Anonymous Authors¹

Abstract

1. Introduction

000

007

009

010

014

015

016

018

019

020

021

022

024

025

026

028

029

030

032

033

034

035

036

037

038

039

041

043

045

046

047

049

050

051

053

Multidimensional array data, a.k.a. a tensor, appears in a huge variety of applications including recommendation systems (Kutty et al., 2012; Adomavicius et al., 2008; Sun et al., 2015), social networks (Sun et al., 2009; Nickel et al., 2011), genomics (Wang et al., 2017), neuroimaging (EEG, fMRI) (Miwakeichi et al.) and signalprocessing (Sidiropoulos et al., 2000; Cichocki et al., 2015). Instead of unfolding those data tensors into matrices where many analysis methods have been proposed, we preserve its inherent multi-modal structure. Studying tensor data while respecting the structure allows us to examine complex interactions among tensor entries. Thereby we can provide extra more interpretation that cannot be addressed by traditional matrix analysis. Also, It has been shown that the tensor preserving analysis improves performance (Zare et al., 2018; Wang & Li, 2018). With those reasons, there is a growing need to develop dimension reduction method without losing multi modal structure. In the line of the attempts, a number of tensor decomposition methods have been proposed in many applications. CANDECOMP/PARAFAC (CP) decomposition was first introduced (Hitchcock, 1928) and revitalized in psychometrics (Harshman, 1970) and in linguistics (Smilde et al., 2004). The Tucker decomposition was proposed in psychometrics (Tucker, 1964; 1966).

Classical tensor completion with those decompositions has treated the entries of data as real-valued. In many cases, However, we encounter data of which the entries are not real-valued but discrete or quantized i.e. binary-valued or multiple-valued. For example, many survey data takes the integer values. To be specific, the data in the Netflix problem has the 3 modes of 'user', 'movie' and 'date of grade'. The entries of the data are the gradings from the users which take integer value from 1 to 5. Also, there are many cases

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute. that the data are quantized in real application. In signal processing, the data are frequently rounded or truncated so that only integer values are available. Also, in graph theory, adjacency matrix can be labeled as from 1 to 3 taking 3 when pairs of vertexes have strong connection and giving 1 when two vertexes have the weak connection according to a given threshold. If we add one more mode on adjacency matrix such as 'context' or 'individual', the data turns into tensor data with 3 integer values.

Therefore, performance improvement can be achieved when the observations are treated as discrete value not as continuous value. In matrix case, there has been many achievements to complete matrix for discrete cases: Models for the case of binary or 1-bit were introduced and studied (Davenport et al., 2012; Bhaskar & Javanmard, 2015). Furthermore, Bhaskar (2016) suggested matrix completion method for general ordinal observations. In tensor case, however, only binary tensor has gotten an attention and achieved performance improvement using binary tensor decomposition methods (Hore et al., 2016; Wang & Li, 2018; Hong et al., 2018; Hu et al., 2018). Accordingly, a general method for the data which has more than 2 ordered label is needed.

We organized this paper as follows. In section 2, we discuss detailed assumptions and descriptions about our probabilistic model in (??). Also, we suggest the estimation method for the latent parameters and related algorithm. In section 3, we provide the statistical properties of the upper, lower bounds and the phase-transition. We then provide the numerical experiments. Our model is applied to real-world data to check validity and performance in section 5. Finally, we wrap up the paper with a discussion.

2. Preliminaries

Let $\mathcal{Y} \in \mathbb{R}^{d_1 \times \cdots \times d_K}$ denote an order-K (d_1,\ldots,d_K) -dimensional tensor. We use y_ω to denote the tensor entry indexed by ω , where $\omega \in [d_1] \times \cdots \times [d_K]$. The Frobenius norm of \mathcal{Y} is defined as $\|\mathcal{Y}\|_F = \sum_\omega y_\omega^2$ and the infinity norm of \mathcal{Y} is defined as $\|\mathcal{Y}\|_\infty = \max_\omega |y_\omega|$. We use $\mathcal{Y}_{(k)}$ to denote the unfolded matrix of size d_k -by- $\prod_{i \neq k} d_k$, obtained by reshaping the tensor along the mode-k. The Tucker rank of \mathcal{Y} is defined as a length-K vector $\mathbf{r} = (r_1,\ldots,r_K)$, where r_k is the rank of matrix $\mathcal{Y}_{(k)}$ for all $k \in [L]$. We say that an event A occurs "with very high

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

066

067

068

069

070

074

075

089

096

100

probability" if $\mathbb{P}(A)$ tends to 1 faster than any polynomial of tensor dimension $d_{\min} = \min\{d_1, \dots, d_K\}$ as $d_{\min} \to \infty$.

We use lower-case letters (e.g., a, b, c) for scalars/vectors, upper-case boldface letters (e.g., A, B, C) for matrices, and calligraphy letters (e.g., A, B, C) for tensors of order three or greater. For ease of notation, we allow the basic arithmetic operators (e.g., \leq , +, -) to be applied to pairs of vectors in an element-wise manner. We use the shorthand [n] to denote the n-set $\{1,\ldots,n\}$ for $n\in N_+$.

3. Model

3.1. Low-rank cumulative model

For the K-mode ordinal tensor $\mathcal{Y} = \llbracket y_\omega \rrbracket \in [L]^{d_1 \times \cdots \times d_K}$, we assume that its entries are realizations of independent multinomial random variables, such that

$$\mathbb{P}(y_{\omega} = l) = f_l(\theta_{\omega}), \quad \omega \in [d_1] \times \cdots \times [d_K].$$
 (1)

In this model, a twice differentiable function $f_l: \mathbb{R} \to [0, 1]$ with $l \in [L]$ is strictly increasing and satisfies $\sum_l f_l(\theta) = 1$ for a fixed θ . The tensor $\Theta = [\![\theta_\omega]\!] \in \mathbb{R}^{d_1 \times \cdots \times d_K}$ is a hidden parameter which we are interested in. We assume the parameter tensor Θ is continuous value and admits a rank $\mathbf{r} = (r_1, \cdots, r_K)$ Tucker decomposition,

$$\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \cdots \times_K \mathbf{M}_N, \tag{2}$$

where $\mathcal{C} \in \mathbb{R}^{r_1 \times \cdots \times r_K}$ is a core tensor and $M_k \in \mathbb{R}^{d_k \times r_k}$, for $n \in [L]$ is a factor matrix. We can get information about the influence of each mode by checking factor matrices.

The tensor model (1) has an equivalent representation as a latent model with K-level quantization. (Davenport et al., 2012; Lan et al., 2014; Bhaskar & Javanmard, 2015; Cai & Zhou, 2013) where $\mathcal{Y} = \llbracket y_{\omega} \rrbracket$ is a quantized value such that,

$$y_{\omega} = \mathcal{Q}(\theta_{\omega} + \epsilon_{\omega}), \quad \omega \in [d_1] \times \cdots [d_K],$$
 (3)

where $\mathcal{E} = \llbracket \epsilon_{\omega} \rrbracket$ is a noise tensor of *i.i.d.* from cumulative distribution function $\Phi(\cdot)$ and the function $\mathcal{Q}: \mathbb{R} \to [L]$ is a quantizer having the following rule.

$$Q(x) = l, \quad \text{if } b_{l-1} < x \le b_l, \quad l \in [L],$$

where $\boldsymbol{b} = (b_1, \dots, b_{L-1})$ is a cutpoints vector such that $-\infty = b_0 < b_1 < \cdots < b_L = \infty$. That is, the entries of observed tensor \mathcal{Y} fall in category l when the associated entries of the latent tensor $\Theta + \mathcal{E}$ fall in the *l*th interval of values. Then we have

$$f_l(\theta_\omega) = \mathbb{P}(y_\omega = l)$$

= $\Phi(b_l - \theta_\omega) - \Phi(b_{l-1} - \theta_\omega).$ (4)

We can diversify our model by the choices of $\Phi(\cdot)$, or equivalently the distribution of \mathcal{E} with given \boldsymbol{b} . The followings are 2 common choices of $\Phi(\cdot)$.

Example 1 (Logistic link/Logistic noise). The logistic model is represented by (1) with $f_l(\theta) = \Phi_{log}(\frac{b_l - \theta}{\sigma})$ – $\Phi_{log}(\frac{b_{l-1}-\theta}{\sigma})$ where $\Phi_{log}(x/\sigma)=(1+e^{-x/\sigma})$. Equivalently, the noise ϵ_{ω} in (3) follows i.i.d. logistic distribution with the scale parameter σ .

Example 2 (Probit link/Gaussian noise). The probit model is represented by (1) with $f_l(\theta) = \Phi_{norm}(\frac{b_l - \theta}{\sigma})$ – $\Phi_{norm}(\frac{b_{l-1}-\theta}{\sigma})$ where Φ_{norm} is the cumulative distribution function of the standard Gaussian. Equivalently, the noise ϵ_{ω} in (3) follows i.i.d. $N(0, \sigma^2)$.

3.2. Rank-constrained likelihood-based estimation

Our goal is to estimate unknown parameter tensor Θ and a cutpoints vector b from observed tensor y using a constrained likelihood approach. With a little abuse of notation, we use Ω to denote either the full index set $\Omega =$ $[d_1] \times \cdots \times [d_K]$ or a random subset indeuced from the subsampling distribution. The log-likelihood function for

$$\mathcal{L}_{\mathcal{Y},\Omega}(\Theta, \boldsymbol{b}) = \sum_{\omega \in \Omega} \Big[\sum_{l \in [L]} \log(f_l(\theta_\omega)) \mathbb{1}_{\{y_\omega = l\}} \Big],$$

where the cutpoints vector \boldsymbol{b} is implicitly contained in the function f_l . Considering the Tucker structure in (2), we have the following constrained optimization problem.

$$\max_{(\Theta, \boldsymbol{b}) \in \mathcal{D} \times \mathcal{B}} \mathcal{L}_{\mathcal{Y}, \Omega}(\Theta, \boldsymbol{b}), \text{ where}$$

$$\mathcal{D} = \{ \Theta \in \mathbb{R}^{d_1 \times \dots \times d_K} : \operatorname{rank}(\Theta) = \boldsymbol{r} \text{ and } \|\Theta\|_{\infty} \le \alpha \}$$

$$\mathcal{B} = \{ \boldsymbol{b} \in \mathbb{R}^{L-1} : b_1 < \dots < b_{L-1} \},$$
(5)

for a given rank $r \in \mathbb{N}_+^K$ and a bound $\alpha \in \mathbb{R}_+$. The search space \mathcal{D} has two constraints on unknown parameter Θ . The first constraint ensures that the unknown parameter Θ admits the Tucker decomposition with rank r. The second constraint makes the entries of Θ bounded by a constant α . This bound condition is a technical assumption to help to recover Θ in the noiseless case. Similar conditions has been imposed in many literatures for the matrix case (Davenport et al., 2012; Bhaskar & Javanmard, 2015; Cai & Zhou, 2013; Bhaskar, 2016). The search space \mathcal{B} makes sure that the probability function f_l in (4) is positive.

3.3. Optimization

In this section, we describe the algorithm to seek the optimizer of (5). The objective function $\mathcal{L}_{\mathcal{V},\Omega}(\Theta, \boldsymbol{b})$ is concave in (Θ, \mathbf{b}) whenever $\Phi(x)$ is log-concave (McCullagh, 1980; Burridge, 1981). However, the feasible set \mathcal{D} is not a convex set, which makes the optimization (5) a non-convex problem. One approach to handle this problem is utilizing the Tucker decomposition and converting optimization into a block-wise convex problem. Let us denote the objective

Algorithm 1 Ordinal tensor decomposition Input: Ordinal tensor $\mathcal{Y} \in [L]^{d_1 \times \cdots \times d_K}$, Rank $r \in \mathbb{N}_+^{K-1}$, Entry-wise bound $\alpha \in \mathbb{R}_+$. Output: $(\hat{\Theta}, \hat{b}) = \arg\max_{(\Theta, b) \in \mathcal{D} \times \mathcal{B}} \mathcal{L}_{\mathcal{Y}, \Omega}(\Theta, b)$. Initialize Core tensor $\mathcal{C}^{(0)}$, Factor matrices $\{M_1^{(0)}, \cdots, M_K^{(0)}\}$, Cut-off points $b^{(0)}$. for $t = 1, 2, \cdots, do$ for $k = 1, 2, \cdots, K$ do Update M_k fixing other blocks. $M_k^{(t+1)} \leftarrow \arg\max_{M_k \in \mathbb{R}^{d_k \times r_k}} \mathcal{L}_{\mathcal{Y}, \Omega}(M_k)$ s.t. $\|\mathcal{C}^{(t)} \times_1 M_1^{(t+1)} \cdots \times_k M_k \cdots \times_K M_K^{(t)}\|_{\infty} \le \alpha$ end for Update \mathcal{C} fixing other blocks. $\mathcal{C}^{(t+1)} \leftarrow \arg\max_{\mathcal{C} \in \mathbb{R}^{r_1 \times \cdots \times r_k}} \mathcal{L}_{\mathcal{Y}, \Omega}(\mathcal{C})$ s.t. $\|\mathcal{C} \times_1 M_1^{(t+1)} \cdots \times_K M_K^{(t+1)}\|_{\infty} \le \alpha$ $\Theta^{(t+1)} \leftarrow \mathcal{C}^{(t+1)} \times_1 M_1^{(t+1)} \cdots \times_K M_K^{(t+1)}$ Update b fixing $\Theta^{(t+1)}$. $b^{(t+1)} \leftarrow \arg\max_{b \in \mathbb{R}^{L-1}} \mathcal{L}_{\mathcal{Y}, \Omega}(\Theta^{(t+1)}, b)$ end for return Θ, b

function as

$$\mathcal{L}(\mathcal{C}, \mathbf{M}_1, \cdots, \mathbf{M}_K, \mathbf{b}) = \mathcal{L}_{\mathcal{Y}, \Omega}(\mathcal{C} \times_1 \mathbf{M}_1 \cdots \times_K \mathbf{M}_K, \mathbf{b}).$$
(6)

From (6), we have K+2 blocks of variables in the objective function, one for the cutpoints vector \boldsymbol{b} , one for the core tensor \mathcal{C} and K for the factor matrices \boldsymbol{M}_K 's. We can change the optimization problem to simple convex problem if any K+1 out of the K+2 blocks being fixed. Therefore, we can alternatively update one block at a time while other blocks being fixed. The algorithm 1 gives the full description.

3.4. Rank selection

Algorithm 1 takes the rank r as an input variable. In practice, the rank r is hardly known. Estimating an appropriate rank r from a given tensor is an important issue. We suggest to use Bayesian Information Criterion(BIC) to choose the rank.

$$\hat{\boldsymbol{r}} = \underset{\boldsymbol{r} \in \mathbb{N}_{+}^{K}}{\min} BIC(\boldsymbol{r})$$

$$= \underset{\boldsymbol{r} \in \mathbb{N}_{+}^{K}}{\min} \left[-2\mathcal{L}_{\mathcal{Y}}(\hat{\boldsymbol{\Theta}}(\boldsymbol{r})) + \left(\prod_{k \in [K]} r_{k} + \sum_{k \in [K]} (d_{k}r_{k} - r_{k}^{2}) \right) \log(\prod_{k \in [K]} d_{k}) \right],$$

where $\hat{\Theta}(r)$ is a maximum likelihood estimater given the rank r. The second part of the summation in (7) is the

number of independent parameters in the model. We select the rank \hat{r} that minimizes BIC value through the grid search method.

4. Real-world Data Applications

In this section, we apply our ordinal tensor decomposition method to two real-world datasets of ordinal tensors. In the first application, we use our model to analyze an ordinal tensor consisting of structural connectivity patterns among 68 brain regions for 136 individuals from Human Connectome Project (HCP) (Geddes, 2016). In the second application, we perform tensor completion from the data with missing values. The data tensor records the ratings from scale 1 to 5 of 42 users to 139 songs on 26 contexts (Baltrunas et al., 2011).

4.1. Human Connectome Project (HCP)

The human connectome project (HCP) is a $68 \times 68 \times 136$ tensor where the first two modes have 68 indices representing brain regions and the last mode has 136 indices meaning individuals. All the individual images were preprocessed following a standard pipeline (Zhang et al., 2018), and the brain was parcellated to 68 regions of interest following the Desikan atlas (Desikan et al., 2006). The tensor entries consist of $\{1,2,3\}$, the strength of fiber connections between 68 brain regions for each of the 136 individuals. First, we apply our ordinal tensor decomposition method with a logistic link function to the HCP data and identify similar brain nodes based on latent parameters. The BIC result suggests r = (23, 23, 8) with $\mathcal{L}_{\mathcal{V}}(\hat{\Theta}, \hat{\boldsymbol{b}}) = -216645.8$.

Also, we compare our ordinal tensor decomposition method with a logistic link function with other methods: the continuous tensor decomposition method, the 1 bit-tensor completion method with probit link function and the 1 bit-tensor completion method with logistic link function. We use 5 folded cross validation method for the comparison. Specifically, we randomly split the tensor entries into 5 similar sized pieces and alternatively use each piece of entries as a test data using other 80% entries as a training data. The entries in the test data are encoded as missing and then predicted based on each method from the training data. We set the rank of the tensor as r = (23, 23, 8) for the tensor decomposition methods which minimizes BIC value. Table 4.1 shows RMSE, MAE and error(false prediction rate), averaged over 5 test set results. We can check that our ordinal tensor decomposition model outperforms the other methods in all criteria.

4.2. Music data with missing values

InCarMusic is a mobile application that offers music recommendation to passengers of cars based on contexts (Bal178

179

191

197

198

199

200

206

214

215

216

217

218

219

METHOD RMSE MAE **ERROR** ORIDNAL TENSOR 0.1503711 0.1502662 0.1502151 DECOMPOSITION CONTI TENSOR 0.1603685 0.1600219 0.1598499 DECOMPOSITION 1BIT-COMPLETION 0.3658390 0.3632757 0.3619940 (PROBIT) 1BIT-COMPLETION 0.45192530.42185130.4068143 (LOGISTIC)

Table 1. Results of comparisons among 4 methods on the HCP data predicting the test data. Four methods are the ordinal tensor decomposition algorithm, the continuous tensor decomposition algorithm and the 1 bit tensor completion method with probit link function and logistic link function. Each method is evaluated by RMSE, MAE, error(false prediction rate).

trunas et al., 2011). Our goal is to get the tensor completion from the $42 \times 139 \times 26$ tensor using the ordinal tensor decomposition and thereby we can offer context-specific music recommendation to users. The tensor entries consist of $\{1,2,3,4,5\}$, the ratings of 42 users to 139 songs on 26 contexts and are encoded as -1 for missing values. The number of missing values is 148904 and the number of available values is 2884.

References

- Adomavicius, G., Mobasher, B., Ricci, F., and Tuzhilin, A. Context-aware recommender systems. *AI Magazine*, 32: 67–80, 2008.
- Baltrunas, L., Kaminskas, M., Ludwig, B., Moling, O., Ricci, F., Aydin, A., Lüke, K.-H., and Schwaiger, R. Incarmusic: Context-aware music recommendations in a car. In *EC-Web*, 2011.
- Bhaskar, S. A. Probabilistic low-rank matrix completion from quantized measurements. *J. Mach. Learn. Res.*, 17: 60:1–60:34, 2016.
- Bhaskar, S. A. and Javanmard, A. 1-bit matrix completion under exact low-rank constraint. 2015 49th Annual Conference on Information Sciences and Systems (CISS), pp. 1–6, 2015.
- Burridge, J. A note on maximum likelihood estimation for regression models using grouped data. 1981.
- Cai, T. and Zhou, W.-X. A max-norm constrained minimization approach to 1-bit matrix completion. *J. Mach. Learn. Res.*, 14:3619–3647, 2013.

- Cichocki, A., Mandic, D. P., Lathauwer, L. D., Zhou, G., Zhao, Q., Caiafa, C. F., and Phan, A. H. Tensor decompositions for signal processing applications: From two-way to multiway component analysis. *IEEE Signal Processing Magazine*, 32:145–163, 2015.
- Davenport, M. A., Plan, Y., van den Berg, E., and Wootters, M. 1-bit matrix completion. *ArXiv*, abs/1209.3672, 2012.
- Desikan, R. S., Ségonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner, R. L., Dale, A. M., Maguire, R. P., Hyman, B. T., Albert, M. S., and Killiany, R. J. An automated labeling system for subdividing the human cerebral cortex on mri scans into gyral based regions of interest. *NeuroImage*, 31:968–980, 2006.
- Geddes, L. Human brain mapped in unprecedented detail. 2016.
- Harshman, R. A. Foundations of the parafac procedure: Models and conditions for an "explanatory" multi-model factor analysis. 1970.
- Hitchcock, F. L. Multiple invariants and generalized rank of a p-way matrix or tensor. *Journal of Mathematics and Physics*, 7(1-4):39–79, 1928. doi: 10.1002/sapm19287139. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/sapm19287139.
- Hong, D., Kolda, T. G., and Duersch, J. A. Generalized canonical polyadic tensor decomposition. *ArXiv*, abs/1808.07452, 2018.
- Hore, V., Viñuela, A., Buil, A., Knight, J. C., McCarthy, M. I., Small, K. S., and Marchini, J. Tensor decomposition for multi-tissue gene expression experiments. In *Nature Genetics*, 2016.
- Hu, Q., Li, G., Wang, P., Zhang, Y., and Cheng, J. Training binary weight networks via semi-binary decomposition. In *ECCV*, 2018.
- Kutty, S., Chen, L., and Nayak, R. A people-to-people recommendation system using tensor space models. In *SAC '12*, 2012.
- Lan, A. S., Studer, C., and Baraniuk, R. G. Matrix recovery from quantized and corrupted measurements. *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4973–4977, 2014.
- McCullagh, P. Regression models for ordinal data. 1980.
- Miwakeichi, F., Martínez-montes, E., Valdés-sosa, B. P. A., Nishiyama, N., Mizuhara, A. H., and A, Y. Y. Decomposing eeg data into space-time-frequency components using parallel factor analysis. *Neuroimage*, pp. 2004.

- Nickel, M., Tresp, V., and Kriegel, H.-P. A three-way model for collective learning on multi-relational data. In *ICML*, 222 2011.
- Sidiropoulos, N. D., Bro, R., and Giannakis, G. B. Parallel factor analysis in sensor array processing. *IEEE Trans. Signal Processing*, 48:2377–2388, 2000.

- Smilde, A. K., Bro, R., and Geladi, P. Multi-way analysis with applications in the chemical sciences. 2004.
- Sun, J., Papadimitriou, S., Lin, C.-Y., Cao, N., Liu, S., and Qian, W. Multivis: Content-based social network exploration through multi-way visual analysis. In *SDM*, 2009.
- Sun, W. W., Lu, J., Liu, H., and Cheng, G. Provable sparse tensor decomposition. 2015.
- Tucker, L. The extension of factor analysis to three-dimensional matrices. In Gulliksen, H. and Frederiksen, N. (eds.), *Contributions to mathematical psychology.*, pp. 110–127. Holt, Rinehart and Winston, New York, 1964.
- Tucker, L. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966. URL https://EconPapers.repec.org/RePEc: spr:psycho:v:31:y:1966:i:3:p:279-311.
- Wang, M. and Li, L. Learning from binary multiway data: Probabilistic tensor decomposition and its statistical optimality. *ArXiv*, abs/1811.05076, 2018.
- Wang, M., Fischer, J., and Song, Y. S. Three-way clustering of multi-tissue multi-individual gene expression data using constrained tensor decomposition. 2017.
- Zare, A., Ozdemir, A., Iwen, M. A., and Aviyente, S. Extension of pca to higher order data structures: An introduction to tensors, tensor decompositions, and tensor pca. *Proceedings of the IEEE*, 106:1341–1358, 2018.
- Zhang, Z., Descoteaux, M., Zhang, J., Girard, G., Chamberland, M., Dunson, D. B., Srivastava, A., and Zhu, H. Mapping population-based structural connectomes. *NeuroImage*, 172:130–145, 2018.