

# Low Rank Joint Probability Tensor Recovery from 1D Marginals

Jian Vora

Department of Electrical Engineering, IIT Bombay

`jianvora@iitb.ac.in`

December 8, 2020

## 1 Introduction

## 2 Probability Tensor Recovery

- Low Rank of the Joint Probability Tensor
- Canonical Polyadic Decomposition of a Tensor

## 3 Tomography in Density Estimation

- Radon Transform
- Density Estimation from Random Projections

## 4 High-Dimensional Density Estimation

- Prior Literature
- Algorithm for Density Estimation from Marginals
- Empirical Results on Synthetic Data

## 5 References

# Problem Statement

The ever classical problem of given  $N$  random variables  $X_1, X_2, \dots, X_N$ , and their realisations as samples, we need to estimate  $p(X_1, X_2, \dots, X_N)$

Lots of work in this domain like Gaussian mixture models, Bayesian Nets, Markov Random Fields, Variational Autoencoders ...

Can we estimate the joint density from the marginals? Does transforming the data help in the above goal? Do we have any prior knowledge of naturally occurring probability densities?

# Low Rank of the Joint Probability

Q. Do we have any prior knowledge of joint probability densities? Yes!

- For real world data, random variables are reasonably (in)dependent
- For example, in the case of images, we know neighbouring pixels values are highly correlated while value of far off pixels are almost independent.

This implies that the joint probability density  $p(X_1, X_2, \dots, X_N)$  has a *low-rank*. But what does it mean for a tensor to have a low-rank?

We shall now explore a model which helps in capturing this low-rankness of a tensor which will be used to model the joint density throughout this talk.

# Canonical Polyadic Decomposition of a Tensor

$N$ -way tensor  $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3 \times \dots \times I_N}$  admits a Canonical Polyadic Decomposition (CPD) if it can be decomposed as a sum of  $F$  rank-1 tensors. Let  $\mathbf{a} \otimes \mathbf{b}$  denote the outerproduct of two vectors, then formally the CPD model can be stated as follows:

$$\mathbf{X} = \sum_{f=1}^{f=F} \lambda(f) \mathbf{A}_1(:, f) \otimes \mathbf{A}_2(:, f) \otimes \mathbf{A}_3(:, f) \otimes \dots \otimes \mathbf{A}_N(:, f)$$

where  $\lambda \in \mathbb{R}^F$  and  $\mathbf{A}_i \in \mathbb{R}^{I_i \times F}$ .

The joint probability density  $p(X_1, X_2, \dots, X_N)$  is an  $N$ -way tensor (referred to as  $\mathbf{X}$  henceforth), more specifically

$$\mathbf{X}(x_1, x_2, \dots, x_N) = p(X_1 = x_1, X_2 = x_2, \dots, X_N = x_N)$$

# Canonical Polyadic Decomposition of a Tensor(Notation)

- $\{\mathbf{A}_i\}_{i=1}^N$  are referred to as mode latent factors
- Elements of  $\boldsymbol{\lambda}$  are referred to as component mixing weights
- $\mathbf{X}$  is the core pmf tensor which we aim to recover
- $F$  is the rank of the tensor

For a tensor  $\mathbf{X}$  permitting a Canonical Polyadic Decomposition, we refer  $\mathbf{X} = [\boldsymbol{\lambda}, \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_N]$ . Thus there exists a bijection between  $\mathbf{X}$  and the set  $[\boldsymbol{\lambda}, \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_N]$  and recovering the core tensor is equivalent to recovering the mode latent factors and the mixing weights.

This drastically reduces the number of parameters to be estimated which were previously exponential in  $N$ !

# Canonical Polyadic Decomposition of a Tensor

The CPD model can be considered as a naive bayes model with the hidden state  $H$  taking a bounded number of states(= the rank  $F$ )

$$Pr(X_1 = i_1, X_2 = i_2, .., X_N = i_n) = \sum_{f=1}^{f=F} Pr(H = f) \prod_{n=1}^{n=N} Pr(X_n = i_n | H = f)$$

with  $\lambda(f) = Pr(H = f)$  and  $\mathbf{A}_n(i_n, f) = Pr(X_n = i_n | H = f)$

## Uniqueness of the CPD Model

For a tensor  $\mathbf{X}$  of rank  $F$ , we say that a decomposition  $\mathbf{X} = [\lambda, \mathbf{A}_1, \mathbf{A}_2, .., \mathbf{A}_N]$  is essentially unique if the factors are unique up to a common permutation and scaling / counter-scaling of columns of the mode latent factors.

# Uniqueness of the CPD

We clearly do not have the scaling ambiguity if we the non-negativity and sum to one constraints on the columns of the latent factors. Let  $\mathbf{X} = [\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3]$ , where  $\mathbf{A}_1 \in \mathbb{R}^{l_1 \times F}$ ,  $\mathbf{A}_2 \in \mathbb{R}^{l_2 \times F}$ , and  $\mathbf{A}_3 \in \mathbb{R}^{l_3 \times F}$  with  $l_1 \leq l_2 \leq l_3$  without loss of generality, then:

## Lemma on Uniqueness

If  $\min(l_1, l_2) \geq 3$  and  $F \leq l_3$ , then,  $\text{rank}(\mathbf{X}) = F$  and the decomposition of  $\mathbf{X}$  is essentially unique, almost surely, if and only if  $F \leq (l_1 - 1)(l_2 - 1)$

Thus, instead of  $F$  being  $l_1 l_2 l_3$  which would imply a full rank of the core pmf tensor, we can control the low-rankness of the joint density using a single tunable hyperparameter  $F$

[Kargas et.al, IEEE TSP 2018]



# Random Linear Projections

For a given vector  $X \in \mathbb{R}^N$ , we define the following operation as random linear projections:

$$Y = \Phi X$$

where  $\Phi \in \mathbb{R}^{M \times N}$  is composed of entries drawn i.i.d from a standard normal Gaussian or a bernoulli  $\{-1, 1\}$  distribution.

Used a lot in compressed sensing, dimensionality reduction, GMM learning

# Radon and Inverse Radon Transform

The Radon transform of a 2-D function  $f(x, y)$  is defined as:

$$R(r, \theta)[f] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f(x, y) \delta(r - x \cos \theta - y \sin \theta) dx dy$$

The above operation is invertible and is referred to as the inverse radon transform. Filtered backprojection(FBP) is one of the common algorithms used for inverse radon transform defined as:

$$\hat{f}(x, y) = \int_0^\pi \int_{-\infty}^{+\infty} R(v, \theta) H(v) dv d\theta$$

$H(v)$  is the filter transfer function with some common choices being Ram-Lak, Cosine, Shepp Logan.  $\theta$  is the key parameter which determines the direction of projection.

# Density Estimation from Random Projections

Consider a random vector  $X \in \mathbb{R}^N$  and another random vector  $\theta \in \mathbb{R}^N$

$\theta^T X$  is a scalar variable and hence we can estimate  $p(\theta^T X)$  which is a 1D density estimate and is fairly easy to do using either histogramming or kernel density estimate.

We can do this for various angles  $\theta$  and hence we can accumulate a set of 1D densities of the form  $p(\theta_i^T X)$ . But do these mean anything in our original goal of estimating the joint density  $p(X_1, X_2, \dots, X_N)$ ?

We can in fact relate these 1D densities of random projections to the core pmf tensor  $\mathbf{X}$ . It is tomography which comes to our rescue!

# Density Estimation from Random Projections

Let  $\theta = [\theta_1, \theta_2, \dots, \theta_N]$  and  $X = [X_1, X_2, \dots, X_N]$ , then

$$p(\theta^T X = t) = p\left(\sum_i \theta_i X_i = t\right)$$

$$p(\theta^T X = t) = \sum_{a_1} \sum_{a_2} \dots \sum_{a_{N-1}} p(\theta_1 X_1 = a_1, \theta_2 X_2 = a_2, \dots, \theta_N X_N = t - \sum_{i=1}^{N-1} a_i)$$

Thus the 1D density estimate of a projection of a data vector is infact the radon transform of the joint probability tensor taken at an angle  $\theta$ .

Thus performing an inverse radon transform approach on these densities projected along various angles gives us the original core tensor!

[Sullivan et.al, Journal of the Royal Statistical Society, 1993]

We perform preliminary experiments on recovering just 2 dimensional densities from these 1D projections. Inverting the radon transform becomes costly as the dimensionality increases.

**Evaluation Metric:** Jensen-Shannon Divergence

$$JS(P, Q) = \frac{KL(P || (P + Q)/2) + KL(Q || (P + Q)/2)}{2}$$

**Data Generation Process :** We take a truncated 2D gaussian (from -2 to +2 in either directions) with an identity covariance matrix and 0 mean. This is put into 10 bins as cumulative interval measures to get the PDF matrix. PDF of  $\theta_i X$  is generated using **histograms**.

# Empirical Results

Number of samples	JSD Radon	JSD Histogramming
100	0.0608	0.1902
500	0.0478	0.0529
1000	0.0221	0.0203
10000	0.0107	0.0019

**Table:** Reconstruction Comparisons with 1000 projections for Radon

Clearly, the radon approach for density estimation does a better job in the low-sample regime. For the higher samples, we incur some error as the inverse radon operator is lossy in nature.

We now look at reconstruction errors for a variety of Gaussian mixture models as a function of both the number of projections and samples.

# Empirical Results on GMMs

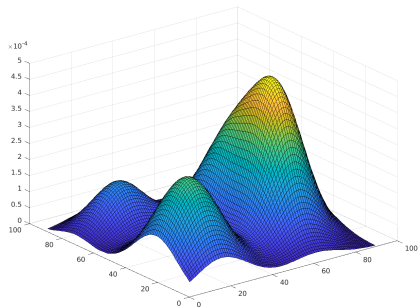


Figure: GMM 1 (4 components)

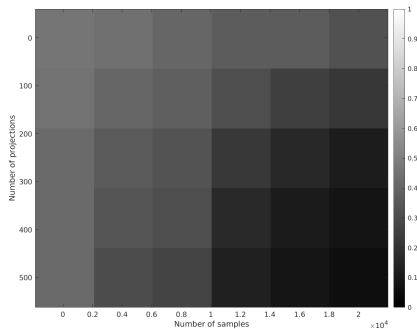


Figure: GMM 1 Errors

# Empirical Results on GMMs

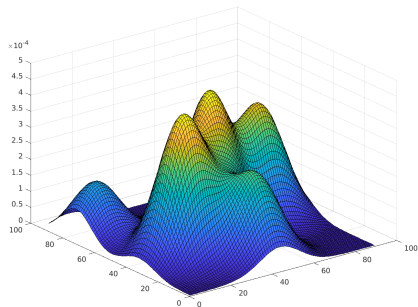


Figure: GMM 2 (9 components)

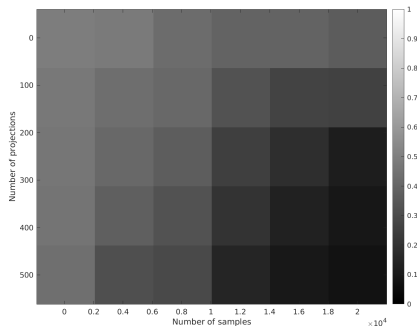


Figure: GMM 2 Errors



# Empirical Results on GMMs

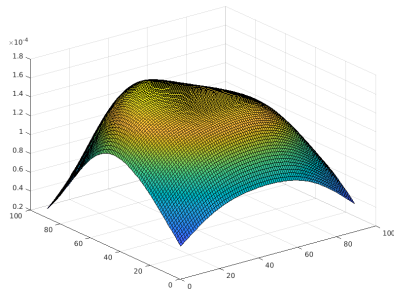


Figure: GMM 3 (4 components)

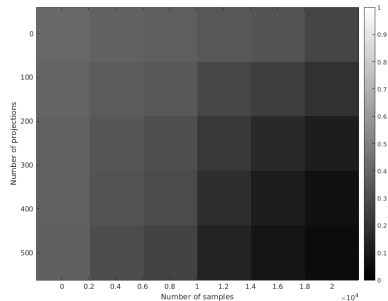


Figure: GMM 3 Errors

# Joint Recovery from 3-way marginals

This was the first work concerning density estimation from marginals by [Kargas et.al, IEEE TSP 2018]

## Procedure: Joint PMF Recovery From Triples

1. Estimate  $\mathbf{X}_{j,k,l}$  from the data
2. Jointly factor  $\mathbf{X}_{j,k,l} = [\lambda, \mathbf{A}_j, \mathbf{A}_k, \mathbf{A}_l]$  using the CPD model of rank  $F$
3. Once all the latent factors are available using the joint factorisation, reconstruct using

$$\mathbf{X} = \sum_{f=1}^{f=F} \lambda(f) \mathbf{A}_1(:, f) \otimes \mathbf{A}_2(:, f) \otimes \mathbf{A}_3(:, f) \otimes \dots \otimes \mathbf{A}_N(:, f)$$

For solving the joint factorisation problem, we use the following least squares optimization problem using alternating minimization:

$$\min_{\{\mathbf{A}_n\}, \lambda} \sum_j \sum_{k>j} \sum_{l>k} \|\mathbf{X}_{j,k,l} - [\lambda, \mathbf{A}_j, \mathbf{A}_k, \mathbf{A}_l]\|_F^2$$

subject to appropriate normalization of  $\lambda$  and  $\mathbf{A}_n, n = 1, 2, 3, \dots, N$

# Joint Recovery using Expectation Maximization

Work by [Yeredor et.al., Signal Process. Letters, 2019] tries to find the parameters by maximising the likelihood of observing the data

$$\log P(y[1], y[2], \dots, y[T]) = \sum_{t=1}^{t=T} \log \sum_{f=1}^{f=F} \lambda_f \prod_{n=1}^{n=N} \mathbf{A}_n(y[t], f)$$

Thus, we formulate the following optimization problem for minimizing the negative log-likelihood:

$$\min_{\{\mathbf{A}_n\}, \lambda} - \sum_{t=1}^{t=T} \log \sum_{f=1}^{f=F} \lambda_f \prod_{n=1}^{n=N} \mathbf{A}_n(y[t], f)$$

subject to appropriate normalization of  $\lambda$  and  $\mathbf{A}_n$ ,  $n = 1, 2, 3, \dots, N$

# Joint Recovery using 2-way marginals

Work by [Ibrahim et.al., <https://arxiv.org/abs/2006.16912>]

Construct the 2-way marginals  $\mathbf{X}_{j,k} = \mathbf{A}_j \mathbf{D}(\lambda) \mathbf{A}_k^T$  using histogramming

Split the  $N$  variables into 2 sets with  $M$  being the index of the split. Thus we have 2 splits namely  $\{1, 2, \dots, M\}$  and  $\{M+1, M+2, \dots, N\}$

Arrange the 2-way marginals in the form a matrix  $\mathbf{A} \in \mathbb{R}^{MI \times (N-M)I}$  where the block  $\mathbf{A}_{i,j} = \mathbf{X}_{i,M+i}$  where NMF can be applied to get the factors

Use the SPA algorithm for NMF, followed by refining the estimates by running EM with the initial conditions as the output of the SPA algorithm

# Algorithm for high-dimensional density estimation

- 1 For  $i, j$  such that  $1 \leq i < j \leq N$ , take random linear projections of data with variables  $X_i$  and  $X_j$ . For each such pair, if we take  $M$  projections, then in total we will have  $M.N(N-1)/2$  scalar numbers.
- 2 Using the above projections, estimate densities of the form  $p(X_i, X_j)$  using the inverse radon transform method.
- 3 Once we have all the pair wise marginals, then reconstruct the core probability tensor using the SPA algorithm.
- 4 Refine the estimates of the mode latent factors by running the EM algorithm on the estimates obtained.

# Experimental Conditions

Mode latent factors were generated with values of each column from a sine wave to ensure the underlying core pmf tensor is smooth. The amplitude, frequency and phase were chosen at random for each column of the mode latent factors.

The mixing ratios were chosen randomly with no such smoothness prior.

200 projections were taken for the radon approach.

$$\text{MSE} = \|\lambda - \hat{\lambda}\|_2^2 / F + \sum_i \|A_i - \hat{A}_i\|_2^2 / NF$$

# Algorithms compared for performance

CNMF-SPA: Using the SPA algorithm to predict the mode latent factors [Ibrahim et.al.]

SPA-EM: Running EM after the CNMF-SPA algorithm to refine the mode latent factors [Ibrahim et.al.]

CNMF-SPA-R: Using the SPA algorithm to predict the mode latent factors but with the initial 2-way marginals coming from the radon transform approach [Our Approach]

SPA-EM-R: Running EM after the CNMF-SPA-R algorithm to refine the mode latent factors [Our Approach]

RAND-EM: Not estimating and marginals, run EM directly with random initializations of the mode latent factors and  $\lambda$  [Yeredor et.al.]

# Selected Empirical Results

	1000	5000	10000	30000
SPA	0.0911	0.0543	0.0459	0.0497
SPA-EM	0.0737	0.0721	0.0358	<b>0.0422</b>
SPA-R	<b>0.0607</b>	<b>0.0469</b>	0.0376	0.0433
SPA-R-EM	0.0668	0.0579	<b>0.0353</b>	0.0452
RAND-EM	0.0713	0.0473	0.0537	0.0598

Table: MSE Comparison for  $F = 10$ ,  $I = 50$ ,  $N = 5$

	1000	5000	10000	30000
SPA	0.0821	0.0963	0.0654	0.0497
SPA-EM	0.0814	0.0518	0.0452	0.0473
SPA-R	0.0611	0.0543	<b>0.0442</b>	0.0475
SPA-R-EM	0.0642	<b>0.0512</b>	0.0494	<b>0.0468</b>
RAND-EM	<b>0.0473</b>	0.0521	0.0489	0.0719

Table: MSE Comparison for  $F = 10$ ,  $I = 200$ ,  $N = 4$



# Selected Empirical Results

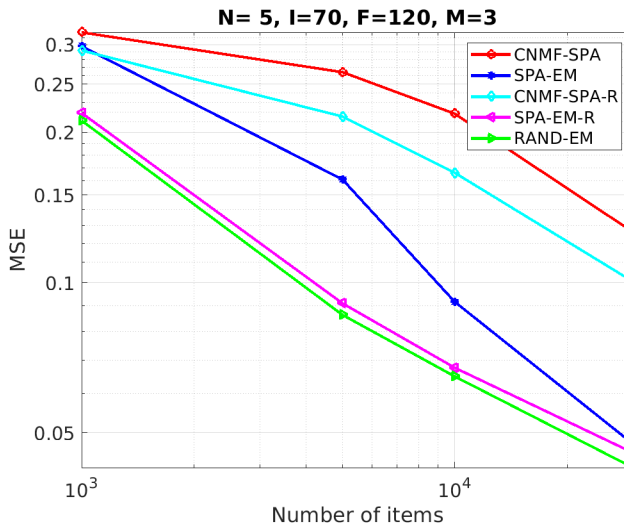


Figure: Errors averaged over 5 different randomly chosen densities

# Selected Empirical Results

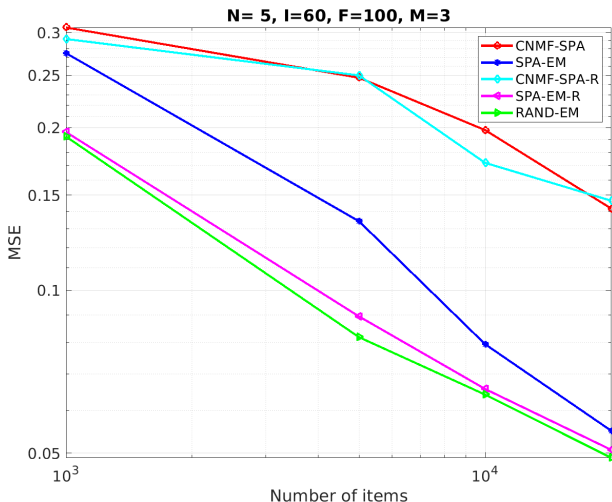


Figure: Errors compared in a different setting

# Empirical Results for Missing Data

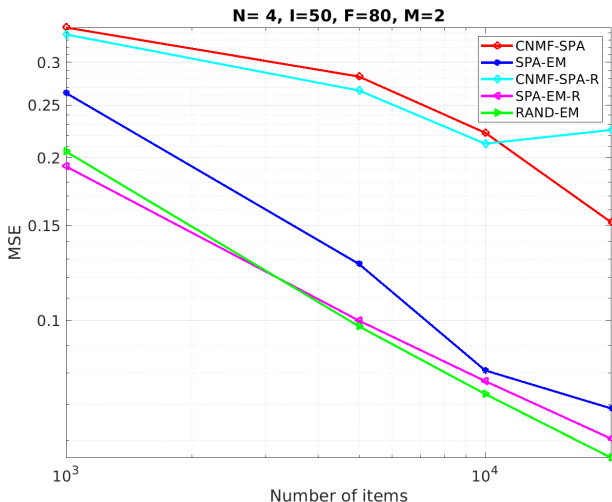


Figure: Errors with **missing data** with probability 0.2

# Empirical Results for Higher Dimensions

	1000	10000	50000
SPA	0.182	0.158	0.083
SPA-EM	0.217	0.126	0.085
SPA-R	0.142	0.117	0.103
SPA-R-EM	<b>0.102</b>	<b>0.081</b>	0.078
RAND-EM	0.113	0.084	<b>0.076</b>

Table: MSE Comparison for  $F = 200$ ,  $I = 90$ ,  $N = 12$

# Future Work

- 1 Try the algorithm on real-world datasets for doing some probabilistic inferences like classification or regression
- 2 Remove the constraint of the rank of the tensor being strictly  $F$  and rather pose the problem as a low-rank tensor recovery where we penalize some metric to penalise the rank of the tensor. This would in general lead to a more expressive model and reduce a tunable parameter from the system.
- 3 Scale the algorithm to very high dimensional data sets (for example 500) which is not possible in the current framework because of the limitations of the SPA algorithm
- 4 Try estimating the actual continuous density by spline interpolation on the mode latent factors.

# 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 lowrank probability mass tensor from partial observations,” IEEE Signal Process. Lett., vol. 26, no. 10, pp. 1551–1555, Oct 2019
- ④ S. Ibrahim, X. Fu, Recovering Joint Probability of Discrete Random Variables from Pairwise Marginals - <https://arxiv.org/abs/2006.16912>
- ⑤ Tensors and Probability: An Intriguing Union - N. Sidiropoulos, N. Kargas, X. Fu, GlobalSIP 2018 Keynote
- ⑥ Finbarr O’Sullivan and Yudi Pawitan, Multidimensional Density Estimation by Tomography, Journal of the Royal Statistical Society. Series B (Methodological) , 1993, Vol. 55, No. 2 (1993), pp. 509-521

# Questions?