
Exponential Family Tensor Regression with Covariates on Multiple Modes

Abstract

Higher-order tensors have recently received increasing attention in many fields across science and engineering. Here, we present an exponential family of tensor-response regression models that incorporate covariates on multiple modes. Such problems are common in neuroimaging, network modeling, and spatial-temporal analysis. We propose a rank-constrained estimator and establish the theoretical accuracy guarantees. Unlike earlier methods, our approach allows covariates from multiple tensor modes whenever available. An efficient alternating updating algorithm is further developed. Our proposal handles a broad range of data types, including continuous, count, and binary observations. We apply the method to multi-relational social network data and diffusion tensor imaging data from human connection project. Our approach identifies the key global connectivity pattern and pinpoints the local regions that are associated with covariates.

1. Introduction

Many contemporary scientific and engineering studies collect multi-way array data, a.k.a. tensors, accompanied by additional covariates. One example is neuroimaging analysis (??), in which the brain connectivity networks are collected from a sample of individuals. Researchers are often interested in identifying connection edges that are affected by individual characteristics such as age, gender, and disease status (see Figure 1a). Another example is in the field of network analysis (??). A typical social network consists of nodes that represent people and edges that represent friendships. In addition, features on nodes and edges are often available, such as people's personality and demographic location. It is of keen scientific interest to identify the variation in the connection patterns (e.g., transitivity, community) that can be attributable to the node features.

This paper presents a general treatment to these seemingly different problems. We formulate the learning task as a regression problem, with tensor observation serving as a response, and the node features and/or their interactions forming the predictor. Figure 1b illustrates the general set-up we consider. The regression approach allows the identification

of variation in the data tensor that is explained by the covariates. In contrast to earlier work (??), our method allows the covariates from multiple modes, whenever available. We utilize a low-rank constraint in the regression coefficient to encourage the sharing among tensor entries. The statistical convergence of our estimator is established, and we quantify the gain in predictive power by taking multiple covariates into account.

A secondary contribution is that our method allows a broad range of tensor types, including continuous, count, and binary observations. While previous tensor regression methods (??) are able to analyze Gaussian responses, none of them is suitable for exponential distribution family of tensors. We develop a generalized tensor regression framework, and as a by product, our models allows heteroscedasticity by relating the variance of tensor entry to its mean. This flexibility is particularly important in practice, because social network, brain imaging, or gene expression datasets are often non-Gaussian.

Related work. Our work is closely related to but also clearly distinctive from several lines of previous work. The first is a class of *unsupervised* tensor decomposition (???) that aims to find a low-rank representation of a data tensor. In contrast, our model can be viewed a *supervised* tensor learning, which aims to identify the association between a data tensor and covariates. The second related line (??) tackles tensor regression where the response is a scalar and the *predictor* is a tensor. Our proposal is orthogonal to theirs because we treat the tensor as a *response*. The tensor-response model is appealing for high-dimensional analysis when both the response and the covariate dimensions grow. The last line of work studies the network-response model (??). The earlier development of this model focuses mostly on binary data in the presence of dyadic covariates (?). We will demonstrate the enhanced accuracy as the order of data grows, and establish the general theory for exponential family which is arguably better suited to various data types.

2. Preliminaries

We begin by reviewing the basic properties about tensors (?). We use $\mathcal{Y} = \llbracket y_{i_1, \dots, i_K} \rrbracket \in \mathbb{R}^{d_1 \times \dots \times d_K}$ to denote an order- K (d_1, \dots, d_K)-dimensional tensor. The multilinear multiplication of a tensor $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$ by matrices

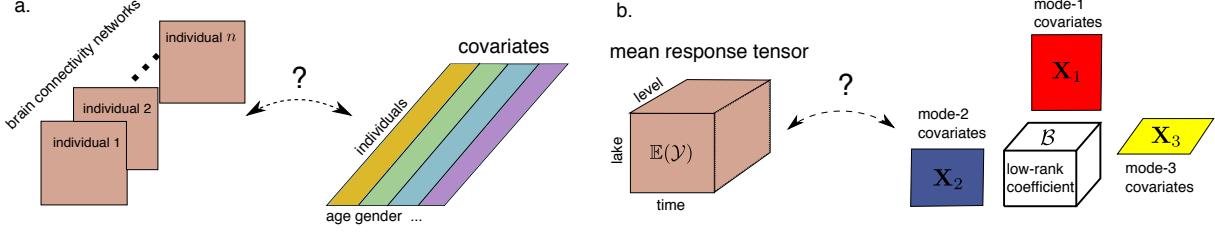


Figure 1. Examples of tensor response regression model with covariates on multiple modes. (a) Network population model. (b) Spatial-temporal growth model.

$\mathbf{X}_k = \llbracket x_{i_k, j_k}^{(k)} \rrbracket \in \mathbb{R}^{p_k \times d_k}$ is defined as

$$\mathcal{Y} \times_1 \mathbf{X}_1 \dots \times_K \mathbf{X}_K = \llbracket \sum_{i_1, \dots, i_K} y_{i_1, \dots, i_K} x_{j_1, i_1}^{(1)} \dots x_{j_K, i_K}^{(K)} \rrbracket,$$

which results in an order- K (p_1, \dots, p_K)-dimensional tensor. For ease of presentation, we use shorthand notion $\mathcal{Y} \times \{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ to denote the tensor-by-matrix product. For any two tensors $\mathcal{Y} = \llbracket y_{i_1, \dots, i_K} \rrbracket$, $\mathcal{Y}' = \llbracket y'_{i_1, \dots, i_K} \rrbracket$ of identical order and dimensions, their inner product is defined as $\langle \mathcal{Y}, \mathcal{Y}' \rangle = \sum_{i_1, \dots, i_K} y_{i_1, \dots, i_K} y'_{i_1, \dots, i_K}$. The Frobenius norm of tensor \mathcal{Y} is defined as $\|\mathcal{Y}\|_F = \langle \mathcal{Y}, \mathcal{Y} \rangle^{1/2}$. A higher-order tensor can be reshaped into a lower-order object (?). We use $\text{vec}(\cdot)$ to denote the operation that reshapes the tensor into a vector, and $\text{Unfold}_k(\cdot)$ the operation that reshapes the tensor along mode- k into a matrix of size d_k -by- $\prod_{i \neq k} d_i$. The Tucker rank of an order- K tensor \mathcal{Y} is defined as a length- K vector $\mathbf{r} = (r_1, \dots, r_K)$, where r_k is the rank of matrix $\text{Unfold}_k(\mathcal{Y})$, $k = 1, \dots, K$. We use lower-case letters (e.g., a, b, c) for scalars/vectors, upper-case boldface letters (e.g., $\mathbf{A}, \mathbf{B}, \mathbf{C}$) for matrices, and calligraphy letters (e.g., $\mathcal{A}, \mathcal{B}, \mathcal{C}$) for tensors of order three or greater. We let \mathbf{I}_d denote the $d \times d$ identity matrix, $[d]$ denote the d -set $\{1, \dots, d\}$, and allow an $\mathbb{R} \rightarrow \mathbb{R}$ function to be applied to tensors in an element-wise manner.

3. Motivation and model

Let $\mathcal{Y} = \llbracket y_{i_1, \dots, i_K} \rrbracket \in \mathbb{R}^{d_1 \times \dots \times d_K}$ denote an order- K data tensor. Suppose we observe covariates on some of the K modes. Let $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$ denote the available covariates on the mode k , where $p_k \leq d_k$. We propose a multilinear structure on the conditional expectation of the tensor. Specifically,

$$\mathbb{E}(\mathcal{Y} | \mathbf{X}_1, \dots, \mathbf{X}_K) = f(\Theta), \text{ with} \quad (1)$$

$$\Theta = \mathcal{B} \times \{\mathbf{X}_1, \dots, \mathbf{X}_K\},$$

where $f(\cdot)$ is a known link function, $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$ is the linear predictor, $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$ is the parameter tensor of interest, and \times denotes the tensor Tucker product. The choice of link function depends on the distribution of the response data. Some common choices are identity link for

Data type	Gaussian	Poisson	Bernoulli
Domain \mathbb{Y}	\mathbb{R}	\mathbb{N}	$\{0, 1\}$
$b(\theta)$	$\theta^2/2$	$\exp(\theta)$	$\log(1 + \exp(\theta))$
link $f(\theta)$	θ	$\exp(\theta)$	$(1 + \exp(-\theta))^{-1}$

Table 1. Canonical links for common distributions.

Gaussian tensor, logistic link for binary tensor, and $\exp(\cdot)$ link for Poisson tensor (see Table 1).

We give three concrete examples of tensor regression that arise in practice.

Example 1 (Spatio-temporal growth model). Let $\mathcal{Y} = \llbracket y_{ijk} \rrbracket \in \mathbb{R}^{d \times m \times n}$ denote the pH measurements of d lakes at m levels of depth and for n time points. Suppose the sampled lakes belong to p types, with q lakes in each type. Let $\{\ell_j\}_{j \in [m]}$ denote the sampled depth levels and $\{t_k\}_{k \in [n]}$ the time points. Assume that the expected pH trend in depth is a polynomial of order r and that the expected trend in time is a polynomial of order s . Then, the spatio-temporal growth model can be represented as

$$\mathbb{E}(\mathcal{Y} | \mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3) = \mathcal{B} \times \{\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3\}, \quad (2)$$

where $\mathcal{B} \in \mathbb{R}^{p \times (r+1) \times (s+1)}$ is the coefficient tensor of interest, $\mathbf{X}_1 = \text{blockdiag}\{\mathbf{1}_q, \dots, \mathbf{1}_q\} \in \{0, 1\}^{d \times p}$ is the design matrix for lake types,

$$\mathbf{X}_2 = \begin{pmatrix} 1 & \ell_1 & \dots & \ell_1^r \\ 1 & \ell_2 & \dots & \ell_2^r \\ \vdots & \ddots & & \vdots \\ 1 & \ell_m & \dots & \ell_m^r \end{pmatrix}, \quad \mathbf{X}_3 = \begin{pmatrix} 1 & t_1 & \dots & t_1^s \\ 1 & t_2 & \dots & t_2^s \\ \vdots & \vdots & \ddots & \vdots \\ 1 & t_n & \dots & t_n^s \end{pmatrix}$$

are the design matrices for spatial and temporal effects, respectively. The model (2) is a higher-order extension of the “growth curve” model originally proposed for matrix data (???). Clearly, the spatial-temporal model is a special case of our tensor regression model, with covariates available on each of the three modes.

Example 2 (Network population model). Network response model is recently developed in the context of neuroimaging analysis. The goal is to study the relationship between network-valued response and the individual covariates. Suppose we observe n i.i.d. observations $\{(\mathbf{Y}_i, \mathbf{x}_i) : i =$

110 $1, \dots, n\}$, where $\mathbf{Y}_i \in \{0, 1\}^{d \times d}$ is the brain connectivity
 111 network on the i -th individual, and $\mathbf{x}_i \in \mathbb{R}^p$ is the individual covariate such as age, gender, cognition, etc. The
 112 network-response model (??) has the form
 113

$$115 \text{logit}(\mathbb{E}(\mathbf{Y}_i | \mathbf{x}_i)) = \mathcal{B} \times_3 \mathbf{x}_i, \quad \text{for } i = 1, \dots, n \quad (3)$$

116 where $\mathcal{B} \in \mathbb{R}^{d \times d \times p}$ is the coefficient tensor of interest.
 117

118 The model (3) is a special case of our tensor-response
 119 model, with covariates on the last mode of the tensor.
 120 Specifically, stacking $\{\mathbf{Y}_i\}$ together yields an order-3 re-
 121 sponse tensor $\mathcal{Y} \in \{0, 1\}^{d \times d \times n}$, along with covariate ma-
 122 trix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times p}$. Then, the model (3) can
 123 be written as
 124

$$125 \text{logit}(\mathbb{E}(\mathcal{Y} | \mathbf{X})) = \mathcal{B} \times_3 \mathbf{X} = \mathcal{B} \times \{\mathbf{I}_d, \mathbf{I}_d, \mathbf{X}\}.$$

126 **Example 3** (Dyadic data with node attributes). Dyadic
 127 dataset consists of measurements on pairs of objects or
 128 under a pair of conditions. Common examples include net-
 129 works and graphs. Let $\mathcal{G} = (V, E)$ denote a network, where
 130 $V = [d]$ is the node set of the graph, and $E \subset V \times V$ is the
 131 edge set. Suppose that we also observe covariate $\mathbf{x}_i \in \mathbb{R}^p$
 132 associated to each $i \in V$. A probabilistic model on the
 133 graph $\mathcal{G} = (V, E)$ can be described by the following matrix
 134 regression. The edge connects the two vertices i and j inde-
 135 pendently of other pairs, and the probability of connection
 136 is modeled as
 137

$$138 \text{logit}(\mathbb{P}((i, j) \in E) = \mathbf{x}_i^T \mathbf{B} \mathbf{x}_j = \langle \mathbf{B}, \mathbf{x}_i^T \mathbf{x}_j \rangle. \quad (4)$$

139 The above model has demonstrated its success in mod-
 140 eling transitivity, balance, and communities in the net-
 141 works (?). We show that our tensor regression model (1)
 142 also incorporates the graph model as a special case. Let
 143 $\mathcal{Y} = [\![y_{ij}]\!]$ be a binary matrix where $y_{ij} = \mathbb{1}_{(i,j) \in E}$. Define
 144 $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times p}$. Then, the graph model (4)
 145 can be expressed as
 146

$$147 \text{logit}(\mathbb{E}(\mathbf{Y} | \mathbf{X})) = \mathcal{B} \times \{\mathbf{X}, \mathbf{X}\}.$$

148 In the above three examples and many other studies, re-
 149 searchers are interested in uncovering the variation in the
 150 data tensor that can be explained by the covariates. The
 151 regression coefficient \mathcal{B} in our model model (1) serves this
 152 goal by collecting the effects of covariates and the inter-
 153 action thereof. To encourage the sharing among effects, we
 154 assume that the coefficient tensor \mathcal{B} lies in a low-dimensional
 155 parameter space:
 156

$$157 \mathcal{P}_{r_1, \dots, r_K} = \{\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K} : r_k(\mathcal{B}) \leq r_k \text{ for all } k \in [K]\},$$

158 where $r_k(\mathcal{B}) \leq p_k$ is the Tucker rank at mode k of the tensor.
 159 The low-rank assumption is plausible in many scientific
 160 applications. In brain imaging analysis, for instance, it is
 161

often believed that the brain nodes can be grouped into fewer
 162 communities, and the numbers of communities are much
 163 smaller than the number of nodes. The low-rank structure
 164 encourages the shared information across tensor entries,
 165 thereby greatly improving the estimation stability. When
 166 no confusion arises, we drop the subscript (r_1, \dots, r_K) and
 167 write \mathcal{P} for simplicity.

Our tensor regression model is able to incorporate covari-
 168 ates on any subset of modes, whenever available. Without
 169 loss of generality, we denote by $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ the
 170 covariates in all modes and treat $\mathbf{X}_k = \mathbf{I}_{d_k}$ if the mode- k
 171 has no (informative) covariate. Then, the final form of our
 172 tensor regression model can be written as:
 173

$$174 \mathbb{E}(\mathcal{Y} | \mathcal{X}) = f(\Theta), \quad \Theta = \mathcal{B} \times \{\mathbf{X}_1, \dots, \mathbf{X}_K\}, \\ 175 \text{where } \text{rank}(\mathcal{B}) \leq (r_1, \dots, r_K), \quad (5)$$

176 where the entries of \mathcal{Y} are independent r.v.'s conditional on
 177 \mathcal{X} , and $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$ is the low-rank coefficient tensor
 178 of interest. We comment that other forms of tensor low-
 179 rankness are also possible, and here we choose Tucker rank
 180 just for parsimony. Similar models can be derived using var-
 181 ious notions of low-rankness based on CP decomposition (?)
 182 and train decomposition (?).

4. Rank-constrained likelihood-based estimation

We develop a likelihood-based procedure to estimate the coefficient tensor \mathcal{B} in (5). We adopt the exponential family as a flexible framework for different data types. In a classical generalized linear model (GLM) with a scalar response y and covariate \mathbf{x} , the density is expressed as:

$$p(y | \mathbf{x}, \boldsymbol{\beta}) = c(y, \phi) \exp \left(\frac{y\theta - b(\theta)}{\phi} \right) \text{ with } \theta = \boldsymbol{\beta}^T \mathbf{x},$$

where $b(\cdot)$ is a known function, θ is the linear predictor, $\phi > 0$ is the dispersion parameter, and $c(\cdot)$ is a known normalizing function. The choice of link functions depends on the data types and on the observation domain of y , denoted \mathbb{Y} . For example, the observation domain is $\mathbb{Y} = \mathbb{R}$ for continuous data, $\mathbb{Y} = \mathbb{N}$ for count data, and $\mathbb{Y} = \{0, 1\}$ for binary data. Note that the canonical link function f is chosen to be $f(\cdot) = b'(\cdot)$. Table 1 summarizes the canonical link functions for common types of distributions.

We model the entries in the response tensor y_{ijk} conditional on θ_{ijk} as independent draws from an exponential family. The quasi log-likelihood of (5) is equal (ignoring constant) to Bregman distance between \mathcal{Y} and $b'(\Theta)$:

$$\mathcal{L}_{\mathcal{Y}}(\mathcal{B}) = \langle \mathcal{Y}, \Theta \rangle - \sum_{i_1, \dots, i_K} b(\theta_{i_1, \dots, i_K}), \\ \text{where } \Theta = \mathcal{B} \times \{\mathbf{X}_1, \dots, \mathbf{X}_K\}.$$

We assume that we have an additional information on an upper bound $\alpha > 0$ such that $\|\Theta\|_\infty \leq \alpha$. This is the case for many applications we have in mind such as brain network analysis where fiber connections are bounded. We propose a constrained maximum likelihood estimator (MLE) for the coefficient tensor:

$$\hat{\mathcal{B}} = \arg \max_{\text{rank}(\mathcal{B}) \leq \mathbf{r}, \|\Theta(\mathcal{B})\|_\infty \leq \alpha} \mathcal{L}_{\mathcal{Y}}(\mathcal{B}). \quad (6)$$

In the following theoretical analysis, we assume the rank $\mathbf{r} = (r_1, \dots, r_K)$ is known and fixed. The adaptation of unknown \mathbf{r} will be addressed in Section 5.2.

4.1. Statistical properties

We assess the estimation accuracy using the deviation in the Frobenius norm. For the true coefficient tensor $\mathcal{B}_{\text{true}}$ and its estimator $\hat{\mathcal{B}}$, define

$$\text{Loss}(\mathcal{B}_{\text{true}}, \hat{\mathcal{B}}) = \|\mathcal{B}_{\text{true}} - \hat{\mathcal{B}}\|_F^2.$$

In modern applications, the response tensor and covariates are often large-scale. We are particularly interested in the high-dimensional region in which both d_k and p_k diverge; i.e. $d_k \rightarrow \infty$ and $p_k \rightarrow \infty$, while $\frac{p_k}{d_k} \rightarrow \gamma_k \in [0, 1]$. As the size of problem grows, and so does the number of unknown parameters. As such, the classical MLE theory does not directly apply. We leverage the recent development in random tensor theory and high-dimensional statistics to establish the error bounds of the estimation.

Assumption 1. *We make the following assumptions:*

- A1. *There exist two positive constants $c_1, c_2 > 0$ such that $c_1 \leq \sigma_{\min}(\mathbf{X}_k) \leq \sigma_{\max}(\mathbf{X}_k) \leq c_2$ for all $k \in [K]$. Here $\sigma_{\min}(\cdot)$ and $\sigma_{\max}(\cdot)$ denotes the smallest and largest singular values, respectively.*
- A2. *There exist positive constants $L, U > 0$ such that $L\phi \leq \text{Var}(y_{i_1, \dots, i_K} | \theta_{i_1, \dots, i_K}) \leq U\phi$ for all $|\theta_{i_1, \dots, i_K}| \leq \alpha$.*
- A2'. *Equivalently, there exists two positive constants $L, U > 0$ such that $L \leq b''(\theta) \leq U$ for all $|\theta| \leq \alpha$, where α is the upper bound of the linear predictor.*

The assumptions are fairly mild. Assumption A1 guarantees the non-singularity of the covariates, and Assumption A2 ensures the log-likelihood $\mathcal{Y}(\Theta)$ is strictly concave in the linear predictor Θ . Assumption A2 and A2' are equivalent, because $\text{Var}(y_{i_1, \dots, i_K} | \mathcal{X}, \mathcal{B}) = \phi b''(\theta_{i_1, \dots, i_K})$ when y_{i_1, \dots, i_K} belongs to an exponential family (?).

Theorem 4.1 (Statistical convergence). *Consider a generalized tensor regression model with covariates on multiple modes $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_K\}$. Suppose the entries in \mathcal{Y} are independent realizations of an exponential family distribution, and $\mathbb{E}(\mathcal{Y} | \mathcal{X})$ follows the low-rank tensor regression model (5). Under Assumption 1, there exist two*

constants $C_1, C_2 > 0$, such that, with probability at least $1 - \exp(-C_1 \sum_k p_k)$,

$$\text{Loss}(\mathcal{B}_{\text{true}}, \hat{\mathcal{B}}) \leq C_2 \sum_k p_k. \quad (7)$$

Here, $C_2 = C_2(\mathbf{r}, \alpha, K) > 0$ is a constant that does not depend on the dimensions $\{d_k\}$ and $\{p_k\}$.

To gain further insight on the bound (7), we consider a special case when tensor dimensions are equal at each of the modes, i.e., $d_k = d$, $p_k = \gamma d$, $\gamma \in [0, 1)$ for all $k \in [K]$, and the covariates \mathbf{X}_k are Gaussian design matrices with i.i.d. $N(0, 1)$ entries. To put the context in the framework of Theorem 4.1, we rescale the covariates into $\check{\mathbf{X}}_k = \frac{1}{\sqrt{d}} \mathbf{X}_k$ so that the singular values of $\check{\mathbf{X}}_k$ are bounded by $1 \pm \sqrt{\gamma}$. The result in (7) implies that the estimated coefficient has a convergence rate $\mathcal{O}(\frac{p}{d^K})$ in the scale of the original covariates $\{\mathbf{X}_k\}$. Therefore, our estimation is consistent as the dimension grows, and the convergence becomes especially favorably as the order of tensor data increases.

As immediate applications, we obtain the convergence rate for the three examples mentioned in Section 3. Without loss of generality, we assume that the singular values of the d_k -by- p_k covariate matrix \mathbf{X}_k are bounded by $\sqrt{d_k}$.

Example 4 (Spatio-temporal growth model). The estimated type-by-time-by-space coefficient tensor converges at the rate $\mathcal{O}(\frac{p+r+s}{d^{mn}})$ where $p \leq d$, $r \leq m$ and $s \leq n$. The estimation achieves consistency as long as the dimension grows in either of the three modes.

Example 5 (Network population model). The estimated node-by-node-by-covariate tensor converges at the rate $\mathcal{O}(\frac{2d+p}{d^2n})$ where $p \leq n$. The estimation achieves consistency as the number of individuals or the number of nodes grows.

Example 6 (Dyadic data with node attributes). The estimated covariate-by-covariate matrix converges at the rate $\mathcal{O}(\frac{p}{d^2})$ where $p \leq d$. Again, our estimation achieves consistency as the number of nodes grows.

We conclude this section by providing the prediction accuracy, measured in KL divergence, for the response distribution.

Theorem 4.2 (Prediction error). *Assume the same set-up as in Theorem 4.1. Let $\mathbb{P}_{\mathcal{Y}_{\text{true}}}$ and $\mathbb{P}_{\hat{\mathcal{Y}}}$ denote the distributions of \mathcal{Y} given the true parameter $\mathcal{B}_{\text{true}}$ and estimated parameter $\hat{\mathcal{B}}$, respectively. Then, we have, with probability at least $1 - \exp(C_1 \sum_k p_k)$,*

$$KL(\mathbb{P}_{\mathcal{Y}_{\text{true}}}, \mathbb{P}_{\hat{\mathcal{Y}}}) \leq C_4 \sum_k p_k,$$

where $C_4 = C_4(\mathbf{r}, \alpha, K) > 0$ is a constant that do not depend on the dimensions $\{d_k\}$ and $\{p_k\}$.

5. Numerical implementation

5.1. Alternating optimization

In this section, we introduce an efficient algorithm to solve (6). The objective function $\mathcal{L}_Y(\mathcal{B})$ is concave in \mathcal{B} when the link f is the canonical link function. However, the feasible set \mathcal{P} is non-convex, and thus the optimization (6) is a non-convex problem. We utilize a Tucker factor representation of the coefficient tensor \mathcal{B} and turn the optimization into a block-wise convex problem.

Specifically, write the rank- r decomposition of coefficient tensor \mathcal{B} as

$$\mathcal{B} = \mathcal{C} \times \{\mathbf{M}_1, \dots, \mathbf{M}_K\}, \quad (8)$$

where $\mathcal{C} \in \mathbb{R}^{r_1 \times \dots \times r_K}$ is a full-rank core tensor, $\mathbf{M}_k \in \mathbb{R}^{p_k \times r_k}$ are factor matrices with orthogonal columns. Estimating \mathcal{B} amounts to finding both the core tensor \mathcal{C} and the factor matrices \mathbf{M}_k 's. The optimization (6) can be written as $(\hat{\mathcal{C}}, \{\hat{\mathbf{M}}_k\}) = \arg \max \mathcal{L}_Y(\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K)$, where

$$\mathcal{L}_Y(\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K) = \langle \mathcal{Y}, \Theta \rangle - \sum_{i_1, \dots, i_K} b(\theta_{i_1, \dots, i_K}),$$

with $\Theta = \mathcal{C} \times \{\mathbf{M}_1 \mathbf{X}_1, \dots, \mathbf{M}_K \mathbf{X}_K\}$.

The decision variables in the above objective function consist of $K + 1$ blocks of variables, one for the core tensor \mathcal{C} and K for the factor matrices \mathbf{M}_k 's. We notice that, if any K out of the $K + 1$ blocks of variables are known, then the optimization with respect to the last block of variables reduced to a simple GLM. We therefore choose to iteratively update one block at a time while keeping others fixed. We leverage on a block relaxation algorithm for optimization, and the classical (local) convergence for block algorithm applies. Although a non-convex optimization of this type usually has no guarantee on global optimality, our numerical experiments have suggested high-quality solutions (see Section 6). The full algorithm is described in Algorithm 1.

5.2. Rank selection

Algorithm 1 takes the rank r as an input. Estimating an appropriate rank given the data is of practical importance. We propose to use Bayesian information criterion (BIC) and choose the rank that minimizes BIC; i.e.

$$\begin{aligned} \hat{r} &= \arg \min_{\mathbf{r}=(r_1, \dots, r_K)} \text{BIC}(\mathbf{r}) \\ &= \arg \min_{\mathbf{r}=(r_1, \dots, r_K)} [-2\mathcal{L}_Y(\hat{\mathcal{B}}) + p_e(\mathbf{r}) \log(\prod_k d_k)], \end{aligned} \quad (9)$$

where $p_e(\mathbf{r}) \stackrel{\text{def}}{=} \sum_k (p_k - r_k) r_k + \prod_k r_k$ is the effective number of parameters in the model. We choose \hat{r} that minimizes $\text{BIC}(\mathbf{r})$ via grid search. Our choice of BIC aims to balance between the goodness-of-fit for the data and the

degree of freedom in the population model. We test its empirical performance in Section 6.

6. Simulation

We evaluate the empirical performance of our generalized tensor regression through simulations. We consider order-3 tensors with a range of distribution types. The coefficient tensor \mathcal{B} is generated using the factorization form (8) where both the core and factor matrices are drawn i.i.d. from Uniform[-1,1]. The linear predictor is then simulated from $\mathcal{U} = \mathcal{B} \times \{\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3\}$, where \mathbf{X}_k is either an identity matrix (i.e. no covariate available) or Gaussian random matrix with i.i.d. entries from $N(0, \sigma_k^2)$. We set $\sigma_k = d_k^{-1/2}$ to ensure the singular values of \mathbf{X}_k are bounded as d_k increases. The \mathcal{U} is scaled such that $\|\mathcal{U}\|_\infty = 1$. Conditional on the linear predictor $\mathcal{U} = [u_{ijk}]$, the entries in the tensor $\mathcal{Y} = [y_{ijk}]$ are drawn independently according to one of the following three probabilistic models:

(a) (Gaussian). Continuous entries $y_{ijk} \sim N(\alpha u_{ijk}, 1)$.

(b) (Poisson). Count entries $y_{ijk} \sim \text{Poi}(e^{\alpha u_{ijk}})$.

(c) (Bernoulli). Binary entries $y_{ijk} \sim \text{Ber}\left(\frac{e^{\alpha u_{ijk}}}{1+e^{\alpha u_{ijk}}}\right)$.

Here $\alpha > 0$ is a scalar controlling the magnitude of the effect size. In each simulation study, we report the mean squared error (MSE) for the coefficient tensor averaged across $n_{\text{sim}} = 30$ replications.

6.1. Finite-sample performance

The experiment I assesses the selection accuracy of our BIC criterion (9). We consider the balanced situation where $d_k = d$, $p_k = 0.4d_k$ for $k = 1, 2, 3$. We set $\alpha = 10$ and consider various combinations of dimension d and rank $\mathbf{r} = (r_1, r_2, r_3)$. For each combination, we simulate tensor data following Gaussian, Bernoulli, and Poisson models. We then minimize BIC using a grid search over three dimensions. The hyper-parameter α is set to infinity in the fitting, which essentially imposes no prior on the coefficient magnitude. Table 2 reports the selected rank averaged over $n_{\text{sim}} = 30$ replicates for Gaussian and Poisson models. We found that when $d = 20$, the selected rank is slightly smaller than the true rank, and the accuracy improves immediately when the dimension increases to $d = 40$. This agrees with our expectation, as in tensor regression, the sample size is related to the number of entries. A larger d implies a larger sample size, so the BIC selection becomes more accurate.

The experiment II evaluates the accuracy when covariates are available on all modes. We set $\alpha = 10$, $d_k = d$, $p_k = 0.4d_k$, $r_k = r \in \{2, 4, 6\}$ and increase d from 25 to 50. Our theoretical analysis suggests that $\hat{\mathcal{B}}$ has a convergence rate $\mathcal{O}(d^{-2})$ in this setting. Figure 1 plots the estimation error versus the “effective sample size”, d^2 , under three different

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

Algorithm 1 Generalized tensor response regression with covariates on multiple modes

Input: Response tensor $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$, covariate matrices $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$ for $k = 1, \dots, K$, target Tucker rank $\mathbf{r} = (r_1, \dots, r_K)$, link function f , infinity norm bound α

Output: Low-rank estimation for the coefficient tensor $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$.

- 1: Calculate $\tilde{\mathcal{B}} = \mathcal{Y} \times_1 [(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T] \times_2 \dots \times_K [(\mathbf{X}_K^T \mathbf{X}_K)^{-1} \mathbf{X}_K^T]$.
- 2: Initialize the iteration index $t = 0$. Initialize the core tensor $\mathcal{C}^{(0)}$ and factor matrices $\mathbf{M}_k^{(0)} \in \mathbb{R}^{p_k \times r_k}$ via rank- \mathbf{r} Tucker approximation of $\tilde{\mathcal{B}}$, in the least-square sense.
- 3: **while** the relative increase in objective function $\mathcal{L}_{\mathcal{Y}}(\mathcal{B})$ is less than the tolerance **do**
- 4: Update iteration index $t \leftarrow t + 1$.
- 5: **for** $k = 1$ to K **do**
- 6: Obtain the factor matrix $\mathbf{M}_k^{(t+1)} \in \mathbb{R}^{p_k \times r_k}$ by solving p_k separate GLMs with link function f .
- 7: Update the columns of $\mathbf{M}_k^{(t+1)}$ by Gram-Schmidt orthogonalization.
- 8: **end for**
- 9: Obtain the core tensor $\mathcal{C}^{(t+1)} \in \mathbb{R}^{r_1 \times \dots \times r_K}$ by solving a GLM with $\text{vec}(\mathcal{Y})$ as response, $\odot_{k=1}^K [\mathbf{X}_k \mathbf{M}_k^{(t)}]$ as covariates, and f as link function. Here \odot denotes the Khatri-Rao product of matrices.
- 10: Rescale the core tensor subject to the infinity norm constraint.
- 11: Update $\mathcal{B}^{(t+1)} \leftarrow \mathcal{C}^{(t+1)} \times_1 \mathbf{M}_1^{(t+1)} \times_2 \dots \times_K \mathbf{M}_K^{(t+1)}$.
- 12: **end while**

True Rank r	Dimension (Gaussian tensors)		Dimension (Poisson tensors)	
	$d = 20$	$d = 40$	$d = 20$	$d = 40$
(3, 3, 3)	(2.1, 2.0, 2.0)	(3, 3, 3)	(2.0, 2.2, 2.1)	(3, 3, 3)
(4, 4, 6)	(3.2, 3.1, 5.0)	(4, 4, 6)	(4.0, 4.0, 5.2)	(4, 4, 6)
(6, 8, 8)	(5.1, 7.0, 6.9)	(6, 8, 8)	(5.0, 6.1, 7.1)	(6, 8, 8)

Table 2. Rank selection via BIC. Bold number indicates no significant difference between the estimate and the ground truth, based on a z -test with a level 0.05.

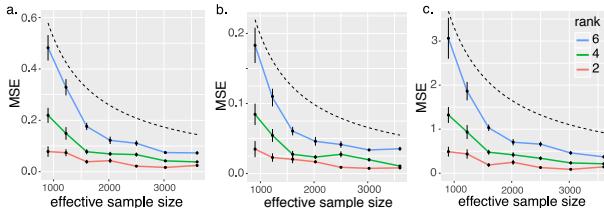


Figure 2. Mean squared error (MSE) against effective sample size. The three panels depict the MSE when the response tensors are generated from (a) Gaussian (b) Poisson and (c) Bernoulli models. The dashed curves correspond to $\mathcal{O}(1/d^2)$.

distribution models. We found that the empirical MSE decreases roughly at the rate of $1/d^2$, which is consistent with our theoretical ascertainment. We also observed that, tensors with higher ranks tend to yield higher estimation errors, as reflected by the upward shift of the curves as r increases. Indeed, a larger r implies a higher model complexity and thus greater difficulty in the estimation. Similar behaviors can be observed in the non-Gaussian data in Figures 2b-c.

The experiment III investigates the capability of our model in handling correlation among coefficients. We mimic the scenario of brain imaging analysis. A sample of $d_3 = 50$ networks are simulated, one for each individual. Each network measures the connections between $d_1 = d_2 = 20$ brain nodes. We simulate $p = 5$ covariates for the each

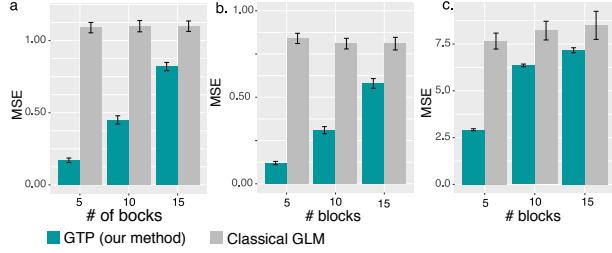


Figure 3. MSE when the networks have block structure. The three panels depict the MSE when the response tensors are generated from (a) Gaussian (b) Poisson and (c) Bernoulli models. The x-axis represents the number of blocks in the networks.

of the 50 individuals. These covariates may represent, for example, age, gender, cognitive score, etc. Recent study (?) has suggested that brain connectivity networks often exhibit community structure represented as a collection of subnetworks, and each subnetwork is comprised of a set of spatially distributed brain nodes. To accommodate this structure, we utilize the stochastic block model (?) to generate the effect size. Specifically, we partition the nodes into r blocks by assigning each node to a block with uniform probability. Edges within a same block are assumed to share the same covariate effects, where the effects are drawn i.i.d. from $N(0, 1)$. We then apply our tensor regression model to the network data using the BIC-selected rank. Note that in this case, the true model rank is unknown; the rank of a r -block matrix is not necessarily equal to r (?).

Figure 3 compares the MSE of our method with a classical GLM approach. A classical GLM is to regress the dyadic edges, one at a time, on the covariates, and this model is repeatedly fitted for each edge. This repeated approach, however, does not account for the correlation among the

330 edges, and may suffer from overfitting. As we can see in
 331 Figure 3, our tensor regression method achieves significant
 332 error reduction in all three models considered. The outer-
 333 performance is significant in the presence of large communi-
 334 ties, and even in the less structured case ($\sim 20/15 = 1.33$
 335 nodes per block), our method still outer-performs GLM.
 336 This is because the low-rankness in our modeling automati-
 337 cally identifies the shared information across entries. By
 338 selecting the rank in a data-driven way, our method is able to
 339 achieve accurate estimation with improved interpretability.
 340

6.2. Comparison with alternative methods

We compare our generalized tensor regression (**GTR**) with three other supervised tensor methods:

- Higher-order low-rank regression (**HOLRR**, ?) is a least-square based tensor regression that allows covariates on a single mode.
- Higher-order partial least square (**HOPLS**, ?) is a dimension-reduction method that jointly models a tensor response and a tensor covariate.
- Subsampled tensor projected gradient (**TPG**, ?) tackles the same question as **HOLRR** but instead uses a different algorithm to solve the problem.

These three methods are the closest algorithms to ours, in that they relate a tensor response to covariates using a low-rank structure. All the three methods allow only Gaussian data, whereas ours is applicable to any exponential family distribution including Gaussian, Bernoulli, Multinomial, etc. For fair comparison, we consider only Gaussian response in the simulation. We measure the accuracy using mean squared prediction error, $MSPE = \sqrt{\sum_k d_k} \|\hat{Y} - \mathbb{E}(Y|X)\|_F$, where \hat{Y} is the fitted value from each of the methods.

The comparison was assessed from three aspects: (a) benefit of incorporating covariates from multiple modes; (b) prediction error with respect to sample size; (c) sensitivity of accuracy with respect to model complexity. We use similar simulation setups as in our experiment II, but consider combinations of rank ($r = (3, 3, 3)$ vs. $(4, 5, 6)$), noise ($\sigma = 1/2$ vs. $1/4$), and dimension (d ranging from 20 to 100 for modes with covariates, $d = 20$ for modes without covariates).

Figure 4 shows the averaged prediction error across 30 replicates. We see that our **GTR** outperforms others, especially in the high-rank high-noise setting. As the number of informative modes (i.e. modes with available covariates) increases, the **GTR** exhibits a reduction in error whereas others have increased errors. This showcases the benefit toward prediction via incorporation of multiple covariates. Note that our method **GTR** is most comparable to **HOLRR** when there is only a single informative mode. In such a case, both methods share a same cost function but have dif-

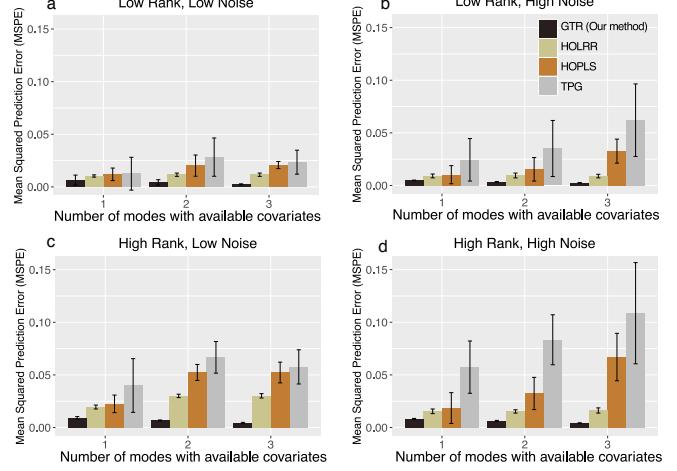


Figure 4. Comparison of MSPE versus the number of modes with covariates. We consider rank $r = (3, 3, 3)$ (low), $r = (4, 5, 6)$ (high), and noise $\sigma = 1/2$ (high), $\sigma = 1/4$ (low).

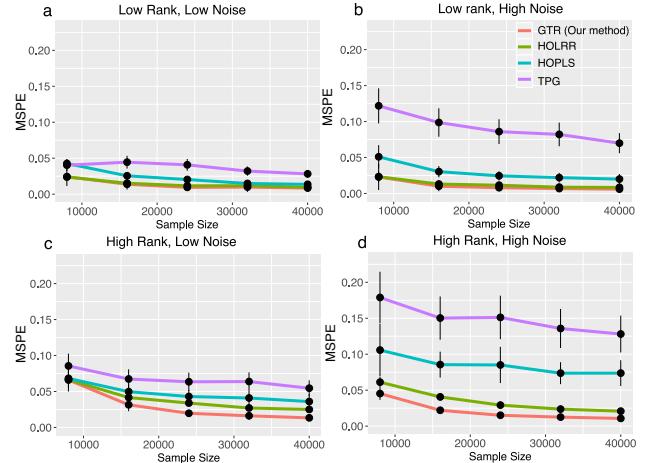


Figure 5. Comparison of MSPE versus sample size. We consider rank $r = (3, 3, 3)$ (low), $r = (4, 5, 6)$ (high), and noise $\sigma = 1/2$ (high), $\sigma = 1/4$ (low).

ferent algorithms. **GTR** alternates between informative and non-informative modes, whereas **HOLRR** approximates the non-informative modes via unfolded response alone. The accuracy gain in Figure 4 demonstrates the benefit of alternating algorithm – having informative modes also improves the estimation along non-informative modes.

Figure 5 compares the prediction error with respect to sample size. The sample size is the total number of entries in the tensor. In the low-rank setting, our method has similar performance as **HOLRR**, and the improvement becomes more pronounced when the rank increases. Neither **HOPLS** nor **TPG** has satisfactory performance in high-rank or high-noise settings. One possible reason is that a higher rank implies a higher inter-mode complexity, and our **GTR** method lends itself well to this context.

Exponential family tensor regression

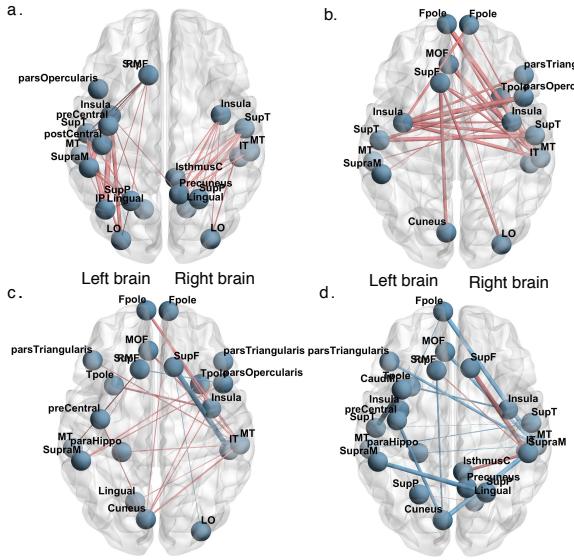


Figure 6. Top edges with large effects. Red edges represent relatively strong connections and blue edges represent relatively weak connections. (a) Global effect; (b) Female effect; (c) Age 22-25; (d) Age 31+.

7. Data analysis

We apply our method to two real datasets. The first application concerns the brain network modeling in response to individual attributes (i.e. covariate on one mode), and the second application focuses on multi-relational network analysis with dyadic attributes (i.e. covariates on two modes).

7.1. Human Connectome Project (HCP)

The Human connectome project (HCP, (?)) aims to build a network map that characterizes the anatomical and functional connectivity within healthy human brains. We take a subset of HCP data that consists of 136 brain structural networks, one for each individual. Each brain network is represented as a 68-by-68 binary matrix, where the entries encode the presence or absence of fiber connections between 68 brain regions. We consider four individual-covariates: gender, age 22-25, age 26-30, and age 31+.

We fit the tensor regression model to the HCP data. The BIC suggests a rank $r = (10, 10, 4)$ with log-likelihood $\mathcal{L}_y = -174654.7$. Figure 6 shows the top edges with high effect size, overlaid on the Desikan atlas brain template (?). We utilize the sum-to-zero contrasts in the effects coding and depict only the top 3% edges whose connections are non-constant across samples. Figure 6a shows that the global connection exhibits clear spatial separation, and that the nodes within each hemisphere are more densely connected with each other. In particular, the superior-temporal (*SupT*), middle-temporal (*MT*) and Insula are the top three popular nodes in the network. Interestingly, female brains display higher inter-hemispheric connectivity, especially

in the frontal, parietal, and temporal lobes (Figure 6b). This is in agreement with a recent study showing that female brains are optimized for inter-hemispheric communication (?). This result demonstrates the applicability of our method in detecting covariates signals.

7.2. Nations data

The second application examines the multi-relational network analysis with node-level attributes. We consider *Nations* dataset (?) which records 56 relations among 14 countries between 1950 and 1965. The multi-relational networks can be organized into a $14 \times 14 \times 56$ binary tensor, with each entry indicating the presence or absence of a connection, such as “sending tourist to”, “export”, “import”, between countries. The 56 relations span the fields of politics, economics, military, religion, etc.

We apply our tensor regression model to the *Nations* data. The BIC criterion suggests a rank $r = (4, 4, 4)$ for the coefficient tensor $\mathcal{B} \in \mathbb{R}^{6 \times 6 \times 56}$. Table ?? shows the K -means clustering of the 56 relations based on the 3rd mode factor $M_3 \in \mathbb{R}^{56 \times 4}$. We find that the relations reflecting the similar aspects of international affairs are grouped together. In particular, cluster I consists of political relations such as *officialvisits*, *intergovorgs*, and *militaryactions*; clusters II and III capture the economical relations such as *economicaid*, *booktranslations*, *tourism*; and Cluster IV represents the Cold War alliance blocs. The annotation similarity among grouped entities indicates the clustering results.

8. Conclusion

We have developed a generalized tensor regression with covariates on multiple modes. A fundamental feature of tensor-valued data is the statistical interdependence among entries. Our proposed rank-constrained estimation achieves high accuracy with sound theoretical guarantees. The estimation accuracy is quantified via deviation in the Frobenius norm and K-L divergence. Other measures of accuracy may also be desirable, such as the spectral norm or the maximum norm of the deviation. Exploiting the properties and benefits of different error quantification warrants future research.

References

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