the low-dimensional feature, and $Z$ is the feature cluster. Two-stage models have distinct models for dimensionality reduction and clustering, whereas HMoGs are a single model of both. PCA-based models share a single standard deviation $\sigma$ amongst all $X_i$, whereas FA-based models have distinct $\sigma_i$ for each $X_i$.

variable $Z$. In this formulation, the components of the mixture are given by the likelihood $p(\mathbf{y} \mid z)$ for each index $z$, the weights $\boldsymbol{\pi}$ of the mixture are the prior probabilities $\pi_i = p(z = i)$, and the weighted sum of the components ("the mixture distribution") is equal to the observable distribution $p(\mathbf{y})$. MoGs, then, are simply the case where each component $p(\mathbf{y} \mid z = i) = \mathrm{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ is a multivariate normal distributions with mean $\boldsymbol{\mu}_i$ and covariance $\boldsymbol{\Sigma}_i$.

FA, PCA, and MoGs can all be fit to data using expectation-maximization (EM) in our probabilistic framework. The two-stage approach to clustering a high-dimensional dataset $\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(d_S)}$ is thus first to fit an appropriate linear Gaussian model (i.e. FA or PCA) $p(\mathbf{x}, \mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{W})$ to the dataset. Then, to fit a MoG $p(\mathbf{y}, z; \boldsymbol{\pi}, \boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_1, \ldots, \boldsymbol{\Sigma}_k)$ to the projected dataset $\mathbf{y}^{(1)}, \ldots, \mathbf{y}^{(d_S)}$, where $\mathbf{y}^{(i)} = \mathbb{E}[Y \mid X = \mathbf{x}^{(i)}]$ (Fig. 1). After training, the cluster membership of an arbitrary datapoint $\mathbf{x}$ is determined by first computing the projection $\mathbf{y} = \mathbb{E}[Y \mid X = \mathbf{x}]$, and then computing (and perhaps taking the maximum of) the index probabilities $p(z \mid \mathbf{y})$.

## 2.2 New exponential family theory enables tractable computation with latent variable models

When analyzing high-dimensional data, one must take care to avoid making complex computations in the high-dimensional space of observations, as this can quickly lead to computational intractability and numerical instability. In order to address this in our present context, we introduce the class of model known as exponential family harmoniums [18, 19], which allows us to unify the mathematics of our latent variable models. We then develop a theory of conjugate priors for

4

harmoniums, which we use to break down complex probabilistic computations with harmoniums (and ultimately HMoGs) into tractable components.

To begin, an exponential family is a statistical model defined by a sufficient statistic $\mathbf{s}$ and base measure $\nu$, and has the form $p(x; \boldsymbol{\theta}) = e^{\mathbf{s}(x) \cdot \boldsymbol{\theta} - \psi(\boldsymbol{\theta})} \nu(x)$, where $\boldsymbol{\theta}$ are the natural parameters, and $\psi$ is the normalizer known as the log-partition function. There are two exponential families that are particularly relevant for our purposes. Firstly, the family of $n$-dimensional multivariate normal distributions is the exponential family with base measure $\nu = (2\pi)^{-\frac{n}{2}}$ and sufficient statistic $\mathbf{s}(\mathbf{x}) = (\mathbf{x}, \mathbf{x} \otimes \mathbf{x})$, where $\otimes$ is the outer product operator. Secondly, the family of categorical distributions over indices $1, \ldots, k$ is the exponential family with base measure $\nu = 1$ sufficient statistic given by a so-called "one-hot" vector, so that $\mathbf{s}(1) = \mathbf{0}$, and $\mathbf{s}(z)$ is a $k - 1$ length vector with all zero elements except for a 1 at element $z - 1$.

An exponential family harmonium is then defined as a kind of product exponential family, which includes various models as special cases [19]. Given two exponential families defined by $\mathbf{s}_X$ and $\nu_X$, and $\mathbf{s}_Y$ and $\nu_Y$, respectively, a harmonium is the model

$$p(x, y; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY}) = e^{\mathbf{s}_X(x) \cdot \boldsymbol{\theta}_X + \mathbf{s}_Y(y) \cdot \boldsymbol{\theta}_Y + \mathbf{s}_X(x) \cdot \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y) - \psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})} \nu_X(x) \nu_Y(y),$$
(1)

which is an exponential family defined by the base measure $\nu_{XY}(x, y) = \nu_X(x) \nu_Y(y)$ and the sufficient statistic $\mathbf{s}_{XY}(x, y) = (\mathbf{s}_x(x), \mathbf{s}_Y(y), \mathbf{s}_X(x) \otimes \mathbf{s}_Y(y))$, with natural parameters $\boldsymbol{\theta}_{XY} = (\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$ and log-partition function $\psi_{XY}$.

To return to our dimensionality-reduction and clustering problem, let us define $\mathbf{s}_X$ and $\nu_X$, $\mathbf{s}_Y$ and $\nu_Y$, and $\mathbf{s}_Z$ and $\nu_Z$ as the sufficient statistics and base measures of the $n$-dimensional multivariate normal family, the $m$-dimensional multivariate normal family, and the categorical family over $k$ indices, respectively. Based on our definitions, we may immediately define a MoG as the harmonium $p(\mathbf{y}, z; \boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{YZ}) \propto e^{\mathbf{s}_Y(\mathbf{y}) \cdot \boldsymbol{\theta}_Y + \mathbf{s}_Z(z) \cdot \boldsymbol{\theta}_Z + \mathbf{s}_Y(\mathbf{y}) \cdot \boldsymbol{\Theta}_{YZ} \cdot \mathbf{s}_Z(z)}$ (where we absorb the constant base measure $\nu_Y \cdot \nu_Z$ into the proportionality relation). We leave the detailed equations of the natural parameters to the Appendix, but sufficed to say the natural parameters $\boldsymbol{\theta}_Y$ and $\boldsymbol{\Theta}_{YZ}$ correspond to the parameters $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_1, \ldots, \boldsymbol{\Sigma}_k$ of the mixture components of the MoG, and the parameters $\boldsymbol{\theta}_Z$ correspond to the mixture weights $\boldsymbol{\pi}$ of the MoG.

On the other hand, we may express a linear Gaussian model in the form $p(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY}) \propto e^{\mathbf{s}_X(\mathbf{x}) \cdot \boldsymbol{\theta}_X + \mathbf{s}_Y(\mathbf{y}) \cdot \boldsymbol{\theta}_Y + \mathbf{x} \cdot \boldsymbol{\Theta}_{XY} \cdot \mathbf{y}}$. This is a restricted version of the harmonium defined by $\mathbf{s}_X$ and $\nu_X$, and $\mathbf{s}_Y$ and $\nu_Y$, where the term $\mathbf{x} \cdot \boldsymbol{\Theta}_{XY} \cdot \mathbf{y}$ in the exponent ensures that the variables $\mathbf{x}$ and $\mathbf{y}$ only interact through their first order statistics, rather than the second order statistics that are part of $\mathbf{s}_X$ and $\mathbf{s}_Y$. We again leave out the details (see Appendix), but the parameters $\boldsymbol{\theta}_X$ correspond to $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, and the parameters $\boldsymbol{\Theta}_{XY}$ correspond to the parameters $\mathbf{W}$ of the PCA and FA models.

To continue, a prior distribution is said to be conjugate to the posterior when both have the same exponential family form, and conjugate priors can greatly simplify statistical inference in exponential family models. Conjugate priors usually arise in the context of inferring the parameters of well known models

such as normal or Poisson distributions, but here we wish to use the concept to infer distributions over our features $Y$ and clusters $Z$. Due to the log-linear form of harmonium, the harmonium posterior is given by

$$p(y \mid x) = e^{\mathbf{s}_Y(y) \cdot (\boldsymbol{\theta}_Y + \mathbf{s}_X(x) \cdot \boldsymbol{\Theta}_{XY}) - \psi_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(x) \cdot \boldsymbol{\Theta}_{XY})} \nu_Y(y), \tag{2}$$

and is therefore always in the exponential family defined by $\mathbf{s}_Y$ and $\nu_Y$. As such we can reduce the question of whether a harmonium prior is conjugate to its posterior to whether the prior is also in the exponential family defined by $\mathbf{s}_Y$ and $\nu_Y$.

**Lemma 1.** *Suppose that $p(x, y; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$ is a harmonium. Then the prior $p(y)$ satisfies $p(y) \propto \nu_Y(y) e^{\mathbf{s}_Y(y) \cdot \boldsymbol{\theta}_Y^*}$ for some $\boldsymbol{\theta}_Y^*$ if and only if there exists a vector $\boldsymbol{\rho}_Y$ and a constant $\rho_0$ such that*

$$\psi_X(\boldsymbol{\theta}_X + \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(\mathbf{y})) = \mathbf{s}_Y(\mathbf{y}) \cdot \boldsymbol{\rho}_Y + \rho_0, \tag{3}$$

*and*

$$\boldsymbol{\theta}_Y^* = \boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y. \tag{4}$$

When Eq. 3 is satisfied, we say that the harmonium $p(x, y)$ is conjugated, and we refer to the parameters $\boldsymbol{\rho}_Y$ and $\rho_0$ as the conjugation parameters. In general, Eq. 3 is a strong constraint, and is not satisfied by most models. Nevertheless, both linear Gaussian models and mixture models are indeed conjugated harmoniums with closed-form expressions for their conjugation parameters. We again leave the out details of their evaluation (see Appendix), but sufficed to say, for PCA and FA, the computational complexity of evaluating the conjugation parameters scales only linearly with the dimension $n$ of the observations $X$.

By applying the theory of conjugation, we can overcome one of our major challenges, namely computing the harmonium observable density $p(x)$. The observable density is given by

$$p(x) = e^{\mathbf{s}_X(x) \cdot \boldsymbol{\theta}_X + \psi_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(x) \cdot \boldsymbol{\Theta}_{XY}) - \psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})} \nu_X(x), \tag{5}$$

which is not in general an exponential family distribution. Typically, the log-partition function $\psi_Y$ will be tractable, which reduces the difficulty of evaluating this density to evaluating the harmonium log-partition function $\psi_{XY}$. Although $\psi_{XY}$ is typically intractable, for conjugated harmoniums we can also reduce its evaluation to the evaluation of $\psi_Y$.

**Corollary 2.** *Suppose that $p(x, y; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$ is a conjugated harmonium, with conjugation parameters $\boldsymbol{\rho}_Y$ and $\rho_0$. Then*

$$\psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY}) = \psi_Y(\boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y) + \rho_0. \tag{6}$$

In the context of our dimensionality reduction problem, Eq. 6 also allows us to evaluate the log-partition function $\psi_{XY}$ of the linear Gaussian model $p(\mathbf{x}, \mathbf{y})$ in terms of the log-partition function $\psi_Y$ on the feature space, and thereby avoid computations in the high-dimensional observable space.

Our other major challenge is fitting our models to data, and thankfully, the exponential family structure of harmoniums also affords a general description of the EM algorithm. This follows from the fundamental property of exponential families that the gradient of the log-partition function is equal to the expected value of the sufficient statistics [20]. Suppose we wish to fit the harmonium $p(x, y; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$ to the sample $x^{(1)}, \ldots, x^{(d_S)}$, and that $\boldsymbol{\tau}_Y$ and $\boldsymbol{\tau}_{XY}$ are the gradients of the log-partition functions $\psi_X$ and $\psi_{XY}$, respectively. Then an iteration of EM may be formulated as

> **for all** $i \in \{1, \ldots, d_S\}$ **do**
>    $\mathbb{E}[\mathbf{s}(Y) \mid X = x^{(i)}] \leftarrow \boldsymbol{\tau}_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(x^{(i)}) \cdot \boldsymbol{\Theta}_{XY})$     ▷ Compute latent expectations
> **end for**
>
> $\boldsymbol{\eta}'_X \leftarrow \frac{1}{d_S} \sum_{i=1}^{d_S} \mathbf{s}_X(x^{(i)})$     ▷ Compute updated average sufficient statistics
> $\boldsymbol{\eta}'_Y \leftarrow \frac{1}{d_S} \sum_{i=1}^{d_S} \mathbb{E}[\mathbf{s}(Y) \mid X = x^{(i)}]$
> $\mathbf{H}'_{XY} \leftarrow \frac{1}{d_S} \sum_{i=1}^{d_S} \mathbf{s}_X(x^{(i)}) \otimes \mathbb{E}[\mathbf{s}(Y) \mid X = x^{(i)}])$
>
> $(\boldsymbol{\theta}'_X, \boldsymbol{\theta}'_Y, \boldsymbol{\Theta}'_{XY}) \leftarrow \boldsymbol{\tau}_{XY}^{-1}(\boldsymbol{\eta}_X, \boldsymbol{\eta}_Y, \mathbf{H}_{XY})$     ▷ Compute updated natural parameters

The tractability of training a harmonium $p(\mathbf{x}, \mathbf{y})$ thus reduces to the tractability of the log-partition gradient $\boldsymbol{\tau}_Y$, and the inverse of the gradient $\boldsymbol{\tau}_{XY}$. In the case of linear Gaussian models and MoGs, both functions are available in closed-form (see Appendix), and there is thus a closed-form EM for training them.

## 2.3 Hierarchical mixtures of Gaussians can be tractably evaluated and fit to data

We exploit the probabilistic formulation of the linear Gaussian model $p(\mathbf{x}, \mathbf{y})$ and MoG $p(\mathbf{y}, z)$ to define an HMoG as the model $p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{x} \mid \mathbf{y})p(\mathbf{y}, z)$ over observations $X$, features $Y$, and clusters $Z$. Intuitively, we define an HMoG by taking a PCA or FA model, and swapping out the standard normal prior for a MoG. This definition is more than a notational trick, as it ensures that fitting the two-stage model (i.e. PCA/FA and MoG separately) maximizes a lower-bound on the marginal log-likelihood $\frac{1}{d_S} \sum_{i=1}^{d_S} \log p(x^{(i)})$ of the HMoG (see Appendix).

Putting the expressions for $p(\mathbf{x} \mid \mathbf{y})$ and $p(\mathbf{y}, z)$ together allows us to write the HMoG density as

$$\log p(\mathbf{x}, \mathbf{y}, z; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{XY}, \boldsymbol{\Theta}_{YZ}) =$$
$$\boldsymbol{\theta}_X \cdot \mathbf{s}_X(\mathbf{x}) + \boldsymbol{\theta}_Y \cdot \mathbf{s}_Y(\mathbf{y}) + \boldsymbol{\theta}_Z \cdot \mathbf{s}_Z(z) + \mathbf{x} \cdot \boldsymbol{\Theta}_{XY} \cdot \mathbf{y} + \mathbf{s}_Y(\mathbf{y}) \cdot \boldsymbol{\Theta}_{YZ} \cdot \mathbf{s}_Z(z)$$
$$+ \log \nu_{XYZ}(\mathbf{x}, \mathbf{y}, z) - \psi_{XYZ}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{XY}, \boldsymbol{\Theta}_{YZ}), \quad (7)$$

where $\nu_{XYZ} = \nu_X \cdot \nu_Y \cdot \nu_Z$ is the HMoG base measure, and $\psi_{XYZ}$ is the log-partition function. This construction ensures that $X$ and $Z$ are conditionally

independent given $Y$, which allows us to depict an HMoG as a hierarchical graphical model (Fig. 1).

By reorganizing terms, the density of an HMoG may expressed as a form of harmonium density, so that we may apply the theory of conjugation to HMoGs. In particular, by applying Corollary 2, we may express the observable density (Eq. 5) of an HMoG as

$$p(x) = \nu_X(x) e^{\mathbf{s}_X(\mathbf{x}) \cdot \boldsymbol{\theta}_X + \psi_Z(\boldsymbol{\theta}_Z + \boldsymbol{\rho}'_Z) - \psi_Z(\boldsymbol{\theta}_Z + \boldsymbol{\rho}^*_Z) + \rho'_1 - \rho^*_1 - \rho_0}, \tag{8}$$

where $\boldsymbol{\rho}_Y$ and $\rho_0$ are the conjugation parameters of the linear Guassian model $p(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$; where $\boldsymbol{\rho}'_Z$ and $\rho'_1$ are the conjugation parameters of the mixture model $p(\mathbf{y}, z; \boldsymbol{\theta}'_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{YZ})$ for $\boldsymbol{\theta}'_Y = \boldsymbol{\theta}_Y + \mathbf{x} \cdot \boldsymbol{\Theta}_{XY}$; and where $\boldsymbol{\rho}^*_Z$ and $\rho^*_1$ are the conjugation parameters of the mixture model $p(\mathbf{y}, z; \boldsymbol{\theta}^*_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{YZ})$ for $\boldsymbol{\theta}^*_Y = \boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y$ (see Appendix for a detailed derivation and the conjugation parameters).

With regards to fitting, computing the expectation step for HMoGs reduces to computing the log-partition gradient $\boldsymbol{\tau}_{YZ}$ of a mixture model, which is available in closed-form. The maximization step, on the other hand, does not afford a closed-form expression. Nevertheless, we can solve the maximization step with gradient ascent as long as we can evaluate the gradient $\boldsymbol{\tau}_{XYZ}$ of the HMoG log-partition function $\psi_{XYZ}$, which is indeed possible for an HMoG through judicious use of Eq. 4 (see Appendix). Our strategy for fitting HMoGs is thus to evaluate the expectation-step in closed-form, and then use the Adam gradient optimizer to approximately solve the maximization step [21].

Finally, once we have trained our HMoG model $p(\mathbf{x}, \mathbf{y}, z)$ we may project and classify new observations. The probabilistic structure of HMoGs suggests that the projection of a data point $\mathbf{x}$ is given by $\mathbb{E}[Y \mid X = \mathbf{x}^{(i)}]$, which mathematically involves marginalizing out the cluster index $Z$. Similarly, the cluster membership of a data point is given by $p(z \mid \mathbf{x})$, which involves marginalizing out the features $Y$. In practice these probabilistic definitions of projection and classification perform very similarly to their two-stage versions, and for consistency we use the two-stage projection and classification algorithms with two-stage models, and the probabilistic versions with unified models. For a more detailed discussion see the Appendix.

## Notes on training parameters and implementation details

When fitting two-stage algorithms, we always run 100 iterations of EM each for the dimensionality-reduction models and the MoG, as training times were negligible. We initialized the PCA/FA parameters by setting the mean equal to the empirical mean, the isotropic/diagonal variance to the corresponding empirical value, and initializing the projection matrix randomly with values between -0.01 and 0.01. After fitting PCA/FA, we initialized the MoG parameters by fitting a multivariate normal to the projected data, and setting the mean and covariance of each cluster to a sample point from the fit normal, and the covariance of the fit normal, respectively. We set the weight distribution to be uniform.

Fitting the HMoG PCA/FA algorithms on the Iris and synthetic datasets was possible using the two-stage initialization schemes and a range of learning parameters. For the RNA-Seq data, we initialized an HMoG by first fitting its parameters with a two-stage algorithm, and then running 800 iterations of unified EM. We used relatively conservative learning parameters: a learning rate of $10^{-4}$ for the Adam optimizer, with 2000 steps per EM iteration. To avoid local maxima during optimization, we always ran 10 simulations in parallel and selected the model that maximized training data performance.

All algorithms were implemented in Haskell and targeted the CPU, and simulations were run on an AMD Ryzen 9 5950X processor. By distributing parallel simulations over cores, a single fit of an HMoG FA model on the RNA-Seq data took about 10 minutes.

# 3   Applications

We next show how HMoG theory is useful for dimensionality-reduction and clustering in three application scenarios: (i) a simple validation of the HMoG theory on the *Iris* flower dataset [22]; (ii) a demonstration on synthetic data of the limitations of standard two-stage algorithms, and how HMoGs can overcome them; and (iii) an application to RNA-Seq data [9, 23] to show how HMoGs enhance performance over two-stage approaches on real data.

On each dataset, we compare four methods for dimensionality-reduction and clustering: we fit (i) probabilistic PCA or (ii) FA to the given dataset with EM, followed by separately fitting a MoG to the projected data with EM; or we combine either (iii) probabilistic PCA or (iv) FA with a MoG into an HMoG, and fit the model with the unified EM algorithm. We refer to these methods as *two-stage PCA*, *two-stage FA*, *HMoG PCA*, and *HMoG FA*, respectively (Fig. 1).

## 3.1   HMoGs afford rigorous comparison of high-dimensional clustering methods

As we have shown, two-stage algorithms maximize the likelihood of an equivalent HMoG, even when they optimize dimensionality-reduction and clustering separately. We demonstrate this by comparing the log-likelihood trajectories of all four methods given the *Iris* flower dataset [22] (Fig. 2**A**). We observed that the log-likelihood increased every iteration for all four training algorithms. For this dataset, we saw little performance difference between the four methods for the fully trained model. Qualitatively, the projections and clusters learned by the two-stage PCA method were indeed sufficient for capturing and clustering features of the data (Fig. 2**B**), and the HMoG FA model did not perform noticeably better (Fig. 2**C**). In conclusion, HMoG theory supports quantitative model comparison through evaluation of the HMoG log-likelihood (Eq. 8).
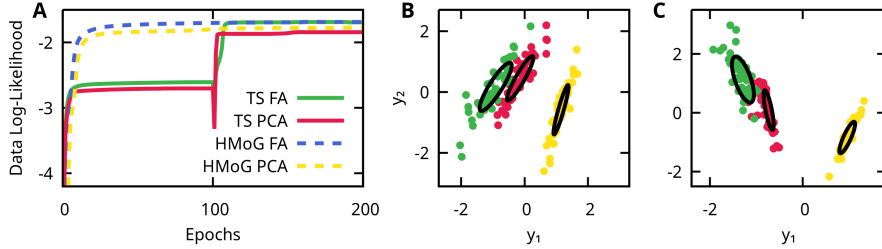
Figure 2: *A simple validation of HMoG theory.* **A:** Ascent of the data log-likelihood $\sum_{i=1}^{d_S} \log p(x^{(i)})$ by two-stage PCA (red) and FA (green), and HMoG PCA (dashed yellow) and HMoG FA (dashed blue) given the *Iris* flower dataset [22]. **B-C:** Projection of the dataset into feature space learned by two-stage PCA (**B**) and HMoG FA (**C**), with colours (red, green, yellow) indicating distinct Iris species, as well as confidence ellipses (black lines) of $p(y \mid z)$ for clusters $z = 0$, 1, and 2.

## 3.2  HMoG EM overcomes limitations of two-stage algorithms

We next compared our methods on a synthetic dataset designed to reveal the limitations of the two-stage approaches (Fig. 3). Data was generated from a ground-truth, HMoG FA model with two clusters, a 1-dimensional feature space, and a 2-dimensional observable space. We designed the ground-truth HMoG so that the dimension along which cluster membership changes in observation space is perpendicular to the direction of maximum variance (Fig. 3**A-C**).

Unsurprisingly, two-stage PCA learned a projection that only followed the direction of maximum variance, and failed to effectively model either the observable density (Fig. 3**A**) or feature density (Fig. 3**D**). Although we do not present the results here, HMoG PCA also failed to effectively model the data in the same way as two-stage PCA, indicating that the limitation is due to the simpler structure of the PCA model with its single $\sigma$ (Fig. 1), rather than the training algorithm used.

On the other hand, two-stage FA performed better than PCA, and learned a projection that captures the feature direction along which cluster membership changes (Fig. 3**B**). This was reflected in the feature distribution learned by the two-stage MoG, which captured the clusters in the feature space of the ground-truth HMoG (Fig. 3**E**). Nevertheless, two-stage FA failed to capture the fine-structure of the observable densities, and in repeated simulations we found that it often learned suboptimal clusterings. In contrast, HMoG FA nearly perfectly modeled the observable density (Fig. 3**C**), and reliably separated the two clusters in feature space (Fig. 3**F**). Overall, we found that FA-based models of dimensionality reduction have important advantages over PCA-based models, which are only fully exploited by training them with the unified, HMoG EM algorithm.
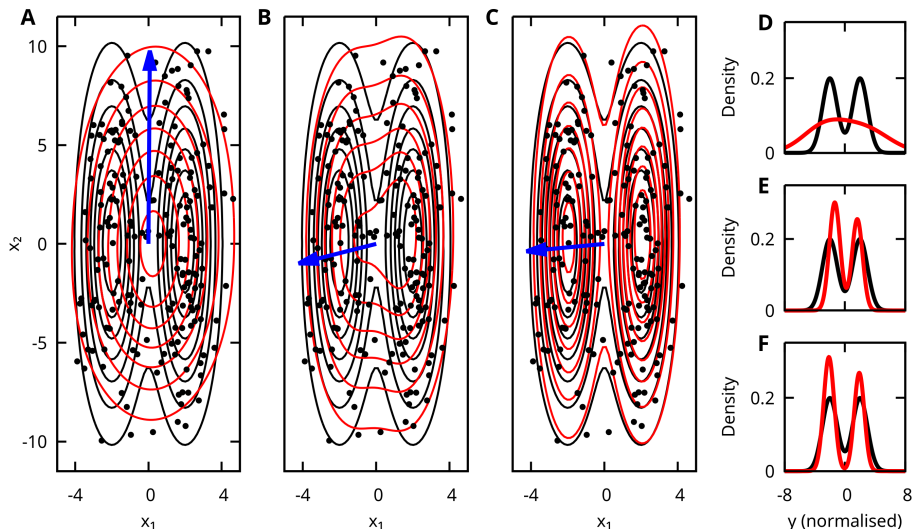
Figure 3: *Limitations of two-stage models.* **A-C:** Contours (black lines) and a sample (black dots) from the observable density $p(x)$ of a ground-truth HMoG, and contours (red lines) and projection direction (blue arrow) of model learned by two-stage PCA (**A**), two-stage FA (**B**), and HMoG FA (**C**). **D-F:** Feature density $p(y)$ of the ground-truth HMoG (black line), compared to feature density (red line) learned by two-stage PCA (**D**), two-stage FA (**E**), and HMoG FA (**F**). Learned density location was normalized by the lengths of the projection vectors.

## 3.3 HMoGs exceed predictive performance of two-stage models on RNA-Seq data

Finally, we applied our methods to a single-cell RNA-Seq dataset from peripheral blood mononuclear cells (PBMCs), a subset of the human immune system. The data contained 3 994 cells with measurements of 15 715 genes [9]. The data was preprocessed by computing analytic Pearson residuals [23] using scanpy 1.9 default settings, and twenty genes with the highest residual variance were selected. The dataset was grouped into eight immune cell subtypes identified by previously known genetic markers by [24], allowing us to compare model classification performance to ground-truth clusters.

We evaluated the predictive performance of all four methods on the RNA-Seq data through 5-fold cross-validation of the log-likelihood, as we varied the number of dimensions and number of clusters in the feature space (Fig. 4**A-C**). In contrast to the Iris flower data, we found that HMoG FA significantly outperformed all the alternative approaches. For example, for four clusters and four features, two-stage PCA, two-stage FA, and HMoG FA achieved predictive log-likelihoods of $-41.27 \pm 0.59$, $-39.33 \pm 0.78$, and $-35.61 \pm 0.37$, respectively. Moreover, HMoG FA could achieve equivalent performance (predictive log-likelihood of $39.12 \pm 0.90$) to the two-stage methods with merely three clusters and three
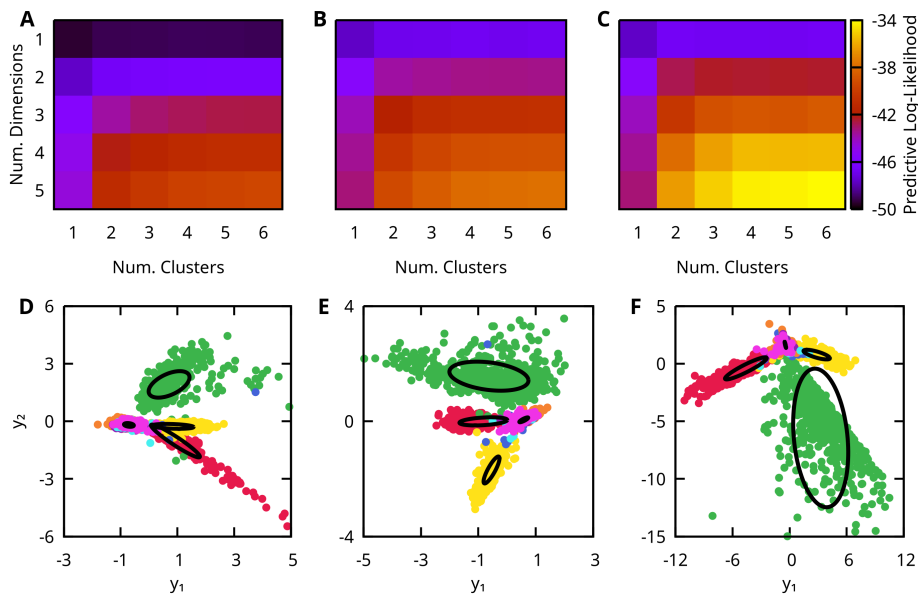
Figure 4: *Dimensionality reduction and clustering with RNA sequencing data.*
**A-C:** Heatmaps of cross-validated average log-likelihood on the data, for the
two-stage PCA (**A**), two-stage FA (**B**), and HMoG FA methods (**C**). **D-F:** Two
representative dimensions of dataset projections into the feature space of two-
stage PCA (**D**), two-stage FA (**E**), and HMoG FA (**F**). Point colours indicate
ground-truth cell subtype identity.

features. We again found that HMoG PCA performance did not significantly
exceed that of two-stage PCA, and so do not present the results here.

To better understand the feature space learned by each model, we refit the
models on the full datasets. For each method, we set the number of clusters and
feature dimensions to four, as we found that predictive performance saturates
around 4 clusters, and chose 4 features for simplicity. We then focused on a
representative plane in the 4-dimensional feature space. Overall we found that
each model successfully separated out three of the cell subtypes, and used a
fourth cluster to represent "everything else" (Fig. 4**D-F**).

Since the feature spaces of the methods were qualitatively similar, we sought
to better understand why the log-likelihood of HMoG FA was so much higher than
that of the other techniques. Due to its hierarchical structure, we can marginalize
out the features in the HMoG density, and then factor the joint distribution
over observations $X$ and cluster indices $Z$ into $p(\mathbf{x}, z) = p(\mathbf{x} \mid z)p(z)$. Given
labelled data $(\mathbf{x}^{(1)}, z^{(1)}, \ldots, \mathbf{x}^{(d_S)}, z^{(d_S)})$, performance may then be quantified as
the combination of classification performance $\sum_{i=1}^{d_S} p(z^{(i)})$ — i.e. how much the
labels under the model match the true labels — and how well the model captures
the data distribution $\sum_{i=1}^{d_S} p(\mathbf{x}^{(i)} \mid z^{(i)})$ Quantitatively, we found no significant

12

difference in classification performance between our four methods when using 4 clusters and 4 features — when we associated each cluster with its most strongly correlated label on the training data, we found that all models achieved held-out classification performance of $\approx 53\%$. When we allowed one cluster to represent multiple labels (to permit one cluster to model the "everything else"), we found each technique achieved held-out classification performance of $\approx 98\%$. Overall, this suggests that each model was able to saturate classification performance given 4 clusters, and that the gains achieved by HMoG FA were entirely due to how well the clusters it learned in the observable space $p(\mathbf{x} \mid z)$ were able to represent the relevant data.

# Discussion

In this paper we have presented a unified model of dimensionality reduction and clustering we call a hierarchical mixture of Gaussians. We showed how the theory of HMoGs generalized existing two-stage algorithms for dimensionality reduction and clustering, allowing us to derive an exact expression for the log-likelihood of these methods, and to enhance their performance with an EM algorithm that trains both stages in parallel. We applied our theory of HMoGs to three datasets, and showed how HMoGs can help existing methods exceed their limitations.

A notable finding in our applications was the performance advantage achieved by FA-based models over PCA-based models. It is well-known that PCA-based dimensionality reduction can fail to capture the dimensions relevant for clustering [10], and it has been suggested that the rescaling properties of FA might help address this [11]. Our applications confirmed the potential advantages of FA over PCA for dimensionality-reduction, both in standard two-stage models and HMoGs. Given the relative dominance of PCA over FA in dimensionality-reduction applications [see 9], our results call for additional investigation of the advantages offered by FA-based methods.

An advantage of our framework is its theoretical simplicity and flexibility, as ultimately we fit HMoGs by maximizing the likelihood using EM, where the maximization step is implemented with a gradient ascent procedure. As such, the basic fitting algorithm we provided could be further enhanced with regularization or sparsity techniques [8, 25] or specialized EM algorithms for large datasets [26].

Finally, our exponential family-based theory of HMoGs opens two directions for further research. On one hand, the theory of conjugation we developed provides a toolbox for designing tractable hierarchical models out of arbitrary exponential families. As such, it should prove useful to researchers who work with data that call for more specialized distributions, such as Poisson distributions for count data. On the other hand, the log-linear structure of our exponential family models makes them highly amenable to GPU-based computation, and should allow optimized implementations of our algorithms to scale up to the massive sizes of modern datasets.

# References

1. Beyer, K., Goldstein, J., Ramakrishnan, R. & Shaft, U. *When Is "Nearest Neighbor" Meaningful?* in *Database Theory — ICDT'99* (Berlin, Heidelberg, 1999).

2. Assent, I. Clustering high dimensional data. *WIREs Data Mining and Knowledge Discovery* **2** (2012).

3. Houdard, A., Bouveyron, C. & Delon, J. High-Dimensional Mixture Models for Unsupervised Image Denoising (HDMI). *SIAM Journal on Imaging Sciences* **11.** Publisher: Society for Industrial and Applied Mathematics (2018).

4. Hertrich, J. *et al.* PCA reduced Gaussian mixture models with applications in superresolution. *Inverse Problems and Imaging* **16.** Company: Inverse Problems and Imaging Distributor: Inverse Problems and Imaging Institution: Inverse Problems and Imaging Label: Inverse Problems and Imaging Publisher: American Institute of Mathematical Sciences (2022).

5. Warren Liao, T. Clustering of time series data—a survey. *Pattern Recognition* **38** (2005).

6. Lewicki, M. S. A review of methods for spike sorting: the detection and classification of neural action potentials. *Network: Computation in Neural Systems* **9.** Publisher: Informa UK Limited (1998).

7. Baden, T. *et al.* The functional diversity of retinal ganglion cells in the mouse. *Nature* **529.** Number: 7586 Publisher: Nature Publishing Group (2016).

8. Witten, D. M. & Tibshirani, R. A Framework for Feature Selection in Clustering. *Journal of the American Statistical Association* **105.** Publisher: Taylor & Francis _eprint: https://doi.org/10.1198/jasa.2010.tm09415 (2010).

9. Duò, A., Robinson, M. D. & Soneson, C. A systematic performance evaluation of clustering methods for single-cell RNA-seq data. *F1000Research* **7** (2020).

10. Chang, W.-C. On Using Principal Components before Separating a Mixture of Two Multivariate Normal Distributions. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* **32.** _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.2307/2347949 (1983).

11. McLachlan, G. J., Lee, S. X. & Rathnayake, S. I. Finite Mixture Models. *Annual Review of Statistics and Its Application* **6** (2019).

12. Ding, C., He, X., Zha, H. & Simon, H. *Adaptive dimension reduction for clustering high dimensional data* in *2002 IEEE International Conference on Data Mining, 2002. Proceedings.* (2002).

13. Fern, X. Z. & Brodley, C. E. *Random projection for high dimensional data clustering: a cluster ensemble approach* in *Proceedings of the Twentieth International Conference on International Conference on Machine Learning* (Washington, DC, USA, 2003).

14. Ghahramani, Z. & Hinton, G. E. *The EM algorithm for mixtures of factor analyzers* tech. rep. (Technical Report CRG-TR-96-1, University of Toronto, 1996).

15. Tipping, M. E. & Bishop, C. M. Mixtures of Probabilistic Principal Component Analyzers. *Neural Computation* **11** (1999).

16. Roweis, S. *EM Algorithms for PCA and SPCA* in *Advances in Neural Information Processing Systems* **10** (1997).

17. Bishop, C. M. *Pattern recognition and machine learning* (New York, 2006).

18. Smolensky, P. *Information Processing in Dynamical Systems: Foundations of Harmony Theory* tech. rep. Section: Technical Reports (COLORADO UNIV AT BOULDER DEPT OF COMPUTER SCIENCE, 1986).

19. Welling, M., Rosen-zvi, M. & Hinton, G. E. in *Advances in Neural Information Processing Systems 17* (2005).

20. Wainwright, M. J. & Jordan, M. I. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning* **1** (2008).

21. Kingma, D. & Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

22. *Iris flower data set* Page Version ID: 1080211522. 2022.

23. Lause, J., Berens, P. & Kobak, D. Analytic Pearson residuals for normalization of single-cell RNA-seq UMI data. *Genome Biology* **22** (2021).

24. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nature Communications* **8.** Number: 1 Publisher: Nature Publishing Group (2017).

25. Yi, X. & Caramanis, C. *Regularized EM Algorithms: A Unified Framework and Statistical Guarantees* in *Advances in Neural Information Processing Systems* **28** (2015).

26. Chen, J., Zhu, J., Teh, Y. W. & Zhang, T. *Stochastic Expectation Maximization with Variance Reduction* in *Advances in Neural Information Processing Systems* **31** (2018).

27. Amari, S.-i. & Nagaoka, H. *Methods of information geometry* (2007).

28. Neal, R. M. & Hinton, G. E. in *Learning in graphical models* (1998).

29. Hinton, G. E., Osindero, S. & Teh, Y. W. A fast learning algorithm for deep belief nets. *Neural computation* **18** (2006).

30. Salakhutdinov, R. Learning Deep Generative Models. *Annual Review of Statistics and Its Application* **2** (2015).

31. Kingma, D. P. & Welling, M. An Introduction to Variational Autoencoders. *Foundations and Trends® in Machine Learning* **12.** arXiv: 1906.02691 (2019).

# A  Appendix

Towards developing and explaining the theory of HMoGs, this appendix is organized into the following sections:

1. A brief introduction to exponential families, largely in order to fix notation.

2. An overview of the theory of exponential family graphical models known as harmoniums.

3. Derivation of the theory of conjugated harmoniums.

4. Formulations of mixture models and linear Gaussian models as conjugated harmoniums.

5. Combining these various components to formalize an HMoG.

The culmination of these theoretical developments is a series of algorithms for computing the log-likelihood and expectation-maximization algorithms for HMoGs.

## A.1  Exponential families

For a thorough development of exponential family theory see Amari & Nagaoka (2007), and Wainwright & Jordan (2008).

Consider a random variable $X \in \mathcal{X}$ on the sample space $\mathcal{X}$, with an unknown distribution $P_X$, and suppose $\{X^{(i)}\}_{i=1}^{n}$ is an independent and identically distributed sample from $P_X$, such that $X^{(i)} \sim P_X$. We may model $P_X$ based on the sample $\{X^{(i)}\}_{i=1}^{n}$ by defining a statistic $\mathbf{s}_X \colon \mathcal{X} \to \mathrm{H}_X$, where $\mathrm{H}_X$ is a space with dimension $d_X$, and looking for a probability distribution $Q_X$ that satisfies $\mathbb{E}_Q[\mathbf{s}_X(X)] = \frac{1}{n} \sum_{i=1}^{n} \mathbf{s}_X(X^{(i)})$, where $\mathbb{E}_Q[f(X)] = \int_{\mathcal{X}} f dQ_X$ is the expected value of $f(X)$ under $Q_X$. On its own this is an under-constrained problem, but if we further assume that $Q_X$ must have maximum entropy, then we may define a family (or manifold) of distributions $\mathcal{M}_X$ that is uniquely described by the possible values of $\frac{1}{n} \sum_{i=1}^{n} \mathbf{s}_X(X^{(i)})$.

A $d_X$-dimensional *exponential family* $\mathcal{M}_X$ is defined by a *sufficient statistic* $\mathbf{s}_X$, as well as a *base measure* $\mu_X$ which helps define integrals and expectations within the family. An exponential family is parameterized by a set of *natural parameters* $\Theta_X$, such that each element of the family $Q_X \in \mathcal{M}_X$ may be identified with some parameters $\boldsymbol{\theta}_X \in \Theta_X$. The density of the distribution $Q_X$ is given by $\log q(x) = \mathbf{s}_X(x) \cdot \boldsymbol{\theta}_X - \psi_X(\boldsymbol{\theta}_X)$, where $\psi_X(\boldsymbol{\theta}_X) = \log \int_{\mathcal{X}} e^{\mathbf{s}_X(x) \cdot \boldsymbol{\theta}_X} \mu_X(dx)$ is the log-partition function. Expectations of any $Q_X \in \mathcal{M}_X$ are then given by $\mathbb{E}_Q[f(X)] = \int_{\mathcal{X}} f dQ_X = \int_{\mathcal{X}} f \cdot q d\mu_X$. Because each $Q_X \in \mathcal{M}_X$ is uniquely defined by $\mathbb{E}_Q[\mathbf{s}_X(X)]$, the means of the sufficient statistic also parameterize $\mathcal{M}_X$. The space of all *mean parameters* is $\mathrm{H}_X$, and we denote them by $\boldsymbol{\eta}_X = \mathbb{E}_Q[\mathbf{s}_X(X)]$. Finally, a sufficient statistic is *minimal* when its component functions are non-constant and linearly independent. If the sufficient statistic $\mathbf{s}_X$ of a given family $\mathcal{M}_X$ is minimal, then $\Theta_X$ and $\mathrm{H}_X$ are isomorphic. Moreover,

the transition functions between them $\boldsymbol{\tau}_X \colon \Theta_X \to \mathrm{H}_X$ and $\boldsymbol{\tau}_X^{-1} \colon \mathrm{H}_X \to \Theta_X$ are given by $\boldsymbol{\tau}_X(\boldsymbol{\theta}_X) = \partial_{\boldsymbol{\theta}_X} \psi_X(\boldsymbol{\theta}_X)$, and $\boldsymbol{\tau}_X^{-1}(\boldsymbol{\eta}_X) = \partial_{\boldsymbol{\eta}_X} \phi_X(\boldsymbol{\eta}_X)$, where $\phi_X(\boldsymbol{\eta}_X) = \mathbb{E}_Q[\log q(X)]$ is the negative entropy of $Q_X$. The transition functions $\boldsymbol{\tau}_X$ and $\boldsymbol{\tau}_X^{-1}$ are also referred to as the forward and backward mappings, respectively.

## A.2  Exponential Family Harmoniums

An *exponential family harmonium* is a kind of product exponential family which includes restricted Boltzmann machines as a special case [18, 19]. We may construct an exponential family harmonium $\mathcal{M}_{XY}$ out of the exponential families $\mathcal{M}_X$ and $\mathcal{M}_Y$ by defining the base measure of $\mathcal{M}_{XY}$ as the product measure $\mu_X \cdot \mu_Y$, and by defining the sufficient statistic of $\mathcal{M}_{XY}$ as the vector which contains the concatenation of all the elements in $\mathbf{s}_X$, $\mathbf{s}_Y$, and the outer product $\mathbf{s}_X \otimes \mathbf{s}_Y$. More concretely, $\mathcal{M}_{XY}$ is the manifold that contains all the distributions $Q_{XY} \in \mathcal{M}_{XY}$ with densities of the form $q(x, y) \propto e^{\mathbf{s}_X(x) \cdot \boldsymbol{\theta}_X + \mathbf{s}_Y(y) \cdot \boldsymbol{\theta}_Y + \mathbf{s}_X(x) \cdot \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y)}$, where $\boldsymbol{\theta}_X$, $\boldsymbol{\theta}_Y$, and $\boldsymbol{\Theta}_{XY}$ are the natural parameters of $Q_{XY}$.

The linear structure of harmoniums affords simple expressions for their training algorithms. On one hand, given a sample $\{X^{(i)}\}_{i=1}^n$ and a harmonium $Q_{XY} \in \mathcal{M}_{XY}$ with parameters $\boldsymbol{\theta}_X$, $\boldsymbol{\theta}_Y$, and $\boldsymbol{\Theta}_{XY}$, an iteration of the expectation-maximization algorithm (EM) may be formulated as:

**Expectation Step:** compute the unobserved means $\boldsymbol{\eta}_{Y,i} = \boldsymbol{\tau}_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(X^{(i)}) \cdot \boldsymbol{\Theta}_{XY})$ for every $i$,

**Maximization Step:** eval. $\boldsymbol{\tau}_{XY}^{-1}(\frac{1}{d_S} \sum_{i=1}^{d_S} \mathbf{s}_X(X^{(i)}), \frac{1}{d_S} \sum_{i=1}^{d_S} \boldsymbol{\eta}_{Y,i}, \frac{1}{d_S} \sum_{i=1}^{d_S} \mathbf{s}_X(X^{(i)}) \otimes \boldsymbol{\eta}_{Y,i})$.

On the other hand, the stochastic log-likelihood gradients of the parameters of $Q_{XY}$ are

$$\partial_{\boldsymbol{\theta}_X} \log q(X^{(i)}) = \mathbf{s}_X(X^{(i)}) - \boldsymbol{\eta}_X,$$
$$\partial_{\boldsymbol{\theta}_Y} \log q(X^{(i)}) = \boldsymbol{\tau}_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(X^{(i)}) \cdot \boldsymbol{\Theta}_{XY}) - \boldsymbol{\eta}_Y,$$
$$\partial_{\boldsymbol{\Theta}_{XY}} \log q(X^{(i)}) = \mathbf{s}_X(X^{(i)}) \otimes \boldsymbol{\tau}_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(X^{(i)}) \cdot \boldsymbol{\Theta}_{XY}) - \mathbf{H}_{XY}.$$

where $(\boldsymbol{\eta}_X, \boldsymbol{\eta}_Y, \mathbf{H}_{XY}) = \boldsymbol{\tau}_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$. It is often the case that we have a closed-form expression for $\boldsymbol{\tau}_{XY}$, but not for $\boldsymbol{\tau}_{XY}^{-1}$, and therefore cannot compute the maximization step in closed-form. Although we could simply train a harmonium with gradient ascent on the log-likelihood, for certain models this can lead to numerical instability, and an effective strategy is rather to evaluate the expectation step, and then approximate the maximization step with gradient descent.

## A.3  Harmoniums and Conjugate Priors

The conditional distributions $Q_{X|Y}$ and $Q_{Y|X}$ of a harmonium $Q_{XY}$ have a simple linear structure, and are themselves always exponential family distributions; in particular, $Q_{X|Y=y} \in \mathcal{M}_X$ and $Q_{Y|X=x} \in \mathcal{M}_Y$ for any $y$ or $x$, respectively. In

general, however, the marginal distributions $Q_X$ and $Q_Y$ are not members $\mathcal{M}_X$ and $\mathcal{M}_Y$, respectively. This is both a blessing and a curse, as on one hand, the marginal distribution $Q_X$ can model much more complex datasets than the simpler elements of $\mathcal{M}_X$. On the other hand, because the prior $Q_Y$ may not be computationally tractable, various computations with harmoniums, such as sampling and inference, may also prove intractable.

Ideally, the prior $Q_Y$ would be in $\mathcal{M}_Y$ to facilitate computability, while the modelled observable distribution $Q_X$ would be more complex than the elements of $\mathcal{M}_X$, and some classes of harmoniums do indeed have this structure. In general, a prior is said to be conjugate to a posterior if both the prior and posterior have the same form for any observation $x$. In the context of harmoniums, since the posterior $Q_{Y|X=x} \in Q_Y$ for any $x$, the prior $Q_Y$ is conjugate to the posterior iff $Q_Y \in \mathcal{M}_Y$. We refer to harmoniums with conjugate priors as *conjugated harmoniums*.

**Lemma 1.** *Suppose that $\mathcal{M}_{XY}$ is a harmonium family defined by the exponential families $\mathcal{M}_X$ and $\mathcal{M}_Y$, and that $Q_{XY} \in \mathcal{M}_{XY}$ has parameters $(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$. Then $Q_{XY}$ is conjugated if and only if there exists a vector $\boldsymbol{\rho}_Y$ and a constant $\rho_0$ such that*

$$\psi_X(\boldsymbol{\theta}_X + \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y)) = \mathbf{s}_Y(y) \cdot \boldsymbol{\rho}_Y + \rho_0, \tag{9}$$

*for any $\mathbf{x} \in \Omega_X$, and where*

$$\boldsymbol{\theta}_Y^* = \boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y. \tag{10}$$

*Proof.* In general, the density of the harmonium prior distribution $Q_Y$ is given by

$$q(y) = e^{\mathbf{s}_Y(y) \cdot \boldsymbol{\theta}_Y + \psi_X(\boldsymbol{\theta}_X + \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y)) - \psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})}. \tag{11}$$

On one hand, if we assume that $Q_{XY}$ is conjugated, then it also holds that $Q_Y \in \mathcal{M}_Y$ with some natural parameters $\boldsymbol{\theta}_Y^*$, and therefore

$$q(y) \propto e^{\boldsymbol{\theta}_Y \cdot \mathbf{s}_Y(x) + \psi_X(\boldsymbol{\theta}_X + \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y))} \propto e^{\boldsymbol{\theta}_Y^* \cdot \mathbf{s}_Y(y)}$$

$$\implies \quad \boldsymbol{\theta}_Y \cdot \mathbf{s}_Y(y) + \psi_X(\boldsymbol{\theta}_X + \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y)) = \boldsymbol{\theta}_Y^* \cdot \mathbf{s}_Y(y) + \rho_0$$

$$\implies \quad \psi_X(\boldsymbol{\theta}_X + \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y)) = \mathbf{s}_Y(y) \cdot \boldsymbol{\rho}_Y + \rho_0.$$

for some $\rho_0$, and $\boldsymbol{\rho}_Y = \boldsymbol{\theta}_Y^* - \boldsymbol{\theta}_Y$.

On the other hand, if we first assume that Eq. 9 holds, then $Q_Y$ is given by

$$q(y) \propto e^{\boldsymbol{\theta}_Y \cdot \mathbf{s}_Y(y) + \psi_X(\boldsymbol{\theta}_X + \boldsymbol{\Theta}_{XY} \cdot \mathbf{s}_Y(y))} \propto e^{(\boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y) \cdot \mathbf{s}_Y(y)}, \tag{12}$$

which implies that $Q_X \in \mathcal{M}_X$ with parameters $\boldsymbol{\theta}_X + \boldsymbol{\rho}_X$. $\qquad\square$

We refer to $\boldsymbol{\rho}_X$ and $\rho_0$ as the conjugation parameters. The value of this lemma is that it allows us to reduce various computations to evaluating the conjugation parameters, and show how these computations are fundamentally the same even for apparently unrelated models. To wit, the following corollary describes how to compute the log-partition function of a conjugated harmonium.

**Corollary 2.** *Suppose that $\mathcal{M}_{XY}$ is a harmonium family defined by the exponential families $\mathcal{M}_X$ and $\mathcal{M}_Y$, and that the harmonium $Q_{XY} \in \mathcal{M}_{XY}$ with parameters $(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})$ is conjugated, with prior $Q_Y \in \mathcal{M}_Y$ and parameters $\boldsymbol{\theta}_Y^*$. Then $Q_{XY}$ satisfies*

$$\psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY}) = \psi_Y(\boldsymbol{\theta}_Y^*) + \rho_0.$$

*Proof.* Lemma 1 implies that the prior $Q_Y$ of the conjugated harmonium $Q_{XY}$ satisfies

$$
\begin{aligned}
q(y) &= e^{\boldsymbol{\theta}_Y^* \cdot \mathbf{s}_Y(y) - \psi_Y(\boldsymbol{\theta}_Y^*)} \\
&= e^{\boldsymbol{\theta}_Y \cdot \mathbf{s}_Y(y) + \boldsymbol{\rho}_Y \cdot \mathbf{s}_Y(y) + \rho_0 - \psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY})} \\
\iff \quad \psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY}) &= (\boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y - \boldsymbol{\theta}_Y^*) \cdot \mathbf{s}_Y(y) + \psi_Y(\boldsymbol{\theta}_Y^*) + \rho_0 \\
&= \psi_Y(\boldsymbol{\theta}_Y^*) + \rho_0.
\end{aligned}
$$

$\square$

Typically, evaluating the log-partition function $\psi_Y$ of $\mathcal{M}_Y$ will be tractable, which means that various computations, including evaluating the density of $Q_X$, are also tractable. In particular, given the conjugated harmonium $Q_{XY} \in \mathcal{M}_{XY}$ with parameters $\boldsymbol{\theta}_X$, $\boldsymbol{\theta}_Y$, and $\boldsymbol{\Theta}_{XY}$

$$
\begin{aligned}
\log q(x) &= \mathbf{s}_X(x) \cdot \boldsymbol{\theta}_X + \psi_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(x) \cdot \boldsymbol{\Theta}_{XY}) - \psi_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\Theta}_{XY}, \boldsymbol{\theta}_Y) \\
&= \mathbf{s}_X(x) \cdot \boldsymbol{\theta}_X + \psi_Y(\boldsymbol{\theta}_Y + \mathbf{s}_X(x) \cdot \boldsymbol{\Theta}_{XY}) - \psi_Y(\boldsymbol{\theta}_Y^*) - \rho_0,
\end{aligned}
$$

where $\boldsymbol{\theta}_Y^* = \boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y$.

## A.4 Mixture Models and Linear Guassian Models

Constructing a hierarchical mixture of Gaussians ultimately requires putting together a mixture model and a linear Gaussian model in the "right" way. To this do we use the theory of conjugated harmoniums, by first defining the exponential families of categorical and normal distributions, and then deriving the conjugation parameters of mixture models and linear Gaussian models.

The $d_Z$-dimensional categorical exponential family $\mathcal{M}_Z$ contains all the distributions over integer values between 0 and $d_Z$. The base measure of $\mathcal{M}_Z$ is the counting measure, and the $k$th element of its sufficient statistic $\mathbf{s}_Z(k)$ at $j$ is 1, and 0 for any other elements. The sufficient statistic $\mathbf{s}_Z$ is thus a vector of all zeroes when $j = 0$, and all zeroes except for the $j$th element when $j > 0$. Finally, the log-partition function is $\psi_Z(\boldsymbol{\theta}_Z) = \log\left(1 + \sum_{i=1}^{d_K} e^{\theta_{Z,i}}\right)$, and the forward mapping is given by $\tau_{Z,i}(\boldsymbol{\theta}_Z) = \frac{e^{\theta_{Z,i}}}{1 + \sum_{i=1}^{d_Z} e^{\theta_{Z,i}}}$.

Now, a mixture distribution is simply a harmonium defined by $\mathcal{M}_X$ and $\mathcal{M}_Y$, where $\mathcal{M}_Y$ is the family of categorical distributions. For a mixture model $Q_{XY}$, the observable distribution $Q_X$ can typically model distributions (e.g. multimodal distributions) outside of $\mathcal{M}_X$, yet the prior $Q_Y$ is indeed in $\mathcal{M}_Y$.

**Algorithm 1** Computing mixture model conjugation parameters
___
**Require:** Mixture model natural parameters $\boldsymbol{\theta}_Y, \boldsymbol{\theta}_{YZ}$
**Ensure:** Mixture model conjugation parameters $\rho_1, \boldsymbol{\rho}_Z$
   **function** MIXTURECONJUGATIONPARAMETERS($\boldsymbol{\theta}_Y, \boldsymbol{\theta}_{YZ}$)
      $\rho_1 \leftarrow \psi_Y(\boldsymbol{\theta}_Y)$
      **for all** $i \in \{1, \ldots, d_Z\}$ **do**
         $\rho_{Z,i} \leftarrow \psi_Y(\boldsymbol{\theta}_Y + \boldsymbol{\Theta}_{YZ} \cdot \mathbf{s}_Z(i)) - \rho_1$
      **end for**
      $\boldsymbol{\rho}_Z = (\rho_{Z,1}, \ldots, \rho_{Z,d_Z})$
      **return** $\rho_1, \boldsymbol{\rho}_Z$
   **end function**
___

We show how to compute the conjugation parameters of a mixture model in Algorithm 1, and how to compute the forward mapping $\boldsymbol{\tau}_{YZ}$ in Algorithm 2.

The $d_X$-dimensional (multivariate) normal exponential family $\mathcal{M}_X$ has base measure $\mu(\mathbf{x}) = \frac{1}{2}d_X \log(2\pi)$, and sufficient statistic $\mathbf{s}_X(\mathbf{x}) = (\mathbf{x}, \mathbf{x} \otimes \mathbf{x})$, where $\otimes$ is the outer product (also to ensure that $\mathbf{s}_X$ is minimal we should only include either the upper or lower triangular part of $\mathbf{x} \otimes \mathbf{x}$, although this detail is often elided in practice).

Linear Gaussian models include various well known constructions from factor analysis and principle component analysis, to the emission and transition distributions of Kalman filters [17]. A linear Gaussian model can be interpreted as a subset of the harmonium family $\mathcal{M}_{XY}$ where both $\mathcal{M}_X$ and $\mathcal{M}_Y$ are families of multivariate normal distributions, restricted so that $\boldsymbol{\Theta}_{XY}$ only captures models interactions between the first-order statistics. More concretely, a linear Gaussian family $\mathcal{M}_{XY}$ contains all distributions of the form

$$p(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{XY}) \propto e^{\mathbf{s}_X(\mathbf{x}) \cdot \boldsymbol{\theta}_X + \mathbf{s}_Y(\mathbf{y}) \cdot \boldsymbol{\theta}_Y + \mathbf{x} \cdot \boldsymbol{\Theta}_{XY} \cdot \mathbf{y}}.$$

The prior of a linear Gaussian model is always a multivariate normal distribution, and is therefore conjugate to the posterior. We may more succinctly express of the conjugation parameters and forward mappings of the linear Gaussian model with the decomposing the parameters $\boldsymbol{\theta}_X$ and $\boldsymbol{\theta}_Y$ into the parameters corresponding to the first and second order statistics, so that $\boldsymbol{\theta}_X = (\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX})$, and $\boldsymbol{\theta}_Y = (\boldsymbol{\theta}_Y^\mu, \boldsymbol{\Theta}_{YY})$, respectively. Using this extra notation, we present the computation of the conjugation parameters in Algorithm 3 algorithm in Algorithm 4.

For a general linear Guassian distribution $Q_{XY}$, the observable distribution $Q_X$ is also in $\mathcal{M}_X$, and so one may wonder what is gained with the linear Gaussian construction. Firstly, the modelling power of $\mathcal{M}_{XY}$ becomes more interesting when one further restricts the observable covariance parameters $\boldsymbol{\Theta}_{XX}$ to have a simpler structure, such as isotropic or diagonal. This makes $\mathcal{M}_{XY}$ essentially equivalent to probabilistic PCA or factor analysis, respectively, and thereby affords modelling high dimensional datasets $\{X^{(i)}\}_{i=1}^{n}$ with a model that better scales with the dimension of the sample points $X^{(i)}$. Secondly, by

**Algorithm 2** The mixture model forward mapping

---

**Require:** Mixture model natural parameters $\boldsymbol{\theta}_{YZ} = (\boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{YZ})$
**Ensure:** Mixture model expectations $(\boldsymbol{\eta}_Y, \boldsymbol{\eta}_Z, \mathbf{H}_{YZ}) = \boldsymbol{\tau}_{YZ}(\boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{YZ})$
  **function** $\boldsymbol{\tau}_{YZ}(\boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{YZ})$
    $\boldsymbol{\rho}_Z, \rho_1 \leftarrow \text{MIXTURECONJUGATIONPARAMETERS}(\boldsymbol{\theta}_Y, \boldsymbol{\theta}_{YZ})$
    $\boldsymbol{\eta}_Z \leftarrow \boldsymbol{\tau}_Z(\boldsymbol{\theta}_Z + \boldsymbol{\rho}_Z)$
    **for** $i = 1, \ldots, d_Z$ **do**
      $\boldsymbol{\eta}_{Y,i} \leftarrow \boldsymbol{\tau}_Y(\boldsymbol{\theta}_Y + \boldsymbol{\Theta}_{YZ} \cdot \mathbf{s}_Z(i))$
      **if** $i = 1$ **then**
        $\boldsymbol{\eta}_Y \leftarrow (1 - \sum_{j=1}^{d_Z-1} \eta_{Z,j})\boldsymbol{\eta}_{Y,1}$
      **else**
        $\boldsymbol{\eta}_Y \leftarrow \boldsymbol{\eta}_Y + \eta_{Z,i-1}\boldsymbol{\eta}_{Y,i}$
      **end if**
    **end for**
    $\mathbf{H}_{YZ} \leftarrow \begin{bmatrix} \boldsymbol{\eta}_{Y,2} \\ \vdots \\ \boldsymbol{\eta}_{Y,d_Z} \end{bmatrix}$
    **return** $(\boldsymbol{\eta}_Y, \boldsymbol{\eta}_Z, \mathbf{H}_{YZ})$
  **end function**

---

expressing mixture models and linear Guassian models in the same language, we can combine them to express the theory of hierarchical mixtures of Gaussians.

## A.5 Hierarchical Mixtures of Gaussians

Let us now begin constructing an HMoG, and more formally impose the PCA or FA structure. Suppose that $\mathcal{M}_X$ is the multivariate normal family with isotropic or diagonal covariance matrices, respectively, that $\mathcal{M}_Y$ is the family of multivariate normals with full covariances, and that $\mathcal{M}_Z$ is the family of categorical distributions. Let $\mathcal{M}_{XY}$ be the harmonium family defined by $\mathcal{M}_X$ and $\mathcal{M}_Y$, and let $\mathcal{M}_{YZ}$ be the harmonium family defined by $\mathcal{M}_Y$ and $\mathcal{M}_Z$. We then define the HMoG family as the harmonium family $\mathcal{M}_{XYZ}$ defined by $\mathcal{M}_X$ and $\mathcal{M}_{YZ}$, restricted to distributions $Q_{XYZ}$ with densities given by

$$q(\mathbf{x}, \mathbf{y}, z) \propto e^{\boldsymbol{\theta}_X \cdot \mathbf{s}_X(\mathbf{x}) + \boldsymbol{\theta}_Y \cdot \mathbf{s}_Y(\mathbf{y}) + \boldsymbol{\theta}_Z \cdot \mathbf{s}_Z(z) + \mathbf{x} \cdot \boldsymbol{\Theta}_{XY} \cdot \mathbf{y} + \mathbf{s}_Y(\mathbf{y}) \cdot \boldsymbol{\Theta}_{YZ} \cdot \mathbf{s}_Z(z)}. \quad (13)$$

To begin our development of HMoG techniques, we first show that the two-stage training algorithm raises a lower bound on the HMoG log-likelihood. Suppose $X^{(1)}, \ldots, X^{(n)}$ is a sample, $Q_{X|Y}$ is an arbitrary linear Gaussian likelihood from a distribution $Q_{XY} \in \mathcal{M}_{XY}$ with a standard normal prior $Q_Y$, and that the MoG distribution $Q_{YZ} \in \mathcal{M}_{YZ}$ is equal to a standard normal (e.g. by defining all components by $q(\mathbf{y} \mid z) = \mathrm{N}(\mathbf{0}, \mathbf{I})$). Then then $Q_{XYZ} \in \mathcal{M}_{XYZ}$ with density $q(\mathbf{x}, \mathbf{y}, z) = q(\mathbf{x} \mid \mathbf{y})q(\mathbf{y}, z)$ is equal to the linear Gaussian distribution $Q_{XY}$. Since the first stage of the two-stage algorithm is to fit $Q_{XY}$ with EM,

---
**Algorithm 3** Computing linear Gaussian model conjugation parameters
---
**Require:** Linear Gaussian model natural parameters $\boldsymbol{\theta}_X = (\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX})$, and $\boldsymbol{\Theta}_{XY}$

**Ensure:** Linear Gaussian model conjugation parameters $\rho_0$, and $\boldsymbol{\rho}_Y = (\boldsymbol{\rho}_Y^\mu, \mathbf{P}_{YY})$

    **function** GAUSSIANCONJUGATIONPARAMETERS($\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\Theta}_{XY}$)

        $\rho_0 \leftarrow -\frac{1}{4}\boldsymbol{\theta}_X^\mu \cdot \boldsymbol{\Theta}_{XX}^{-1} \cdot \boldsymbol{\theta}_X^\mu - \frac{1}{2}\log|-2\boldsymbol{\Theta}_{XX}|$

        $\boldsymbol{\rho}_Y^\mu \leftarrow -\frac{1}{2}\boldsymbol{\Theta}_{YX} \cdot \boldsymbol{\Theta}_{XX}^{-1} \cdot \boldsymbol{\theta}_X^\mu$

        $\mathbf{P}_{YY} \leftarrow -\frac{1}{4}\boldsymbol{\Theta}_{YX} \cdot \boldsymbol{\Theta}_{XX}^{-1} \cdot \boldsymbol{\Theta}_{XY}$

        $\boldsymbol{\rho}_Y \leftarrow (\boldsymbol{\rho}_Y^\mu, \mathbf{P}_{YY})$

        **return** $\rho_0, \boldsymbol{\rho}_Y$

    **end function**
---

---
**Algorithm 4** Computing the Gaussian model forward mapping
---
**Require:** Linear Gaussian model natural parameters $(\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\theta}_Y^\mu, \boldsymbol{\Theta}_{YY}, \boldsymbol{\Theta}_{XY})$

**Ensure:** Expectations $(\boldsymbol{\eta}_X^\mu, \mathbf{H}_{XX}, \boldsymbol{\eta}_Y^\mu, \mathbf{H}_{YY}, \mathbf{H}_{XY}) = \boldsymbol{\tau}_{XY}(\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\theta}_Y^\mu, \boldsymbol{\Theta}_{YY}, \boldsymbol{\Theta}_{XY})$

    **function** $\boldsymbol{\tau}_{XY}(\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\theta}_Y^\mu, \boldsymbol{\Theta}_{YY}, \boldsymbol{\Theta}_{XY})$

$$\boldsymbol{\Sigma} \leftarrow -\frac{1}{2}\begin{bmatrix} \boldsymbol{\Theta}_{XX} & \boldsymbol{\Theta}_{XY} \\ \boldsymbol{\Theta}_{YX} & \boldsymbol{\Theta}_{YY} \end{bmatrix}^{-1}$$

        $(\boldsymbol{\eta}_X^\mu, \boldsymbol{\eta}_Y^\mu) \leftarrow \boldsymbol{\Sigma} \cdot (\boldsymbol{\theta}_X^\mu, \boldsymbol{\theta}_Y^\mu)$

$$\begin{bmatrix} \mathbf{H}_{XX} & \mathbf{H}_{XY} \\ \mathbf{H}_{YX} & \mathbf{H}_{YY} \end{bmatrix} \leftarrow \boldsymbol{\Sigma} + (\boldsymbol{\eta}_X^\mu, \boldsymbol{\eta}_Y^\mu) \otimes (\boldsymbol{\eta}_X^\mu, \boldsymbol{\eta}_Y^\mu)$$

        **return** $(\boldsymbol{\eta}_X^\mu, \mathbf{H}_{XX}, \boldsymbol{\eta}_Y^\mu, \mathbf{H}_{YY}, \mathbf{H}_{XY})$

    **end function**
---

the first stage is equivalent to maximizing the log-likelihood $Q_{XYZ}$ with a fixed standard normal prior, given the sample to $X^{(1)}, \ldots, X^{(n)}$.

Given an HMoG $Q_{XYZ}$ trained by the first stage, the second stage is then to fit a MoG to the projected dataset $\mathbf{y}^{(1)}, \ldots, \mathbf{y}^{(n)}$ defined by $\mathbf{y}^{(i)} = \mathbb{E}_Q[Y \mid X = X^{(i)}]$. While this does not maximize the likelihood of $Q_{XYZ}$ directly, it does approximately increase the evidence lower-bound on $Q_{XYZ}$ [28–31]. What this bound tells us is that we may increase a lower bound on the data log-likelihood by maximizing the log-likelihood of a sample $Y^{(1)}, \ldots, Y^{(n)}$ in the feature space, where each $Y^{(i)} \sim Q_{Y|X=X^{(i)}}$. Although sampled features $Y^{(i)}$ and projected datapoints $\mathbf{y}^{(i)}$ are not in general equivalent, in the Gaussian case where the projection is linear, they are fairly interchangeable. Similarly, this equivalence accounts for why the cluster probabilities $Q_{Z|X}$ of a data point $X$ can be evaluated by either marginalizing out the feature space, or instead projecting $X$ into feature space via $Q_{Y|X}$, and then and evaluating the cluster via probabilities $Q_{Z|Y}$.

To apply HMoGs, we provide Algorithms 5, 6, and 7, for evaluating the

**Algorithm 5** Computing the HMoG log-density function

**Require:** HMoG natural parameters $(\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{XY}, \boldsymbol{\Theta}_{YZ})$
**Require:** Sample point $\mathbf{x}$
**Ensure:** HMoG log-density $\log p(\mathbf{x})$

$\quad \boldsymbol{\theta}_Y' \leftarrow \boldsymbol{\theta}_Y + \mathbf{x} \cdot \boldsymbol{\Theta}_{XY}$
$\quad \rho_1', \boldsymbol{\rho}_Z' \leftarrow \text{MixtureConjugationParameters}(\boldsymbol{\theta}_Y', \boldsymbol{\theta}_{YZ})$
$\quad \boldsymbol{\theta}_Z' \leftarrow \boldsymbol{\theta}_Z + \boldsymbol{\rho}_Z'$
$\quad \rho_0, \boldsymbol{\rho}_Y \leftarrow \text{GaussianConjugationParameters}(\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\Theta}_{XY})$
$\quad \boldsymbol{\theta}_Y^* \leftarrow \boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y$
$\quad \rho_1^*, \boldsymbol{\rho}_Z^* \leftarrow \text{MixtureConjugationParameters}(\boldsymbol{\theta}_Y^*, \boldsymbol{\theta}_{YZ})$
$\quad \boldsymbol{\theta}_Z^* \leftarrow \boldsymbol{\theta}_Z + \boldsymbol{\rho}_Z^*$
$\quad \textbf{return } \mathbf{x} \cdot \boldsymbol{\theta}_X^\mu + \mathbf{x} \cdot \boldsymbol{\Theta}_{XX} \cdot \mathbf{x} + \psi_Z(\boldsymbol{\theta}_Z') + \rho_1' - \psi_Z(\boldsymbol{\theta}_Y^*) - \rho_1^* - \rho_0$

HMoG density, the HMoG forward mapping, and the HMoG EM algorithm, respectively.

---
**Algorithm 6** Computing the HMoG forward mapping

---
**Require:** HMoG natural parameters $(\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{XY}, \boldsymbol{\Theta}_{YZ})$
**Ensure:** HMoG expectations $(\boldsymbol{\eta}_X^\mu, \mathbf{H}_{XX}, \boldsymbol{\eta}_Y, \boldsymbol{\eta}_Z, \mathbf{H}_{XY}, \mathbf{H}_{YZ})$
   **function** $\boldsymbol{\tau}_{XYZ}(\boldsymbol{\theta}_X, \boldsymbol{\Theta}_{XX}, \boldsymbol{\Theta}_{XY}, \boldsymbol{\theta}_Y, \boldsymbol{\Theta}_{YZ}, \boldsymbol{\theta}_Z)$
      $\rho_0, \boldsymbol{\rho}_Y \leftarrow \textsc{GaussianConjugationParameters}(\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\Theta}_{XY})$
      $\boldsymbol{\theta}_Y^* \leftarrow \boldsymbol{\theta}_Y + \boldsymbol{\rho}_Y$
      $\rho_1^*, \boldsymbol{\rho}_Z^* \leftarrow \textsc{MixtureConjugationParameters}(\boldsymbol{\theta}_Y^*, \boldsymbol{\Theta}_{YZ})$
      $\boldsymbol{\eta}_Z \leftarrow \boldsymbol{\tau}_Z(\boldsymbol{\theta}_Z + \boldsymbol{\rho}_Z^*)$
      **for** $i = 1, \ldots, d_Z$ **do**
         $\boldsymbol{\eta}_{XY,i} \leftarrow \boldsymbol{\tau}_{XY}(\boldsymbol{\theta}_X, \boldsymbol{\theta}_{XX}, \boldsymbol{\Theta}_{XY}, \boldsymbol{\theta}_Y + \boldsymbol{\Theta}_{YZ} \cdot \mathbf{s}_Z(i))$
         **if** $i = 1$ **then**
            $\boldsymbol{\eta}_{XY} \leftarrow (1 - \sum_{j=1}^{d_Z-1} \eta_{Z,j}) \boldsymbol{\eta}_{XY,i}$
         **else**
            $\boldsymbol{\eta}_{XY} \leftarrow \boldsymbol{\eta}_Y + \eta_{Z,i-1} \boldsymbol{\eta}_{XY,i}$
         **end if**
      **end for**
      $(\boldsymbol{\eta}_X^\mu, \mathbf{H}_{XX}, \boldsymbol{\eta}_Y, \mathbf{H}_{XY}) \leftarrow \boldsymbol{\eta}_{XY}$
      **return** $(\boldsymbol{\eta}_X^\mu, \mathbf{H}_{XX}, \boldsymbol{\eta}_Y, \boldsymbol{\eta}_Z, \mathbf{H}_{XY}, \mathbf{H}_{YZ})$
   **end function**

---

 

---
**Algorithm 7** Expectation-maximization for HMoGs

---
[t]
**Require:** Sample $X^{(1)}, \ldots, X^{(d_S)}$
**Require:** HMoG parameters $\boldsymbol{\theta}_{XYZ} = (\boldsymbol{\theta}_X^\mu, \boldsymbol{\Theta}_{XX}, \boldsymbol{\theta}_Y, \boldsymbol{\theta}_Z, \boldsymbol{\Theta}_{XY}, \boldsymbol{\Theta}_{YZ})$
**Require:** Adam learning parameters $\alpha, \epsilon, \beta_1, \beta_2$
**Ensure:** Updated HMoG parameters $\boldsymbol{\theta}_{XYZ}'$
   **for all** $i \in \{1, \ldots, d_S\}$ **do**
      $\boldsymbol{\theta}_{Y,i} \leftarrow \boldsymbol{\theta}_Y + X^{(i)} \cdot \boldsymbol{\Theta}_{XY}$
      $(\boldsymbol{\eta}_{Y,i}', \boldsymbol{\eta}_{Z,i}', \mathbf{H}_{YZ,i}') \leftarrow \boldsymbol{\tau}_{YZ}(\boldsymbol{\theta}_{Y,i}, \boldsymbol{\Theta}_Z, \boldsymbol{\Theta}_{YZ})$
      $\mathbf{H}_{XY,i}' \leftarrow X^{(i)} \otimes \boldsymbol{\eta}_{Y,i}'$
      $\boldsymbol{\eta}_{X,i}^{\mu'} \leftarrow X^{(i)}$
      $\mathbf{H}_{XX,i}' \leftarrow X^{(i)} \otimes X^{(i)}$
      $\boldsymbol{\eta}_{XYZ,i}' \leftarrow (\boldsymbol{\eta}_{X,i}^{\mu'}, \mathbf{H}_{XX,i}', \mathbf{H}_{XY,i}', \boldsymbol{\eta}_{Y,i}', \mathbf{H}_{YZ,i}', \boldsymbol{\eta}_{Z,i}')$
   **end for**
   $\boldsymbol{\eta}_{XYZ} \leftarrow \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{\eta}_{XYZ,i}'$
   **function** $\nabla\mathcal{L}(\boldsymbol{\theta}_{XYZ})$
      $\boldsymbol{\eta}_{XYZ} \leftarrow \boldsymbol{\tau}_{XYZ}(\boldsymbol{\theta}_{XYZ})$
      **return** $\boldsymbol{\eta}_{XYZ} - \boldsymbol{\eta}_{XYZ}'$
   **end function**
   **return** $\textsc{AdamOptimizer}(\nabla\mathcal{L}, \boldsymbol{\theta}_{XYZ}, \alpha, \epsilon, \beta_1, \beta_2)$

---