

---

# Efficient Learning for Time Series Models by Non-Negative Moment Matrix Factorization

---

## Abstract

In recent years, method-of-moments (MoM) based algorithms for latent variable models have been popular in the machine learning community, as a computationally cheaper alternatives to more conventional maximum likelihood based algorithms. Although there are MoM algorithms for the Hidden Markov Model (HMM), Gaussian Mixture Model (GMM), and Latent Dirichlet Allocation (LDA), it is unclear how to extend these methods to more complex models. In this paper, we develop a MoM-based approach using Non-negative Matrix Factorization (NMF) for learning several time series models. These include the Mixture of HMMs (MHMM), Switching HMM (SHMM) and Factorial HMM (FHMM). We show experimentally that NMF is considerably faster than conventional maximum-likelihood approaches and gives comparable estimation accuracies.

## 1. Introduction

In recent years, method-of-moments (MoM) learning algorithms with unique solution guarantees for several latent variable models have been developed in the machine learning literature. Examples include the Hidden Markov Model (HMM) (Hsu et al., 2009; Anandkumar et al., 2012c;b), Gaussian Mixture Model (GMM) (Hsu & Kakade, 2013), Latent Dirichlet Allocation (LDA) (Anandkumar et al., 2012a), and Independent Component Analysis (ICA) (Anandkumar et al., 2012b). These algorithms provide computationally cheaper alternatives to conventional Expectation-Maximization (EM) (Dempster et al., 1977) based approaches that are free of locally optimal solutions.

While there are MoM algorithms for a handful of models, a general framework that is as widely applicable as EM does not exist. Therefore, it is not clear how to derive MoM learning algorithms with unique solution guarantees for more sophisticated models such as the Mixture

of Hidden Markov Models (MHMM) (Smyth, 1997), the Switching HMM (SHMM) (Murphy, 2002), or the Factorial HMM (FHMM) (Ghahramani & Jordan, 1997).

There exists Non-negative Matrix Factorization (NMF) (Lee & Seung, 1999) based MoM algorithms for learning the parameters of an HMM (Cybenko & Crespi, 2011; Lakshminarayanan & Raich, 2010; Shashanka, 2011). This approach amounts to doing alternating least squares or having an iterative multiplicative update scheme. However, it has not been applied to more sophisticated models.

In this paper, we propose a general NMF-MoM framework to learn the parameters of several time series models. We do this by converting each to an HMM with a specific transition structure and developing multiplicative update rules (Lee & Seung, 2001) that iteratively update the parameter matrices.

Although one loses the uniqueness properties of the spectral learning algorithms with NMF, the computational advantage of NMF-MoM is significant. Learning of time series models is typically done with maximum-likelihood training using EM (ML-EM). An important drawback of ML-EM is that one has to perform inference for the latent variables at each iteration. For example, for MHMM, one has to do forward-backward message-passing for each sequence at each iteration. Therefore, in abundance of data, learning the parameters of the model may be computationally expensive. However, learning the model with a NMF-MoM approach amounts to the estimation of some observable moments and the factorization of a moment matrix (or tensor) with NMF, where both of which can be executed quite fast.

We perform experiments on synthetic data to compare the accuracy and speed of the proposed learning methods with EM. We observe that our approach is substantially faster than ML-EM when data is plentiful. In addition, we achieve comparable estimation accuracies. We also conduct experiments with real datasets to show the validity of our approach.

## 2. Notation

We used square brackets  $[-]$  for the indicator function. We use MATLAB notation to select a column or row of a matrix. For example,  $O(:, k)$ ,  $O(k, :)$  and  $O(k, j)$  denote the  $k$ 'th column and row and  $(i, j)$ 'th entry of the matrix  $O$ , respectively. We use boldface for matrices when doing matrix calculations. The outer product operator  $\otimes$  is defined as  $(a \otimes b)(i, j) = a(i)b(j)$ . We use  $\odot$  to denote element-wise matrix multiplication. We denote a sequence of length  $T$  as  $x_{1:T} = \{x_1, x_2, \dots, x_T\}$ . We use boldface to denote a sequence so that  $x_{1:T, n} = \mathbf{x}_n$  and consequently we use  $\mathbf{x}_{1:N} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  to denote a dataset of  $N$  sequences. Note that  $x_t \in \mathbb{R}^L$  in general. The estimate for a parameter  $\theta$  is shown as  $\theta^*$ . Finally, we use subscript to denote parameters belonging to slices of a tensor, e.g.  $O_k$  matrix is the  $k$ 'th slice of  $O$  tensor.

Some important indices, variables and parameters that are used throughout the paper are as follows.  $n \in \{1, \dots, N\}$ : sequence index,  $t \in \{1, \dots, T\}$ : time index,  $k \in \{1, \dots, K\}$ : cluster index,  $h \in \{1, \dots, K\}$ : cluster/regime indicator,  $r_t \in \{1, \dots, J\}$ : state indicator,  $x_t \in \{1, \dots, L\}$ : observation,  $\mathbf{O}_k \in \mathbb{R}^{L \times J}$ : cluster emission matrix,  $\tilde{\mathbf{O}} = [\mathbf{O}_1 \ \mathbf{O}_2 \ \dots \ \mathbf{O}_K] \in \mathbb{R}^{L \times KJ}$ : stacked emission matrix,  $\mathbf{B} \in \mathbb{R}^{K \times K}$ : regime transition matrix,  $\mathbf{A}_k \in \mathbb{R}^{J \times J}$ : cluster transition matrix,  $\nu_k \in \mathbb{R}^J$ : cluster initial state distribution,  $\hat{\mathbf{A}}_k = \mathbf{A}_k \text{diag}(\nu_k)$ ,  $\tilde{\mathbf{A}} \in \mathbb{R}^{JK \times JK}$ : global joint state matrix,  $\boldsymbol{\pi} \in \mathbb{R}^K$ : prior vector,  $\mathbf{1}_{ab} \in \mathbb{R}^{a \times b}$ : ones matrix.

## 3. Model Definitions

In this section, we give the definitions of the models used in this paper, and the corresponding variables and parameters. We also describe the modeling choices used throughout the paper.

### 3.1. Hidden Markov Model

In Hidden Markov Model (HMM), an observed sequence  $\mathbf{x} = x_{1:T}$  is generated conditioned on a latent Markov chain  $r_{1:T}$ , with  $r_t \in \{1, \dots, J\}$ . Given the model parameters  $\theta = (O, A, \nu)$ , the likelihood  $p(x_{1:T}|\theta)$  of an observation sequence  $x_{1:T}$  of length  $T$  is defined as follows:

$$\begin{aligned} p(x_{1:T}|\theta) &= \sum_{r_{1:T}} p(x_{1:T}, r_{1:T}|\theta) \\ &= \sum_{r_{1:T}} \prod_{t=1}^T p(x_t|r_t)p(r_t|r_{t-1}) \\ &= \sum_{r_{1:T}} \prod_{t=1}^T O(x_t, r_t)A(r_t, r_{t-1}) \end{aligned} \quad (1)$$

The model parameters are defined as follows:

- $\nu(u) = p(r_1 = u|r_0) = p(r_1 = u)$ , the initial latent state distribution.
- $A(u, v) = p(r_t = u|r_{t-1} = v)$ ,  $t \geq 2$ , latent state transition matrix.
- $O(:, u) = \mathbb{E}[x_t|r_t = u]$ , emission matrix.

where,  $\nu \in \mathbb{R}^J$ ,  $A \in \mathbb{R}^{J \times J}$  and  $O \in \mathbb{R}^{L \times J}$ . A column  $O(:, u)$  of the observation matrix  $O$  is defined as the expectation of  $x_t$ , conditioned on the latent state  $r_t$ . The choice of  $p(x_t|r_t)$ , determines what the columns of  $O$  correspond to. Some frequently used choices are:

- Gaussian:  $p(x_t|r_t = u) = \mathcal{N}(x_t; \mu_u, \sigma^2)$   
then,  $O(:, u) = \mathbb{E}[x_t|r_t = u] = \mu_u$ .
- Poisson:  $p(x_t|r_t = u) = \mathcal{PO}(x_t; \lambda_u)$   
then,  $O(:, u) = \mathbb{E}[x_t|r_t = u] = \lambda_u$ .
- Multinomial:  $p(x_t|r_t = u) = \text{Multinomial}(p_u, S)$   
then,  $O(:, u) = \mathbb{E}[x_t|r_t = u] = p_u$ .

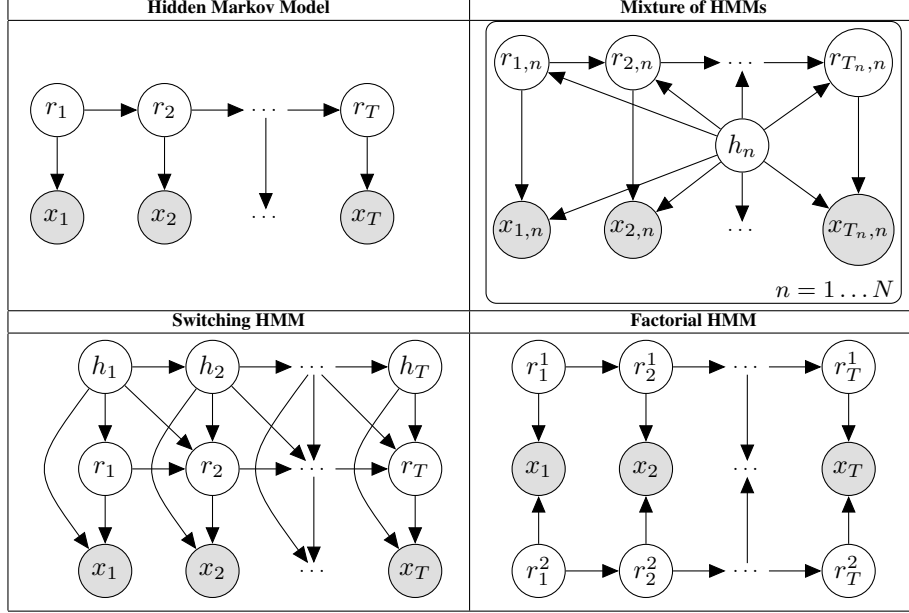
The first choice is multivariate, isotropic Gaussian with mean  $\mu_u \in \mathbb{R}^L$ . The second distribution is Poisson with intensity parameter  $\lambda_u \in \mathbb{R}^L$ . This choice is particularly useful for counts data. The last density is a multinomial distribution, with parameter  $p_u \in \mathbb{R}^L$ , and number of draws  $S$ . In this work, we use the Multinomial distribution in all of the models. In audio modeling applications, for example,  $S$  is taken to be a large number so that the observations describe the discrete Fourier transform (DFT) magnitude of the signal over time (Mysore et al., 2010).

### 3.2. Mixture of HMMs

Mixture of HMMs is a useful model for clustering sequences where each one is modeled with one of the  $K$  HMMs. Given model parameters  $\theta = (O_{1:K}, A_{1:K}, \nu_{1:K}, \pi)$ , the likelihood  $p(\mathbf{x}_n|\theta)$  of an observation sequence  $\mathbf{x}_n = \{x_{1,n}, x_{2,n}, \dots, x_{T,n}\}$  of length  $T_n$  is computed as convex combination of the likelihood of  $K$  hidden Markov models:

165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219

Table 1. Probabilistic graphical models of the time series models.



$$p(\mathbf{x}_n|\theta) = \sum_{k=1}^K p(h_n = k)p(\mathbf{x}_n|h_n = k) \quad (2)$$

$$\begin{aligned} &= \sum_{k=1}^K \pi_k \sum_{\mathbf{r}_n} p(\mathbf{x}_n, \mathbf{r}_n|h_n = k) \\ &= \sum_{k=1}^K \pi_k \sum_{r_{1:T_n,n}} \prod_{t=1}^{T_n} p(x_{t,n}|r_{t,n}, h_n = k)p(r_{t,n}|r_{t-1,n}, h_n = k) \\ &= \sum_{k=1}^K \pi_k \sum_{r_{1:T_n,n}} \prod_{t=1}^{T_n} O_k(x_{t,n}, r_{t,n})A_k(r_{t,n}, r_{t-1,n}) \end{aligned}$$

where,  $h_n \in \{1, 2, \dots, K\}$  is the latent cluster indicator, and  $\mathbf{r}_n = \{r_{1,n}, r_{2,n}, \dots, r_{T_n,n}\}$  is the latent state sequence of the observed sequence  $\mathbf{x}_n$ .  $\pi \in \mathbb{R}^K$  is defined as the cluster prior probability vector. Note that, if a sequence is assigned to  $k$ 'th cluster (meaning  $h_n = k$ ), then the transition matrix  $A_k$  and the observation matrix  $O_k$  are used to generate that particular sequence. That is to say, we have a mixture of HMMs where each mixture component is specified by the model parameters  $\theta_k$ .

### 3.3. Switching HMM

Switching state space HMM is a generalization of Mixture of HMMs, where the cluster indicator variable  $h$  is allowed to take differing values within a single sequence, as illustrated in the graphical model in Table 1. Allowing  $h$  to change within a sequence enables us to use this model for segmentation. The emission and transition matrices used to

generate the observations vary as  $h$  changes throughout the sequence.

Given the model parameters  $\theta = (O_{1:K}, A_{1:K}, \nu_{1:K}, B, \pi)$ , the likelihood  $p(x_{1:T}|\theta)$  of a given sequence  $x_{1:T}$  is defined as follows:

$$\begin{aligned} p(x_{1:T}|\theta) &= \sum_{h_{1:T}, r_{1:T}} p(x_{1:T}, h_{1:T}, r_{1:T}|\theta_{1:K}) \quad (3) \\ &= \sum_{h_{1:T}, r_{1:T}} \prod_{t=1}^T p(x_t|r_t, h_t)p(r_t|r_{t-1}, h_t, h_{t-1})p(h_t|h_{t-1}) \\ &= \sum_{h_{1:T}, r_{1:T}} \prod_{t=1}^T O_{r_t}(x_t, r_t) \left\{ [h_t = h_{t-1}]A_{h_t}(r_t, r_{t-1}) \right. \\ &\quad \left. + [h_t \neq h_{t-1}]A_{h_t, h_{t-1}}(r_t, r_{t-1}) \right\} B(h_t, h_{t-1}) \end{aligned}$$

where  $h_t \in \{1, \dots, K\}$  is the latent regime indicator, and  $r_t \in \{1, \dots, J\}$  is the latent state indicator.  $A_{h_t, h_{t-1}}$  denotes the transition matrix to be used whenever there is a regime change ( $h_t \neq h_{t-1}$ ). If there is no regime change ( $h_t = h_{t-1}$ ), then  $A_{h_t, h_{t-1}} = A_{h_t}$ . Note that  $A_{h_t, h_{t-1}}$  matrices are fixed in learning. The regime transition matrix  $B \in \mathbb{R}^{K \times K}$  is defined as  $B(u, v) = p(h_t = u|h_{t-1} = v)$ , and regime initial distribution  $\pi \in \mathbb{R}^K$  is defined as  $\pi(u) = p(h_1 = u)$ .

### 3.4. Factorial HMM

The factorial HMM is useful for modeling additive mixtures. An observed sequence is modeled as the weighted sum of the outputs of  $K$  HMMs. For  $K = 2$ , given the

275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329

model parameters  $\theta = (O_{1:2}, A_{1:2}, \nu_{1:2})$ , the likelihood  $p(x_{1:T}|\theta)$  of a given sequence  $x_{1:T}$  is defined as follows:

$$\begin{aligned} p(x_{1:T}|\theta) &= \sum_{r_{1:T}^1, r_{1:T}^2} p(x_{1:T}, r_{1:T}^1, r_{1:T}^2|\theta) \quad (4) \\ &= \sum_{r_{1:T}^1, r_{1:T}^2} \prod_{t=1}^T p(x_t|r_t^1, r_t^2) p(r_t^1|r_{t-1}^1) p(r_t^2|r_{t-1}^2) \\ &= \sum_{r_{1:T}^1, r_{1:T}^2} \prod_{t=1}^T \frac{1}{2} \left\{ O_1(x_t, r_t^1) + O_2(x_t, r_t^2) \right\} \\ &\quad \times A_1(r_t^1, r_{t-1}^1) A_2(r_t^2, r_{t-1}^2) \end{aligned}$$

where  $r_{1:T}^1$  is the first, and  $r_{1:T}^2$  is the second latent chain. Conditioned on them, the observation sequence  $x_{1:T}$  is generated using the observation model  $\frac{1}{2} \left\{ O_1(x_t, r_t) + O_2(x_t, r_t) \right\}$ . Note that the mixing weights are set to  $\frac{1}{2}$  in this paper. In general, mixing weights are also parameters to be estimated.

#### 4. MoM-NMF algorithms for learning

A common approach to learn latent variable models like the ones that are defined in Section 3 is to use a maximum likelihood criterion. A parameter set  $\theta^*$  is chosen such that the log-likelihood of the observed data is maximized:

$$\theta^* = \operatorname{argmax}_{\theta} \log p(x|\theta) = \operatorname{argmax}_{\theta} \log \sum_h p(x, h|\theta), \quad (5)$$

where  $x$  is the observed data,  $h$  is a set of hidden variables, and  $\theta$  is the parameter set to be optimized. Because of the presence of a summation over the latent variables  $h$ , which may be hard to compute in general, it is easier to iteratively maximize a lower bound on the log likelihood. The estimate in iteration  $\tau$  is:

$$(\theta^\tau)^* = \operatorname{argmax}_{\theta^\tau} \mathbb{E}_{p(\theta^{\tau-1}|x, h)} [\log p(x, h|\theta^\tau)]. \quad (6)$$

This EM lower bound involves computing the posterior  $p(\theta^{\tau-1}|x, h)$  over the model parameters at each iteration. In the temporal models that we consider, this requires a message passing algorithm. For example, for a mixture of HMMs, we need to compute the posterior for each sequence  $\mathbf{x}_n$  for each cluster  $k$  at each iteration  $\tau$ . This can be computationally prohibitive for large datasets. The number of calculations required per iteration is  $\mathcal{O}(J^2 K T N + JKLTN)$ , where  $N$  is the number of sequences and  $T$  is the average sequence length. Since the dependence is linear in  $N$ , EM is infeasible for large datasets.

A computationally cheaper learning scheme is as follows. We can compute moment estimates from the data and estimate the model parameters via a factorization of the moment matrix/tensor. We note that the 2<sup>nd</sup> order moment matrices for these models all factorize in the form:

$$\mathbf{P} := \mathbb{E}[x_t \otimes x_{t-1}] = \tilde{\mathbf{O}} \tilde{\mathbf{A}} \tilde{\mathbf{O}}^\top, \quad (7)$$

where  $\mathbf{P} \in \mathbb{R}^{L \times L}$  is defined as the true moment matrix, and the block matrix  $\tilde{\mathbf{O}} \in \mathbb{R}_+^{L \times JK}$  contains the HMMs' emission matrices, stacked side-by-side.  $\tilde{\mathbf{A}} \in \mathbb{R}_+^{JK \times JK}$  is the joint state distribution matrix that contains the transition structure of the model embedded in a single, larger HMM. The details of  $\tilde{\mathbf{A}}$  for each model are given in the subsequent sections.

Having made this observation, the learning problem reduces to the following non-negative factorization problem:

$$\theta^* = (\tilde{\mathbf{O}}^*, \tilde{\mathbf{A}}^*) = \operatorname{argmin}_{\tilde{\mathbf{O}} \geq 0, \tilde{\mathbf{A}} \geq 0} d(\mathbf{V} \| \tilde{\mathbf{O}} \tilde{\mathbf{A}} \tilde{\mathbf{O}}^\top), \quad (8)$$

where  $\mathbf{V}$  is an empirical estimate for the moment matrix  $\mathbf{P} = \mathbb{E}[x_t \otimes x_{t-1}]$ , and  $d(\cdot, \|\cdot)$  is an appropriate divergence measure. We use KL divergence as it is commonly used in NMF community. We continue by providing 2<sup>nd</sup> order moment matrix decomposition descriptions for each model. Higher order decompositions for each model are given in the supplemental material.

##### 4.1. Hidden Markov Model

The 2<sup>nd</sup>-order moment matrix of hidden Markov model is decomposed as follows:

$$\begin{aligned} \mathbf{P} &= \mathbb{E}[x_{t+1} \otimes x_t] \quad (9) \\ &= \sum_{r_t} \sum_{r_{t+1}} \mathbb{E}[x_{t+1}|r_{t+1}] p(r_{t+1}|r_t) p(r_t) \mathbb{E}[x_t|r_t] \\ &= \sum_{r_t} \sum_{r_{t+1}} O(x_{t+1}, r_{t+1}) A(r_{t+1}|r_t) \nu(r_t) O(x_t, r_t) \\ &= \mathbf{O} \mathbf{A} \operatorname{diag}(\boldsymbol{\nu}) \mathbf{O}^\top = \mathbf{O} \hat{\mathbf{A}} \mathbf{O}^\top \end{aligned}$$

where  $\mathbf{O}$  and  $\mathbf{A}$  are respectively the emission and transition matrices of the HMM. Note that  $\hat{\mathbf{A}} = p(r_{t+1}, r_t)$  is the joint consecutive state distribution. Thus, the true moment  $\mathbf{P}$  is:

$$\mathbf{P} = \mathbf{O} \hat{\mathbf{A}} \mathbf{O}^\top = \tilde{\mathbf{O}} \tilde{\mathbf{A}} \tilde{\mathbf{O}}^\top. \quad (10)$$

The multiplicative update rules to factorize the empirical moment matrix  $\mathbf{V}$  are:

Table 2. Global joint state distributions  $\tilde{\mathbf{A}}$  for the time series models.

Hidden Markov Model	Mixture of HMMs
$\hat{\mathbf{A}}$	$\begin{bmatrix} \pi_1 \hat{\mathbf{A}}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \pi_2 \hat{\mathbf{A}}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \pi_K \hat{\mathbf{A}}_K \end{bmatrix}$
Switching HMM	Factorial HMM
$\begin{bmatrix} \mathbf{A}_1 B_{11} & \mathbf{A}_{12} B_{12} & \cdots & \mathbf{A}_{1K} B_{1K} \\ \mathbf{A}_{21} B_{21} & \mathbf{A}_2 B_{22} & \cdots & \mathbf{A}_{2K} B_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}_{K1} B_{K1} & \mathbf{A}_{K2} B_{K2} & \cdots & \mathbf{A}_K B_{KK} \end{bmatrix} \text{diag}(\tilde{\pi})$	$\frac{1}{K^2} \begin{bmatrix} \hat{\mathbf{A}}_1 & \hat{\mathbf{A}}_1 \mathbf{1}_{JJ} \hat{\mathbf{A}}_2 & \cdots & \hat{\mathbf{A}}_1 \mathbf{1}_{JJ} \hat{\mathbf{A}}_K \\ \hat{\mathbf{A}}_2 \mathbf{1}_{JJ} \hat{\mathbf{A}}_1 & \hat{\mathbf{A}}_2 & \cdots & \hat{\mathbf{A}}_2 \mathbf{1}_{JJ} \hat{\mathbf{A}}_K \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{A}}_K \mathbf{1}_{JJ} \hat{\mathbf{A}}_1 & \hat{\mathbf{A}}_K \mathbf{1}_{JJ} \hat{\mathbf{A}}_2 & \cdots & \hat{\mathbf{A}}_K \end{bmatrix}$

$$\tilde{\mathbf{O}} \leftarrow \tilde{\mathbf{O}} \odot \frac{\left(\frac{\mathbf{v}_\cdot}{\tilde{\mathbf{V}}}\right) \tilde{\mathbf{O}} \tilde{\mathbf{A}}^\top + \left(\frac{\mathbf{v}_\cdot}{\tilde{\mathbf{V}}}\right)^\top \tilde{\mathbf{O}} \tilde{\mathbf{A}}}{\mathbf{1}_{LL} \tilde{\mathbf{O}} (\tilde{\mathbf{A}}^\top + \tilde{\mathbf{A}})}, \quad (11)$$

$$\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} \odot \frac{\tilde{\mathbf{O}}^\top \left(\frac{\mathbf{v}_\cdot}{\tilde{\mathbf{V}}}\right) \tilde{\mathbf{O}}}{\tilde{\mathbf{O}}^\top \mathbf{1}_{LL} \tilde{\mathbf{O}}}. \quad (12)$$

where  $\hat{\mathbf{V}}$  is the approximation of  $\mathbf{V}$  constructed with  $\hat{\mathbf{A}}$  and  $\tilde{\mathbf{O}}$ . Matrix division is performed element-wise. Iterating (11) and (12) is guaranteed to converge to a local minimum of (8) (Lee & Seung, 2001).

## 4.2. Mixture of Hidden Markov Models

The 2<sup>nd</sup>-order moment matrix of the mixture of HMMs is decomposed as follows:

$$\begin{aligned} \mathbf{P} &= \mathbb{E}[x_{t+1} \otimes x_t] = \sum_{h=1}^K p(h) \cdots \\ & \sum_{r_t} \sum_{r_{t+1}} \mathbb{E}[x_{t+1} | r_{t+1}, h] p(r_t, r_{t+1} | h) \mathbb{E}[x_t | r_t, h] \\ &= \sum_{h=1}^K \pi(h) \sum_{r_t} \sum_{r_{t+1}} O_h(x_{t+1}, r_{t+1}) A_h(r_{t+1} | r_t) \\ & \quad \cdots \nu_k(r_t) O_h(x_t, r_t) \\ &= \sum_{k=1}^K \pi_k \mathbf{O}_k \hat{\mathbf{A}}_k \mathbf{O}_k^\top = \tilde{\mathbf{O}} \tilde{\mathbf{A}} \tilde{\mathbf{O}}^\top \end{aligned} \quad (13)$$

We have dropped the sequence index  $n$  to avoid clutter and defined a block-diagonal  $\tilde{\mathbf{A}}$  for the entire MHMM (see Table 2). Since the factorization is identical in form to that of the single HMM case, the NMF updates are identical to (11) and (12). A similar approach was taken to derive update rules based on the 3<sup>rd</sup> and 4<sup>th</sup> moments (see supplemental materials). Note that if  $\tilde{\mathbf{A}}$  is initialized as block-

diagonal, it will remain this way throughout the multiplicative updates.

## 4.3. Switching Hidden Markov Model

The 2<sup>nd</sup>-order moment matrix for the switching HMM is decomposed as follows:

$$\begin{aligned} \mathbf{P} &= \mathbb{E}[x_{t+1} \otimes x_t] = \sum_{h_t} \sum_{h_{t+1}} p(h_t, h_{t+1}) \sum_{r_t} \sum_{r_{t+1}} \cdots \\ & \mathbb{E}[x_{t+1} | r_{t+1}, h_{t+1}] p(r_{t+1} | r_t, h_t, h_{t+1}) p(r_t | h_t) \mathbb{E}[x_t | r_t, h_t] \\ &= \sum_{h_t} \sum_{h_{t+1}} B(h_t, h_{t+1}) \sum_{r_t} \sum_{r_{t+1}} O_{h_{t+1}}(x_{t+1}, r_{t+1}) \cdots \\ & \quad A_{h_{t+1}, h_t}(r_{t+1}, r_t) \nu_{h_t}(r_t) O_{h_t}(x_t, r_t) \\ &= \sum_{k=1}^K \sum_{k'=1}^K \mathbf{O}_k \tilde{\mathbf{A}}_{(kk')} \mathbf{O}_{k'}^\top = \tilde{\mathbf{O}} \tilde{\mathbf{A}} \tilde{\mathbf{O}}^\top \end{aligned} \quad (14)$$

where  $\tilde{\mathbf{A}}_{(kk')}$  denotes the  $(k, k')$ <sup>th</sup> block of the  $\tilde{\mathbf{A}}$  matrix defined for the entire SHMM (see Table 2) and  $\mathbf{B}$  is the regime transition matrix. The blocks of  $\tilde{\mathbf{A}}$  encode the state transition information given the regime indicators  $h_t$  and  $h_{t+1}$ . The multiplicative update for  $\tilde{\mathbf{O}}$  is identical to (11). The other updates are:

$$\mathbf{A}_{kk'} \leftarrow \mathbf{A}_{kk'} \odot \frac{\mathbf{O}_k^\top \left(\frac{\mathbf{v}_\cdot}{\tilde{\mathbf{V}}}\right) \mathbf{O}_{k'}}{\mathbf{O}_k^\top \mathbf{1}_{LL} \mathbf{O}_{k'}}, \quad (15)$$

$$B_{kk'} \leftarrow B_{kk'} \odot \frac{\mathbf{1}_{1L} \left[ \left(\frac{\mathbf{v}_\cdot}{\tilde{\mathbf{V}}}\right) \odot \mathbf{O}_k \tilde{\mathbf{A}}_{(kk')} \mathbf{O}_{k'}^\top \right] \mathbf{1}_{L1}}{\mathbf{1}_{1L} \left[ \mathbf{O}_k \tilde{\mathbf{A}}_{(kk')} \mathbf{O}_{k'}^\top \right] \mathbf{1}_{L1}}, \quad (16)$$

$$\tilde{\pi} \leftarrow \tilde{\pi} \odot \frac{\left[ \tilde{\mathbf{A}}^\top \tilde{\mathbf{O}}^\top \left(\frac{\mathbf{v}_\cdot}{\tilde{\mathbf{V}}}\right) \odot \tilde{\mathbf{O}}^\top \right] \mathbf{1}_{L1}}{\tilde{\mathbf{A}}^\top \tilde{\mathbf{O}}^\top \mathbf{1}_{L1} \odot \tilde{\mathbf{O}}^\top \mathbf{1}_{L1}}, \quad (17)$$

where:

440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494

495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549



$$\tilde{\boldsymbol{\pi}} = [\pi_1 \boldsymbol{\nu}_1^\top \quad \pi_2 \boldsymbol{\nu}_2^\top \quad \cdots \quad \pi_K \boldsymbol{\nu}_K^\top]^\top. \quad (18)$$

We can fix the off-diagonal blocks in  $\tilde{\mathbf{A}}$  for simplicity. For example, we may know that when a regime change occurs (i.e.  $h_t \neq h_{t+1}$ ), all states are equally likely at time  $t + 1$ . In this case,  $\mathbf{A}_{kk'} = \frac{1}{J} \mathbf{1}_{JJ}$ ,  $k \neq k'$ . Alternatively, we may impose that  $r_{t+1} = 1$  (e.g. for left-to-right models used for speech or gesture datasets), in which case we have that  $\mathbf{A}_{kk'} = [1 \ 0 \ \cdots \ 0]^\top \mathbf{1}_{1J}$ ,  $k \neq k'$ .

#### 4.4. Factorial Hidden Markov Model

In order for our approach to be tractable, an emission model  $O(x_t, r_t^{1:K})$  is assumed:

$$O(x_t, r_t^{1:K}) = \frac{1}{K} \sum_{k=1}^K O_k(x_t, r_t^k). \quad (19)$$

where  $r_t^{1:K} = \{r_t^1, \dots, r_t^K\}$  is the set of HMM states at time  $t$ . Since we average over the entire sequence to compute moments, this is equivalent to assuming that the mixing proportions  $p(r_t^k)$  are uniform *on average*. This is reasonable in many cases.

The 2<sup>nd</sup>-order moment matrix for the factorial HMM is decomposed as follows:

$$\begin{aligned} \mathbf{P} &= \mathbb{E}[x_{t+1} \otimes x_t] \\ &= \sum_{r_t^{1:K}} \sum_{r_{t+1}^{1:K}} \mathbb{E}[x_{t+1} | r_{t+1}^{1:K}] p(r_{t+1}^{1:K}, r_t^{1:K}) \mathbb{E}[x_t | r_t^{1:K}] \\ &= \frac{1}{K^2} \sum_{k=1}^K \left[ \mathbf{O}_k \hat{\mathbf{A}}_k \mathbf{O}_k^\top + \sum_{k' \neq k} \mathbf{O}_k \hat{\mathbf{A}}_k \mathbf{1}_{JJ} \hat{\mathbf{A}}_{k'} \mathbf{O}_{k'}^\top \right] \\ &= \tilde{\mathbf{O}} \tilde{\mathbf{A}} \tilde{\mathbf{O}}^\top. \end{aligned} \quad (20)$$

where we have defined an  $\tilde{\mathbf{A}}$  for the entire FHMM (see Table 2). This formulation can easily be generalized to the case where the HMMs have differing numbers of states  $J$ .

The NMF update for  $\tilde{\mathbf{O}}$  is identical to (11). The updates for the transition matrices can be written as:

$$\mathbf{A}_k \leftarrow \mathbf{A}_k \odot \frac{\mathbf{E}_k \mathbf{R}_1 \left( \tilde{\mathbf{O}}^\top \left( \frac{\mathbf{V}_k}{\tilde{\mathbf{V}}} \right) \tilde{\mathbf{O}} \odot \mathbf{Q}_k \right) \mathbf{R}_2 \mathbf{E}_k^\top}{\mathbf{E}_k \mathbf{R}_1 \left( \tilde{\mathbf{O}}^\top \mathbf{1}_{LL} \tilde{\mathbf{O}} \odot \mathbf{Q}_k \right) \mathbf{R}_2 \mathbf{E}_k^\top}. \quad (21)$$

where, for example, when  $K = 2$ :

$$\begin{aligned} \mathbf{Q}_1 &= \begin{bmatrix} \mathbf{1}_{JJ} & \mathbf{1}_{JJ} \\ \mathbf{1}_{JJ} & \mathbf{0} \end{bmatrix}, \quad \mathbf{Q}_2 = \begin{bmatrix} \mathbf{0} & \mathbf{1}_{JJ} \\ \mathbf{1}_{JJ} & \mathbf{1}_{JJ} \end{bmatrix}, \\ \mathbf{R}_1 &= \begin{bmatrix} \mathbf{I} & \mathbf{1}_{JJ} \mathbf{A}_2^\top \\ \mathbf{1}_{JJ} \mathbf{A}_1^\top & \mathbf{I} \end{bmatrix}, \quad \mathbf{R}_2 = \begin{bmatrix} \mathbf{I} & \mathbf{A}_1^\top \mathbf{1}_{JJ} \\ \mathbf{A}_2^\top \mathbf{1}_{JJ} & \mathbf{I} \end{bmatrix}, \\ \mathbf{E}_1 &= [\mathbf{I} \ \mathbf{0}], \quad \mathbf{E}_2 = [\mathbf{0} \ \mathbf{I}]. \end{aligned} \quad (22)$$

We also derived updates based on the factorization of the 3<sup>rd</sup> moment tensor (see supplemental materials).

#### 4.5. Computational complexity

The computation complexities of the ML-EM and NMF-MoM algorithms are given in Table 3. The length  $T$  of the data sequences (and the number of sequences  $N$ , in the case of the MHMM) appears as a multiplicative factor for the ML-EM algorithms and as an additive term for the NMF algorithms. Although ML-EM typically converges in fewer iterations, message-passing in the E step is very costly. This shows the computational advantage of the MoM approach. We also note that for the factorial HMM, the EM algorithm runs in time that is exponential in the number of latent state chains  $K$ . However, this dependence is at worst cubic in the MoM methods. Thus, the computational advantage of the proposed approach may be substantial for complicated models.

## 5. Experiments

In all the synthetic experiments described in this section, parameter matrices were generated in the same way. The columns of the  $\mathbf{O}$  matrices were generated from a Dirichlet distribution with parameters  $\alpha_l = 0.15$ ,  $l = 1, \dots, L$ . This ensures that they are relatively sparse. The transition matrices were generated as  $\mathbf{A} = \beta \mathbf{I} + (1 - \beta) \mathbf{I}_s$ , where  $0.6 < \beta < 1$  and  $\mathbf{I}_s$  is the identity matrix with a downward circular shift applied to the rows. Thus, the  $\mathbf{A}$  matrices are ergodic but relatively sparse. Prior vectors were generated from a uniform distribution over the probability simplex. For the NMF-MoM algorithms, parameters are initialized with random non-negative values.

### 5.1. Mixture of HMMs

#### 5.1.1. SYNTHETIC DATA

The goal of the experiments with synthetic data is to understand the impact of dataset size on sequence clustering accuracy and computation time. We generated 100 sets of MHMM data consisting of 3 clusters, each of which contains 10 sequences and  $J = 3$ ,  $L = 9$ . We report the clustering accuracies and estimation time for sequences of

Table 3. Computational complexities of ML-EM and NMF-MoM algorithms per iteration.

Model	EM	2 <sup>nd</sup> order	3 <sup>rd</sup> order	4 <sup>th</sup> order
MHMM	$\mathcal{O}(J^2 KTN + JKLTN)$	$\mathcal{O}(J^2 KL + JKL^2)$	$\mathcal{O}(J^2 KL^2 + JKL^3)$	$\mathcal{O}(J^2 KL^3 + JKL^4)$
SHMM	$\mathcal{O}(J^2 K^2 T + JKLT)$	$\mathcal{O}(J^2 K^2 L + JKL^2)$	$\mathcal{O}(J^2 K^2 L^2 + JKL^3)$	
FHMM	$\mathcal{O}(J^2 K^2 T + J^K LT)$	$\mathcal{O}(JKL^2 + J^2 K^2 L + J^3 K^3)$	$\mathcal{O}(J^2 K^3 L^2 + JK^3 L^3 + J^3 K^3 L)$	

length  $L \in [10, 10,000]$  (see supplemental materials for pseudocode).

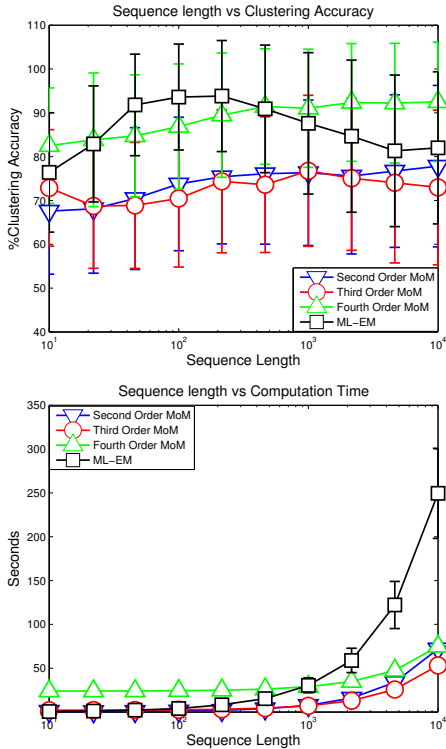


Figure 1. **Top:** Sequence Length vs Clustering accuracy Experiment, **Bottom:** Sequence Length vs Computation Time Experiment

We see in Figure 1 that the 4<sup>th</sup> order MoM algorithm performs nearly as well as the maximum-likelihood based EM algorithm with short sequences. With long sequences, the MoM algorithm performs slightly better. The savings in computation with NMF-MoM is significant when the sequences are long. The increase in computation time is due to the assignment step for each sequence after learning (a single forward pass for each sequence). However, the situation is much worse for ML-EM since we need to do forward-backward passes for each sequence in every iteration. The EM iterations were stopped automatically as soon as the clustering became stable.

5.1.2. REAL DATA

We used the 3<sup>rd</sup> order decomposition algorithm for the mixture of HMMs on the "Character Trajectories" dataset from

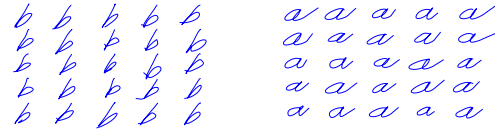


Figure 2. Clustering "a" and "b" character trajectories.

the UCI machine learning repository (Bache & Lichman, 2013). The data consists of handwritten character trajectories written on a PDA screen. The data in its original format is a 3-dimensional, continuous time series. The first two dimensions are the derivative of the spatial coordinates and the third is the pen pressure as a function of time. We discretized the data using K-means clustering with 14 clusters. The original dataset contains 2,858 sequences in 20 clusters. When we use the first 160 data items containing 97 "a" letters and 63 "b" letters, we obtain a 98% clustering accuracy. The clustering (with moment estimation) takes 7.5 seconds. Examples of clustered characters are given in Figure 2. NMF-MoM is able to cluster the trajectories correctly despite within-cluster variations.

5.2. Switching HMM

We ran 100 trials with synthetic data in which parameter matrices for a 2-regime SHMM were generated and a training sequence was sampled from the resulting model. The regime transition matrix was strongly diagonal. Each regime had 3 states. We also generated 10 test sequences of length 1,000 and dimensionality 10. Observation vectors were sampled from the appropriate multinomial distribution with 10 draws. The off-diagonal transition matrix blocks in  $\tilde{A}$  were set to be uniform. We learned the parameters with 2<sup>nd</sup> and 3<sup>rd</sup> order NMF-MoM algorithms as well as ML-EM. The NMF-MoM algorithms were run for 10,000 iterations each and the ML-EM algorithm was terminated when the log likelihood improved by less than 0.01%. The Viterbi algorithm was run on the test sequences and the decoding accuracy was measured via regime labeling. Thus, the specific state reported within a regime is irrelevant.

Accuracy and training time results for these experiments are shown in Figure 3. The NMF-MoM algorithms perform better than ML-EM in terms of regime labeling. They also don't suffer from a rapid increase in time complexity with increasing sequence length that is incurred due to the

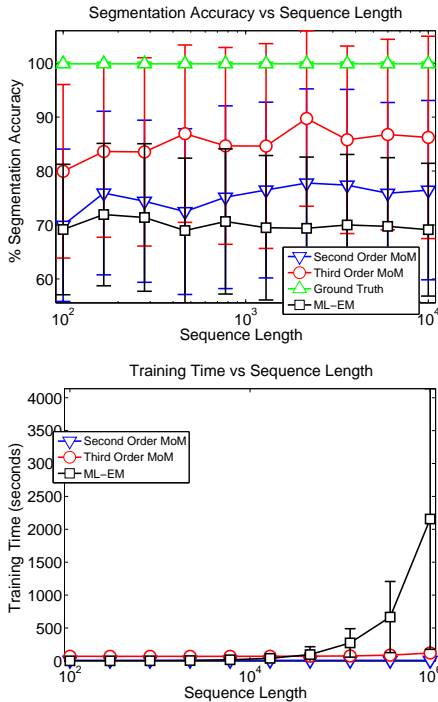


Figure 3. Results of synthetic trials for the switching HMM. **Top:** Segmentation accuracy. **Bottom:** Training time.

forward and backward passes of EM. The 3<sup>rd</sup> order NMF-MoM algorithm performs the best since it is able to take advantage of higher-order moment information.

### 5.3. Factorial HMM

The same synthetic trials were run for the factorial HMM as were run for the switching HMM. The 2<sup>nd</sup> and 3<sup>rd</sup> order NMF-MoM algorithms were run on the training sequence for 1,000 iterations each. The Viterbi algorithm was then applied to find the most likely state sequences for the test data.

The decoding accuracy and training time of the algorithms are shown in Figure 4. The 3<sup>rd</sup> order algorithm performs better than the 2<sup>nd</sup> order one, but it also takes significantly longer to run. This is due both to the increased amount of moment information and to the increased complexity of the updates.

We note that the ML-EM algorithm corresponding to the multinomial FHMM is non-trivial to implement since a closed-form expression for the mean updates cannot be found. These updates can be replaced with a gradient descent procedure. However, evaluating the objective function in this optimization step (i.e. data log likelihood) requires a forward pass along the observation sequence. In addition, we must consider  $J^K$  possible state pairs in both

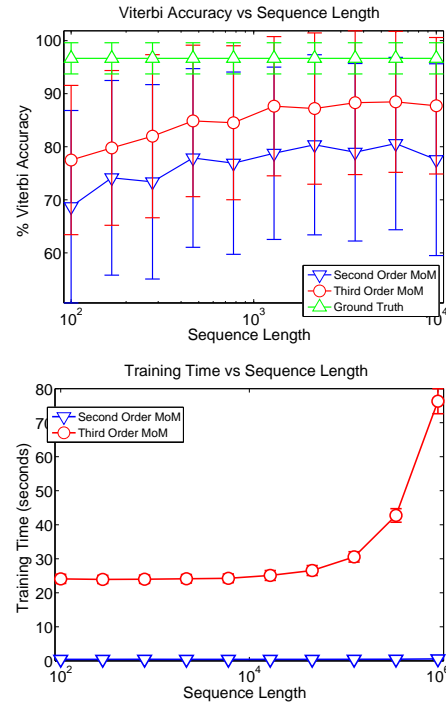


Figure 4. Results of synthetic trials for the factorial HMM. **Top:** Viterbi decoding accuracy. **Bottom:** Training time.

the E and M steps. These factors make the ML-EM algorithm extremely computationally prohibitive and so we have omitted a numerical comparison with it.

## 6. Conclusions

In this paper, we introduced a learning framework for time series models such as the mixture of HMMs, switching HMM, and factorial HMM. This involves computing observable moments and factorizing a non-negative moment matrix/tensor in terms of the model parameters. This is a general approach since we only require that the factorization exists. The main advantage of the proposed NMF-MoM algorithms is computational. Once moments are computed, the learning procedure is independent of the amount of data. We showed that for large datasets, this advantage is substantial. This suggests that this approach may be beneficial for on-line learning. At the same time, we observed that NMF-MoM often performs just as well as ML, if not better. As future work, it may also be beneficial to consider other divergence measures besides the KL criterion used in this work. Finally, we note that NMF-MoM may reach local optima. It is important to explore clever initialization strategies and methods for finding solutions with uniqueness guarantees, as has been done in the spectral learning community.



880	<b>References</b>	
881	Anandkumar, A., Foster, D. P., Hsu, D., Kakade, S.M., and	Mysore, Gautham J., Smaragdis, Paris, and Raj, Bhiksha.
882	Liu, Y.K. Two svds suffice: Spectral decompositions for	Non-negative Hidden Markov Modeling of Audio with
883	probabilistic topic modeling and latent dirichlet alloca-	Application to Source Separation. <i>In Proceedings of</i>
884	tion. In <i>Neural Information Processing Systems</i> , 2012a.	<i>the International Conference on Latent Variable Anal-</i>
885		<i>ysis and Signal Separation (LVA / ICA)</i> , 2010.
886	Anandkumar, A., Ge, R., Hsu, D., Kakade, S.M., and Tel-	935
887	garsky, M. Tensor decompositions for learning latent	936
888	variable models. <i>arXiv:1210.7559v2</i> , 2012b.	937
889		938
890	Anandkumar, A., Hsu, D., and Kakade, S.M. A method of	939
891	moments for mixture models and hidden markov mod-	940
892	els. In <i>COLT</i> , 2012c.	941
893		942
894	Bache, K. and Lichman, M. UCI machine learning repos-	943
895	itory, 2013. URL <a href="http://archive.ics.uci.edu/ml">http://archive.ics.uci.</a>	944
896	<a href="http://archive.ics.uci.edu/ml">edu/ml</a> .	945
897		946
898	Cybenko, G. and Crespi, V. Learning Hidden Markov Mod-	947
899	els using nonnegative matrix factorization. <i>IEEE Trans-</i>	948
900	<i>actions on Information Theory</i> , 57(6):3963–3970, 2011.	949
901		950
902	Dempster, A. P., Laird, N. M., and Rubin, D. B. Maximum	951
903	likelihood from incomplete data via the EM algorithm.	952
904	<i>Journal of the Royal Statistical Society, Series B</i> , 39(1):	953
905	1–38, 1977.	954
906		955
907	Ghahramani, Z. and Jordan, M. I. Factorial hidden Markov	956
908	models. <i>Journal of Machine Learning</i> , 29(2-3):245–273,	957
909	1997.	958
910		959
911	Hsu, Daniel and Kakade, Sham M. Learning mixtures	960
912	of spherical gaussians: moment methods and spectral	961
913	decompositions. In <i>Fourth Innovations in Theoretical</i>	962
914	<i>Computer Science</i> , 2013.	963
915		964
916	Hsu, Daniel, Kakade, Sham M., and Zhang, Tong. A	965
917	spectral algorithm for learning hidden markov models a	966
918	spectral algorithm for learning hidden markov models.	967
919	<i>Journal of Computer and System Sciences</i> , (1460-1480),	968
920	2009.	969
921		970
922	Lakshminarayanan, Balaji and Raich, Raviv. Non-negative	971
923	matrix factorization for parameter estimation in hidden	972
924	markov models. <i>Machine Learning for Signal Process-</i>	973
925	<i>ing (MLSP)</i> , 2010 <i>IEEE International Workshop on</i> , pp.	974
926	89–94, 2010.	975
927		976
928	Lee, D. and Seung, H. S. Learning the parts of objects by	977
929	non-negative matrix factorization. <i>Nature</i> , 401:788–791,	978
930	1999.	979
931		980
932	Lee, D and Seung, H S. Algorithms for non-negative matrix	981
933	factorization. <i>Conference on Advances in Neural Infor-</i>	982
934	<i>mation Processing Systems</i> , pp. 556–562, 2001.	983
		984
		985
		986
		987
		988
		989