

Recovering Joint PMF from Pairwise Marginals

Shahana Ibrahim

*School of Electrical Engineering and Computer Science
Oregon State University
Corvallis, OR, USA
ibrahish@oregonstate.edu*

Xiao Fu

*School of Electrical Engineering and Computer Science
Oregon State University
Corvallis, OR, USA
xiao.fu@oregonstate.edu*

Abstract—To overcome the curse of dimensionality in joint probability learning, recent work has proposed to recover the joint probability mass function (PMF) of an arbitrary number of random variables (RVs) from three-dimensional marginals, exploiting the uniqueness of tensor decomposition and the (unknown) dependence among the RVs. Nonetheless, accurately estimating three-dimensional marginals is still costly in terms of sample complexity. Tensor decomposition also poses a computationally intensive optimization problem. This work puts forth a new framework that learns the joint PMF using pairwise marginals that are relatively easy to acquire. The method is built upon nonnegative matrix factorization (NMF) theory, and features a Gram–Schmidt-like economical algorithm that works provably well under realistic conditions. Theoretical analysis of a recently proposed expectation maximization (EM) algorithm for joint PMF recovery is also presented. In particular, the EM algorithm is shown to provably improve upon the proposed pairwise marginal-based approach. Synthetic and real-data experiments are employed to showcase the effectiveness of the proposed approach.

Index Terms—joint probability learning, nonnegative matrix factorization, probability tensors, two-dimensional marginals

I. INTRODUCTION

Many learning and inference tasks in high-dimensional statistics boil down to estimating/approximating the joint probability of a set of random variables (RVs). However, in the high-dimensional regime, directly estimating the joint probability via “structure-free” methods such as sample averaging is considered not viable—due to the need of a huge amount of data. Many workarounds, e.g., linear estimators, kernels, and neural networks, have been proposed for combating this curse of dimensionality [1]. However, the fundamental challenge of estimating the joint probability from limited data remains.

Very recently, Kargas *et al.* proposed a new framework for *blindly* estimating the *joint probability mass function* (PMF) of N discrete finite-alphabet RVs [2] by modelling the N -dimensional joint PMF as an N th-order tensor. The work in [2] shows that if the RVs are “reasonably dependent”, the joint PMF can be recovered via jointly decomposing the three-dimensional marginal PMFs (which are third-order tensors). The approach does not use any *a priori* structural information of the RVs, and the recoverability of the joint PMF is provably guaranteed [2]. The work in [2] has shown

promising results, but a couple of major challenges remain. First, estimating third-order marginals accurately is not a trivial task, since real-life high-dimensional data are often very sparse. Second, computing coupled third-order tensor decomposition poses a challenging and resource-consuming optimization problem, which in general does not have known polynomial-time solvers.

Instead of working with a large number of third-order marginals, a recent work in [3] offers an *expectation maximization* (EM) algorithm that directly estimates the latent factors of the N th-order probabilistic tensor. The EM algorithm is well-motivated, since it tackles the *maximum likelihood estimator* (MLE). In addition, it admits simple and economical updates, and thus is quite scalable. However, unlike the tensor-based approach in [2], it is unclear if the recoverability of the joint PMF can be guaranteed using EM and the MLE. In addition, since the ML estimation problem is nonconvex, convergence properties of the EM algorithm is unclear.

In this work, we propose a new framework that offers provable recoverability of the joint PMF and at the same time enjoys low sample and computational complexities. To be specific, we propose an approach that utilizes only pairwise marginals to recover the joint PMF of an arbitrary number of discrete finite-alphabet RVs. This way, the sample complexity is substantially reduced relative to the three-dimensional marginal-based approach in [2]. We propose a pragmatic and easy-to-implement joint PMF estimation procedure, which is based on performing a simple and scalable Gram–Schmidt (GS)-like algorithm (namely, the *successive projection algorithm* [4]) on a carefully constructed “virtual nonnegative matrix factorization (NMF)” model. We also show that the EM algorithm in [3] can provably recover the joint PMF tensor if initialized properly, e.g., using our GS-like algorithm. We illustrate the effectiveness of the proposed approach using a number of synthetic and real-data experiments.

II. PROBLEM STATEMENT

Consider a set of discrete and finite-alphabet RVs, i.e., Z_1, \dots, Z_N . We will use $\Pr(i_1, \dots, i_N)$ as the shorthand notation to represent $\Pr(Z_1 = z_1^{(i_1)}, \dots, Z_N = z_N^{(i_N)})$ in the sequel, where $\{z_n^{(1)}, \dots, z_n^{(I_n)}\}$ denotes the alphabet of Z_n . The work in [2] shows a connection between joint PMFs and low-rank tensors under the *canonical polyadic decompo-*

This work is supported by the National Science Foundation under NSF-ECCS 1608961 and NSF-ECCS 1808159 and the Army Research Office under ARO W911NF-19-1-0247 and W911NF-19-1-0407.

sition (CPD) model. To be specific, if an N -th order tensor $\underline{\mathbf{X}} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ has CP rank F , it can be written as:

$$\underline{\mathbf{X}}(i_1, i_2, \dots, i_N) = \sum_{f=1}^F \lambda(f) \prod_{n=1}^N \mathbf{A}_n(i_n, f), \quad (1)$$

where $\mathbf{A}_n \in \mathbb{R}^{I_n \times F}$ is called the mode- n latent factor. In the above, $\boldsymbol{\lambda} = [\lambda(1), \dots, \lambda(F)]^T$ with $\|\boldsymbol{\lambda}\|_0 = F$ is employed to “absorb” the norms of columns. The work in [2] shows that any joint PMF admits a naive Bayes (NB) model representation. It follows that the joint PMF of $\{Z_n\}_{n=1}^N$ can always be decomposed as

$$\Pr(i_1, i_2, \dots, i_N) = \sum_{f=1}^F \Pr(f) \prod_{n=1}^N \Pr(i_n|f), \quad (2)$$

where $\Pr(f) := \Pr(H = f)$ is the prior distribution of a latent variable H and $\Pr(i_n|f) := \Pr(Z_n = z_n^{(i_n)} | H = f)$ are the conditional distributions. Consequently, one can represent any joint PMF as an N th-order tensor by letting $\underline{\mathbf{X}}(i_1, \dots, i_N) = \Pr(i_1, \dots, i_N)$ and $\mathbf{A}_n(i_n, f) = \Pr(i_n|f)$, $\lambda(f) = \Pr(f)$ (see more details in [2]).

The approach in [2] showed that if one has access to three-dimensional marginals, i.e., $\Pr(i_j, i_k, i_\ell)$ for different j, k, ℓ , then the joint PMF can be provably recovered through a latent factor-coupled tensor decomposition approach—if the tensor rank F is small (meaning that if the N RVs are reasonably dependent). However, estimating $\Pr(i_j, i_k, i_\ell)$ is still not easy, since one needs many co-realizations of three RVs. In addition, jointly decomposing a large number of third-order tensors poses a challenging optimization problem, whose global optimality is not guaranteed. The more recent work by Yeredor and Haardt in [3] takes an ML perspective and directly estimates the model parameters in (2) using an EM algorithm. The algorithm admits economical updates, and is effective if carefully initialized, e.g., using a coupled tensor decomposition algorithm. Nonetheless, the ML formulation’s recoverability properties are unclear. Since the ML estimator is a nonconvex optimization criterion, it is also unclear if the EM algorithm converges to the desired latent factors or not.

III. PROPOSED APPROACH

Our idea is to utilize pairwise marginals instead of the three-dimensional marginals. Under the naive Bayes model, the pairwise marginals can be expressed as $\Pr(i_j, i_k) = \sum_{f=1}^F \Pr(f) \Pr(i_j|f) \Pr(i_k|f)$, or, equivalently,

$$\mathbf{X}_{jk} = \mathbf{A}_j \mathbf{D}(\boldsymbol{\lambda}) \mathbf{A}_k^T, \quad \mathbf{X}_{jk}(i_j, i_k) = \Pr(i_j, i_k),$$

where $\mathbf{D}(\boldsymbol{\lambda}) = \text{Diag}(\boldsymbol{\lambda})$ and $\{\mathbf{A}_n\}_{n=1}^N$ and $\boldsymbol{\lambda}$ are defined as before. It is readily seen that if \mathbf{A}_n ’s and $\boldsymbol{\lambda}$ can be identified from the marginals, $\Pr(i_1, \dots, i_N)$ can be recovered by (2).

In practice, the pairwise marginals \mathbf{X}_{jk} ’s are estimated from realizations of the joint PMF. Consider a set of realizations (data samples) of $\Pr(Z_1, \dots, Z_N)$, denoted as $\{\mathbf{d}_s \in \mathbb{R}^N\}_{s=1}^S$. Assuming that there is no missing observations, the following sample averaging estimator can be employed:

$$\widehat{\mathbf{X}}_{jk}(i_j, i_k) = \frac{1}{S} \sum_{s=1}^S \mathbb{1}[\mathbf{d}_s(j) = z_j^{(i_j)}, \mathbf{d}_s(k) = z_k^{(i_k)}],$$

where $\mathbf{d}_s(n)$ denotes the realization of Z_n in the s -th data sample and $\mathbb{1}[E] = 1$ if the event E happens and $\mathbb{1}[E] = 0$ otherwise. Note that using such sample averaging schemes, the pair-wise marginals can be estimated to a much higher accuracy compared to the three-dimensional ones, under the same amount of data sample; see, e.g., [5].

However, *recoverability* of the joint PMF using pairwise marginals is nontrivial to establish. The marginal distributions, i.e., $\mathbf{X}_{jk} = \mathbf{A}_j \mathbf{D}(\boldsymbol{\lambda}) \mathbf{A}_k^T$ for all j, k , are matrices, and low-rank matrix decomposition is in general *nonunique*—while uniqueness of the employed factorization model was the key stepping stone in [2] to establish recoverability for the joint PMF from the three-dimensional marginals. A natural thought to handle the identifiability problem would be employing certain NMF tools [6], [7], since the latent factors are all nonnegative, per their physical interpretations. However, the identifiability of NMF models holds only if $F \leq \min\{I_j, I_k\}$ (and preferably $F \ll \min\{I_j, I_k\}$). The pairs $\mathbf{X}_{jk} = \mathbf{A}_j \mathbf{D}(\boldsymbol{\lambda}) \mathbf{A}_k^T \in \mathbb{R}^{I_j \times I_k}$ inherit the inner dimension F (i.e., the column dimension of \mathbf{A}_j) from the joint PMF of all the variables, which is the tensor rank of an N th-order tensor. Note that the tensor rank F could be much larger than the I_j ’s [8]. Hence, one may not directly use the available NMF uniqueness results on individual \mathbf{X}_{jk} ’s to argue for joint PMF recoverability.

A. A Virtual NMF-based Approach

To see how we approach these challenges, consider a splitting of the indices of the N variables, i.e., $\mathcal{S}_1 = \{\ell_1, \dots, \ell_M\}$ and $\mathcal{S}_2 = \{\ell_{M+1}, \dots, \ell_N\}$ such that $\mathcal{S}_1 \cup \mathcal{S}_2 = \{1, \dots, N\}$ and $\mathcal{S}_1 \cap \mathcal{S}_2 = \emptyset$, where \emptyset denotes the empty set. Then, we construct the following matrix:

$$\widetilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X}_{\ell_1 \ell_{M+1}} & \dots & \mathbf{X}_{\ell_1 \ell_N} \\ \vdots & \vdots & \vdots \\ \mathbf{X}_{\ell_M \ell_{M+1}} & \dots & \mathbf{X}_{\ell_M \ell_N} \end{bmatrix} = \mathbf{W} \mathbf{H}^T, \quad (3)$$

where $\mathbf{W} = [\mathbf{A}_{\ell_1}^T, \dots, \mathbf{A}_{\ell_M}^T]^T \in \mathbb{R}^{MI \times F}$ and $\mathbf{H} = \mathbf{D}[\mathbf{A}_{\ell_{M+1}}^T, \dots, \mathbf{A}_{\ell_N}^T]^T \in \mathbb{R}^{(N-M)I \times F}$, respectively, if $I_1 = \dots = I_N = I$. Note that \mathbf{W} and \mathbf{H} are both non-negative and thus (3) is an NMF model. By constructing $\widetilde{\mathbf{X}}$ such that $F \leq \min\{MI, (N-M)I\}$, the identifiability of \mathbf{W} and \mathbf{H} can be established using certain NMF tools. One celebrated condition is the so-called *separability* [6]:

Definition 1 (Separability) If $\mathbf{H} \geq \mathbf{0}$, and $\boldsymbol{\Lambda} = \{\ell_1, \dots, \ell_F\}$ such that $\mathbf{H}(\boldsymbol{\Lambda}, \cdot) = \boldsymbol{\Sigma}$ holds, where $\boldsymbol{\Sigma} = \text{Diag}(\alpha_1, \dots, \alpha_F)$ and $\alpha_f > 0$, then, \mathbf{H} satisfies the separability condition.

Under the exact separability condition on \mathbf{H} , we have $\mathbf{H}(\boldsymbol{\Lambda}, \cdot) = \boldsymbol{\Sigma} = \text{Diag}(\alpha_1, \dots, \alpha_F)$ and $\mathbf{W} \boldsymbol{\Sigma} = \mathbf{X}(\boldsymbol{\Lambda}, \cdot)$. Hence, the coupled NMF task boils down to identifying the index set $\boldsymbol{\Lambda}$. The *successive projection algorithm* (SPA) from the NMF literature [4] can be employed for this purpose. Notably, this algorithm admits Gram–Schmidt-like economical and scalable updates and is provably robust to noise.

Once \mathbf{W} is identified, we can recover $\mathbf{A}_{\ell_n} \in \mathbb{R}^{I_{\ell_n} \times F}$ for $\ell_n \in \mathcal{S}_1$ up to identical column permutations, by extracting the corresponding rows of \mathbf{W} . Unlike general NMF models, since we know every column of \mathbf{A}_n is a conditional PMF, there is no scaling ambiguity. The \mathbf{H} matrix can be estimated using (constrained) least squares, and \mathbf{A}_{ℓ_n} for $\ell_n \in \mathcal{S}_2$ can then be extracted in the similar way. Denote (3) as $\widetilde{\mathbf{X}} = \mathbf{W}\mathbf{D}(\boldsymbol{\lambda})\widetilde{\mathbf{H}}^\top$, where $\widetilde{\mathbf{H}} = [\mathbf{A}_{\ell_{M+1}}^\top \dots \mathbf{A}_{\ell_N}^\top]^\top$. Then, the PMF of the latent variable can be estimated via $\boldsymbol{\lambda} = (\widetilde{\mathbf{H}} \odot \mathbf{W})^\dagger \text{vec}(\widetilde{\mathbf{X}})$ where we have used the fact that the Khatri-Rao product $\widetilde{\mathbf{H}} \odot \mathbf{W}$ has full column rank since both \mathbf{W} and $\widetilde{\mathbf{H}}$ have full column rank. Note that the permutation ambiguity across \mathbf{A}_n 's and $\boldsymbol{\lambda}$ are identical. Hence, the existence of column permutations does not affect the ‘‘assembling’’ of $\Pr(i_n|f)$ and $\Pr(f)$ to recover $\Pr(i_1, \dots, i_N)$. We refer to this procedure as *coupled NMF via SPA* (CNMF-SPA); see Algorithm 1.

Algorithm 1: CNMF-SPA

input : data samples $\{\mathbf{d}_s\}_{s=1}^S$ and M
1 estimate second order statistics $\widehat{\mathbf{X}}_{jk}$;
2 split $\{1, \dots, N\}$ into $\mathcal{S}_1 = \{1, \dots, M\}$ and $\mathcal{S}_2 = \{M+1, \dots, N\}$;
3 Construct $\widetilde{\mathbf{X}}$;
4 Estimate $\widehat{\mathbf{W}}$ using the SPA algorithm [4] to select \mathbf{A} ;
5 **for** $n = 1$ **to** M **do**
6 $\widehat{\mathbf{A}}_n \leftarrow \widehat{\mathbf{W}}((n-1)I_n + 1 : nI_n, :)$;
7 normalize columns of $\widehat{\mathbf{A}}_n$ with respect to ℓ_1 norm;
8 **end**
9 $\widehat{\mathbf{H}} \leftarrow \arg \min_{\mathbf{H} \geq 0} \|\widetilde{\mathbf{X}} - \widehat{\mathbf{W}}\mathbf{H}^\top\|_F^2$;
10 **for** $n = M+1$ **to** N **do**
11 $\widehat{\mathbf{A}}_n \leftarrow \widehat{\mathbf{H}}((n-1)I_n + 1 : nI_n, :)$;
12 normalize columns of $\widehat{\mathbf{A}}_n$ with respect to ℓ_1 norm;
13 **end**
14 $\widehat{\mathbf{W}}^\top \leftarrow [\widehat{\mathbf{A}}_1^\top, \dots, \widehat{\mathbf{A}}_M^\top]$;
15 $\widehat{\mathbf{H}}^\top \leftarrow [\widehat{\mathbf{A}}_{M+1}^\top \dots \widehat{\mathbf{A}}_N^\top]^\top$;
16 $\widehat{\boldsymbol{\lambda}} \leftarrow (\widehat{\mathbf{H}} \odot \widehat{\mathbf{W}})^\dagger \text{vec}(\widetilde{\mathbf{X}})$;
output: estimates $\{\widehat{\mathbf{A}}_n\}_{n=1}^N, \widehat{\boldsymbol{\lambda}}$.

B. Performance Analysis

The CNMF-SPA procedure looks simple, but several caveats exist. First, SPA only works if one can construct \mathcal{S}_1 and \mathcal{S}_2 such that \mathbf{H} in (3) satisfies the separability condition. Testing all combinations of \mathcal{S}_1 and \mathcal{S}_2 gives a rise to a combinatorial problem, which is apparently impossible. Second, the pairwise marginals \mathbf{X}_{ij} are estimated through a finite number of data samples, and typically there are many missing values in different data samples (i.e., not all the realizations of Z_1, \dots, Z_n are observed in any \mathbf{d}_s)—which both make the estimated pairwise marginals very noisy. It is unclear how the performance of CNMF-SPA is affected.

In practice, we observe that using the ‘‘naive’’ construction $\mathcal{S}_1 = \{1, \dots, M\}$ and $\mathcal{S}_2 = \{M+1, \dots, N\}$ seems to work reasonably well, even if the number of samples is finite and many missing values exist. To understand such effectiveness, we assume that S realizations of $\Pr(Z_1, \dots, Z_N)$ are available. In each realization, every variable is observed with probability

p —which determines how much is missing in the dataset. For simplicity, we also assume that $I_n = I$ for all n and utilize the following generative model for the nonnegative \mathbf{A}_m 's:

Assumption 1 Assume that the rows of \mathbf{A}_m 's are generated from the $(F-1)$ -probability simplex uniformly at random and then positively scaled, so that $\mathbf{1}^\top \mathbf{A}_m = 1$ is respected.

Under the above settings, we show that the following holds:

Theorem 1 Assume that $\|\widehat{\mathbf{X}}_{ij}(:, q)\|_1 \geq \eta > 0$ for any q, i, j . Also, assume that $M \geq F/I$, $p \geq (\frac{8}{S} \log(4/\delta))^{1/2}$,

$$S = \Omega\left(\frac{M^2 I \log(1/\delta)}{\sigma_{\max}^2(\mathbf{W}) \eta^2 \varepsilon^2 p^2}\right),$$

$$N = M + \Omega\left(\frac{\varepsilon^{-2(F-1)}}{IF} \log\left(\frac{F}{\delta}\right)\right),$$

where $\varepsilon = \frac{M \min(\frac{1}{2\sqrt{F-1}}, \frac{1}{4})}{2\kappa(\mathbf{W})(1+80\kappa^2(\mathbf{W}))}$. Then, under the defined S, p and Assumption 1, CNMF-SPA outputs $\widehat{\mathbf{A}}_m$'s such that

$$\min_{\Pi: \text{permutation}} \|\widehat{\mathbf{A}}_m \Pi - \mathbf{A}_m\|_2 = O\left(\kappa^2(\mathbf{W}) \sqrt{F} \zeta\right) \quad (4)$$

for $m \in \mathcal{S}_1$ with a probability greater than or equal to $1 - \delta$, where $\zeta = \max(\sigma_{\max}(\mathbf{W})\varepsilon, M\sqrt{I \log(1/\delta)}/\eta p \sqrt{S})$.

The proof can be found in a long version of the paper in the appendix. Note that if \mathbf{A}_m for all $m \in \mathcal{S}_1$ can be accurately estimated, the estimation accuracy of \mathbf{A}_n for all $n \in \mathcal{S}_2$ and $\boldsymbol{\lambda}$ can also be guaranteed and quantified, following standard sensitivity analyses of least squares. We leave this part out of the work for conciseness.

Remark 1 Theorem 1 is not entirely surprising. The insight behind is to model the finite sample-induced noise and the violation of separability as combined virtual noise, and then utilize the robustness of SPA [4]. The challenge lies in quantifying this virtual noise, for which we leverage concentration theorems and Assumption 1. We would like to remark that Assumption 1 is a working assumption for us to understand the effectiveness of CNMF-SPA under the naive \mathcal{S}_1 and \mathcal{S}_2 . The key fact is that when $|\mathcal{S}_2|$ grows, \mathbf{H} in (3) has a good chance to attain the separability condition approximately [9] under Assumption 1. In principle, this fact holds if the rows of \mathbf{A}_m are drawn from any joint continuous distribution, but using the uniform distribution assumed in Assumption 1 helps simplify the analysis.

C. Refinement via EM and Optimality Guarantees

In [3], Yeredor and Haardt proposed an EM algorithm to handle the ML estimator for the naive Bayes model in (2). Using MLE exhibits promising performance for joint PMF recovery—after all the MLE is a ‘‘gold standard’’ for statistical learning. However, the EM algorithm often converges to undesired solutions, if randomly initialized. This is perhaps due to the nonconvex nature of the MLE problem. As observed

in [3], using good initialization may improve the performance of EM. However, quantification for the required quality of the initial estimates has been elusive.

In this work, we offer theoretical supports for the EM algorithm in [3]. To be specific, we show that, with good initial estimations for \mathbf{A}_n and $\boldsymbol{\lambda}$ (e.g., those output by CNMF-SPA), the EM algorithm improves the solution towards the ground truth. To proceed, we make the following assumption:

Assumption 2 Define \bar{D}_1 and \bar{D}_2 as follows:

$$\bar{D}_1 = \min_{f \neq f'} \frac{1}{N} \sum_{n=1}^N p \mathbb{D}_{\text{KL}}(\mathbf{A}_n(:, f), \mathbf{A}_n(:, f')),$$

$$\bar{D}_2 = \frac{2}{N} \min_{f \neq f'} \log(\boldsymbol{\lambda}(f)/\boldsymbol{\lambda}(f')).$$

and $\bar{D} = (\bar{D}_1 + \bar{D}_2)/2$. Assume that $\mathbf{A}_n, \boldsymbol{\lambda}$ and the initial estimates $\hat{\mathbf{A}}_n^0, \hat{\boldsymbol{\lambda}}^0$ satisfy $|\hat{\mathbf{A}}_n^0(i, f) - \mathbf{A}_n(i, f)| \leq \delta_1 := \frac{4}{\rho_1(4+\bar{D})}$, $\mathbf{A}_n(i, f) \geq \rho_1$, $|\hat{\boldsymbol{\lambda}}^0(f) - \boldsymbol{\lambda}(f)| \leq \delta_2 := \frac{4}{\rho_2(4+N\bar{D})}$ and $\boldsymbol{\lambda}(f) \geq \rho_2$ for all n, i, f .

Note that \bar{D}_1 characterizes the ‘‘conditioning’’ of \mathbf{A}_n under the Kullback-Leibler (KL) divergence (denoted by $\mathbb{D}_{\text{KL}}(\cdot, \cdot)$) sense, and $|\bar{D}_2|$ measures how far $\boldsymbol{\lambda}$ is away from the uniform distribution. Under the above assumption, we show that the EM algorithm improves upon the initialization:

Theorem 2 Let $\delta_{\min} = \min(\delta_1, \delta_2)$. Assume that the following hold:

$$N \geq \max \left(\frac{33 \log(3SF/\mu)}{\rho_1 \bar{D}_1}, \frac{4 \log(4SF^2/(3p\rho_2\mu))}{\bar{D}} \right),$$

$$S \geq \frac{192F^2 \log(12NFI/\mu)}{p^2 \rho_2^2 \delta_{\min}^2}, \quad \bar{D} \geq \max \left\{ \frac{8 - 4\rho_1^2}{\rho_1^2}, \frac{8 - 4\rho_2^2}{N\rho_2^2} \right\}.$$

Then, under Assumption 2, the EM algorithm in [3] outputs $\hat{\mathbf{A}}_n(i, f), \hat{\boldsymbol{\lambda}}(f)$ that satisfy the following with a probability greater than or equal to $1 - \mu$:

$$|\hat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f)|^2 \leq \frac{48 \log(12NFI/\mu)}{Sp\boldsymbol{\lambda}(f)} \leq \delta_1^2,$$

$$|\hat{\boldsymbol{\lambda}}(f) - \boldsymbol{\lambda}(f)|^2 \leq \frac{192F^2 \boldsymbol{\lambda}(f) \log(12NFI/\mu)}{S} \leq \delta_2^2.$$

The proof of the theorem can be found in a longer version of the paper in the appendix. Our proof extends the analysis of a different EM algorithm proposed in [10] that is designed for learning the Dawid-Skene model in crowdsourcing. The EM algorithm there effectively learns a naive Bayes model when the latent variable is uniform. Our analysis covers the more recent EM algorithm in [3] that can handle $\boldsymbol{\lambda}$'s who have general PMFs beyond the uniform distribution.

Since CNMF-SPA is a natural economical initialization for Yeredor and Haardt's EM, we combine the two algorithms together and refer to this procedure as CNMF-SPA-EM.

IV. EXPERIMENTS

In this section, we present experiments to showcase the effectiveness of the proposed framework.

TABLE I: MSE & MRE for $N = 5, F = 5, I = 10, p = 0.5$

Algorithms	Metric	$S = 10^3$	$S = 10^4$	$S = 10^5$	$S = 10^6$
CNMF-SPA [Proposed]	MSE	0.0702	0.0257	0.0211	0.0204
CNMF-SPA-EM [Proposed]	MSE	0.0560	0.0230	0.0207	0.0204
RAND-EM [3]	MSE	0.0855	0.0405	0.0298	0.0502
CTD [2]	MSE	0.1589	0.0260	0.0211	0.0205
CTD-EM [3]	MSE	0.1196	0.0233	0.0208	0.0204
CNMF-SPA [Proposed]	MRE	0.8084	0.3228	0.1137	0.0356
CNMF-SPA-EM [Proposed]	MRE	0.6922	0.2077	0.0682	0.0219
RAND-EM [3]	MRE	0.8285	0.3399	0.2226	0.3931
CTD [2]	MRE	0.9237	0.3081	0.0955	0.0309
CTD-EM [3]	MRE	0.8312	0.2180	0.0681	0.0220

TABLE II: MovieLens Action Movies set

Algorithm	RMSE	MAE	Time (s)
CNMF-SPA [Proposed]	0.8536±0.0071	0.6679±0.0061	0.029
CNMF-SPA-EM [Proposed]	0.7761±0.0039	0.5936±0.0039	2.554
CTD [2]	0.8792±0.0137	0.6640±0.0104	17.808
CTD-EM [3]	0.7872±0.0053	0.6038±0.0076	20.578
BMF [11]	0.8020±0.0015	0.6268±0.0017	45.967
Global Average	0.9471±0.0010	0.6954±0.0010	-
User Average	0.8960±0.0012	0.6834±0.0007	-
Movie Average	0.8844±0.0009	0.6979±0.0008	-

A. Synthetic Data

We consider $N = 5$ RV's where each variable takes $I = 10$ discrete values. The columns of the conditional PMF matrices (factor matrices) $\mathbf{A}_n \in \mathbb{R}^{I_n \times F}$ and the prior probability vector $\boldsymbol{\lambda} \in \mathbb{R}^F$ are generated with $F = 5$. The so-called ε -separability condition on \mathbf{H} in (3) holds with $\varepsilon = 0.1$; see the definition of ε -separability in [9]. We generate S realizations of the joint PMF and randomly hide each variable's realization with probability $p = 0.5$. We fix $M = 3$. The mean squared error (MSE) of the factors and the mean relative error (MRE) of the recovered joint PMFs (see [2]) are evaluated. MRE is more preferred for evaluation, but it is hard to compute (due to memory issues) for large N . The results are averaged from 20 random trials. We benchmark our method using the EM algorithm by Yeredor and Haardt [3] initialized by random guesses (denoted as RAND-EM) and an alternating optimization (AO) algorithm-based coupled tensor decomposition (CTD) (denoted as CTD-EM), respectively.

Table I shows that the proposed approaches CNMF-SPA and CNMF-SPA-EM exhibit promising performance. In particular, CNMF-SPA is effective for initializing the EM algorithm whereas RAND-EM sometimes struggles to attain good performance. CTD-EM also works well when S is large, perhaps because the CTD stage needs a large S to estimate the three-dimensional marginals accurately.

B. Real Data : Recommender Systems

We test the approaches using the MovieLens 20M dataset [12], which has many missing values. We first round the ratings to the closest integers so that every movie's rating resides in $\{1, 2, \dots, 5\}$. We choose select a subset of movies from the action movie genre. This way, Z_i for $i = 1, \dots, 30$ represent the ratings of movie i , and all Z_i 's alphabets are $\{1, 2, \dots, 5\}$. We predict the rating for a movie (e.g., movie N) by user k via computing $\mathbb{E}[i_N | r_k(1), \dots, r_k(N-1)]$ (i.e., using the MMSE estimator), where $r_k(i)$ denotes the rating of movie i by user k . This can be done via estimating $\Pr(i_1, \dots, i_N)$.

TABLE III: UCI Dataset Car

Algorithm	Accuracy (%)	Time (s)
CNMF-SPA [Proposed]	69.26±2.28	0.007
CNMF-SPA-EM [Proposed]	86.61±1.76	0.018
CTD [2]	83.47±2.34	0.845
CTD-EM [3]	85.72±1.88	0.955
SVM	83.65±1.58	0.147
Linear Regression	80.68±1.61	0.029
Neural Net	85.00±3.22	0.193
SVM-RBF	76.22±3.93	0.793
Naive Bayes	83.42±2.15	0.026

We create the validation and testing sets by randomly hiding 20% and 30% of the dataset for each trial. The remaining 50% is used for training (learning joint PMF in our approach). In this task, we also use one of the popular recommender system algorithms, *biased matrix factorization* (BMF) method [11] as a baseline. The rank F for all the methods (ranging from 5 to 25) and the number of iterations needed for EM are chosen using the validation set. The results are taken from 20 random trials. We report the *root mean squared error* (RMSE) and *mean absolute error* (MAE) of the predicted ratings.

From Table II, one can see that the proposed methods are promising in terms of prediction accuracy and runtime. Note BMF is specialized for recommender systems, while the proposed approaches are for generic joint PMF recovery. The fact that our methods perform better suggests that the underlying joint PMF is well captured by the proposed CNMF approach. Another important observation is that the CPD method in [2] does not perform as well compared to the proposed pairwise marginals based methods. This may be because of the noisy estimation for three-dimensional marginals, due to the sparse nature of the user-movie datasets. In particular, CNMF-SPA is fast with acceptable prediction accuracy. In addition, CNMF-SPA-EM presents a good accuracy and speed tradeoff—it exhibits the lowest RMSEs and MAEs, and is 8 times faster than the state-of-the-art algorithm, i.e., CTD-EM.

C. Real Data : Data Classification

We consider UCI datasets for classification tasks. We split each dataset into training, validation and testing sets in the ratio of 50 : 20 : 30. For our approach, we estimate the joint PMF of the features and the label using the training set, and then predict the labels on the testing data by constructing a *Maximum A Posterior* (MAP) predictor (i.e., predicting the labels conditioned on the features). For each dataset, we perform 20 trials with randomly partitioned training/testing/validation sets and take average of the results.

Tables III and IV show results on the UCI datasets Car and Mushroom. Car has 1,728 data samples from 4 classes, and Mushroom has 8,124 data samples belonging to two classes. We have $N = 7$ and $N = 22$ for the two datasets, respectively, and the average I 's are 4 and 6, respectively. We set $M = 5$ in all experiments, and select rank F as before. In all the cases, one can see that the proposed combination CNMF-SPA-EM gives the most promising results. In particular, the solely using CNMF-SPA is not as promising, perhaps because the N and I are not large enough to ensure a high-accuracy performance

TABLE IV: UCI Dataset Mushroom

Algorithm	Accuracy (%)	Time (sec.)
CNMF-SPA [Proposed]	92.23+/-6.15	0.025
CNMF-SPA-EM [Proposed]	99.47+/-0.80	0.242
CTD [2]	96.40+/-0.59	13.695
CTD-EM [3]	97.18+/-1.21	13.931
SVM	97.47+/-0.46	37.213
Linear Regression	93.38+/-0.59	0.040
Neural Net	98.98+/-1.97	1.036
SVM-RBF	98.89+/-0.34	2.291
Naive Bayes	94.84+/-0.55	0.048

of the CNMF-SPA stage. However, when initialized using CNMF-SPA, the CNMF-SPA-EM outputs the best classification accuracy and is at least 50 times faster than CTD-EM that is suggested in [3]. This suggests that even under critical scenarios, CNMF-SPA still offers useful initialization for EM.

V. CONCLUSION

We proposed a new framework for recovering joint PMF of any number of discrete RVs from marginal distributions. Unlike a recent approach that relies on three-dimensional marginals, our approach only uses two-dimensional marginals, which naturally has reduced-sample complexity and lighter computational burden. We proposed a virtual NMF framework and employed a Gram-Schmidt-like scalable algorithm for handling our formulation. We showed that the proposed framework is effectiveness under realistic conditions, e.g., finite samples. We also showed that an existing EM algorithm can provably improve the output of our NMF approach, using theoretical analysis and experimental validation. The combined NMF and EM approach admits economical updates and exhibits appealing joint PMF recovery accuracy.

REFERENCES

- [1] K. P. Murphy, *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [2] N. Kargas, N. D. Sidiropoulos, and X. Fu, "Tensors, learning, and 'Kolmogorov extension' for finite-alphabet random vectors," *IEEE Trans. Signal Process.*, vol. 66, no. 18, pp. 4854–4868, 2018.
- [3] A. Yeredor and M. Haardt, "Maximum likelihood estimation of a low-rank probability mass tensor from partial observations," *IEEE Signal Process. Lett.*, vol. 26, no. 10, pp. 1551–1555, Oct 2019.
- [4] N. Gillis and S. Vavasis, "Fast and robust recursive algorithms for separable nonnegative matrix factorization," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 4, pp. 698–714, April 2014.
- [5] Y. Han, J. Jiao, and T. Weissman, "Minimax estimation of discrete distributions under ℓ_1 loss," *IEEE Transactions on Information Theory*, vol. 61, no. 11, pp. 6343–6354, 2015.
- [6] X. Fu, K. Huang, N. D. Sidiropoulos, and W.-K. Ma, "Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications," *IEEE Signal Process. Mag.*, vol. 36, no. 2, pp. 59–80, March 2019.
- [7] N. Gillis, "The why and how of nonnegative matrix factorization," *Regularization, Optimization, Kernels, and Support Vector Machines*, vol. 12, p. 257, 2014.
- [8] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos, "Tensor decomposition for signal processing and machine learning," *IEEE Trans. Signal Process.*, vol. 65, no. 13, pp. 3551–3582, 2017.
- [9] S. Ibrahim, X. Fu, N. Kargas, and K. Huang, "Crowdsourcing via pairwise co-occurrences: Identifiability and algorithms," *arXiv preprint arXiv:1909.12325*, 2019.
- [10] Y. Zhang, X. Chen, D. Zhou, and M. I. Jordan, "Spectral methods meet em: A provably optimal algorithm for crowdsourcing," in *Advances in Neural Information Processing Systems*, 2014, pp. 1260–1268.

- [11] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, no. 8, pp. 30–37, 2009.
- [12] F. M. Harper and J. A. Konstan, "The movielens datasets: History and context," *ACM Trans. Interact. Intell. Syst.*, vol. 5, no. 4, pp. 19:1–19:19, Dec. 2015.
- [13] A. Anandkumar, D. Hsu, and S. Kakade, "A method of moments for mixture models and hidden markov models," *Journal of Machine Learning Research*, vol. 23, 03 2012.

APPENDIX A
PROOF OF THEOREM 1

Consider the noisy matrix factorization model as below:

$$\widetilde{\mathbf{X}} = \mathbf{W}\mathbf{H}^\top + \mathbf{N}, \quad (5)$$

where $\mathbf{W} \in \mathbb{R}^{L \times F}$, $\mathbf{H} \in \mathbb{R}^{K \times F}$, $\mathbf{W} \geq \mathbf{0}$ and $\mathbf{H} \geq \mathbf{0}$, and $\mathbf{N} \in \mathbb{R}^{L \times K}$ represents the noise. Also assume that $\text{rank}(\mathbf{W}) = F$ and $\mathbf{H}\mathbf{1} = \mathbf{1}$, $\mathbf{H} = \mathbf{\Pi} \begin{bmatrix} \mathbf{I}_F \\ \mathbf{H}^* \end{bmatrix}$ where \mathbf{I}_F is the identity matrix of size F and $\mathbf{\Pi}$ is the permutation matrix. This implies that \mathbf{H} satisfies *seperability* condition and that there exists $\mathbf{A} = \{l_1, \dots, l_F\}$ such that $\mathbf{H}(\mathbf{A}, :) = \mathbf{I}_F$. Gillis and Vavasis [4] have shown that under the model in (5), SPA is provably robust to noise in estimating the factor matrix \mathbf{W} (see Theorem 3).

First, we characterize the noise in our data model given by (3) by re-expressing it as a ‘virtual’ separable NMF as in (5). Then, we utilize Theorem 3 to characterize the error in estimating \mathbf{W} via SPA algorithm.

Consider the pairwise marginals \mathbf{X}_{jk} ’s used to construct the matrix $\widetilde{\mathbf{X}}$ in (3). \mathbf{X}_{jk} ’s are estimated by sample averaging of a finite number of realizations and thus the estimated \mathbf{X}_{jk} (denoted as $\widehat{\mathbf{X}}_{jk}$) is always noisy; i.e., we have

$$\widehat{\mathbf{X}}_{jk} = \mathbf{X}_{jk} + \mathbf{N}_{jk}, \quad (6)$$

where the noise matrix $\mathbf{N}_{jk} \in \mathbb{R}^{I \times I}$, assuming $I_n = I$ for all $n \in \{1, \dots, N\}$.

In order to characterize the estimation accuracy of $\widehat{\mathbf{X}}_{jk}$ using finite number of realizations, we have the following proposition:

Proposition 1 *Let $p \in (0, 1]$ be the probability that an RV is observed. Let S be the number of available realizations of N RVs. Assume that $p \geq (\frac{8}{9} \log(2/\delta))^{1/2}$. Then, with probability at least $1 - \delta$,*

$$\|\mathbf{X}_{jk} - \widehat{\mathbf{X}}_{jk}\|_F = \|\mathbf{N}_{jk}\|_F \leq \phi,$$

holds for any distinct j, k where $\phi = \frac{\sqrt{2}(1 + \sqrt{\log(2/\delta)})}{(p\sqrt{S})}$.

The proof of Proposition 1 is given in Sec. C.

By the definition of the Frobenius norm, we have

$$\sum_{c=1}^I \|\mathbf{N}_{ij}(:, c)\|_2^2 = \|\mathbf{N}_{ij}\|_F^2 \leq \phi^2.$$

Applying norm equivalence $\frac{\|\mathbf{N}_{ij}(:, c)\|_1}{\sqrt{I}} \leq \|\mathbf{N}_{ij}(:, c)\|_2$, we get

$$\sum_{c=1}^I \|\mathbf{N}_{ij}(:, c)\|_1^2 \leq I\phi^2. \quad (7)$$

Eq. (7) implies that, for all $c \in \{1, \dots, I\}$ and any i, j where $i \neq j$,

$$\begin{aligned} \|\mathbf{N}_{ij}(:, c)\|_1^2 &\leq I\phi^2, \\ \implies \|\mathbf{N}_{ij}(:, c)\|_1 &\leq \sqrt{I}\phi. \end{aligned} \quad (8)$$

By using the estimates $\widehat{\mathbf{X}}_{jk}$, the model given by (3) can be represented as

$$\begin{aligned} \widehat{\mathbf{X}} &= \begin{bmatrix} \widehat{\mathbf{X}}_{\ell_1 \ell_{M+1}} & \cdots & \widehat{\mathbf{X}}_{\ell_1 \ell_N} \\ \vdots & & \vdots \\ \widehat{\mathbf{X}}_{\ell_M \ell_{M+1}} & \cdots & \widehat{\mathbf{X}}_{\ell_M \ell_N} \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} \mathbf{A}_{\ell_1} \\ \vdots \\ \mathbf{A}_{\ell_M} \end{bmatrix}}_{\mathbf{W}} \underbrace{D(\lambda)[\mathbf{A}_{\ell_{M+1}}^\top, \dots, \mathbf{A}_{\ell_N}^\top]}_{\mathbf{H}^\top} + \widetilde{\mathbf{N}} \\ &= \widetilde{\mathbf{X}} + \widetilde{\mathbf{N}}. \end{aligned} \quad (9)$$

Note that $\widehat{\mathbf{X}}$, $\widetilde{\mathbf{X}}$ and $\widetilde{\mathbf{N}}$ all have the same size of $L \times K$. Assuming $I_n = I$ for all $n \in \{1, \dots, N\}$, we have $L = MI$ and $K = (N - M)I$. Also note that \mathbf{W} has a size of $L \times F$ and \mathbf{H} has a size of $K \times F$.

Since any column of $\widetilde{\mathbf{N}}$ formed from the columns of M number of \mathbf{N}_{ij} 's, we have

$$\|\widetilde{\mathbf{N}}(:, q)\|_1 \leq M\sqrt{I}\phi, \quad (10)$$

where the last inequality is obtained by using triangle inequality and (8).

Next, we consider estimating \mathbf{W} given in (9) using SPA algorithm. Before performing SPA to the data matrix $\widehat{\mathbf{X}}$, the columns of $\widehat{\mathbf{X}}$ are normalized with respect to the ℓ_1 -norm. Let us denote the normalized data model as follows:

$$\overline{\mathbf{X}} = \overline{\mathbf{W}}\overline{\mathbf{H}}^\top + \overline{\mathbf{N}}, \quad (11)$$

where $\overline{\mathbf{X}}$ and $\overline{\mathbf{W}}$ are column normalized versions (with respect to the ℓ_1 norm) of $\widehat{\mathbf{X}}$ and \mathbf{W} , respectively, and $\overline{\mathbf{H}}$ is row normalized version of \mathbf{H} .

Since the matrix $\widehat{\mathbf{X}}$ is noisy, the effect of normalization on $\widetilde{\mathbf{N}}$ can be characterized by Lemma 6. From the assumption $\|\widehat{\mathbf{X}}_{ij}(:, c)\|_1 \geq \eta$ for any $i \neq j$ and $c \in \{1, \dots, I\}$, we get $\|\widehat{\mathbf{X}}(:, q)\|_1 \geq M\eta$ for any q . Combining Lemma 6 and Eq. (10), we get

$$\|\overline{\mathbf{N}}(:, q)\|_1 \leq \frac{2\sqrt{I}\phi}{\eta}. \quad (12)$$

Applying norm equivalence, we further have $\|\overline{\mathbf{N}}(:, q)\|_2 \leq \|\overline{\mathbf{N}}(:, q)\|_1$ and hence we get

$$\|\overline{\mathbf{N}}(:, q)\|_2 \leq \frac{2\sqrt{I}\phi}{\eta}. \quad (13)$$

Lemma 1 Assume that $\|\overline{\mathbf{N}}(:, q)\|_2 \leq \varphi$ for any q and that $\overline{\mathbf{H}}$ satisfies ε -separability assumption in the model (11). Suppose

$$(\sigma_{\max}(\overline{\mathbf{W}})\varepsilon + \varphi) \leq \sigma_{\min}(\overline{\mathbf{W}}) \min\left(\frac{1}{2\sqrt{F}-1}, \frac{1}{4}\right) (1 + 80\kappa^2(\overline{\mathbf{W}}))^{-1}.$$

Then, SPA identifies an index set $\widehat{\mathbf{A}} = \{\widehat{l}_1, \dots, \widehat{l}_F\}$ such that

$$\max_{1 \leq f \leq F} \min_{\widehat{l}_f \in \widehat{\mathbf{A}}} \left\| \overline{\mathbf{W}}(:, f) - \overline{\mathbf{X}}(:, \widehat{l}_f) \right\|_2 \leq (\sigma_{\max}(\overline{\mathbf{W}})\varepsilon + \varphi) (1 + 80\kappa^2(\overline{\mathbf{W}})), \quad (14)$$

where $\kappa(\overline{\mathbf{W}}) = \frac{\sigma_{\max}(\overline{\mathbf{W}})}{\sigma_{\min}(\overline{\mathbf{W}})}$ is the condition number of $\overline{\mathbf{W}}$.

The proof of Lemma 1 is given in Sec. D.

The right hand side of (14) can be written as

$$\begin{aligned} \max_{1 \leq f \leq F} \min_{\widehat{l}_f \in \widehat{\mathbf{A}}} \left\| \overline{\mathbf{W}}(:, f) - \overline{\mathbf{X}}(:, \widehat{l}_f) \right\|_2^2 &= \frac{1}{M} \sum_{m=1}^M \max_{1 \leq f \leq F} \min_{\widehat{l}_f \in \widehat{\mathbf{A}}} \left\| \mathbf{A}_{\ell_m}(:, f) - \widehat{\mathbf{A}}_{\ell_m}(:, f) \right\|_2^2 \\ &\geq \frac{1}{M} \max_{1 \leq f \leq F} \min_{\widehat{l}_f \in \widehat{\mathbf{A}}} \left\| \mathbf{A}_{\ell_m}(:, f) - \widehat{\mathbf{A}}_{\ell_m}(:, f) \right\|_2^2, \end{aligned} \quad (15)$$

for any $m \in \{1, \dots, M\}$, where the first equality is due to $\overline{\mathbf{W}} = [\mathbf{A}_{\ell_1}^\top, \dots, \mathbf{A}_{\ell_M}^\top]^\top / M$, in which $\widehat{\mathbf{A}}_{\ell_m}$ denotes the corresponding estimate of \mathbf{A}_{ℓ_m} .

Since $\|\mathbf{W}(:, f)\|_1 = M$ for any f , $\overline{\mathbf{W}} = \mathbf{W}/M$. Therefore, we have

$$\sigma_{\max}(\overline{\mathbf{W}}) = \sigma_{\max}(\mathbf{W})/M, \quad \sigma_{\min}(\overline{\mathbf{W}}) = \sigma_{\min}(\mathbf{W})/M, \quad \kappa(\overline{\mathbf{W}}) = \frac{\sigma_{\max}(\overline{\mathbf{W}})}{\sigma_{\min}(\overline{\mathbf{W}})} = \kappa(\mathbf{W}). \quad (16)$$

Therefore, by combining (13),(15),(16) and Lemma 1, SPA estimates for any $m \in \{1, \dots, M\}$ such that

$$\max_{1 \leq f \leq F} \min_{\widehat{l}_f \in \widehat{\mathbf{A}}} \left\| \mathbf{A}_{\ell_m}(:, f) - \widehat{\mathbf{A}}_{\ell_m}(:, f) \right\|_2 \leq \left(\sigma_{\max}(\mathbf{W})\varepsilon + \frac{2M\sqrt{I}\phi}{\eta} \right) (1 + 80\kappa^2(\mathbf{W})), \quad (17)$$

if the below condition is satisfied:

$$\frac{\sigma_{\max}(\mathbf{W})}{M}\varepsilon + \frac{2\sqrt{I}\phi}{\eta} \leq \sigma_{\min}(\mathbf{W}) \min\left(\frac{1}{2\sqrt{F}-1}, \frac{1}{4}\right) (1 + 80\kappa^2(\mathbf{W}))^{-1}. \quad (18)$$

Letting $\varepsilon = \frac{M \min\left(\frac{1}{2\sqrt{F-1}}, \frac{1}{4}\right)}{2\kappa(\mathbf{W})(1+80\kappa^2(\mathbf{W}))}$, from (18), we get the condition on ϕ as follows:

$$\phi \leq \frac{\eta \sigma_{\max}(\mathbf{W}) \varepsilon}{4M\sqrt{I}}. \quad (19)$$

From Proposition 1, we have $\phi = \frac{\sqrt{2(1+\sqrt{\log(2/\delta)})}}{p\sqrt{S}}$ with probability greater than $1 - \delta$. By substituting ϕ on the left hand side of (19), we get the number of realizations S required to get the estimation error bound (17) as below:

$$S \geq \frac{32M^2I(1 + \sqrt{\log(2/\delta)})^2}{\sigma_{\max}^2(\mathbf{W})\eta^2\varepsilon^2p^2}.$$

Note that Lemma 1 holds if $\overline{\mathbf{H}}$ satisfies ε -separability condition. By combining Assumption 1 and Lemma 7, we get that $\overline{\mathbf{H}}$ satisfies ε -separability assumption with probability greater than $1 - \rho$, if

$$(N - M)I = \Omega\left(\frac{\varepsilon^{-2(F-1)}}{F} \log\left(\frac{F}{\rho}\right)\right). \quad (20)$$

By substituting ϕ in (17) and using the fact that for any matrix $\mathbf{A} \in \mathbb{R}^{I \times F}$, the matrix 2-norm $\|\mathbf{A}\|_2 \leq \sqrt{F} \max_{1 \leq f \leq F} \|\mathbf{A}(:, f)\|_2$, we get the result (4) in the theorem.

Finally, we combine the probabilities involved in the results used in our proof. For the concentration bound in Proposition 1 and ε -separability condition on $\overline{\mathbf{H}}$ given by Lemma 7 to jointly occur with probability greater than $1 - 2\delta$, we can assign $\rho = \delta$ in (20).

This completes the proof.

APPENDIX B PROOF OF THEOREM 2

The EM algorithm proposed in [3] is as follows. Let $\mathbf{d}_s \in \mathbb{R}^N$ be the s -th joint realization of N RVs. Let $f_s \in \{1, \dots, F\}$ be the realization of the ‘latent variable’ H in the s th realization. Suppose $\{z_n^{(1)}, z_n^{(2)}, \dots, z_n^{(I_n)}\}$ denotes the alphabet set of n -th RV, then n -th entry in \mathbf{d}_s denoted as $\mathbf{d}_s(n) \in \{z_n^{(1)}, z_n^{(2)}, \dots, z_n^{(I_n)}\}$. The expectation maximization algorithm proposed in [3] has the following E-step and M-step which are executed alternatively until convergence:

E-step: The parameter $\hat{q}_{s,f}$ is updated for all s, f using the current estimates of $\hat{\mathbf{A}}_n$ ’s and $\hat{\lambda}$

$$\hat{q}_{s,f} = \frac{\exp(\log(\hat{\lambda}(f)) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log(\hat{\mathbf{A}}_n(i, f)))}{\sum_{f'=1}^F \exp(\log(\hat{\lambda}(f')) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log(\hat{\mathbf{A}}_n(i, f')))} \quad (21)$$

M-step: Using the estimated $\hat{q}_{s,f}$, $\hat{\mathbf{A}}_n$ and $\hat{\lambda}$ are updated as:

$$\hat{\mathbf{A}}_n(i, f) \leftarrow \frac{\sum_{s=1}^S \hat{q}_{s,f} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)})}{\sum_{i'=1}^{I_n} \sum_{s=1}^S \hat{q}_{s,f} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i')})}, \quad \forall i, f \quad (22a)$$

$$\hat{\lambda}(f) \leftarrow \frac{\sum_{s=1}^S \hat{q}_{s,f}}{\sum_{f'=1}^F \sum_{s=1}^S \hat{q}_{s,f'}}, \quad \forall f. \quad (22b)$$

The proof of convergence for the above defined iterates is inspired by the convergence proof of an EM algorithm handling the crowdsourcing problem in [10]. There, the EM algorithm assumes a uniform latent distribution, i.e., $\lambda(f) = 1/F$ while formulating the maximum likelihood function. In our case, we do not assume uniform prior for λ .

First, we define certain events as below:

$$\mathcal{E}_1 : \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log\left(\frac{\mathbf{A}_n(i, f_s)}{\mathbf{A}_n(i, f)}\right) \geq N\overline{D}_1/2, \quad \text{for all } s \text{ and } f \neq f_s$$

$$\mathcal{E}_2 : \left| \sum_{s=1}^S \mathbb{1}(f_s = f) \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) - S\lambda(f)p\mathbf{A}_n(i, f) \right| \leq St_{nif}, \quad \text{for all } n, i, f,$$

$$\mathcal{E}_3 : \left| \sum_{s=1}^S \mathbb{1}(f_s = f) \mathbb{1}(\mathbf{d}_s(n) \neq 0) - S\lambda(f)p \right| \leq St_{nif}, \quad \text{for all } n, i, f,$$

$$\mathcal{E}_4 : \left| \sum_{s=1}^S \mathbb{1}(f_s = f) - S\lambda(f) \right| \leq Sc_f, \quad \text{for all } f,$$

where $\mathbf{d}_s(n) \neq 0$ represents that n -th RV is observed with any value from its alphabet set in the s -th sample, and $t_{nif} > 0, c_f > 0$ are scalars which will be assigned specific value later in the proof.

First, we consider the E-step update given in (21). The parameter $\hat{q}_{s,f}$ can be bounded using the following lemma:

Lemma 2 Assume that the event \mathcal{E}_1 happens and also assume that $\mathbf{A}_n, \boldsymbol{\lambda}$ and the initial estimates satisfy $|\hat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f)| \leq \delta_1$, $\mathbf{A}_n(i, f) \geq \rho_1$, $|\hat{\boldsymbol{\lambda}}(f) - \boldsymbol{\lambda}(f)| \leq \delta_2$ and $\boldsymbol{\lambda}(f) \geq \rho_2$ for all n, i, f . Then, if $\hat{q}_{s,f}$ is updated by (21), the below holds:

$$|\hat{q}_{s,f} - \mathbb{1}(f_s = f)| \leq \exp\left(-\left(N\bar{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N\left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)}\right)\right) + \log(F)\right), \forall f, s. \quad (23)$$

The proof of Lemma 2 is given in Sec. E Next lemma shows that once $\hat{q}_{s,f}$ updated in E-step is bounded, the subsequent M-step updates to \mathbf{A}_n 's and $\boldsymbol{\lambda}$ are bounded.

Lemma 3 Assume that $\mathcal{E}_2 \cap \mathcal{E}_3 \cap \mathcal{E}_4$ holds. Suppose $\hat{q}_{s,f}$ updated by (21) satisfies the following:

$$|\hat{q}_{s,f} - \mathbb{1}(f_s = f)| \leq \beta, \forall f, s, \quad (24)$$

where $\beta > 0$ is a scalar. Then $\hat{\mathbf{A}}_n$ and $\hat{\boldsymbol{\lambda}}$ updated by (22) are bounded by:

$$|\hat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f)| \leq \frac{2St_{nif} + 2S\beta}{S\boldsymbol{\lambda}(f)p - St_{nif} - S\beta}, \quad (25a)$$

$$|\hat{\boldsymbol{\lambda}}(f) - \boldsymbol{\lambda}(f)| \leq \frac{Sc_f + S\beta + SF\beta}{S - SF\beta}. \quad (25b)$$

The proof of Lemma 3 is given in Sec. F

Next, we show that the estimation accuracy bounds for $\hat{\mathbf{A}}_n$ and $\hat{\boldsymbol{\lambda}}$ given in Lemma 3 are less than or equal to the initial estimation accuracy. For this, we have the following lemma:

Lemma 4 Assume that $\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3 \cap \mathcal{E}_4$ happens. Also assume that $\mathbf{A}_n, \boldsymbol{\lambda}$ and the initial estimates satisfy $|\hat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f)| \leq \delta_1 := \frac{4}{\rho_1(4+\bar{D})}$, $\mathbf{A}_n(i, f) \geq \rho_1$, $|\hat{\boldsymbol{\lambda}}(f) - \boldsymbol{\lambda}(f)| \leq \delta_2 := \frac{4}{\rho_2(4+N\bar{D})}$, $\boldsymbol{\lambda}(f) \geq \rho_2$ for all n, i, f and $\bar{D} \geq \max\left\{\frac{8-4\rho_1^2}{\rho_1^2}, \frac{8-4\rho_2^2}{N\rho_2^2}\right\}$. Suppose that the following holds $\forall g \in \{\{t_{nif}\}_{n,i,f}, \{c_f\}_f\}$:

$$2 \exp\left(-\frac{N\bar{D}}{2} + \log(F)\right) \leq g \leq \frac{p\rho_2}{8F} \min\left(\frac{4}{\rho_1(4+\bar{D})}, \frac{4}{\rho_2(4+N\bar{D})}\right). \quad (26)$$

Then, by updating the parameters using (21) and (22) at least once (i.e., after running the EM algorithm for at least one iteration), we have the following:

$$\left|\hat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f)\right| \leq \frac{4t_{nif}}{\boldsymbol{\lambda}(f)p} \leq \delta_1, \quad (27a)$$

$$\left|\hat{\boldsymbol{\lambda}}(f) - \boldsymbol{\lambda}(f)\right| \leq 8Fc_f \leq \delta_2. \quad (27b)$$

The proof of Lemma 4 is given in Sec G.

Next step is to find out the probabilities that the bounds in (27) hold true. Specifically, we need to characterize the probability for the event $\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3 \cap \mathcal{E}_4$ to happen and the conditions under which (26) holds. Theorem 4 in [10] characterizes the probabilities of the occurrence of events \mathcal{E}_1 , \mathcal{E}_2 and \mathcal{E}_3 . Specifically, we get the following results:

$$\Pr(\mathcal{E}_1) \geq 1 - SF \exp\left(-\frac{N\bar{D}_1}{33 \log(1/\rho_1)}\right), \quad (28a)$$

$$\Pr(\mathcal{E}_2) \geq 1 - \sum_{n=1}^N \sum_{f=1}^F \sum_{i=1}^{I_n} 2 \exp\left(-\frac{St_{nif}^2}{3p\boldsymbol{\lambda}(f)}\right), \quad (28b)$$

$$\Pr(\mathcal{E}_3) \geq 1 - \sum_{n=1}^N \sum_{f=1}^F \sum_{i=1}^{I_n} 2 \exp\left(-\frac{St_{nif}^2}{3p\boldsymbol{\lambda}(f)}\right). \quad (28c)$$

In order to characterize $\Pr(\mathcal{E}_4)$, we observe that $\sum_{s=1}^S \mathbb{1}(f_s = f)$ is sum of i.i.d. Bernoulli random variables with mean $S\lambda(f)$. Therefore, using the Chernoff bound, we have

$$\Pr\left(\left|\sum_{s=1}^S \mathbb{1}(f_s = f) - S\lambda(f)\right| \geq Sc_f\right) \leq 2 \exp(-Sc_f^2/(3\lambda(f))), \quad \forall f. \quad (29)$$

By taking the union bound over all $f \in \{1, \dots, F\}$, we obtain

$$\Pr(\mathcal{E}_4) \geq 1 - \sum_{f=1}^F 2 \exp(-Sc_f^2/(3\lambda(f))). \quad (30)$$

Summing the probability bounds for $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3$ given by (28) and for \mathcal{E}_4 given by (30), one can see that $\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3 \cap \mathcal{E}_4$ holds with probability at least

$$1 - SI \exp\left(\frac{N\bar{D}_1}{33 \log(1/\rho_1)}\right) - \sum_{n=1}^N \sum_{f=1}^F \sum_{i=1}^{I_n} 4 \exp\left(-\frac{St_{nif}^2}{3p\lambda(f)}\right) - \sum_{f=1}^F 2 \exp\left(-\frac{Sc_f^2}{3\lambda(f)}\right).$$

To ensure that the estimation error bounds for \mathbf{A}_n and λ given by (27) hold with probability greater than $1 - \epsilon$, the following conditions has to be satisfied simultaneously:

$$N \geq \frac{33 \log(1/\rho_1) \log(3SF/\epsilon)}{\bar{D}_1} \quad (31)$$

$$S \geq \frac{3p\lambda(f) \log(12NFI/\epsilon)}{t_{nif}^2} \quad (32)$$

$$S \geq \frac{3\lambda(f) \log(6F/\epsilon)}{c_f^2}. \quad (33)$$

We can assign specific values to t_{nif} and c_f such that the above conditions are satisfied. Let

$$t_{nif} := \sqrt{\frac{3p\lambda(f) \log(12NFI/\epsilon)}{S}}, \quad (34a)$$

$$c_f := \sqrt{\frac{3\lambda(f) \log(12NFI/\epsilon)}{S}}. \quad (34b)$$

By this selection of t_{nif} and c_f , the conditions in (32) and (33) hold. To enforce the condition (26), the following equalities have to hold:

$$\begin{aligned} \sqrt{\frac{3p\lambda(f) \log(12NFI/\epsilon)}{S}} &\geq 2 \exp\left(-\frac{N\bar{D}}{2} + \log(F)\right) \\ \sqrt{\frac{3\lambda(f) \log(12NFI/\epsilon)}{S}} &\leq \frac{p\rho_2}{8F} \min\left(\frac{4}{\rho_1(4 + \bar{D})}, \frac{4}{\rho_2(4 + N\bar{D})}\right) = \frac{p\rho_2\delta_{\min}}{8F}, \end{aligned}$$

where $\delta_{\min} = \min(\delta_1, \delta_2)$. The above can be implied by the following:

$$N \geq \frac{4 \log(2SF^2/(3p\rho_2 \log(12NFI/\epsilon)))}{\bar{D}} \quad (35)$$

$$S \geq \frac{192F^2 \log(12NFI/\epsilon)}{p^2 \rho_2^2 \delta_{\min}^2}, \quad (36)$$

where we have used $1 \geq \lambda(f) \geq \rho_2$.

Using the inequality $\log x > 1 - \frac{1}{x}$, $x > 0$, we can express the condition (35) as

$$N \geq \frac{4 \log(2SF^2/(3p\rho_2(1 - \epsilon/(12NFI))))}{\bar{D}}$$

Using the fact that $1 - \epsilon/(12NFI) > \epsilon/2$, we can further write the above condition as follows:

$$N \geq \frac{4 \log(4SF^2/(3p\rho_2\epsilon))}{\bar{D}}. \quad (37)$$

Combing the two conditions (31) and (37), we have

$$N \geq \max \left(\frac{33 \log(3SF/\epsilon)}{\rho_1 \bar{D}_1}, \frac{4 \log(4SF^2/(3p\rho_2\epsilon))}{\bar{D}} \right), \quad (38)$$

where we have used the fact that $\log(1/\rho_1) \leq (1/\rho_1) - 1 < 1/\rho_1$.

To summarize, if (36) and (38) hold and t_{nif} and c_f are chosen to be as in (34), then, with probability at least $1 - \epsilon$, the following inequalities hold by Lemma 4:

$$\begin{aligned} |\widehat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f)|^2 &\leq \frac{16t_{nif}^2}{p^2 \boldsymbol{\lambda}(f)^2} \leq \frac{48 \log(12NFI/\epsilon)}{Sp \boldsymbol{\lambda}(f)} \\ |\widehat{\boldsymbol{\lambda}}(f) - \boldsymbol{\lambda}(f)|^2 &\leq 64F^2 c_f^2 \leq \frac{192F^2 \boldsymbol{\lambda}(f) \log(12NFI/\epsilon)}{S} \end{aligned}$$

This completes the proof.

APPENDIX C PROOF PROPOSITION 1

Let $\mathbf{d}_s \in \mathbb{R}^N$ denote the s th realization of the joint PMF $\Pr(Z_1, \dots, Z_N)$. Recall that p is the probability of an RV being observed in any realization. Let $\{z_n^{(1)}, \dots, z_n^{(I_n)}\}$ denote the alphabet set of Z_n . For simplicity, we assume that 0 does not belong to the alphabets of Z_1, \dots, Z_N , and we use the notation $\mathbf{d}_s(j) = 0$ to represent that ' Z_j is not observed in the s th realization'.

For S realizations of the joint PMF, i.e., $\{\mathbf{d}_s\}_{s=1}^S$, the sample averaging expressions for estimating \mathbf{X}_{jk} is defined as follows:

$$\widehat{\mathbf{X}}_{jk}(i_j, i_k) = \frac{1}{|\mathcal{S}_{jk}|} \sum_{s \in \mathcal{S}_{jk}} \mathbb{I}[\mathbf{d}_s(j) = z_j^{(i_j)}, \mathbf{d}_s(k) = z_k^{(i_k)}],$$

where the indicator function $\mathbb{I}[E]$ is one if the event E happens and zero otherwise; e.g.,

$$\mathbb{I}[\mathbf{d}_s(j) = z_j^{(i_j)}, \mathbf{d}_s(k) = z_k^{(i_k)}] = \begin{cases} 1, & \mathbf{d}_s(j) = z_j^{(i_j)}, \mathbf{d}_s(k) = z_k^{(i_k)} \\ 0, & \text{o.w.} \end{cases}$$

In addition, $\mathcal{S}_{jk} = \{s \mid \mathbb{I}[\mathbf{d}_s(j) \neq 0, \mathbf{d}_s(k) \neq 0]\}$.

Let us construct a random variable $V_{j,s}$, where $V_{j,s} = 1$ if Z_j is observed in \mathbf{d}_s ; otherwise $V_{j,s} = 0$. With this definition, we can see that the parameter S_{jk} in Lemma 5 is the sum of S i.i.d. Bernoulli random variable given by

$$S_{jk} = \sum_{s=1}^S \mathbb{I}[V_{j,s} = 1 \text{ and } V_{k,s} = 1], \quad (39)$$

with mean $\mathbb{E}[S_{jk}] = Sp^2$, since $V_{j,s}$ and $V_{k,s}$ are independent.

In order to characterize the random variable S_{jk} , we can use Chernoff lower tail bound such that for $0 < t < 1$,

$$\Pr(S_{jk} \leq (1-t)Sp^2) \leq e^{-Sp^2 t^2/2}. \quad (40)$$

Eq. (40) also implies that

$$\Pr(S_{jk} \geq (1-t)Sp^2) \geq 1 - e^{-Sp^2 t^2/2}. \quad (41)$$

Combining Lemma 5 and (41), we have

$$\begin{aligned} \Pr \left[\|\widehat{\mathbf{X}}_{jk} - \mathbf{X}_{jk}\|_F \leq \frac{(1 + \sqrt{\log(1/\delta)})}{\sqrt{(1-t)Sp^2}} \right] \\ = \Pr \left[\|\widehat{\mathbf{X}}_{jk} - \mathbf{X}_{jk}\|_F \leq \frac{(1 + \sqrt{\log(1/\delta)})}{\sqrt{S_{jk}}}, S_{jk} \geq (1-t)Sp^2 \right] \\ \geq 1 - \delta - e^{-Sp^2 t^2/2}, \end{aligned}$$

where we have applied the De Morgan's law and the union bound to obtain the last inequality.

Setting $t = 1/2$, we have

$$\begin{aligned} \Pr \left[\|\widehat{\mathbf{X}}_{jk} - \mathbf{X}_{jk}\|_F \leq \frac{\sqrt{2}(1 + \sqrt{\log(1/\delta)})}{p\sqrt{S}} \right] \\ \geq 1 - \delta - e^{-Sp^2/8}. \end{aligned} \quad (42)$$

It follows that if $p^2 \geq \frac{8}{S} \log(1/\delta)$, the right hand side of (42) is greater than $1 - 2\delta$.

APPENDIX D
PROOF OF LEMMA 1

From the assumption that $\overline{\mathbf{H}}$ satisfies ε -separability, there exists a set of indices $\Lambda = \{l_1, \dots, l_F\}$ such that

$$\overline{\mathbf{H}}(\Lambda, :) = \mathbf{I}_F + \mathbf{E},$$

$\mathbf{E} \in \mathbb{R}^{F \times F}$ is the error matrix with $\|\mathbf{E}(l, :)\|_2 \leq \varepsilon$. and \mathbf{I}_F is the identity matrix of size $F \times F$

Now we can write the normalized data model given in (11) as

$$\begin{aligned} \overline{\mathbf{X}} &= \overline{\mathbf{W}}\overline{\mathbf{H}}^\top + \overline{\mathbf{N}} \\ &= \overline{\mathbf{W}}[\mathbf{I}_F + \mathbf{E}^\top, (\mathbf{H}^*)^\top] + \overline{\mathbf{N}} \\ &= \overline{\mathbf{W}}[\mathbf{I}_F, (\mathbf{H}^*)^\top] + [\overline{\mathbf{W}}\mathbf{E}^\top, \mathbf{0}] + \overline{\mathbf{N}}, \end{aligned}$$

where the zero matrix $\mathbf{0}$ has the same dimension as that of \mathbf{H}^* . By defining the noise matrix $\mathbf{N} \in \mathbb{R}^{L \times K}$ such that $\mathbf{N} := [\overline{\mathbf{W}}\mathbf{E}^\top, \mathbf{0}] + \overline{\mathbf{N}}$, we have $\overline{\mathbf{X}} = \overline{\mathbf{W}}[\mathbf{I}_F \quad (\mathbf{H}^*)^\top] + \mathbf{N}$. Then, for any $q \in \{1, \dots, K\}$, the following inequality holds:

$$\begin{aligned} \|\mathbf{N}(:, q)\|_2 &\leq \|\overline{\mathbf{W}}\|_2 \|\mathbf{E}(q, :)\|_2 + \|\overline{\mathbf{N}}(:, q)\|_2 \\ &\leq \sigma_{\max}(\overline{\mathbf{W}})\varepsilon + \varphi, \end{aligned} \quad (43)$$

where the first inequality is by the Cauchy-Schwartz inequality and by the assumptions in the lemma.

Then, we invoke Lemma 3 to characterize the estimation accuracy of $\overline{\mathbf{W}}$. Combining (43) and Lemma 3, we get the final result of the lemma given by (14) if

$$(\sigma_{\max}(\overline{\mathbf{W}})\varepsilon + \varphi) \leq \sigma_{\min}(\overline{\mathbf{W}}) \min\left(\frac{1}{2\sqrt{F-1}}, \frac{1}{4}\right) (1 + 80\kappa^2(\overline{\mathbf{W}}))^{-1}.$$

APPENDIX E
PROOF OF LEMMA 2

Consider the update to $\hat{q}_{s,f}$ in the E-step given by (21). For any $f \neq f_s$,

$$\begin{aligned} \hat{q}_{s,f} &\leq \frac{\exp(\log(\hat{\lambda}(f)) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log(\hat{\mathbf{A}}_n(i, f)))}{\exp(\log(\hat{\lambda}(f_s)) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log(\hat{\mathbf{A}}_n(i, f_s)))} \\ &= \frac{\hat{\lambda}(f) \prod_{n=1}^N \prod_{i=1}^{I_n} (\hat{\mathbf{A}}_n(i, f))^{\mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)})}}{\hat{\lambda}(f_s) \prod_{n=1}^N \prod_{i=1}^{I_n} (\hat{\mathbf{A}}_n(i, f_s))^{\mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)})}} \\ &= \frac{\hat{\lambda}(f)}{\hat{\lambda}(f_s)} \prod_{n=1}^N \prod_{i=1}^{I_n} \left(\frac{\hat{\mathbf{A}}_n(i, f)}{\hat{\mathbf{A}}_n(i, f_s)} \right)^{\mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)})} \\ &= 1 / \exp \left(\underbrace{\log \left(\hat{\lambda}(f_s) / \hat{\lambda}(f) \right) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log \left(\hat{\mathbf{A}}_n(i, f_s) / \hat{\mathbf{A}}_n(i, f) \right)}_{B_f} \right). \end{aligned} \quad (44)$$

Then it follows that

$$\begin{aligned} B_f &= \log \left(\hat{\lambda}(f_s) / \hat{\lambda}(f) \right) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log \left(\hat{\mathbf{A}}_n(i, f_s) / \hat{\mathbf{A}}_n(i, f) \right) \\ &= \log \left(\frac{\lambda(f_s)}{\lambda(f)} \frac{\hat{\lambda}(f_s)}{\hat{\lambda}(f)} \frac{\lambda(f)}{\lambda(f_s)} \right) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log \left(\frac{\mathbf{A}_n(i, f_s) \hat{\mathbf{A}}_n(i, f)}{\mathbf{A}_n(i, f) \hat{\mathbf{A}}_n(i, f_s)} \right) \\ &= \log \left(\frac{\lambda(f_s)}{\lambda(f)} \right) + \log \left(\frac{\hat{\lambda}(f_s)}{\hat{\lambda}(f)} \right) - \log \left(\frac{\lambda(f)}{\lambda(f_s)} \right) + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \log \left(\frac{\mathbf{A}_n(i, f_s)}{\mathbf{A}_n(i, f)} \right) \\ &\quad + \sum_{n=1}^N \sum_{i=1}^{I_n} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \left[\log \left(\frac{\hat{\mathbf{A}}_n(i, f_s)}{\hat{\mathbf{A}}_n(i, f)} \right) - \log \left(\frac{\hat{\mathbf{A}}_n(i, f)}{\hat{\mathbf{A}}_n(i, f_s)} \right) \right]. \end{aligned} \quad (45)$$

To proceed, we can bound all the terms in (45). First we have

$$\begin{aligned} \log \left(\frac{\widehat{\boldsymbol{\lambda}}(f_s)}{\boldsymbol{\lambda}(f_s)} \right) - \log \left(\frac{\widehat{\boldsymbol{\lambda}}(f)}{\boldsymbol{\lambda}(f)} \right) &\geq \log(\rho_2 - \delta_2) + \log(\rho_2) \\ &= \log(\rho_2(\rho_2 - \delta_2)) \\ &\geq 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)}, \end{aligned} \quad (46)$$

where the first inequality uses the results $\widehat{\boldsymbol{\lambda}}(f) \geq \rho_2 - \delta_2$, $\boldsymbol{\lambda}(f) \geq \rho_2$ and the facts $\widehat{\boldsymbol{\lambda}}(f), \boldsymbol{\lambda}(f) \leq 1$. The last inequality is due to the fact that $\log(x) > 1 - \frac{1}{x}$ for $x > 0$.

Similarly, we can bound

$$\log \left(\frac{\widehat{\mathbf{A}}_n(i, f_s)}{\mathbf{A}_n(i, f_s)} \right) - \log \left(\frac{\widehat{\mathbf{A}}_n(i, f)}{\mathbf{A}_n(i, f)} \right) \geq 1 - \frac{1}{\rho_1(\rho_1 - \delta_1)}. \quad (47)$$

Assuming that \mathcal{E}_1 happens, we have

$$\begin{aligned} B_f &\geq \frac{N\overline{D}_2}{2} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + \frac{N\overline{D}_1}{2} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right) \\ &= \frac{N(\overline{D}_1 + \overline{D}_2)}{2} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right), \end{aligned}$$

where the first inequality is obtained by using the definitions of \overline{D}_2 and event \mathcal{E}_1 , equations (46) and (47). Defining $\overline{D} := \frac{\overline{D}_1 + \overline{D}_2}{2}$, we can have

$$B_f \geq N\overline{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right). \quad (48)$$

Combining (48) with (44), we have for every $f \neq f_s$,

$$\widehat{q}_{s,f} \leq 1/\exp B_f \leq \exp \left(- \left(N\overline{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right) \right) \right). \quad (49)$$

Using (49) and $\sum_{f=1}^F \widehat{q}_{s,f} = 1$, we have

$$\widehat{q}_{s,f_s} = 1 - \sum_{f \neq f_s} \widehat{q}_{s,f} \geq 1 - F \exp \left(- \left(N\overline{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right) \right) \right). \quad (50)$$

The inequalities in (49) and (50) can be summarized as follows:

$$|\widehat{q}_{s,f} - \mathbb{1}(f_s = f)| \leq \exp \left(- \left(N\overline{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right) \right) \right) + \log(F), \quad \forall f, s, \quad (51)$$

where we have used the fact that $\exp(\cdot)$ is an increasing function.

APPENDIX F PROOF OF LEMMA 3

The first result given in (27a) follows similar steps as in Lemma 9 from [10] which is detailed below using the notations in our case:

According to the M-step update (22), we can write

$$\widehat{\mathbf{A}}_n(i, f) = \frac{A}{B},$$

where $A := \sum_{s=1}^S \widehat{q}_{s,f} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)})$ and $B := \sum_{i'=1}^{I_n} \sum_{s=1}^S \widehat{q}_{s,f} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i')})$.

Assuming that the event \mathcal{E}_2 holds true, we can have

$$\begin{aligned} |A - S\boldsymbol{\lambda}(f)p\mathbf{A}_n(i, f)| &\leq \left| \sum_{s=1}^S \mathbb{1}(f_s = f) \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) - S\boldsymbol{\lambda}(f)p\mathbf{A}_n(i, f) \right| \\ &\quad + \left| \sum_{s=1}^S \widehat{q}_{s,f} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) - \sum_{s=1}^S \mathbb{1}(f_s = f) \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i)}) \right| \\ &\leq St_{ni}f + S\beta, \end{aligned} \quad (52)$$

where the last result is obtained from the definition of \mathcal{E}_2 and (24).

Assuming that the event \mathcal{E}_3 holds true, we can have

$$\begin{aligned}
|B - S\lambda(f)p| &\leq \left| \sum_{s=1}^S \mathbb{1}(f_s = f) \mathbb{1}(\mathbf{d}_s(n) \neq 0) - S\lambda(f)p \right| \\
&\quad + \left| \sum_{i'=1}^{I_n} \sum_{s=1}^S \hat{q}_{s,f} \mathbb{1}(\mathbf{d}_s(n) = z_n^{(i')}) - \sum_{s=1}^S \mathbb{1}(f_s = f) \mathbb{1}(\mathbf{d}_s(n) \neq 0) \right| \\
&\leq St_{nif} + S\beta
\end{aligned} \tag{53}$$

where the last inequality is obtained by assuming that event \mathcal{E}_3 holds true and using (24).

Combining the bounds for A and B, we can get

$$\begin{aligned}
\left| \hat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f) \right| &= \left| \frac{A}{B} - \mathbf{A}_n(i, f) \right| \\
&= \left| \frac{S\lambda(f)p\mathbf{A}_n(i, f) + A - S\lambda(f)p\mathbf{A}_n(i, f)}{S\lambda(f)p + B - S\lambda(f)p} - \mathbf{A}_n(i, f) \right| \\
&= \left| \frac{A - S\lambda(f)p\mathbf{A}_n(i, f) - \mathbf{A}_n(i, f)(B - S\lambda(f)p)}{S\lambda(f)p + B - S\lambda(f)p} \right| \\
&\leq \frac{|A - S\lambda(f)p\mathbf{A}_n(i, f)| + \mathbf{A}_n(i, f) |B - S\lambda(f)p|}{|S\lambda(f)p + B - S\lambda(f)p|} \\
&\leq \frac{2St_{nif} + 2S\beta}{S\lambda(f)p - St_{nif} - S\beta},
\end{aligned}$$

where the last inequality is by the fact that $\mathbf{A}_n(i, f) \leq 1$ and the bounds from (52) and (53).

For the second result, consider the M-step update for $\hat{\lambda}$ given by (22). One can write $\hat{\lambda}(f) = C/D$ where

$$C = \sum_{s=1}^S \hat{q}_{s,f}, \quad D = \sum_{f'=1}^F \sum_{s=1}^S \hat{q}_{s,f'}.$$

Assume that the event \mathcal{E}_4 happens, using (24), we have

$$\begin{aligned}
|C - S\lambda(f)| &\leq \left| \sum_{s=1}^S \mathbb{1}(f_s = f) - S\lambda(f) \right| + \left| \sum_{s=1}^S \hat{q}_{s,f} - \sum_{s=1}^S \mathbb{1}(f_s = f) \right| \\
&\leq Sc_f + S\beta.
\end{aligned}$$

In addition, we have

$$|D - S| \leq \left| \sum_{f'=1}^F \sum_{s=1}^S \mathbb{1}(f_s = f') - S \right| + \left| \sum_{f'=1}^F \sum_{s=1}^S \hat{q}_{s,f'} - \sum_{f'=1}^F \sum_{s=1}^S \mathbb{1}(f_s = f') \right| \leq SF\beta.$$

Combining the bounds for C and D , we obtain

$$\begin{aligned}
|\hat{\lambda}(f) - \lambda(f)| &= \left| \frac{C}{D} - \lambda(f) \right| \\
&= \left| \frac{(C - S\lambda(f)) + S\lambda(f)}{(D - S) + S} - \lambda(f) \right| \\
&= \left| \frac{(C - S\lambda(f)) - \lambda(f)(D - S)}{(D - S) + S} \right| \\
&\leq \frac{Sc_f + S\beta + SF\beta}{S - SF\beta},
\end{aligned}$$

where the last inequality is by using triangle inequality and the fact that $\lambda(f) \leq 1$.

APPENDIX G
PROOF OF LEMMA 4

Consider the below term in (23) from Lemma 2:

$$N\bar{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right) =$$

$$N \left(\underbrace{\frac{\bar{D}}{2} + 1 - \frac{1}{\rho_1(\rho_1 - \delta_1)}}_E \right) + \underbrace{\frac{N\bar{D}}{2} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)}}_F.$$

In order to bound the term E as below

$$E := \frac{\bar{D}}{2} + 1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \geq \frac{\bar{D}}{4},$$

δ_1 has to be bounded such that

$$\begin{aligned} \delta_1 &\leq \rho_1 - \frac{4}{\rho_1(4 + \bar{D})} \\ &= \frac{4\rho_1^2 + \bar{D}\rho_1^2 - 4}{\rho_1(4 + \bar{D})}. \end{aligned} \tag{54}$$

Similarly, in order to bound F as below

$$F := \frac{N\bar{D}}{2} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} \geq \frac{N\bar{D}}{4},$$

δ_2 has to be bounded such that

$$\begin{aligned} \delta_2 &\leq \rho_2 - \frac{4}{\rho_2(4 + N\bar{D})} \\ &= \frac{4\rho_2^2 + N\bar{D}\rho_2^2 - 4}{\rho_2(4 + N\bar{D})}. \end{aligned} \tag{55}$$

Without loss of generality, we can fix values to δ_1 and δ_2 such that the conditions (54) and (55) get satisfied respectively. Since $\bar{D} \geq \max\left\{\frac{8-4\rho_1^2}{\rho_1^2}, \frac{8-4\rho_2^2}{N\rho_2^2}\right\}$, $4\rho_1^2 + \bar{D}\rho_1^2 \geq 8$ and $4\rho_2^2 + N\bar{D}\rho_2^2 \geq 8$ hold true. Therefore

$$\delta_1 := \frac{4}{\rho_1(4 + \bar{D})}, \quad \delta_2 := \frac{4}{\rho_2(4 + N\bar{D})} \tag{56}$$

satisfy the conditions (54) and (55). Note that since $\bar{D} \geq \frac{8-4\rho_1^2}{\rho_1^2}$, $\delta_1 := \frac{4}{\rho_1(4 + \bar{D})} \leq \frac{\rho_1}{2} \leq 1$ and is a valid assignment. Similarly, we can observe that $\delta_2 := \frac{4}{\rho_2(4 + N\bar{D})} \leq \frac{\rho_2}{2} \leq 1$. Therefore, using the assigned values of δ_1 and δ_2 , we have

$$N\bar{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right) \geq \frac{N\bar{D}}{2}. \tag{57}$$

Then, we have

$$\begin{aligned} |\hat{q}_{s,f} - \mathbb{1}(f_s = f)| &\leq \exp \left(- \left(N\bar{D} + 1 - \frac{1}{\rho_2(\rho_2 - \delta_2)} + N \left(1 - \frac{1}{\rho_1(\rho_1 - \delta_1)} \right) \right) + \log(F) \right) \\ &\leq \exp \left(- \frac{N\bar{D}}{2} + \log(F) \right) \\ &\leq t_{\min}/2, \end{aligned} \tag{58}$$

where $t_{\min} = \min\{t_{nif}, \{c_f\}\}$ and the first inequality is by Lemma 2, the second inequality is by using (57) and the third inequality is from the condition (26).

From (26), we have $t_{nif} \leq \frac{p\rho_2}{8} \leq \frac{p\lambda(f)}{8}$. Combining this with Lemma 3, we get

$$\left| \hat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f) \right| \leq \frac{2St_{nif} + 2S\beta}{(7/8)S\lambda(f)p - S\beta},$$

where β is defined such that $|\widehat{q}_{s,f} - \mathbb{1}(f_s = f)| \leq \beta$, $\forall f, s$.

From (58), the scalar β can be assigned a value such that $\beta = t_{\min}/2$. Then,

$$\begin{aligned} \left| \widehat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f) \right| &\leq \frac{2St_{nif} + St_{\min}}{(7/8)S\lambda(f)p - St_{\min}/2} \\ &\leq \frac{3St_{nif}}{(7/8)S\lambda(f)p - St_{\min}/2}, \end{aligned}$$

where the last step is obtained by using $t_{\min} \leq t_{nif}$. By using the condition $t_{\min} \leq \frac{p\rho_2}{8} \leq \frac{p\lambda(f)}{8}$ from (26), we get

$$\left| \widehat{\mathbf{A}}_n(i, f) - \mathbf{A}_n(i, f) \right| \leq \frac{3St_{nif}}{(7/8)S\lambda(f)p - S\lambda(f)p/16} \leq \frac{4t_{nif}}{\lambda(f)p}.$$

Since

$$\frac{4t_{nif}}{\lambda(f)p} \leq \frac{4t_{nif}}{p\rho_2} \leq \frac{1}{2F} \min\left(\frac{4}{\rho_1(4 + \overline{D})}, \frac{4}{\rho_2(4 + N\overline{D})}\right) \leq \frac{4}{\rho_1(4 + \overline{D})} = \delta_1,$$

the estimation error of the newly updated \mathbf{A}_n as given in (22) is at least no worse than the initial estimation error δ_1 .

Next, consider the result (27b) from Lemma 3. Assigning $\beta = t_{\min}/2$,

$$\begin{aligned} |\widehat{\lambda}(f) - \lambda(f)| &\leq \frac{Sc_f + S\beta + SF\beta}{S - SF\beta} \leq \frac{Sc_f + St_{\min}/2 + SFt_{\min}/2}{S - SFt_{\min}/2} \\ &\leq \frac{2Sc_f + St_{\min} + SFt_{\min}}{2S - SFt_{\min}} \leq \frac{2Sc_f + Sc_f + SFc_f}{2S - SFt_{\min}} \\ &\leq \frac{4SFc_f}{2S - SFt_{\min}} \leq \frac{4SFc_f}{S - SFt_{\min}} \leq 8Fc_f, \end{aligned}$$

where we have used the fact that $t_{\min} \leq 1/2F$ according to (26).

The above inequality also implies that

$$|\widehat{\lambda}(f) - \lambda(f)| \leq 8Fc_f \leq p\rho_2 \min\left(\frac{4}{\rho_1(4 + \overline{D})}, \frac{4}{\rho_2(4 + N\overline{D})}\right) \leq \frac{4}{\rho_2(4 + N\overline{D})} = \delta_2.$$

That is, the estimation error of the newly updated λ as given in (22) is at least no worse than the initial estimation error δ_2 .

APPENDIX H LEMMATA

In this section, we present a collection of lemmata that are used in our proofs.

Lemma 5 [13] *Let $\delta \in (0, 1]$ and let $\widehat{\mathbf{X}}_{jk}$ be the empirical average of S_{jk} independent co-occurrences of random variables Z_j and Z_k . Then the following holds*

$$\Pr\left[\|\widehat{\mathbf{X}}_{jk} - \mathbf{X}_{jk}\|_F \leq \frac{1 + \sqrt{\log(1/\delta)}}{\sqrt{S_{jk}}}\right] \geq 1 - \delta, \quad (59)$$

Theorem 3 [4] *Under the described NMF model in Eq. (5), assume that $\|\mathbf{N}(:, q)\|_2 \leq \epsilon$ for all $q \in \{1, \dots, K\}$, if the below holds:*

$$\epsilon \leq \sigma_{\min}(\mathbf{W}) \min\left(\frac{1}{2\sqrt{F-1}}, \frac{1}{4}\right) (1 + 80\kappa^2(\mathbf{W}))^{-1},$$

then SPA identifies an index set $\widehat{\Lambda} = \{\widehat{l}_1, \dots, \widehat{l}_F\}$ such that

$$\max_{1 \leq f \leq F} \min_{\widehat{l}_f \in \widehat{\Lambda}} \left\| \mathbf{W}(:, f) - \widetilde{\mathbf{X}}(:, \widehat{l}_f) \right\|_2 \leq \epsilon(1 + 80\kappa^2(\mathbf{W})), \quad (60)$$

where $\kappa(\mathbf{W}) = \frac{\sigma_{\max}(\mathbf{W})}{\sigma_{\min}(\mathbf{W})}$ is the condition number of \mathbf{W} .

Lemma 6 [9] *Consider a vector $\mathbf{x} \in \mathbb{R}^L$ and the corresponding estimate of the vector $\widehat{\mathbf{x}}$ such that $\widehat{\mathbf{x}} = \mathbf{x} + \mathbf{n}$ where \mathbf{n} represents the noise vector and $\mathbf{x}, \widehat{\mathbf{x}} \geq 0$. Assume that $\|\widehat{\mathbf{x}}\|_1 \geq \eta$ where $\eta > 0$ and $\|\mathbf{n}\|_1 < \|\mathbf{x}\|_1$. Suppose, the vector $\widehat{\mathbf{x}}$ is normalized with respect to its ℓ_1 norm. The normalized version can be represented as*

$$\frac{\widehat{\mathbf{x}}}{\|\widehat{\mathbf{x}}\|_1} = \frac{\mathbf{x}}{\|\mathbf{x}\|_1} + \widetilde{\mathbf{n}},$$

where $\|\tilde{\mathbf{n}}\|_1 \leq \frac{2\|\mathbf{n}\|_1}{\eta}$.

Lemma 7 [9] Let $\rho > 0, \varepsilon > 0$, and assume that the rows of $\mathbf{H} \in \mathbb{R}^{K \times F}$ are generated within the $(F - 1)$ -probability simplex uniformly at random (and then nonnegatively scaled). If the number of rows satisfies

$$K = \Omega \left(\frac{\varepsilon^{-2(F-1)}}{F} \log \left(\frac{F}{\rho} \right) \right), \quad (61)$$

then, with probability greater than or equal to $1 - \rho$, there exist rows of \mathbf{H} indexed by l_1, \dots, l_F such that

$$\|\mathbf{H}(l_f, :) - \mathbf{e}_f^\top\|_2 \leq \varepsilon, \quad f = 1, \dots, F.$$