

A Spectral Algorithm for Latent Dirichlet Allocation

Anima Anandkumar · Dean P. Foster · Daniel Hsu ·
Sham M. Kakade · Yi-Kai Liu

Received: 1 October 2013 / Accepted: 12 June 2014 / Published online: 3 July 2014
© Springer Science+Business Media New York 2014

Abstract Topic modeling is a generalization of clustering that posits that observations (words in a document) are generated by *multiple* latent factors (topics), as opposed to just one. The increased representational power comes at the cost of a more challenging unsupervised learning problem for estimating the topic-word distributions when only words are observed, and the topics are hidden. This work provides a simple and efficient learning procedure that is guaranteed to recover the parameters for a wide class of multi-view models and topic models, including latent Dirichlet allocation (LDA). For LDA, the procedure correctly recovers both the topic-word distributions and the parameters of the Dirichlet prior over the topic mixtures, using only trigram statistics (i.e., third order moments, which may be estimated with documents containing just

Preliminary versions of this article appeared as [3,4].

A. Anandkumar
University of California, Irvine, Irvine, CA, USA
e-mail: a.anandkumar@uci.edu

D. P. Foster
Yahoo! Labs, New York, NY, USA
e-mail: dean@foster.net

D. Hsu (✉)
Columbia University, New York, NY, USA
e-mail: djhsu@cs.columbia.edu

S. M. Kakade
Microsoft Research, Cambridge, MA, USA
e-mail: skakade@microsoft.com

Y.-K. Liu
National Institute of Standards and Technology, Gaithersburg, MD, USA
e-mail: yi-kai.liu@nist.gov

three words). The method is based on an efficiently computable orthogonal tensor decomposition of low-order moments.

Keywords Topic models · Mixture models · Method of moments · Latent Dirichlet allocation

1 Introduction

Topic models use latent variables to explain the observed (co-)occurrences of words in documents. They posit that each document is associated with a (possibly sparse) mixture of active topics, and that each word in the document is accounted for (in fact, generated) by one of these active topics. In latent Dirichlet allocation (LDA) [12], a Dirichlet prior gives the distribution of active topics in documents. LDA and related models possess a rich representational power because they allow for documents to be comprised of words from several topics, rather than just a single topic. This increased representational power comes at the cost of a more challenging unsupervised estimation problem, when only the words are observed and the corresponding topics are hidden.

In practice, the most common unsupervised estimation procedures for topic models are based on finding maximum likelihood estimates, through either local search or sampling based methods, e.g., Expectation-Maximization [43], Gibbs sampling [22], and variational approaches [12]. Another body of tools is based on matrix factorization [26, 36]. For document modeling, a typical goal is to form a sparse decomposition of a term-document matrix (which represents the word counts in each document) into two parts: one which specifies the active topics in each document and the other which specifies the distributions of words under each topic.

This work provides an alternative approach to parameter recovery based on the method of moments [42], which attempts to match the observed moments with those posited by the model. Our approach does this efficiently through a particular decomposition of the low-order observable moments, which can be extracted using an orthogonal tensor decomposition. This method is simple and efficient to implement, and is guaranteed to recover the parameters of a wide class of topic models, including the LDA model. We exploit exchangeability of the observed variables and, more generally, the availability of multiple views drawn independently from the same hidden component.

1.1 Summary of Contributions

We present a spectral approach to topic model estimation based on decomposing the low-order (cross) moments of observed variables. The approach differs from other spectral methods (e.g., those based on Principal Component Analysis and Canonical Correlation Analysis) in that it is based on an orthogonal tensor decomposition of a $k \times k \times k$ third-order moment tensor, where k is the number of latent factors (topics). In many applications, k is typically much smaller than the dimension of the observed space d (number of words).

The method is applicable to a wide class of latent variable models including exchangeable and multi-view models. We first consider the class of exchangeable variables with independent latent factors. We show that the low-order moment tensors possess a decomposition known as a canonical polyadic decomposition [34] (Lemma 1 and Lemma 2 in Sect. 3.1). We then consider LDA and show that, even though it does not directly possess an independent latent factor, it *nearly* does so in a rigorous sense, and hence a simple combination of lower-order moments has the required decomposition just as in the previous case (Lemma 3 in Sect. 3.2). Given these moments from either the independent latent factors model or LDA, a simple and computationally efficient algebraic procedure [23] recovers the model parameters exactly under a mild rank condition (Theorem 2 in Sect. 3.3). Informally, we have the following.

Theorem 1 (Informal) *Given low-order (i.e., order ≤ 3 or ≤ 4) moments from either the independent latent factors model or LDA satisfying a rank condition, there is a polynomial time randomized algorithm that returns the model parameters up to scaling and permutation.*

We describe the parameter recovery procedures assuming exact moments as input, but it is straightforward to write down analogous “plug-in” estimators that use estimates of these moments based on sampled data. We do just this using a particular tensor decomposition procedure from [5], and analyze the error of the resulting parameter estimates in terms of the errors in the moment estimates (Theorem 3).

1.2 Related Work

Under the assumption that a single active topic occurs in each document, the work of [41] provides the first theoretical guarantees for recovering the topic distributions (i.e., the distribution of words under each topic), albeit with a rather stringent separation condition (where the words in each topic are essentially non-overlapping). Understanding what separation conditions permit efficient learning is a natural question; in the clustering literature, a line of work has focussed on understanding the relationship between the separation of the mixture components and the complexity of learning. For clustering, the first learnability result [19] was under a rather strong separation condition; subsequent results relaxed [1, 10, 16, 17, 20, 33, 45] or removed these conditions [11, 28, 32, 38]; roughly speaking, learning under a weaker separation condition is more challenging, both computationally and statistically. For the topic modeling problem in which only a single active topic is present per document, [6] provides an algorithm for learning topics with no separation requirement, but under a certain full rank assumption on the topic probability matrix.

For the case of LDA (where each document may be about multiple topics), the recent work of [8] provides the first theoretical result under a natural separation condition. The condition requires that each topic be associated with “anchor words” that only occur in documents about that topic. This is a significantly milder assumption than the one in [41]. Under this assumption, [8] provides the first rigorously analyzed algorithm for learning the topic distributions. Their work also justifies the use of non-negative matrix (NMF) as a provable procedure for this problem (the original motivation for NMF was

as a topic modeling technique, though, prior to this work, formal guarantees as such were rather limited). Furthermore, [8] provides results for certain correlated topic models. Our approach makes further progress on this problem by relaxing the need for this separation condition and providing a simpler parameter estimation procedure.

The underlying approach we take is a certain diagonalization technique of the observed moments. We know of at least three different settings which use this idea for parameter estimation.

The algebraic tensor decomposition technique of [23] (see also [35,37])¹ was rediscovered by [15] for parameter estimation in discrete Markov models. The idea has been extended to other discrete mixture models such as discrete hidden Markov models (HMMs) and mixture models with a single active topic in each document (see [6,29,39]). For such single topic models, the work in [6] demonstrates the generality of the eigenvector method and the irrelevance of the noise model for the observations, making it applicable to both discrete models like HMMs as well as certain Gaussian mixture models (see also [28]).

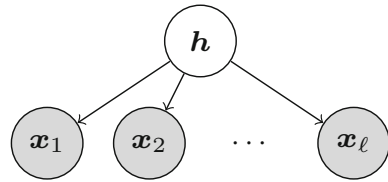
Another set of related techniques is the body of algebraic methods used for the problem of blind source separation [14]. These approaches are tailored for independent source separation with additive noise (usually Gaussian) [18]. Much of the literature focuses on understanding the effects of measurement noise, which often requires more sophisticated algebraic tools (typically, knowledge of noise statistics or the availability of multiple views of the latent factors is not assumed). These ideas are also used by [21,40] for learning a linear transformation (in a noiseless setting) and provides a different algorithm, based on a certain ascent algorithm (rather than joint diagonalization approach, as in [14]), and a provably correct algorithm for the noisy case was only recently obtained [9].

The models we consider, including LDA, are distinguished by the presence of exchangeable (or multi-view) variables (e.g., multiple words in a document), drawn independently conditioned on the same hidden state. This allows us to exploit ideas from the aforementioned works [14,15,23]. In particular, we show that the topic modeling problem exhibits a similarly simple algebraic solution as in previous works. Furthermore, the exchangeability assumption permits us to have an *arbitrary* noise model (rather than an additive Gaussian noise, which is not appropriate for multinomial and other discrete distributions). A key technical contribution is that we show how the basic diagonalization approach can be adapted for Dirichlet models, through a rather careful construction. This construction bridges the gap between mixture models and independent latent factors models.

The multi-view approach has been exploited in previous works for semi-supervised learning and for learning mixtures of well-separated distributions (e.g., [7,16,17,31]). These previous works essentially use variants of canonical correlation analysis [27] between the two views. This present work follows [6] in showing that having a third view of the data permits rather simple estimation procedures for parameter recovery without separation conditions.

¹ The technique of [23] is actually attributed to Robert Jennrich.

Fig. 1 Directed graphical model representation of the multi-view models



A preliminary version of this work [3,4] proposed a method based on using two singular value decompositions—essentially a symmetrized version of the algorithm from [6]. However, that method is not particularly robust, as it depends on randomization to create separation between eigenvalues of a particular matrix. In this article, we propose to use the robust orthogonal tensor decomposition method from [5], which is more robust than the previous method and still computationally efficient.

2 The Multi-view Models

Let $\mathbf{h} = (h_1, h_2, \dots, h_k) \in \mathbb{R}^k$ be a random vector specifying the latent factors (i.e., \mathbf{h} is the hidden state) in a model, where h_i is the value of the i th factor. Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell \in \mathbb{R}^d$ be random vectors which we take to be observable. Throughout, we assume $\ell \geq 3$, and the reason will become clear in the sequel. The primary modeling assumption is that the observable random vectors and latent factors satisfy the following conditional independence and linearity condition; we assume this condition holds the remainder of this paper.

Condition 1 [Multi-view model] *The observations $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell$ are conditionally independent given \mathbf{h} . Moreover, for each $v \in [\ell] := \{1, 2, \dots, \ell\}$, there exists a matrix $\mathbf{O}^{(v)} \in \mathbb{R}^{d \times k}$ such that*

$$\mathbb{E}[\mathbf{x}_v | \mathbf{h}] = \mathbf{O}^{(v)} \mathbf{h}.$$

The conditional independence assumption from Condition 1 is depicted in the directed graphical model in Fig. 1. This model generalizes the multi-view model from [6] in that \mathbf{h} is *not* assumed to only take values in $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_k\}$, where $\mathbf{e}_i \in \{0, 1\}^k$ is the i th coordinate basis vector. We note that at this stage, we have not made any assumptions on the noise model; it need not be additive nor even independent of \mathbf{h} .²

We also assume throughout that the following rank condition on the matrices $\mathbf{O}^{(v)}$ from Condition 1.

Condition 2 (Rank condition) *Each $\mathbf{O}^{(v)}$ has full column rank.*

Condition 2 allows for the identifiability of the columns of $\mathbf{O}^{(v)}$ and was used in previous works on parameter estimation [6, 8, 9, 15, 21, 28, 39, 40].

As shown in [6] (for the case where \mathbf{h} has support only on k points), the multi-view model captures a number of interesting statistical models, including certain Gaussian

² By additive noise, we mean a model in which $\mathbf{x}_v = \mathbf{O}^{(v)} \mathbf{h} + \boldsymbol{\eta}_v$ where $\boldsymbol{\eta}_v$ is a zero-mean random vector independent of \mathbf{h} .

mixture models and HMMs. In our setting, \mathbf{h} is allowed to have a more general distribution, thus enhancing the flexibility of the model. Our goal in this work is to given estimators of the matrices $\mathbf{O}^{(v)}$ (and possibly parameters related to the latent factor \mathbf{h}), solely from repeated observations (independent copies of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell$).

We now consider two specializations of this multi-view model in which the different views $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell$ are naturally *exchangeable*. For these cases, the conditional means of the different views are the same. That is, the following condition holds.

Condition 3 $\mathbf{O}^{(v)} = \mathbf{O}$ for all $v \in [\ell]$.

In the context of topic models, the common matrix \mathbf{O} is referred to as the *topic matrix*, as it specifies parameters associated with each topic in the model. Borrowing this terminology, we generally refer to all $\mathbf{O}^{(v)}$ as topic matrices.

2.1 Independent Latent Factors Model

In the independent latent factors model, we assume \mathbf{h} has a product distribution, i.e., h_1, h_2, \dots, h_k are independent. In the case where the \mathbf{x}_v are deterministic linear functions of \mathbf{h} , (i.e., $\mathbf{x}_v = \mathbf{O}\mathbf{h}$), the model reduces to the noiseless ICA model of [30], which was reinterpreted as a parallelepiped learning problem in [21,40].

Two important examples of this setting with noise are as follows.

Mixture of Gaussians Suppose $\mathbf{x}_v = \mathbf{O}\mathbf{h} + \boldsymbol{\eta}$, where $\boldsymbol{\eta}$ is zero-mean Gaussian noise and \mathbf{h} has a product distribution over $\{0, 1\}^k$. Here, the i th column \mathbf{o}_i of \mathbf{O} can be considered to be the mean of the i th Gaussian component. This generalizes the classic mixture of k Gaussians, as the model now permits any number of Gaussians to be responsible for generating an observation (i.e., \mathbf{h} is permitted to be any of the 2^k vectors on the hypercube, while in the classic mixture problem, only one component is responsible). Alternatively, the model can be viewed as a classical mixture of 2^k Gaussians, where the mean of a component $S \in 2^{[k]}$ is $\sum_{i \in S} \mathbf{o}_i$ and its mixing weight is $\Pr(h_i = 1 \forall i \in S)$ [9]. Note that $\boldsymbol{\eta}$ is allowed to be heteroskedastic (i.e., the noise may depend on \mathbf{h}), so the Gaussians need not have the same covariance.

Poisson topic model This is a model proposed by Canny [13]. Suppose $[\mathbf{O}\mathbf{h}]_j$ specifies the Poisson rate of counts for $[\mathbf{x}_v]_j$. For example, \mathbf{x}_v could be a vector of word counts in the v th sentence of a document. Here, \mathbf{O} would be a matrix with positive entries, and h_i would scale the rate at which topic i generates words in a sentence (as specified by the i th column of \mathbf{O}). The linearity assumption in Condition 1 is satisfied as $\mathbb{E}[\mathbf{x}_v | \mathbf{h}] = \mathbf{O}\mathbf{h}$ (note the noise is not additive in this case, in contrast to the mixture of Gaussians example). Here, multiple topics may be responsible for generating the words in each sentence. This model similar to LDA, except for the fact that here, \mathbf{h} has a product distribution (whereas in LDA, \mathbf{h} is a probability vector).

2.2 The Dirichlet Model

Now suppose the hidden state \mathbf{h} is a probability distribution itself, with a density specified by the Dirichlet distribution with parameter vector $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_k) \in$

$\mathbb{R}_{>0}^k$ ($\boldsymbol{\alpha}$ is a strictly positive real vector); this assumption is written as $\mathbf{h} \sim \text{Dir}(\boldsymbol{\alpha})$. We think of $\mathbf{h} \in \Delta^{k-1}$ as a distribution over topics; here, Δ^{k-1} denotes the probability simplex of distributions over k outcomes. The density of \mathbf{h} , the Dirichlet density, is

$$p_{\boldsymbol{\alpha}}(\mathbf{h}) = \frac{1}{Z(\boldsymbol{\alpha})} \prod_{i=1}^k h_i^{\alpha_i - 1}$$

where $Z(\boldsymbol{\alpha}) := \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\alpha_0)}$ and $\alpha_0 := \alpha_1 + \alpha_2 + \dots + \alpha_k$. Intuitively, α_0 (the sum of the “pseudo-counts”) characterizes the concentration of the distribution. As $\alpha_0 \rightarrow 0$, the distribution degenerates to one over pure topics (i.e., the limiting density is one in which, almost surely, exactly one coordinate of \mathbf{h} is 1, and the rest are 0).

If the \mathbf{x}_v are deterministic linear functions of \mathbf{h} (i.e., $\mathbf{x}_v = \mathbf{O}\mathbf{h}$), then the model can be viewed as the problem of learning a certain class of distributions over a simplex with vertices $\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_k$. The special case of a uniform distribution over a simplex (a problem suggested in [21]) is obtained when $\boldsymbol{\alpha} = (1, 1, \dots, 1)$.

Latent Dirichlet Allocation LDA makes the further assumption that each random vector $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell$ takes on discrete values out of d outcomes (e.g., \mathbf{x}_v represents what the v th word in a document is, so d represents the number of words in a vocabulary). The i th column \mathbf{o}_i of \mathbf{O} is a probability vector representing the distribution over words for the i th topic. The sampling process for a document is as follows. First, the topic mixture \mathbf{h} is drawn from the Dirichlet distribution. Then, the v th word in the document (for $v \in [\ell] := \{1, 2, \dots, \ell\}$) is generated by: (i) drawing $t \in [k]$ according to the discrete distribution specified by h , then (ii) drawing \mathbf{x}_v according to the discrete distribution specified by \mathbf{o}_t . Note that \mathbf{x}_v is independent of \mathbf{h} given t . For this model to fit in our setting, we use the “one-hot” encoding for \mathbf{x}_v (also used by [6] in this context): $\mathbf{x}_v \in \{0, 1\}^d$ with

$$[\mathbf{x}_v]_j = 1 \iff \text{the } v\text{th word in the document is the } j\text{th vocabulary word.}$$

Observe that

$$\mathbb{E}[\mathbf{x}_v | \mathbf{h}] = \sum_{i=1}^k \Pr[t = i | \mathbf{h}] \cdot \mathbb{E}[\mathbf{x}_v | t = i, \mathbf{h}] = \sum_{i=1}^k h_i \cdot \mathbf{o}_i = \mathbf{O}\mathbf{h}$$

as required in Condition 1. Again, note that the noise model is not additive.

3 Moment Structure in Multi-view Models

In this section, we describe the structure of moments for the multi-view models, focusing primarily on the exchangeable models (the independent latent factors model and the Dirichlet model). Recall that for these exchangeable models, we assume $\mathbf{O}^{(v)} = \mathbf{O}$ for all $v \in [\ell]$ (Condition 3). In all cases, the p th-order (not necessarily raw) moments are of the form

$$\sum_{i=1}^k c_{i,p} \underbrace{\mathbf{o}_i \otimes \mathbf{o}_i \otimes \cdots \otimes \mathbf{o}_i}_{p \text{ times}},$$

where $c_{1,p}, c_{2,p}, \dots, c_{k,p} \in \mathbb{R}$ are scalars possibly depending on p . This is known as a symmetric canonical polyadic decomposition [24, 25, 34] (or a Tucker decomposition with a diagonal core tensor [44]). We show that if the $c_{i,p}$ are non-zero, then the \mathbf{o}_i (up to some scaling and permutation) can be extracted from these moments using simple algebraic techniques such as that from [23]. We also give a reduction from the general multi-view setting to the exchangeable setting.

3.1 Moments of Skewed and Kurtotic Independent Factors

Let $m_{i,p} := \mathbb{E}[(h_i - \mathbb{E}[h_i])^p]$ denote the p th central moment of h_i . The variance, skewness, and kurtosis of h_i are given by $\sigma_i^2 := m_{i,2}$, $\gamma_i := m_{i,3}/\sigma_i^3$, and $\kappa_i := m_{i,4}/\sigma_i^4 - 3$, respectively.

Define the following moments of the observable random vectors:

$$\begin{aligned} \boldsymbol{\mu} &:= \mathbb{E}[\mathbf{x}_1] \in \mathbb{R}^d, \\ \mathbf{Pairs} &:= \mathbb{E}[(\mathbf{x}_1 - \boldsymbol{\mu}) \otimes (\mathbf{x}_2 - \boldsymbol{\mu})] \in \mathbb{R}^{d \times d}, \\ \mathbf{Triples} &:= \mathbb{E}[(\mathbf{x}_1 - \boldsymbol{\mu}) \otimes (\mathbf{x}_2 - \boldsymbol{\mu}) \otimes (\mathbf{x}_3 - \boldsymbol{\mu})] \in \mathbb{R}^{d \times d \times d}. \end{aligned} \tag{1}$$

Here, we use \otimes to denote the tensor product. For instance, given vectors $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{R}^d$, the tensor product $\mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} \in \mathbb{R}^{d \times d \times d}$ is the third-order tensor whose (i, j, k) th entry is $u_i v_j w_k$.

Lemma 1 (Independent latent factors moments) *Assume Conditions 1 and 3 and that \mathbf{h} has a product distribution. Then*

$$\boldsymbol{\mu} = \sum_{i=1}^k m_{i,1} \mathbf{o}_i, \quad \mathbf{Pairs} = \sum_{i=1}^k m_{i,2} \mathbf{o}_i \otimes \mathbf{o}_i, \quad \mathbf{Triples} = \sum_{i=1}^k m_{i,3} \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i.$$

Proof Since $\mathbb{E}[\mathbf{x}_v | \mathbf{h}] = \mathbf{O}\mathbf{h}$ (by Conditions 1 and 3) we have

$$\boldsymbol{\mu} = \mathbf{O}\mathbb{E}[\mathbf{h}] = \sum_{i=1}^k m_{i,1} \mathbf{o}_i.$$

This also implies

$$\mathbb{E}[(\mathbf{x}_v - \boldsymbol{\mu}) | \mathbf{h}] = \mathbf{O}(\mathbf{h} - \mathbb{E}[\mathbf{h}]) = \sum_{i=1}^k (h_i - \mathbb{E}[h_i]) \mathbf{o}_i.$$

Since $\mathbf{x}_1 - \boldsymbol{\mu}$ and $\mathbf{x}_2 - \boldsymbol{\mu}$ are conditionally independent given \mathbf{h} (under Condition 1), we may write **Pairs** as the expectation of a tensor product:

$$\mathbf{Pairs} = \mathbb{E} \left[\left(\sum_{i=1}^k (h_i - \mathbb{E}[h_i]) \mathbf{o}_i \right) \otimes \left(\sum_{i=1}^k (h_i - \mathbb{E}[h_i]) \mathbf{o}_i \right) \right].$$

Expanding out the tensor product and applying linearity leaves only the diagonal terms with $\mathbb{E}[(h_i - \mathbb{E}[h_i])^2]$. This is because \mathbf{h} has a product distribution by assumption, and thus $\mathbb{E}[(h_i - \mathbb{E}[h_i])(h_j - \mathbb{E}[h_j])] = 0$ for $i \neq j$. With only the diagonal terms in the product remaining, we have

$$\mathbf{Pairs} = \sum_{i=1}^k \mathbb{E}[(h_i - \mathbb{E}[h_i])^2] \mathbf{o}_i \otimes \mathbf{o}_i = \sum_{i=1}^k m_{i,2} \mathbf{o}_i \otimes \mathbf{o}_i.$$

An analogous argument gives the claim for **Triples**. □

Lemma 1 shows that $\boldsymbol{\mu}$, **Pairs**, and **Triples** possess the structure needed to apply the tensor decomposition technique of [5] to extract the \mathbf{o}_i . However, this is only possible if the scalar factors are non-zero; in Lemma 1, the scalar factors are the central moments $m_{i,p}$. If \mathbf{h} has a distribution symmetric about its mean, the third central moment is zero, and hence **Triples** cannot be used. One recourse comes from the literature on independent component analysis; if the distribution of \mathbf{h} is kurtotic (i.e., $\kappa_i \neq 0$), then one may use fourth-order moments in the form of the fourth cumulant tensor [14]. In the next lemma, we show that this can also be applied with multi-view/exchangeable models.

Define

$$\mathbf{Quadruples} := \mathbb{E}[(\mathbf{x}_1 - \boldsymbol{\mu}) \otimes (\mathbf{x}_2 - \boldsymbol{\mu}) \otimes (\mathbf{x}_3 - \boldsymbol{\mu}) \otimes (\mathbf{x}_4 - \boldsymbol{\mu})] - \mathbf{T} \quad (2)$$

where $\mathbf{T} \in \mathbb{R}^{d \times d \times d \times d}$ is the fourth-order tensor whose (i, j, m, n) th entry is

$$T_{i,j,m,n} = [\mathbf{Pairs}]_{i,j} [\mathbf{Pairs}]_{m,n} + [\mathbf{Pairs}]_{i,m} [\mathbf{Pairs}]_{j,n} + [\mathbf{Pairs}]_{i,n} [\mathbf{Pairs}]_{j,m}.$$

Lemma 2 (Independent latent factors moments, fourth-order) *Under the same setting as Lemma 1,*

$$\mathbf{Quadruples} = \sum_{i=1}^k (m_{i,4} - 3m_{i,2}^2) \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i.$$

Proof By the same argument as in Lemma 1, we have that

$$\mathbb{E}[(\mathbf{x}_1 - \boldsymbol{\mu}) \otimes (\mathbf{x}_2 - \boldsymbol{\mu}) \otimes (\mathbf{x}_3 - \boldsymbol{\mu}) \otimes (\mathbf{x}_4 - \boldsymbol{\mu})] = \mathbb{E}[\mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v}]$$

where

$$\mathbf{v} := \sum_{i=1}^k (h_i - \mathbb{E}[h_i]) \mathbf{o}_i.$$

Expanding the tensor product and applying linearity of expectation leaves only terms involving $\mathbb{E}[(h_i - \mathbb{E}[h_i])^4]$ and $\mathbb{E}[(h_i - \mathbb{E}[h_i])^2 (h_j - \mathbb{E}[h_j])^2]$, again due to the product distribution of \mathbf{h} . Thus

$$\begin{aligned} & \mathbb{E}[\mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v}] \\ &= \sum_{i=1}^k m_{i,4} \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \\ & \quad + \sum_{i \neq j} m_{i,2} m_{j,2} (\mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_j \otimes \mathbf{o}_j + \mathbf{o}_i \otimes \mathbf{o}_j \otimes \mathbf{o}_i \otimes \mathbf{o}_j + \mathbf{o}_i \otimes \mathbf{o}_j \otimes \mathbf{o}_j \otimes \mathbf{o}_i) \\ &= \sum_{i=1}^k (m_{i,4} - 3m_{i,2}^2) \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \\ & \quad + \sum_{i,j} m_{i,2} m_{j,2} (\mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_j \otimes \mathbf{o}_j + \mathbf{o}_i \otimes \mathbf{o}_j \otimes \mathbf{o}_i \otimes \mathbf{o}_j + \mathbf{o}_i \otimes \mathbf{o}_j \otimes \mathbf{o}_j \otimes \mathbf{o}_i) \\ &= \sum_{i=1}^k (m_{i,4} - m_{i,2}^2) \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i + \mathbf{T} \end{aligned}$$

where the step uses the identity for **Pairs** from Lemma 1. Rearranging the terms gives the desired identity. \square

We note that if \mathbf{h} is an isotropic Gaussian random vector, then both $m_{i,3}$ and $m_{i,4} - 3m_{i,2}^2$ are zero for all $i \in [k]$, causing both **Triples** and **Quadruples** to vanish. As expected, these higher-order moments cannot help with the identification of \mathbf{O} without further assumptions: this is simply because $\mathbf{g} := \mathbf{U}^\top \mathbf{h}$ has the same distribution as \mathbf{h} when \mathbf{U} is orthogonal, and $\mathbb{E}[\mathbf{x}_v | \mathbf{h}] = \mathbf{O}\mathbf{h} = (\mathbf{O}\mathbf{U})\mathbf{g}$.

3.2 Moments of Dirichlet Factors

The Dirichlet distribution is not a product distribution, and therefore Lemma 1 does not apply to models where $\mathbf{h} \sim \text{Dir}(\boldsymbol{\alpha})$. However, it is *nearly* a product distribution. Indeed, if $\mathbf{h} \sim \text{Dir}(\boldsymbol{\alpha})$, then \mathbf{h} has the same distribution $\mathbf{x} / \sum_{i=1}^k x_i$, where $x_i \sim \text{Gamma}(\alpha_i, 1)$ (i.e., x_i follows a Gamma distribution with shape parameter α_i and scale parameter 1), and x_1, x_2, \dots, x_k are independent. Thus, the essential reason why the Dirichlet distribution is not a product distribution is because it is restricted to the probability simplex.

Lemma 3 (below) nevertheless shows that if the concentration parameter α_0 is known, the correlations introduced by restricting to the probability simplex can be

accounted for by slightly tweaking the moments; this tweaking causes the same kind of structure as in Lemma 1 to manifest.

Define the following moments of the observable random vectors:

$$\begin{aligned}
 \boldsymbol{\mu} &:= \mathbb{E}[\mathbf{x}_1] \in \mathbb{R}^d, \\
 \mathbf{Pairs}_{\alpha_0} &:= \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2] - \frac{\alpha_0}{\alpha_0 + 1} \boldsymbol{\mu} \otimes \boldsymbol{\mu} \in \mathbb{R}^{d \times d}, \\
 \mathbf{Triples}_{\alpha_0} &:= \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3] \\
 &\quad - \frac{\alpha_0}{\alpha_0 + 2} \left(\mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \boldsymbol{\mu}] + \mathbb{E}[\mathbf{x}_1 \otimes \boldsymbol{\mu} \otimes \mathbf{x}_3] + \mathbb{E}[\boldsymbol{\mu} \otimes \mathbf{x}_2 \otimes \mathbf{x}_3] \right) \\
 &\quad + \frac{2\alpha_0^2}{(\alpha_0 + 2)(\alpha_0 + 1)} \boldsymbol{\mu} \otimes \boldsymbol{\mu} \otimes \boldsymbol{\mu} \in \mathbb{R}^{d \times d \times d}. \tag{3}
 \end{aligned}$$

Lemma 3 (Dirichlet factors moments) *Assume Conditions 1 and 3 and that $\mathbf{h} \sim \text{Dir}(\boldsymbol{\alpha})$. Then*

$$\begin{aligned}
 \boldsymbol{\mu} &= \sum_{i=1}^k \frac{\alpha_i}{\alpha_0} \mathbf{o}_i, \quad \mathbf{Pairs}_{\alpha_0} = \sum_{i=1}^k \frac{\alpha_i}{(\alpha_0 + 1)\alpha_0} \mathbf{o}_i \otimes \mathbf{o}_i, \\
 \mathbf{Triples}_{\alpha_0} &= \sum_{i=1}^k \frac{2\alpha_i}{(\alpha_0 + 2)(\alpha_0 + 1)\alpha_0} \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i.
 \end{aligned}$$

Proof We directly calculate univariate, bivariate, and trivariate moments of \mathbf{h} : for any distinct $i, j, l \in [k]$,

$$\begin{aligned}
 \mathbb{E}[h_i] &= \frac{\alpha_i}{\alpha_0}, \quad \mathbb{E}[h_i^2] = \frac{(\alpha_i + 1)\alpha_i}{(\alpha_0 + 1)\alpha_0}, \quad \mathbb{E}[h_i^3] = \frac{(\alpha_i + 2)(\alpha_i + 1)\alpha_i}{(\alpha_0 + 2)(\alpha_0 + 1)\alpha_0}, \\
 \mathbb{E}[h_i h_j] &= \frac{\alpha_i \alpha_j}{(\alpha_0 + 1)\alpha_0}, \quad \mathbb{E}[h_i^2 h_j] = \frac{(\alpha_i + 1)\alpha_i \alpha_j}{(\alpha_0 + 2)(\alpha_0 + 1)\alpha_0}, \\
 \mathbb{E}[h_i h_j h_l] &= \frac{\alpha_i \alpha_j \alpha_l}{(\alpha_0 + 2)(\alpha_0 + 1)\alpha_0}.
 \end{aligned}$$

Putting these in vector, matrix, and third-order tensor form,

$$\begin{aligned}
 \mathbb{E}[\mathbf{h}] &= \frac{1}{\alpha_0} \sum_{i=1}^k \alpha_i \mathbf{e}_i \\
 \mathbb{E}[\mathbf{h} \otimes \mathbf{h}] &= \frac{1}{(\alpha_0 + 1)\alpha_0} \left(\sum_{i=1}^k \alpha_i \mathbf{e}_i \otimes \mathbf{e}_i + \boldsymbol{\alpha} \otimes \boldsymbol{\alpha} \right) \\
 \mathbb{E}[\mathbf{h} \otimes \mathbf{h} \otimes \mathbf{h}] &= \frac{1}{(\alpha_0 + 2)(\alpha_0 + 1)\alpha_0} \left(2 \sum_{i=1}^k \alpha_i (\mathbf{e}_i \otimes \mathbf{e}_i \otimes \mathbf{e}_i) \right)
 \end{aligned}$$

$$\begin{aligned}
 & + \sum_{i=1}^k \alpha_i (\mathbf{e}_i \otimes \mathbf{e}_i \otimes \boldsymbol{\alpha}) + \sum_{i=1}^k \alpha_i (\boldsymbol{\alpha} \otimes \mathbf{e}_i \otimes \mathbf{e}_i) \\
 & + \sum_{i=1}^k \alpha_i (\mathbf{e}_i \otimes \boldsymbol{\alpha} \otimes \mathbf{e}_i) + \boldsymbol{\alpha} \otimes \boldsymbol{\alpha} \otimes \boldsymbol{\alpha} \Big).
 \end{aligned}$$

Since $\mathbb{E}[\mathbf{x}_v | \mathbf{h}] = \mathbf{O}\mathbf{h}$ by Conditions 1 and 3 (as in the proof of Lemma 1), and because the \mathbf{x}_v are conditionally independent given \mathbf{h} , the claim follows by linearity and rearranging. \square

3.3 Identifiability from Low-order Moments

We now provide a simple argument that establishes the identifiability of the columns of the topic matrix \mathbf{O} (up to scaling and permutation) from the low-order moments considered in Lemmas 1, 2, and 3. Here, Condition 2 plays an essential role, as does the non-Gaussianity of \mathbf{h} in the case of the independent latent factors model.

Theorem 2 (Identifiability from low-order moments) *Assume Conditions 1, 2, and 3 hold.*

1. *If \mathbf{h} has a product distribution and each h_i has positive variance σ_i^2 and non-zero skew γ_i , then there is a randomized algorithm that, given **Pairs** and **Triples** from (1), returns the set of vectors $\{(s_i \sigma_i \mathbf{o}_i, s_i \gamma_i) : i \in [k]\}$ for some $\{s_1, s_2, \dots, s_k\} \subseteq \{\pm 1\}$.*
2. *If \mathbf{h} has a product distribution and each h_i has positive variance σ_i^2 and non-zero kurtosis κ_i , then there is a randomized algorithm that, given **Pairs** and **Quadruples** from (1) and (2), returns $\{(s_i \sigma_i \mathbf{o}_i, s_i \kappa_i) : i \in [k]\}$ for some $\{s_1, s_2, \dots, s_k\} \subseteq \{\pm 1\}$.*
3. *If $\mathbf{h} \sim \text{Dir}(\boldsymbol{\alpha})$ for some $\boldsymbol{\alpha} \in \mathbb{R}_{>0}^k$, then there is a randomized algorithm that, given **Pairs** $_{\alpha_0}$ and **Triples** $_{\alpha_0}$ from (3), returns $\{(s_i \sqrt{c_{i,2}} \mathbf{o}_i, s_i c_{i,3} / c_{i,2}^{3/2}) : i \in [k]\}$ for some $\{s_1, s_2, \dots, s_k\} \subseteq \{\pm 1\}$, where $c_{i,2} = \alpha_i / ((\alpha_0 + 1)\alpha_0)$ and $c_{i,3} = 2\alpha_i / ((\alpha_0 + 2)(\alpha_0 + 1)\alpha_0)$.*

Proof The theorem follows from the structural results from Lemma 1, Lemma 2, and Lemma 3, combined with Lemma 4 (below). \square

Lemma 4 is a variant of a result from [23], which establishes the uniqueness of (symmetric) canonical polyadic decompositions under a rank condition. Our variant uses a symmetrized version of procedure from [23], and is proved using the following observations. First, if the second-order moment matrix (called \mathbf{P} in Lemma 4) has rank k , then it defines an inner product system under which the \mathbf{o}_i are orthogonal. Moreover, under this inner product, the \mathbf{o}_i are orthogonal eigenvectors of a generic flattening of higher-order moment tensors (\mathbf{T} or \mathbf{Q} in Lemma 4).

Below, for a vector $\mathbf{v} = (v_1, v_2, \dots, v_d) \in \mathbb{R}^d$, let $\mathbf{T}\{\mathbf{v}\} \in \mathbb{R}^{d \times d}$ denote the matrix whose (i, j) th entry is $\sum_{l=1}^d v_l [\mathbf{T}]_{i,j,l}$.

Lemma 4 *Let $\mathbf{O} := [\mathbf{o}_1 | \mathbf{o}_2 | \dots | \mathbf{o}_k] \in \mathbb{R}^{d \times k}$, $\mathbf{P} := \sum_{i=1}^k c_{i,2} \mathbf{o}_i \otimes \mathbf{o}_i \in \mathbb{R}^{d \times d}$, $\mathbf{T} := \sum_{i=1}^k c_{i,3} \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \in \mathbb{R}^{d \times d \times d}$, and $\mathbf{Q} := \sum_{i=1}^k c_{i,4} \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \otimes \mathbf{o}_i \in \mathbb{R}^{d \times d \times d \times d}$.*

Algorithm 1 Identification from exact moments

- input** positive integer $k \in \mathbb{N}$, $\mathbf{P} \in \mathbb{R}^{d \times d}$ and $\mathbf{T} \in \mathbb{R}^{d \times d \times d}$ satisfying conditions in Lemma 4.
output $\{(\mathbf{v}_i, \lambda_i) : i \in [k]\}$.
- 1: Let $\sum_{i=1}^k \eta_i \boldsymbol{\xi}_i \otimes \boldsymbol{\xi}_i$ be an eigendecomposition of \mathbf{P} , where $\{\boldsymbol{\xi}_i : i \in [k]\}$ are orthonormal eigenvectors, and $\{\eta_i : i \in [k]\}$ are corresponding positive eigenvalues.
 - 2: Set $\mathbf{W} := [\boldsymbol{\zeta}_1 | \boldsymbol{\zeta}_2 | \dots | \boldsymbol{\zeta}_k] \mathbf{diag}(1/\sqrt{\eta_1}, 1/\sqrt{\eta_2}, \dots, 1/\sqrt{\eta_k})$.
 - 3: Draw $\boldsymbol{\theta}$ uniformly at random from the unit sphere in \mathbb{R}^k .
 - 4: Let $\sum_{i=1}^k \delta_i \boldsymbol{\xi}_i \otimes \boldsymbol{\xi}_i$ be an eigendecomposition of $\mathbf{W}^\top \mathbf{T} \{\mathbf{W}\boldsymbol{\theta}\} \mathbf{W}$, where $\{\boldsymbol{\theta}_i : i \in [k]\}$ are orthonormal eigenvectors, and $\{\delta_i : i \in [k]\}$ are corresponding non-zero eigenvalues.
 - 5: Let $\mathbf{v}_i := (\mathbf{W}^\top)^\dagger \boldsymbol{\xi}_i$ and $\lambda_i := (\mathbf{W}\boldsymbol{\xi}_i)^\top \mathbf{T} \{\mathbf{W}\boldsymbol{\xi}_i\} \mathbf{W}\boldsymbol{\xi}_i$ for each $i \in [k]$.

Assume that \mathbf{O} has full column rank and that $c_{i,2} > 0$ for all $i \in [k]$. If $c_{i,3} \neq 0$ for all $i \in [k]$, then there is a randomized algorithm that, given \mathbf{P} and \mathbf{T} as input, returns $\{(s_i c_{i,2}^{1/2} \mathbf{o}_i, s_i c_{i,3}/c_{i,2}^{3/2}) : i \in [k]\}$ for some $\{s_1, s_2, \dots, s_k\} \subseteq \{\pm 1\}$. If $c_{i,4} \neq 0$ for all $i \in [k]$, then the same holds with \mathbf{Q} in place of \mathbf{T} and $c_{i,4}/c_{i,2}^2$ in place of $c_{i,3}/c_{i,2}^{3/2}$.

Proof Under the assumptions on \mathbf{O} and the $c_{i,2}$, the matrix \mathbf{P} is positive definite. By rescaling the columns of \mathbf{O} (replacing \mathbf{o}_i by $\sqrt{c_{i,2}} \mathbf{o}_i$), we may assume that $c_{i,2} = 1$ for all $i \in [k]$. Therefore there exists a matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$, which can be obtained from the singular value decomposition of \mathbf{P} , such that $\mathbf{W}^\top \mathbf{P} \mathbf{W} = (\mathbf{W}^\top \mathbf{O})(\mathbf{W}^\top \mathbf{O})^\top = \mathbf{I}_k$. This in turn implies that $\mathbf{M} := \mathbf{W}^\top \mathbf{O}$ is orthogonal. Then, for any vector $\boldsymbol{\theta} \in \mathbb{R}^k$,

$$\begin{aligned} \mathbf{W}^\top \mathbf{T} \{\mathbf{W}\boldsymbol{\theta}\} \mathbf{W} &= \mathbf{W}^\top \mathbf{O} \mathbf{diag}(\mathbf{diag}(c_{1,3}, c_{2,3}, \dots, c_{k,3}) \mathbf{O}^\top \mathbf{W}\boldsymbol{\theta}) \mathbf{O}^\top \mathbf{W} \\ &= \mathbf{M} \mathbf{diag}(\mathbf{diag}(c_{1,3}, c_{2,3}, \dots, c_{k,3}) \mathbf{M}^\top \boldsymbol{\theta}) \mathbf{M}^\top. \end{aligned}$$

Since \mathbf{M} is orthogonal, the above displayed equation gives an eigendecomposition of $\mathbf{W}^\top \mathbf{T} \{\mathbf{W}\boldsymbol{\theta}\} \mathbf{W}$, and the set of eigenvectors corresponding to non-repeated eigenvalues are uniquely defined up to sign. Each such unit-norm eigenvector $\boldsymbol{\xi}_i$ (after appropriate reordering) is of the form $s_i \mathbf{M} \mathbf{e}_i = s_i \mathbf{W}^\top \mathbf{o}_i$ for some $s_i \in \{\pm 1\}$. Therefore $\mathbf{v}_i := (\mathbf{W}^\top)^\dagger \boldsymbol{\xi}_i = s_i \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{o}_i = s_i \mathbf{o}_i$, since $\text{range}(\mathbf{W}) = \text{range}(\mathbf{O})$. Finally,

$$\begin{aligned} \lambda_i &:= (\mathbf{W}\boldsymbol{\xi}_i)^\top \mathbf{T} \{\mathbf{W}\boldsymbol{\xi}_i\} \mathbf{W}\boldsymbol{\xi}_i = \sum_{j=1}^k c_{j,3} (s_i \mathbf{o}_j^\top \mathbf{W} \mathbf{W}^\top \mathbf{o}_i)^3 \\ &= s_i \sum_{j=1}^k c_{j,3} (\mathbf{e}_j^\top \mathbf{O}^\top \mathbf{W} \mathbf{W}^\top \mathbf{O} \mathbf{e}_i)^3 = s_i c_{i,3}. \end{aligned}$$

Assume that $c_{i,3} \neq 0$ for all $i \in [k]$. If $\boldsymbol{\theta}$ is drawn uniformly at random from \mathbb{S}^{k-1} , then so is $\mathbf{M}^\top \boldsymbol{\theta}$. In this case, almost surely, the entries of $\mathbf{diag}(\mathbf{diag}(c_{1,3}, c_{2,3}, \dots, c_{k,3}) \mathbf{M}^\top \boldsymbol{\theta})$ are unique. Hence, no eigenvalue of $\mathbf{W}^\top \mathbf{T} \{\mathbf{W}\boldsymbol{\theta}\} \mathbf{W}$ is repeated, so the set $\{(\mathbf{v}_i, \lambda_i) : i \in [k]\}$ has the property claimed by the lemma and is returned by Algorithm 1.

An analogous argument proves the claim assuming $c_{i,4} \neq 0$ for all $i \in [k]$ using \mathbf{Q} in place of \mathbf{T} . □

The critical aspect in Theorem 2 is that only low-order observable moments are needed to identify the columns of \mathbf{O} . The observability implies that they can be estimated from data, and used in a plug-in estimator. The low-order nature of the moments mean that they can be estimated reliably (relative to higher-order moments).

From Theorem 2, we may also conclude that since $\boldsymbol{\mu}$, $\mathbf{Pairs}_{\alpha_0}$, and $\mathbf{Triples}_{\alpha_0}$ involve only \mathbf{x}_1 , \mathbf{x}_2 , and \mathbf{x}_3 , a random document under the LDA model need only have three words (conditionally independent given \mathbf{h}). This is the same conclusion as in the mixture of multinomial model for documents studied in [6], even though in that model, the words in documents were assumed to only be generated from a single topic.

It is worth noting that $\mathbf{Pairs}_{\alpha_0}$ and $\mathbf{Triples}_{\alpha_0}$ depend on the value $\alpha_0 = \sum_{i=1}^k \alpha_i$, so it must be known in order to form the appropriate moment estimates. In some applications, one may have prior knowledge of α_0 , as it characterizes the concentration of the Dirichlet distribution (and, indeed, having prior knowledge of α_0 is indeed much weaker than knowing the entire parameter vector $\boldsymbol{\alpha}$). The particular case where $\boldsymbol{\alpha} = (1, 1, \dots, 1)$ is essentially the problem of learning the vertices of a simplex when given access to (noisy) samples drawn from the uniform distribution over this simplex; therefore our result resolves this open problem posed by [21].

In the case of the independent latent factors model, the third- and fourth-order moments are exploited to take advantage of the non-Gaussianity of \mathbf{h} (which is not possible with only first- and second-order moments, without further assumptions).

3.4 Reducing the General Multi-view Setting to the Exchangeable Setting

To conclude this section, we show how to reduce the general multi-view setting (where the different $\mathbf{O}^{(v)}$ are not necessarily the same) to the case where the $\mathbf{O}^{(v)} = \mathbf{O}$ for all $v \in [\ell]$. Here, we must assume that the number of views is at least three (i.e., $\ell \geq 3$, just as in the independent latent factors and Dirichlet models, since $\mathbf{Triples}$, $\mathbf{Quadruples}$, and $\mathbf{Triples}_{\alpha_0}$ all depend on at least \mathbf{x}_1 , \mathbf{x}_2 , and \mathbf{x}_3). This reduction is based on a proof technique used in [2].

Define $\mathbf{Pairs}_{u,v} := \mathbb{E}[\mathbf{x}_u \otimes \mathbf{x}_v]$ for $u, v \in [\ell]$. Fix any $v_0 \in [\ell]$, which we call the *target view*. Then, for each $v \in [\ell] \setminus \{v_0\}$, pick any $u = u(v) \in [\ell] \setminus \{v, v_0\}$ (which is possible since $\ell \geq 3$) and define

$$\mathbf{C}_{v \rightarrow v_0} := \mathbf{Pairs}_{v_0,u} \mathbf{Pairs}_{v,u}^\dagger, \quad \mathbf{C}_{v_0 \rightarrow v} := \mathbf{Pairs}_{v,u} \mathbf{Pairs}_{v_0,u}^\dagger.$$

Proposition 1 (View symmetrization) *Assume Conditions 1 and 2 hold, and also that $\mathbb{E}[\mathbf{h} \otimes \mathbf{h}]$ is invertible. Fix a target view $v_0 \in [\ell]$, and let $\tilde{\mathbf{O}} := \mathbf{O}^{(v_0)}$. Let $\mathbf{C}_{v_0 \rightarrow v_0} := \mathbf{I}_d$, and for each view $v \in [\ell]$, let $\tilde{\mathbf{x}}_v := \mathbf{C}_{v \rightarrow v_0} \mathbf{x}_v$ where $\mathbf{C}_{v \rightarrow v_0}$ is defined above for $v \neq v_0$. For each $v \in [\ell]$,*

$$\mathbb{E}[\tilde{\mathbf{x}}_v | \mathbf{h}] = \tilde{\mathbf{O}} \mathbf{h}, \quad \mathbf{O}^{(v)} = \mathbf{C}_{v_0 \rightarrow v} \tilde{\mathbf{O}}.$$

Proof Assume without loss of generality that $v_0 := 1$, $v := 2$, and $u = u(v) := 3$ (in the definition of $\mathbf{C}_{2 \rightarrow 1}$). By Condition 1, $\mathbf{Pairs}_{1,3} = \mathbf{O}^{(1)} \mathbb{E}[\mathbf{h} \otimes \mathbf{h}] \mathbf{O}^{(3)\top}$ and $\mathbf{Pairs}_{2,3} = \mathbf{O}^{(2)} \mathbb{E}[\mathbf{h} \otimes \mathbf{h}] \mathbf{O}^{(3)\top}$. By Condition 2 and the assumption that

$\mathbb{E}[\mathbf{h} \otimes \mathbf{h}]$ is invertible, both $\mathbf{Pairs}_{1,3}$ and $\mathbf{Pairs}_{2,3}$ have rank k . Moreover, $\mathbf{Pairs}_{2,3}^\dagger = (\mathbf{O}^{(3)\top})^\dagger \mathbb{E}[\mathbf{h} \otimes \mathbf{h}]^{-1} \mathbf{O}^{(2)\dagger}$, so

$$\mathbf{C}_{2 \rightarrow 1} = \left(\mathbf{O}^{(1)} \mathbb{E}[\mathbf{h} \otimes \mathbf{h}] \mathbf{O}^{(3)\top} \right) \left((\mathbf{O}^{(3)\top})^\dagger \mathbb{E}[\mathbf{h} \otimes \mathbf{h}]^{-1} \mathbf{O}^{(2)\dagger} \right) = \mathbf{O}^{(1)} \mathbf{O}^{(2)\dagger}.$$

Since $\mathbf{O}^{(2)}$ has full column rank (by Condition 2),

$$\mathbb{E}[\tilde{\mathbf{x}}_2 | \mathbf{h}] = \mathbf{C}_{2 \rightarrow 1} \mathbb{E}[\mathbf{x}_2 | \mathbf{h}] = \mathbf{O}^{(1)} \mathbf{O}^{(2)\dagger} \mathbf{O}^{(2)} \mathbf{h} = \mathbf{O}^{(1)} \mathbf{h} = \tilde{\mathbf{O}} \mathbf{h}.$$

By the same argument, $\mathbf{C}_{1 \rightarrow 2} = \mathbf{O}^{(2)} \mathbf{O}^{(1)\dagger} = \mathbf{O}^{(2)} \tilde{\mathbf{O}}^\dagger$; therefore $\mathbf{C}_{1 \rightarrow 2} \tilde{\mathbf{O}} = \mathbf{O}^{(2)} \tilde{\mathbf{O}}^\dagger \tilde{\mathbf{O}} = \mathbf{O}^{(2)}$ since $\tilde{\mathbf{O}}$ has full column rank. \square

4 Estimation from Moments via Orthogonal Tensor Decomposition

In addition to establishing identifiability, Theorem 2 provides an efficient randomized algorithm for recovering the columns of \mathbf{O} up to permutation and scaling, and it indeed can be shown to work using only estimates of the required moments using techniques similar to [6, 39]. However, the resulting sample complexity is rather high (i.e., a high-degree polynomial) on account of the randomization technique used to ensure the uniqueness of an eigendecomposition. Indeed, the randomization collapses moment tensors into matrices, but the resulting spacing between eigenvalues is small; hence, such procedures may not be very robust to errors in the moment estimates. It turns out this is only an artifact of the algorithmic technique, as a different technique based on orthogonal tensor decomposition has both a better analysis and better empirical support. In this section, we recall the tensor decomposition technique from [5] and show how it can be applied to the estimation problem for the multi-view models considered in this work.

4.1 Multi-way Arrays, Multi-linear Forms, and Tensor Algebras

Recall that low-order moments from Sect. 3 have the form $\sum_{i=1}^k c_{i,p} \mathbf{o}_i^{\otimes p}$, where $\mathbf{v}^{\otimes p} := \mathbf{v} \otimes \mathbf{v} \otimes \dots \otimes \mathbf{v}$ denotes the p -fold tensor product of a vector \mathbf{v} with itself. We assume that

$$\mathbf{P} := \sum_{i=1}^k c_{i,2} \mathbf{o}_i^{\otimes 2} \quad \text{and} \quad \mathbf{T} := \sum_{i=1}^k c_{i,p} \mathbf{o}_i^{\otimes p}$$

for some $p \geq 3$ are given, with $c_{i,2} > 0$ and $c_{i,p} > 0$ for all $i \in [k]$, and $\{\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_k\} \subset \mathbb{R}^d$ linearly independent. We note that it is also possible to handle the case where $c_{i,2} < 0$ and $c_{i,p} < 0$ for some $i \in [k]$ (using techniques from [46]), though we avoid here this case for simplicity. In our examples, we will specialize to the $p = 3$ case.

From the second-order moment matrix \mathbf{P} , as discussed in the proof of Lemma 4, we may extract a so-called *whitening matrix* $\mathbf{W} \in \mathbb{R}^{d \times k}$ using the singular value decomposition of \mathbf{P} such that $\mathbf{W}^\top \mathbf{P} \mathbf{W} = \mathbf{I}_k$. Put another way, \mathbf{P} defines a k -dimensional inner product subspace of \mathbb{R}^d in which $\mathcal{B} := \{\sqrt{c_{1,2}}\mathbf{o}_1, \sqrt{c_{2,2}}\mathbf{o}_2, \dots, \sqrt{c_{k,2}}\mathbf{o}_k\}$ are orthonormal. To see this, we define the inner product by

$$\langle \mathbf{u}, \mathbf{v} \rangle := \mathbf{u}^\top \mathbf{P}^\dagger \mathbf{v}$$

Since the \mathbf{o}_i are linearly independent and $c_{i,2} > 0$,

$$\mathbf{P}^\dagger = (\mathbf{O}^\top)^\dagger \mathbf{diag}(c_{1,2}, c_{2,2}, \dots, c_{k,2})^{-1} \mathbf{O}^\dagger,$$

where $\mathbf{O}^\dagger \mathbf{O} = \mathbf{I}_k$. This implies that

$$\langle \sqrt{c_{i,2}}\mathbf{o}_i, \sqrt{c_{j,2}}\mathbf{o}_j \rangle = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

as claimed.

The p th-order tensor \mathbf{T} can be viewed in a number of different ways. The first is as a p -way array of real numbers T_{i_1, i_2, \dots, i_p} for $i_1, i_2, \dots, i_p \in [d]$. The second is as a p -linear form $\mathbf{T} : \mathbb{R}^d \times \mathbb{R}^d \times \dots \times \mathbb{R}^d \rightarrow \mathbb{R}$ (generalizing the bilinear form view of a matrix): for any $\mathbf{u} = (u_1, u_2, \dots, u_d) \in \mathbb{R}^d$, $\mathbf{v} = (v_1, v_2, \dots, v_d) \in \mathbb{R}^d$, and $\mathbf{w} = (w_1, w_2, \dots, w_d) \in \mathbb{R}^d$,

$$\mathbf{T}(\mathbf{u}, \mathbf{v}, \mathbf{w}) = \sum_{i,j,l} T_{i,j,l} u_i v_j w_l,$$

which connects the p -way array with the p -linear form. We will generally denote, for $\mathbf{U} = [\mathbf{u}_1 | \mathbf{u}_2 | \dots | \mathbf{u}_d] \in \mathbb{R}^{m_1 \times d}$, $\mathbf{V} = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_d] \in \mathbb{R}^{m_2 \times d}$, and $\mathbf{W} = [\mathbf{w}_1 | \mathbf{w}_2 | \dots | \mathbf{w}_d] \in \mathbb{R}^{m_3 \times d}$, the tensor $\mathbf{T}(\mathbf{U}, \mathbf{V}, \mathbf{W}) \in \mathbb{R}^{m_1 \times m_2 \times m_3}$ given by

$$\mathbf{T}(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \sum_{i,j,l} T_{i,j,l} \mathbf{u}_i \otimes \mathbf{v}_j \otimes \mathbf{w}_l.$$

(Note that the notation “ $\mathbf{T}\{\mathbf{v}\}$ ” from Lemma 4 can be written as $\mathbf{T}(\mathbf{I}_d, \mathbf{I}_d, \mathbf{v})$.)

The third view of \mathbf{T} is as a member of $(\mathbb{R}^d)^{\otimes p}$, the p th-order tensor product of \mathbb{R}^d . In this view, we may pick a basis for this vector space, and represent \mathbf{T} in that basis. A natural choice is derived from the standard coordinate basis for \mathbb{R}^d , giving $\{\mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_l : i, j, l \in [d]\}$. In this basis, we identify the p -way array as a tensor algebra element:

$$\mathbf{T} = \sum_{i,j,l} T_{i,j,l} \mathbf{e}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_l.$$

Alternatively, if $\mathcal{B}' := \{\sqrt{c_{1,2}}\mathbf{o}_1, \sqrt{c_{2,2}}\mathbf{o}_2, \dots, \sqrt{c_{d,2}}\mathbf{o}_d\}$ is a completion of \mathcal{B} from above, we may derive a different basis $\{\mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} : \mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathcal{B}'\}$ for $(\mathbb{R}^d)^{\otimes p}$, under

which T has the following *diagonal representation*:

$$T = \sum_{i=1}^k \frac{c_{i,3}}{c_{i,2}^{3/2}} (\sqrt{c_{i,2}} \mathbf{o}_i) \otimes (\sqrt{c_{i,2}} \mathbf{o}_i) \otimes (\sqrt{c_{i,2}} \mathbf{o}_i).$$

Note that T as a p -linear form can be expressed in terms of the standard inner product $\mathbf{u}^\top \mathbf{v} = \sum_{i=1}^d u_i v_i$ giving

$$T(\mathbf{u}, \mathbf{v}, \mathbf{w}) = \sum_{i,j,l} T_{i,j,l} (\mathbf{e}_i^\top \mathbf{u}) (\mathbf{e}_j^\top \mathbf{v}) (\mathbf{e}_l^\top \mathbf{w}).$$

The proposed orthogonal tensor decomposition algorithm from [5] exploits the p -linear form based on the inner product $\langle \cdot, \cdot \rangle$, as it exploits the orthogonality of the $\sqrt{c_{i,2}} \mathbf{o}_i$:

$$\begin{aligned} T(\mathbf{P}^\dagger \mathbf{u}, \mathbf{P}^\dagger \mathbf{v}, \mathbf{P}^\dagger \mathbf{w}) &= \sum_{i,j,l} T_{i,j,l} \langle \mathbf{e}_i, \mathbf{u} \rangle \langle \mathbf{e}_j, \mathbf{v} \rangle \langle \mathbf{e}_l, \mathbf{w} \rangle \\ &= \sum_{i=1}^k \frac{c_{i,3}}{c_{i,2}^{3/2}} \langle \sqrt{c_{i,2}} \mathbf{o}_i, \mathbf{u} \rangle \langle \sqrt{c_{i,2}} \mathbf{o}_i, \mathbf{v} \rangle \langle \sqrt{c_{i,2}} \mathbf{o}_i, \mathbf{w} \rangle \end{aligned}$$

Indeed, it is shown that the map

$$\mathbf{v} \mapsto \frac{T(\mathbf{I}_d, \mathbf{P}^\dagger \mathbf{v}, \mathbf{P}^\dagger \mathbf{v})}{\sqrt{\langle T(\mathbf{I}_d, \mathbf{P}^\dagger \mathbf{v}, \mathbf{P}^\dagger \mathbf{v}), T(\mathbf{I}_d, \mathbf{P}^\dagger \mathbf{v}, \mathbf{P}^\dagger \mathbf{v}) \rangle}} \tag{4}$$

has non-zero stable fixed points only at $\sqrt{c_{i,2}} \mathbf{o}_i$ for $i \in [k]$. Moreover, repeated application of the above map starting from a random $\mathbf{v} \in \text{range}(\mathbf{P}) \cap \{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\|_2 = 1\}$ converges at a quadratic rate to one of the $\sqrt{c_{i,2}} \mathbf{o}_i$ (Lemma 5.1 from [5]). Finally, it is easy to check that $T(\mathbf{P}^\dagger \sqrt{c_{i,2}} \mathbf{o}_i, \mathbf{P}^\dagger \sqrt{c_{i,2}} \mathbf{o}_i, \mathbf{P}^\dagger \sqrt{c_{i,2}} \mathbf{o}_i) = c_{i,3}/c_{i,2}^{3/2}$.

4.2 Efficient Estimation Algorithm

In an actual implementation, we must estimate the moment matrices and tensors \mathbf{P} and T from data. Let $S := \{(\mathbf{x}_1^{(j)}, \mathbf{x}_2^{(j)}, \dots, \mathbf{x}_\ell^{(j)}) : j \in [n]\}$ be an i.i.d. sample comprised of n independent copies of $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell)$. We may form estimates, $\widehat{\mathbf{P}}$ and \widehat{T} , of \mathbf{P} and T , respectively, using empirical averages with respect to S . For instance, in LDA, we may we let

$$\begin{aligned} \widehat{\boldsymbol{\mu}} &:= \frac{1}{n} \sum_{j=1}^n \mathbf{x}_1^{(j)}, \\ \widehat{\text{Pairs}}_{\alpha_0} &:= \frac{1}{n} \sum_{j=1}^n \mathbf{x}_1^{(j)} \otimes \mathbf{x}_2^{(j)} - \frac{\alpha_0}{\alpha_0 + 1} \widehat{\boldsymbol{\mu}} \otimes \widehat{\boldsymbol{\mu}}, \end{aligned}$$

$$\begin{aligned} \widehat{\text{Triples}}_{\alpha_0} &:= \frac{1}{n} \sum_{j=1}^n \left(\mathbf{x}_1^{(j)} \otimes \mathbf{x}_2^{(j)} \otimes \mathbf{x}_3^{(j)} \right. \\ &\quad \left. - \frac{\alpha_0}{\alpha_0 + 2} \left(\mathbf{x}_1^{(j)} \otimes \mathbf{x}_2^{(j)} \otimes \widehat{\boldsymbol{\mu}} + \mathbf{x}_1^{(j)} \otimes \widehat{\boldsymbol{\mu}} \otimes \mathbf{x}_3^{(j)} + \widehat{\boldsymbol{\mu}} \otimes \mathbf{x}_2^{(j)} \otimes \mathbf{x}_3^{(j)} \right) \right) \\ &\quad + \frac{2\alpha_0^2}{(\alpha_0 + 2)(\alpha_0 + 1)} \widehat{\boldsymbol{\mu}} \otimes \widehat{\boldsymbol{\mu}} \otimes \widehat{\boldsymbol{\mu}}, \end{aligned}$$

and set $\widehat{\mathbf{P}} := \widehat{\text{Pairs}}_{\alpha_0}$ and $\widehat{\mathbf{T}} := \widehat{\text{Triples}}_{\alpha_0}$. Note that we may also ensure that $\widehat{\mathbf{P}}$ and $\widehat{\mathbf{T}}$ are symmetric (so $[\mathbf{T}]_{i,j,l} = [\mathbf{T}]_{i,l,j} = [\mathbf{T}]_{j,i,l}$ etc.). This gives consistent estimates of \mathbf{P} and \mathbf{T} .

Using $\widehat{\mathbf{P}}$ and $\widehat{\mathbf{T}}$, we propose a plug-in approach to estimation for \mathbf{O} , given as Algorithm 2 (Lines 5–13 make up the robust tensor power method of [5]). This is essentially a robust version of the iteration given by (4) which finds an approximate orthogonal tensor decomposition of \mathbf{T} given the nearby estimate $\widehat{\mathbf{T}}$. The algorithm from [5] is applied to the tensor $\widehat{\mathbf{T}}(\widehat{\mathbf{W}}, \widehat{\mathbf{W}}, \widehat{\mathbf{W}})$, and the outputs $\widehat{\mathbf{v}}_i$ and $\widehat{\lambda}_i$ should be interpreted as estimates of the $\sqrt{c_{j,2}}\mathbf{o}_j$ and $c_{j,3}/c_{j,2}^{3/2}$ for some $j \in [k]$.

4.3 Analysis of Algorithm 2

In this section, we give a simple error analysis for Algorithm 2. Let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k > 0$ be the non-zero singular values of \mathbf{P} . Also, let $\lambda_i := c_{i,3}/c_{i,2}^{3/2}$, ordered so that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0$. We define the operator norm of a symmetric third-order tensor $\mathbf{E} \in \mathbb{R}^{m \times m \times m}$ by

$$\|\mathbf{E}\|_2 := \max_{\|\mathbf{v}\|_2=1} |\mathbf{E}(\mathbf{v}, \mathbf{v}, \mathbf{v})|.$$

Define $\epsilon_{\mathbf{P}} := \|\widehat{\mathbf{P}} - \mathbf{P}\|_2/\sigma_k$ and $\epsilon_{\mathbf{T}} := \|\widehat{\mathbf{T}} - \mathbf{T}\|_2$. The error bound is given in terms of $\epsilon_{\mathbf{P}}$ and $\epsilon_{\mathbf{T}}$.

Theorem 3 (Error analysis of Algorithm 2) *There exist universal constants $C, C', c, c' > 0$ such that the following holds. Pick any $\delta \in (0, 1)$. If*

$$\begin{aligned} \epsilon_{\mathbf{P}} &\leq c \cdot \frac{\lambda_k/\lambda_1}{k}, \quad \epsilon_{\mathbf{T}} \leq c' \cdot \frac{\lambda_k \sigma_k^{3/2}}{k}, \\ N &\geq C \left(\log(k) + \log \left(\log \left(\frac{\lambda_1 \sigma_k^{3/2}}{\epsilon_{\mathbf{T}}} + \frac{1}{\epsilon_{\mathbf{P}}} \right) \right) \right), \\ L &\geq \text{poly}(k) \log(1/\delta) \end{aligned}$$

(for some fixed polynomial $\text{poly}(k)$ specified in Theorem 5.1 of [5]), then with probability at least $1 - \delta$, Algorithm 2 returns $\{(\widehat{\mathbf{v}}_i, \widehat{\lambda}_i) : i \in [k]\}$ satisfying, after appropriate reordering,

Algorithm 2 Plug-in estimator based on orthogonal tensor decomposition

input positive integer $k \in \mathbb{N}$, symmetric matrix $\widehat{\mathbf{P}} \in \mathbb{R}^{d \times d}$, symmetric tensor $\widehat{\mathbf{T}} \in \mathbb{R}^{d \times d \times d}$, number of iterations L and N .

output $\{(\widehat{\mathbf{v}}_i, \widehat{\lambda}_i) : i \in [k]\}$.

- 1: Compute eigendecomposition of $\widehat{\mathbf{P}}$, let $\eta_1 \geq \eta_2 \geq \dots \geq \eta_k$ be the top k eigenvalues, and let $\xi_1, \xi_2, \dots, \xi_k \in \mathbb{R}^d$ be the corresponding orthonormal eigenvectors. (Halt and signal failure if $\eta_k \leq 0$.)
- 2: Set $\widehat{\mathbf{W}} := [|\xi_1\rangle|\xi_2\rangle \dots |\xi_k\rangle] \mathbf{diag}(1/\sqrt{\eta_1}, 1/\sqrt{\eta_2}, \dots, 1/\sqrt{\eta_k})$. ($\widehat{\mathbf{W}}\widehat{\mathbf{W}}^\top$ is the pseudoinverse of a rank- k approximation to $\widehat{\mathbf{P}}$.)
- 3: Initialize $\widetilde{\mathbf{T}} := \widehat{\mathbf{T}}$.
- 4: **for** $i = 1$ to k **do**
- 5: **for** $\tau = 1$ to L **do**
- 6: Draw $\theta_0^{(\tau)}$ uniformly at random from the unit sphere in \mathbb{R}^k .
- 7: **for** $t = 1$ to N **do**
- 8: Compute power iteration update:

$$\mathbf{u} := \widetilde{\mathbf{T}}(\widehat{\mathbf{W}}, \widehat{\mathbf{W}}\theta_{t-1}^{(\tau)}, \widehat{\mathbf{W}}\theta_{t-1}^{(\tau)}), \quad \theta_t^{(\tau)} := \frac{\mathbf{u}}{\|\mathbf{u}\|_2}. \tag{*}$$

- 9: **end for**
- 10: **end for**
- 11: Let

$$\tau^* := \arg \max_{\tau \in [L]} \widetilde{\mathbf{T}}(\widehat{\mathbf{W}}\theta_N^{(\tau)}, \widehat{\mathbf{W}}\theta_N^{(\tau)}, \widehat{\mathbf{W}}\theta_N^{(\tau)})$$

- and execute N more power iteration updates $(*)$ starting from $\theta_N^{(\tau^*)}$ to obtain $\widehat{\theta}_i$.
- 12: Let

$$\widehat{\mathbf{v}}_i := (\widehat{\mathbf{W}}^\top)^\dagger \widehat{\theta}_i, \quad \widehat{\lambda}_i := \widetilde{\mathbf{T}}(\widehat{\mathbf{W}}\widehat{\theta}_i, \widehat{\mathbf{W}}\widehat{\theta}_i, \widehat{\mathbf{W}}\widehat{\theta}_i).$$

- 13: Deflate $\widetilde{\mathbf{T}}$: set $\widetilde{\mathbf{T}} := \widetilde{\mathbf{T}} - \widehat{\lambda}_i \widehat{\mathbf{v}}_i \otimes \widehat{\mathbf{v}}_i \otimes \widehat{\mathbf{v}}_i$.
- 14: **end for**

$$\begin{aligned} \|\widehat{\mathbf{v}}_i - \mathbf{o}_i\|_2 &\leq C' \cdot \left(\frac{1}{\lambda_i} \cdot \frac{1}{\sigma_k^2} \cdot \epsilon_T + \left(\frac{\lambda_1}{\lambda_i} \cdot \frac{1}{\sqrt{\sigma_k}} + 1 \right) \cdot \epsilon_P \right), \\ |\widehat{\lambda}_i - \lambda_i| &\leq C' \cdot \left(\frac{1}{\sigma_k^{3/2}} \cdot \epsilon_T + \lambda_1 \epsilon_P \right) \end{aligned}$$

for all $i \in [k]$.

Proof We assume without loss of generality that $c_{i,2} = 1$ for all $i \in [k]$, so $\mathbf{P} = \mathbf{O}\mathbf{O}^\top$, and $\lambda_i = c_{i,3}$ for all $i \in [k]$. By Lemma 8 from [28], if $\widehat{\Pi}$ is the orthogonal projection to $\text{span}\{|\xi_1\rangle, |\xi_2\rangle, \dots, |\xi_k\rangle\}$ (w.r.t. the standard inner product), and Π is the orthogonal projection to $\text{range}(\mathbf{O})$ (again w.r.t. the standard inner product), then

$$\|(\mathbf{I}_d - \widehat{\Pi})\Pi\|_2 \leq 1.5\epsilon_P. \tag{5}$$

Note that $(\widehat{\mathbf{W}}^\top)^\dagger = [|\xi_1\rangle|\xi_2\rangle \dots |\xi_k\rangle] \mathbf{diag}(\sqrt{\eta_1}, \sqrt{\eta_2}, \dots, \sqrt{\eta_k})$ and therefore $\widehat{\Pi} = (\widehat{\mathbf{W}}^\top)^\dagger \widehat{\mathbf{W}}^\top$.

Let $\mathbf{W} := \widehat{\mathbf{W}}(\widehat{\mathbf{W}}^\top \mathbf{P} \widehat{\mathbf{W}})^\dagger^{1/2}$. By Lemma 9 and Lemma 11 from [28],

$$\|\widehat{\mathbf{W}}\|_2 \leq \frac{1}{\sqrt{(1 - \epsilon_P)\sigma_k}}, \tag{6}$$

$$\|(\widehat{\mathbf{W}} - \mathbf{W})^\top \mathbf{O}\|_2 \leq \sqrt{1 + 1.5\epsilon_P} 1.5\epsilon_P, \tag{7}$$

$$\|\widehat{\mathbf{T}}(\widehat{\mathbf{W}}, \widehat{\mathbf{W}}, \widehat{\mathbf{W}}) - \mathbf{T}(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_2 \leq \frac{2}{\sigma_k^{3/2}} \cdot \epsilon_T + 6\lambda_1 \cdot \epsilon_P =: \varepsilon. \tag{8}$$

Note that

$$\mathbf{T}(\mathbf{W}, \mathbf{W}, \mathbf{W}) = \sum_{i=1}^k \lambda_i (\mathbf{W}^\top \mathbf{o}_i) \otimes (\mathbf{W}^\top \mathbf{o}_i) \otimes (\mathbf{W}^\top \mathbf{o}_i)$$

where $\{\mathbf{W}^\top \mathbf{o}_i : i \in [k]\}$ are orthonormal. Using (8), Theorem 5.1 from [5] implies the following: if $\varepsilon \leq C_1 \lambda_k / k$, $N \geq C_2 (\log(k) + \log \log(\lambda_1 / \varepsilon))$, and $L \geq \text{poly}(k) \log(1/\delta)$, then with probability at least $1 - \delta$, the robust tensor power method returns $\hat{\boldsymbol{\theta}}_i$ and $\hat{\lambda}_i$ satisfying (after appropriate reordering)

$$\|\hat{\boldsymbol{\theta}}_i - \mathbf{W}^\top \mathbf{o}_i\|_2 \leq \frac{8\varepsilon}{\lambda_i}, \quad |\hat{\lambda}_i - \lambda_i| \leq 5\varepsilon. \tag{9}$$

Therefore,

$$\begin{aligned} \|\hat{\mathbf{v}}_i - \mathbf{o}_i\|_2 &= \left\| ((\widehat{\mathbf{W}}^\top)^\dagger \hat{\boldsymbol{\theta}}_i - (\widehat{\mathbf{W}}^\top)^\dagger \mathbf{W}^\top \mathbf{o}_i) \right. \\ &\quad \left. + ((\widehat{\mathbf{W}}^\top)^\dagger \mathbf{W}^\top \mathbf{o}_i - (\widehat{\mathbf{W}}^\top)^\dagger \widehat{\mathbf{W}}^\top \mathbf{o}_i) + ((\widehat{\mathbf{W}}^\top)^\dagger \widehat{\mathbf{W}}^\top \mathbf{o}_i - \mathbf{o}_i) \right\|_2 \\ &\leq \|(\widehat{\mathbf{W}}^\top)^\dagger \hat{\boldsymbol{\theta}}_i - (\widehat{\mathbf{W}}^\top)^\dagger \mathbf{W}^\top \mathbf{o}_i\|_2 \\ &\quad + \|(\widehat{\mathbf{W}}^\top)^\dagger \mathbf{W}^\top \mathbf{o}_i - (\widehat{\mathbf{W}}^\top)^\dagger \widehat{\mathbf{W}}^\top \mathbf{o}_i\|_2 + \|(\widehat{\mathbf{W}}^\top)^\dagger \widehat{\mathbf{W}}^\top \mathbf{o}_i - \mathbf{o}_i\|_2 \\ &\leq \|(\widehat{\mathbf{W}}^\top)^\dagger\|_2 \left(\|\hat{\boldsymbol{\theta}}_i - \mathbf{W}^\top \mathbf{o}_i\|_2 + \|(\widehat{\mathbf{W}} - \mathbf{W})^\top \mathbf{O}\|_2 \right) \\ &\quad + \|(\widehat{\boldsymbol{\Pi}} - \mathbf{I}_d) \boldsymbol{\Pi} \mathbf{o}_i\|_2 \\ &\leq \frac{1}{\sqrt{(1 - \epsilon_P)\sigma_k}} \left(\frac{8\varepsilon}{\lambda_i} + \sqrt{1 + 1.5\epsilon_P} 1.5\epsilon_P \right) + 1.5\epsilon_P \end{aligned}$$

using (5), (6), (9), and (7). □

We remark that Theorem 3 immediately implies estimation consistency under appropriate assumptions on the σ_i and λ_i , and it is straightforward to obtain finite sample guarantees using concentration arguments to bound $\|\widehat{\mathbf{P}} - \mathbf{P}\|_2$ and $\|\widehat{\mathbf{T}} - \mathbf{T}\|_2$. We leave this as an exercise for the reader.

Acknowledgments We thank Kamalika Chaudhuri, Adam Kalai, Percy Liang, Chris Meek, David Sontag, and Tong Zhang for valuable insights. We also thank Rong Ge for sharing preliminary results (in [8]) and the anonymous reviewers for their comments, suggestions, and pointers to references. Part of this work was completed while DH was a postdoctoral researcher at Microsoft Research New England, and while DPF, YKL, and AA were visiting the same lab. AA is supported in part by Microsoft Faculty Fellowship, NSF Career award CCF-1254106, NSF Award CCF-1219234, NSF BIGDATA IIS-1251267 and ARO YIP Award W911NF-13-1-0084.

References

1. Achlioptas, D., McSherry, F.: On spectral learning of mixtures of distributions. Eighteenth Annual Conference on Learning Theory, pp. 458–469. Springer, Bertinoro (2005)
2. Anandkumar, A., Chaudhuri, K., Hsu, D., Kakade, S.M., Song, L., Zhang, T.: Spectral methods for learning multivariate latent tree structure. *Adv. Neural Inf. Process. Syst.* **24**, 2025–2033 (2011)
3. Anandkumar, A., Foster, D.P., Hsu, D., Kakade, S.M., Liu, Y.K.: A spectral algorithm for latent Dirichlet allocation. *Adv. Neural Inf. Process. Syst.* **25**, 917–925 (2012)
4. Anandkumar, A., Foster, D.P., Hsu, D., Kakade, S.M., Liu, Y.K.: Two SVDs suffice: spectral decompositions for probabilistic topic models and latent Dirichlet allocation (2012). [arXiv:1204.6703v1](https://arxiv.org/abs/1204.6703v1)
5. Anandkumar, A., Ge, R., Hsu, D., Kakade, S.M., Telgarsky, M.: Tensor decompositions for learning latent variable models. *J. Mach. Learn. Res.* (2014). To appear.
6. Anandkumar, A., Hsu, D., Kakade, S.M.: A method of moments for mixture models and hidden Markov models. In: Twenty-Fifth Annual Conference on Learning Theory, vol. 23, pp. 33.1–33.34 (2012)
7. Ando, R., Zhang, T.: Two-view feature generation model for semi-supervised learning. In: Twenty-Fourth International Conference on Machine Learning, pp. 25–32 (2007)
8. Arora, S., Ge, R., Moitra, A.: Learning topic models – going beyond SVD. In: Fifty-Third IEEE Annual Symposium on Foundations of Computer Science, pp. 1–10 (2012)
9. Arora, S., Ge, R., Moitra, A., Sachdeva, S.: Provable ICA with unknown Gaussian noise, with implications for Gaussian mixtures and autoencoders. *Adv. Neural Inf. Process. Syst.* **25**, 2375–2383 (2012)
10. Arora, S., Kannan, R.: Learning mixtures of separated nonspherical Gaussians. *Ann. Appl. Probab.* **15**(1A), 69–92 (2005)
11. Belkin, M., Sinha, K.: Polynomial learning of distribution families. In: Fifty-First Annual IEEE Symposium on Foundations of Computer Science, pp. 103–112 (2010)
12. Blei, D.M., Ng, A., Jordan, M.: Latent Dirichlet allocation. *J. Mach. Learn. Res.* **3**, 993–1022 (2003)
13. Canny, J.: GaP: A factor model for discrete data. In: Proceedings of the Twenty-Seventh Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 122–129 (2004)
14. Cardoso, J.F., Comon, P.: Independent component analysis, a survey of some algebraic methods. In: IEEE International Symposium on Circuits and Systems, pp. 93–96 (1996)
15. Chang, J.T.: Full reconstruction of Markov models on evolutionary trees: identifiability and consistency. *Math. Biosci.* **137**, 51–73 (1996)
16. Chaudhuri, K., Kakade, S.M., Livescu, K., Sridharan, K.: Multi-view clustering via canonical correlation analysis. In: Twenty-Sixth Annual International Conference on Machine Learning, pp. 129–136 (2009)
17. Chaudhuri, K., Rao, S.: Learning mixtures of product distributions using correlations and independence. In: Twenty-First Annual Conference on Learning Theory, pp. 9–20 (2008)
18. Comon, P., Jutten, C.: *Handbook of Blind Source Separation: Independent Component Analysis and Applications*. Academic Press, Waltham (2010)
19. Dasgupta, S.: Learning mixtures of Gaussians. In: Fortieth Annual IEEE Symposium on Foundations of Computer Science, pp. 634–644 (1999)
20. Dasgupta, S., Schulman, L.: A probabilistic analysis of EM for mixtures of separated, spherical Gaussians. *J. Mach. Learn. Res.* **8**, 203–226 (2007)
21. Frieze, A.M., Jerrum, M., Kannan, R.: Learning linear transformations. In: Thirty-Seventh Annual Symposium on Foundations of Computer Science, pp. 359–368 (1996)
22. Griffiths, T.: Gibbs sampling in the generative model of latent Dirichlet allocation. Tech. rep., Stanford University (2002)

23. Harshman, R.: Foundations of the PARAFAC procedure: model and conditions for an ‘explanatory’ multi-mode factor analysis. Tech. rep., UCLA Working Papers in Phonetics (1970)
24. Hitchcock, F.: The expression of a tensor or a polyadic as a sum of products. *J. Math. Phys.* **6**, 164–189 (1927)
25. Hitchcock, F.: Multiple invariants and generalized rank of a p-way matrix or tensor. *J. Math. Phys.* **7**, 39–79 (1927)
26. Hofmann, T.: Probabilistic latent semantic indexing. In: Proceedings of the Twenty-Second Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 50–57 (1999)
27. Hotelling, H.: The most predictable criterion. *J. Educ. Psychol.* **26**(2), 139–142 (1935)
28. Hsu, D., Kakade, S.M.: Learning mixtures of spherical Gaussians: moment methods and spectral decompositions. In: Fourth Innovations in Theoretical Computer Science (2013)
29. Hsu, D., Kakade, S.M., Zhang, T.: A spectral algorithm for learning hidden Markov models. *J. Comput. Syst. Sci.* **78**(5), 1460–1480 (2012). <http://www.sciencedirect.com/science/article/pii/S002000012000244>
30. Jutten, C., Herault, J.: Blind separation of sources, part I: an adaptive algorithm based on neuromimetic architecture. *Signal Process.* **24**, 1–10 (1991)
31. Kakade, S.M., Foster, D.P.: Multi-view regression via canonical correlation analysis. In: Twentieth Annual Conference on Learning Theory, pp. 82–96 (2007)
32. Kalai, A.T., Moitra, A., Valiant, G.: Efficiently learning mixtures of two Gaussians. In: Forty-second ACM Symposium on Theory of Computing, pp. 553–562 (2010)
33. Kannan, R., Salmasian, H., Vempala, S.: The spectral method for general mixture models. *SIAM J. Comput.* **38**(3), 1141–1156 (2008)
34. Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM Rev.* **51**(3), 455–500 (2009)
35. Kruskal, J.B.: Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear Algebra Its Appl.* **18**(2), 95–138 (1977)
36. Lee, D.D., Seung, H.S.: Learning the parts of objects by nonnegative matrix factorization. *Nature* **401**, 788–791 (1999)
37. Leurgans, S., Ross, R., Abel, R.: A decomposition for three-way arrays. *SIAM J. Matrix Anal. Appl.* **14**(4), 1064–1083 (1993)
38. Moitra, A., Valiant, G.: Settling the polynomial learnability of mixtures of Gaussians. In: Fifty-First Annual IEEE Symposium on Foundations of Computer Science, pp. 93–102 (2010)
39. Mossel, E., Roch, S.: Learning nonsingular phylogenies and hidden Markov models. *Ann. Appl. Probab.* **16**(2), 583–614 (2006)
40. Nguyen, P.Q., Regev, O.: Learning a parallelepiped: cryptanalysis of GGH and NTRU signatures. *J. Cryptol.* **22**(2), 139–160 (2009)
41. Papadimitriou, C.H., Raghavan, P., Tamaki, H., Vempala, S.: Latent semantic indexing: a probabilistic analysis. *J. Comput. Syst. Sci.* **61**(2), 217–235 (2000)
42. Pearson, K.: Contributions to the mathematical theory of evolution. *Philos. Trans. R. Soc. Lond. A* **185**, 71–110 (1894)
43. Redner, R.A., Walker, H.F.: Mixture densities, maximum likelihood and the EM algorithm. *SIAM Rev.* **26**(2), 195–239 (1984)
44. Tucker, L.R.: Some mathematical notes on three-mode factor analysis. *Psychometrika* **31**(3), 279–311 (1966)
45. Vempala, S., Wang, G.: A spectral algorithm for learning mixtures models. *J. Comput. Syst. Sci.* **68**(4), 841–860 (2004)
46. Zou, J., Hsu, D., Parkes, D., Adams, R.: Contrastive learning using spectral methods. *Adv. Neural Inf. Process. Syst.* **26**, 2238–2246 (2013)