

Generalized tensor response regression with multi-sided covariates

Anonymous Author 1
Unknown Institution 1

Anonymous Author 2
Unknown Institution 2

Anonymous Author 3
Unknown Institution 3

Abstract

We consider the problem of learning higher-order tensors with side information on a set of modes. Such data problems arise frequently in applications such as neuroimaging, network analysis, and recommendation systems. We propose a new family of tensor response regression models that incorporate covariate information, and obtain the theoretical accuracy guarantees. An efficient alternating updating algorithm is further developed. Our proposal handles a broad range of applications, including modeling brain connection in populations, link prediction in networks, and spatial-temporal growth model. The efficacy of our method is demonstrated through both simulations and analyses of two real-world datasets.

1 Introduction

Many contemporary scientific and engineering studies collect multi-way array data, a.k.a. tensor, accompanied by additional covariates. For example, in neuro-imaging analysis, researchers measure brain connections from a sample of individuals with the goal to identify the brain edges affected by age and gender. In social network analysis, how to explain the connection (e.g. community, transitive, etc.) by attributable of both nodes. In this article, we provide a general treatment to these seemingly different problems.

2 Preliminaries

We begin by reviewing a few basic factors about tensors [1]. 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 $\mathbf{M}_k = \llbracket m_{i_k, j_k}^{(k)} \rrbracket \in \mathbb{R}^{s_k \times d_k}$ is defined as

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

which results in an order- K tensor (s_1, \dots, s_K)-dimensional tensor. For ease of notation, we also write the above Tucker product $\mathcal{Y} \times \{\mathbf{M}_1, \dots, \mathbf{M}_K\}$ for short. 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 $\|\mathcal{Y}\|_F = \langle \mathcal{Y}, \mathcal{Y} \rangle^{1/2}$; it is the Euclidean norm of \mathcal{Y} regarded as an $\prod_k d_k$ -dimensional vector. We use lower-case letters a, b, c, \dots for scalars and vectors, upper-case boldface letters $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$ for matrices, and calligraphy letter $\mathcal{A}, \mathcal{B}, \mathcal{C}, \dots$ for tensors of order 3 or greater. We denote by \mathbf{I}_n the identity matrix of dimension n .

3 Motivation and models

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 of interest. In addition, suppose we observe covariates on a subset of modes. Let $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$ be the available covariates on the mode- k , where $p_k \leq d_k$. We propose the following multilinear structure in the mean of the tensor. Specifically,

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

where $f(\cdot)$ is a known link function, $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$ is called the linear predictor tensor, $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$ is the parameter tensor of interest, and \times denotes the tensor Tucker product. The link function depends on the distribution family of the response. Some common choices are identity link for Gaussian tensor, logistic link for binary tensor, and log link for Poisson tensor. We give three examples of multi-covariates tensor regression model that arises 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

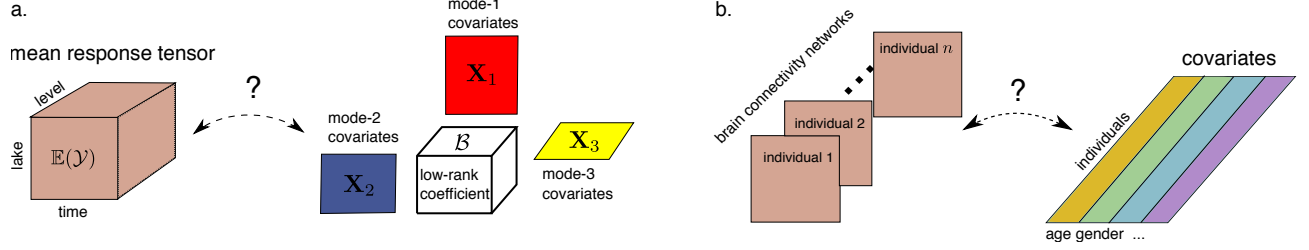


Figure 1: Examples of tensor regression with multi-sided covariates. (a) Spatio-temporal growth model. (b) Network population model.

of d lakes at m levels of depth and for n time points. Suppose the sampled lakes belong to q types, with p 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 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, a classical spatio-temporal growth model can be represented as

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

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

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

are the design matrices for spatial and temporal effects, respectively.

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 the network-valued response with the individual covariates. Suppose we observe n i.i.d. observations $\{(\mathbf{Y}_i, \mathbf{x}_i) : i = 1, \dots, n\}$, where $\mathbf{Y}_i \in \{0, 1\}^{d \times n}$ is the brain connectivity network on the i -th individual and $\mathbf{x}_i \in \mathbb{R}^p$ is the subject covariate such as age, gender. The network-response model has the form

$$\text{logit}(\mathbb{E}(\mathbf{Y}_i|\mathbf{x}_i)) = \mathcal{B} \times_3 \mathbf{x}_i, \quad \text{for } i = 1, \dots, n \quad (2)$$

where $\mathcal{B} \in \mathbb{R}^{d \times d \times p}$ is the coefficient tensor of interest. In fact, the model (2) is a special case of our multilinear tensor-response model. To see this, let $\mathcal{Y} \in \{0, 1\}^{d \times d \times n}$ denote the response tensor by stacking $\{\mathbf{Y}_i\}$ together along the 3rd mode and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$, then model (2) can be expressed as

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

where \mathbf{I}_d denotes the identity matrix of dimension d .

Example 3 (Link model with node attributes). Let $V = [n]$ be a set of vertices and explanatory variable $x_i \in \mathbb{R}^p$ associated to each $i \in V$. The network $G = (V, E)$ is described by the following matrix model. The edge connects the two vertices i and j independently of the others. The probability of connection is modeled as

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

Again, we show that this model is a special case of our tensor regression model. Let $\mathcal{Y} = \llbracket y_{ij} \rrbracket$ where $y_{ij} = \mathbb{1}_{(i, j) \in E}$. Define $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$. Then the above model can be expressed as

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

In the above three example and many other studies, researchers are interested in uncovering the variation in the data tensor that are explained by the covariates.

Without any structure on the coefficient tensor \mathcal{B} : A naive approach is to regress the tensor entry, one at a time, on the covariates, and this model is repeatedly fitted for each tensor element. Though this approach is scalable, it suffers from two drawbacks: (1) ignore the multilinear structure in the tensor (?) and (2) suffers from the multiplicity issue. To allow the structure among ..we further impose a multilinear low-rank structure on the coefficient tensor \mathcal{B}

$$\mathcal{P} = \{\mathcal{B} \in \mathbb{R}^{p_1 \times \cdots \times p_K} : r_k(\mathcal{B}) \leq r_k \text{ for } k \in [K]\}, \quad (3)$$

where $r_k(\mathcal{B}) \leq p_k$ is the Tucker rank of the tensor at mode k . Other low-rankness such as CP rank is also possible.

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

$$\mathbb{E}(\mathcal{Y}|\mathcal{X}) = f(\Theta), \quad \Theta = \mathcal{B} \times \{\mathbf{X}_1, \dots, \mathbf{X}_K\}, \quad \text{where } \text{rank}(\mathcal{B}) = (r_1, \dots, r_K), \quad (4)$$

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 link functions for various distribution types.

where $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$ is the low-rank coefficient tensor of interest, and the entries of \mathcal{Y} are independent realizations from an exponential distribution which are detailed in the next paragraph.

Note that the key assumption in the model is the low-rankness of the tensor form. The low-rank structure in (3) implies that the coefficient tensor can be expressed as $\mathcal{B} = \mathcal{C} \times \{\mathbf{M}_1, \dots, \mathbf{M}_K\}$. Then, our tensor regression model (4) is equivalent to

$$f(\mathbb{E}(\mathcal{Y}|\mathbf{X}_1, \dots, \mathbf{X}_K)) = \mathcal{C} \times \{\mathbf{X}_1 \mathbf{M}_1, \dots, \mathbf{X}_K \mathbf{M}_K\}.$$

The goal is to find a joint dimension reduction of \mathcal{Y} and \mathbf{X}_K such that the unexplained variation in the mean tensor. The factorization is restricted to the space spanned by \mathbf{X}_k . Here $\mathbf{X}_1 \mathbf{M}_1$ can be interpreted as the latent covariates that explains the variation in the response tensor. The core tensor \mathcal{C} collects the interaction effects of latent covariates across the K modes.

4 Rank-constrained likelihood-based estimation

We develop a likelihood-based procedure to estimate \mathcal{B} . The exponential family is a flexible framework for different data types. In a classical GLM with a scalar response y and covariate \mathbf{x} , the density is expressed as:

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

where $b(\cdot)$ is the 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 are tailored to the data types and the observation domain of y , denoted \mathbb{Y} . For example, the observation space for continuous data is $\mathbb{Y} = \mathbb{R}$, and for count data $\mathbb{Y} = \mathbb{N}$, and for binary data, $\mathbb{Y} = \{0, 1\}$. 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.

In our context, the entries in the response tensor y_{ijk} conditional on θ_{ijk} are independent drawn from exponential family. The quasi log-likelihood of (4) is equal (ignoring constant) to Bregman distance between \mathcal{Y}

and $b'(\Theta)$:

$$\mathcal{L}_{\mathcal{Y}}(\mathcal{B}) = \langle \mathcal{Y}, \Theta \rangle - \sum_{i_1, \dots, i_K} b(\theta_{i_1, \dots, i_K}),$$

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 recommendation systems where user ratings are bounded.* We propose a constrained maximum likelihood estimator for the coefficient tensor \mathcal{B} :

$$\hat{\mathcal{B}} = \arg \max_{\text{rank}(\mathcal{B})=\mathbf{r}, \|\Theta(\mathcal{B})\|_{\infty} \leq \alpha} \mathcal{L}_{\mathcal{Y}}(\mathcal{B}). \quad (5)$$

In the theoretical analysis, we assume the rank $\mathbf{r} = (r_1, \dots, r_K)$ is known. 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

$$\text{Loss}(\mathcal{B}_{\text{true}}, \hat{\mathcal{B}}) = \frac{1}{\prod_k p_k} \|\mathcal{B}_{\text{true}} - \hat{\mathcal{B}}\|_F^2.$$

We focus on the high-dimensional region in which both $d_k \rightarrow \infty$ and $p_k \rightarrow \infty$ while $\frac{p_k}{d_k} \rightarrow \gamma_k \in [0, 1]$.

Assumption 1. *We make the following assumptions:*

1. *There exists 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 $k \in [K]$.*
2. *There exist two positive constants $L, U > 0$ such that $L \leq \text{Var}(y_{i_1, \dots, i_K} | \mathcal{X}) \leq U$ uniformly over the parameter space \mathcal{P} .*
- 2'. *Equivalently, $L \leq b''(x) \leq U$ for all $x \leq \alpha$, where $b(\cdot)$ is the known function in the exponential family distribution and α is the upper bound of the linear predictor.*

Theorem 4.1 (Statistical convergence). *Consider a generalized tensor regression model with multi-sided covariates. Let $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$ be the tensor response and $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ the covariates, where $\mathbf{X}_k \in \mathbb{R}^{p_k \times d_k}$ is the covariate matrix on mode- 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 (4). Under Assumption 1, there exist two constants $C_1, C_2 > 0$ such that, with probability at least $1 - \exp(-C_1 \sum_k p_k)$,*

$$\text{Loss}(\mathcal{B}_{\text{true}}, \hat{\mathcal{B}}) \leq \frac{2}{c_1^{2K}} \min \left\{ \frac{C(\mathbf{r}, \alpha) \sum_k p_k}{\prod_k p_k}, 2\alpha^2 \right\},$$

Algorithm 1 Generalized tensor response regression with multi-sided covariates

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 (r_1, \dots, r_K) , link function f , entrywise bound α

Output: Estimated low-rank coefficient tensor $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$.

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

where $C(\alpha, \mathbf{r}) = \frac{1}{b''(\alpha)} \frac{\prod_k r_k}{r_{\max}} > 0$ is a constant that does not depend on dimension $\{d_k\}$ and $\{p_k\}$.

To gain further insight on the bound we consider the special case when dimensions are equal at each of the modes. 1. binary case; 2. large dimension region. $\mathcal{O}\left(\frac{p}{d^k}\right) \leq \mathcal{O}(d^{-(k-1)})$.

Corollary 1 (Spatio-temporal growth model). Our method yields the convergence rate $\mathcal{O}\left(\frac{p+r+s}{dmn}\right)$. Note that $p \leq d$, $r \leq m$ and $s \leq m$, so consistent estimator.

Corollary 2 (Network population model). Our method yields the convergence rate $\mathcal{O}\left(\frac{2d+p}{d^2n}\right)$. Note that $p \leq m$, so this is a consistent estimator. In contrast, a naive repeated glm will give $\mathcal{O}\left(\frac{p}{n}\right)$.

Corollary 3 (Link model with node attributes). Our method yields the convergence rate is $\mathcal{O}\left(\frac{p}{d^2}\right)$. Note that $p \leq m$, so again a consistent estimator. In contrast, a naive repeated glm will give $\mathcal{O}\left(\frac{p}{n}\right)$.

We provide the prediction accuracy for the response tensor.

Theorem 4.2 (Prediction error). Assume the same set-up as in Theorem (4.1). Let $\mathbb{P}_{\mathcal{Y}_{true}}$ the distribution of \mathcal{Y} given the true \mathcal{B}_{true} and $\mathbb{E}(\mathcal{Y}|\mathcal{X})$ the true mean. Let $\mathbb{P}_{\hat{\mathcal{Y}}}$ denote the distribution given the estimated $\hat{\mathcal{B}}$ and $\mathbb{E}(\hat{\mathcal{Y}}|\mathcal{X})$ the predicted mean. We have, with proba-

bility at least $1 - \exp(-C_1 \sum_k d_k)$,

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

$$Loss\left(\mathbb{E}(\mathcal{Y}|\mathcal{X}), \mathbb{E}(\hat{\mathcal{Y}}|\mathcal{X})\right) \leq b''(\alpha) C(\dots) \frac{\sum_k p_k}{\prod_k p_k}.$$

5 Numerical implementation

5.1 Alternating optimization

In this section, we introduce an efficient algorithm to solve (5). The objective function $\mathcal{L}_{\mathcal{Y}}(\mathcal{B})$ is concave in \mathcal{B} when the link f is canonical link function. However, the feasible set \mathcal{P} is non-convex, and thus the optimization (5) is a non-convex problem. We utilize a Tucker factor representation of coefficient tensor \mathcal{B} , and turn the optimization into a block-wise convex problem.

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

$$\mathcal{B} = \mathcal{C} \times \{\mathbf{M}_1, \dots, \mathbf{M}_K\}, \quad (6)$$

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 whose columns are orthogonal. Estimating \mathcal{B} amounts to finding both the core tensor \mathcal{C} and the factor matrices \mathbf{M}_k 's. Then, the optimization (5) can be written as $(\hat{\mathcal{C}}, \{\hat{\mathbf{M}}_k\}) = \arg \max \mathcal{L}_{\mathcal{Y}}(\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K)$, where

$$\mathcal{L}_{\mathcal{Y}}(\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K) = \langle \mathcal{Y}, \Theta \rangle - \sum_{i_1, \dots, i_K} b(\theta_{i_1, \dots, i_K}) \quad (7)$$

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

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. This observation suggests that we can iteratively update one block at a time while keeping others fixed. Specifically, suppose the core tensor and the factor matrix \mathbf{M}_k are known for $k = 1, \dots, K - 1$. It turns out that last factor matrix \mathbf{M}_K can be solved in a row-by-row fashion via d_K separate GLMs. To see this, let $\mathcal{C}^{(t)}$ denote the core tensor at the t -th iteration, $\mathbf{M}_k^{(t)}$ the k th factor matrix for $k \in [K - 1]$ at the t -th iteration, and define

$$\mathbf{X}_{-K}^{(t)} = \mathcal{C}^{(t)} \times \{\mathbf{M}_1^{(t)} \mathbf{X}_1, \dots, \mathbf{M}_{K-1}^{(t)} \mathbf{X}_{K-1}\},$$

Then, the objective (7) implies that the i -th row of $\mathbf{M}_K^{(t)}$ is the "regression coefficient" for a GLM whose response vector is $\text{vec}(\mathcal{Y}(:, :, i)) \in \mathbb{R}^{d_{-K}}$ and covariate matrix is $\text{Unfold}_K(\mathbf{X}_{-K}^{(t)}) \in \mathbb{R}^{d_{-K} \times r_K}$, where $d_{-K} \stackrel{\text{def}}{=} \prod_{k \in [K-1]} d_k$. Here $\text{vec}(\cdot)$ is the operator that reshapes a tensor into a vector, $\text{Unfold}_K(\cdot)$ is the operator that reshapes a tensor along the mode K into a matrix. The property of separation by row allows us to leverage state-of-art GLM solvers and parallel processing to achieve computational efficiency. After each iteration, we rescale the core tensor $\mathcal{C}^{(t+1)}$ subject to the infinity norm constraint. This post-processing in principle may not guarantee the monotonic increase of the objective, but we found that in our experiment this simple post-processing appears to be good enough for a desirable solution. The full algorithm is described in Algorithm 1.

5.2 Rank selection, missing data handling

Before concluding this section, we briefly comment on two implementation details. First, Algorithm 1 takes the rank \mathbf{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.

$$\begin{aligned} \hat{\mathbf{r}} &= \arg \min_{\mathbf{r}=(r_1, \dots, r_K)} \text{BIC}(\mathbf{r}) \\ &= \arg \min_{\mathbf{r}=(r_1, \dots, r_K)} \left[-2\mathcal{L}_{\mathcal{Y}}(\hat{\mathcal{B}}) + p_e(\mathbf{r}) \log \left(\prod_k d_k \right) \right], \end{aligned} \quad (8)$$

where $p_e(\mathbf{r}) \stackrel{\text{def}}{=} \sum_k (p_k - r_k - 1)r_k + \prod_k r_k$ is the effective number of parameters in the model. We choose $\hat{\mathbf{r}}$ 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.

Second, when some response entries y_{i_1, \dots, i_K} are missing, we replace the objective function by $\mathcal{L}_{\mathcal{Y}}$

$= \sum_{(i_1, \dots, i_K) \in \Omega} (y_{i_1, \dots, i_K} \theta_{i_1, \dots, i_K} - b(\theta_{i_1, \dots, i_K}))$, where $\Omega \subset [d_1] \times \dots \times [d_K]$ is the index set of non-missing entries. That is, we model the observed response entries and exclude the missing entries in the fitting. Similar strategy has been used for classical (unsupervised) Tucker and CP tensor decomposition with missing data [2, 3, 4]. In the presence of missing response, we modify line 7 in Algorithm 1 by fitting GLMs to the data for which y_{i_1, \dots, i_K} are observed. This approach requires there are no entirely missing slice, e.g. in the form of $\mathcal{Y}(:, :, k)$ for order-3 tensors. We regard this as a fairly mild condition akin to the coherence condition as in the completion literature [5, 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 (6) where both the core and factor matrices are drawn i.i.d. from Uniform[0,1]. The linear predictor is then simulated as $\mathcal{U} = \mathcal{B} \times \{\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3\}$, where \mathbf{X}_k is either an identity matrix (i.e. no covariate available) or Gaussian random matrix with entries i.i.d. from $N(0, \sigma^2)$. Without loss of generality, we set $\sigma^2 = d_k^{-1/2}$ to ensure the singular values of \mathbf{X}_k are bounded as d_k increases. We rescale \mathcal{U} such that $\|\mathcal{U}\|_{\infty} = 1$. Conditional on the linear predictor $\mathcal{U} = \llbracket u_{ijk} \rrbracket$, the entries in tensor $\mathcal{Y} = \llbracket y_{ijk} \rrbracket$ are drawn independently according to one of the following three probabilistic models:

- (a) (Gaussian). Continuous data $y_{ijk} \sim N(\alpha u_{ijk}, 1)$.
- (b) (Poisson). Count data $y_{ijk} \sim \text{Poi}(e^{\alpha u_{ijk}})$.
- (c) (Bernoulli). Binary data $y_{ijk} \sim \text{Ber}\left(\frac{e^{\alpha u_{ijk}}}{1 + e^{\alpha u_{ijk}}}\right)$.

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

The first experiment assesses the selection accuracy of our BIC criterion (8). We consider the balanced situation where $d_i = d$, $p_i = 0.4d_i$ for $i = 1, 2, 3$. We set $\alpha = 10$, $\alpha = 4$, 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 over \mathbf{r} using a grid search over three dimensions. We set the hyper-parameter α to infinity in the fitting, which essentially poses no prior on the coefficient magnitude. Table 2 reports the selected rank averaged over $n_{\text{sim}} = 30$ replicates for Gaussian and Poisson models. (The results for Bernoulli model is shown in the Supplements.) We found that when $d = 20$ the selected rank

True Rank \mathbf{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: Performance for 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.

is slightly smaller than the true rank, and the accuracy increases 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 second experiment evaluates the accuracy when covariates are available on all modes. We set $\alpha = 10$, $d_i = d$, $r_i = r \in \{2, 4, 6\}$ and increase d from 25 to 50. Our theoretical analysis suggests $\hat{\mathcal{B}}$ 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 distribution models. We found that the empirical RMSE decreases roughly at the rate of $1/d^2$, which is consistent with our theoretical ascertainment. We also observed that coefficients 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 higher model complexity, thus increasing the difficulty of the estimation. Similar behaviors can be observed in the non-Gaussian data in Figure 2b-c.

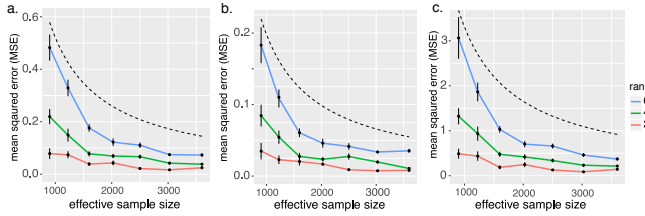


Figure 2: Estimation error against effective sample size. The three panels depicts the MSE when the response tensor is generated from (a) Gaussian (b) Poisson and (c) Bernoulli models. Each solid curve corresponds to a fixed rank. The dashed curve corresponds to $\mathcal{O}(1/d^2)$.

The third experiment investigates our model’s ability 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 50 individuals. These covariates may represent, for example, age, gender, cognitive score, etc. Recent study [7] has suggested that brain con-

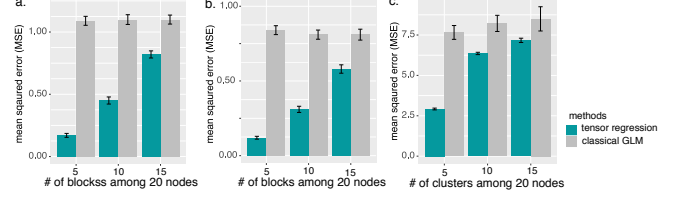


Figure 3: Performance comparison when the population network admit block structure. The three panels depicts the MSE when the response tensor is generated from (a) Gaussian (b) Poisson and (c) Bernoulli models. The x -axis represents the number of blocks in the networks.

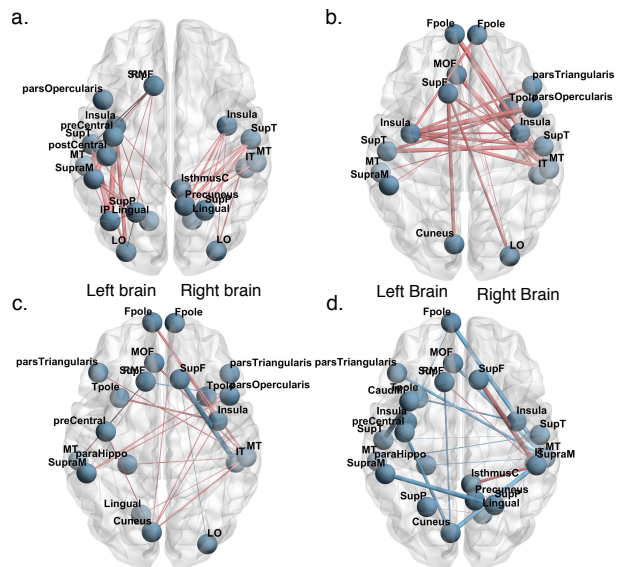
nectivity networks often exhibit community structure represented as a collection of subnetworks, and each subnetwork is comprised of a set of spatially distributed brain nodes. To accommodate this structure, we utilize the popular stochastic block model to generate the effect size. Specifically, we created r blocks among the nodes by randomly assigning each node to a cluster 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 applied 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 bounded by but not necessarily equal to r [8].

Figure 3 compares the MSE of our model 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 approach, however, does not account for the correlation structure among the edges. As we can see in Figure 3, our tensor regression method achieves significant error reduction in all three models considered. The out-performance is more apparent in the presence of large communities, and even in the less structured case ($\sim 20/15 = 1.33$ nodes per block), our method still outperforms GLM. This is because the low-rankness in our modeling automatically identifies the shared information across entries. By a data-driven rank selection, our method is able to achieve accurate estimation with improved interpretability.

7 Data analysis

We apply our tensor regression model to two real datasets. The first application concerns the modeling of brain network population 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).

7.2 Nations data



To investigate how the dyadic attributes affect the connection, we depicted the estimated coefficients

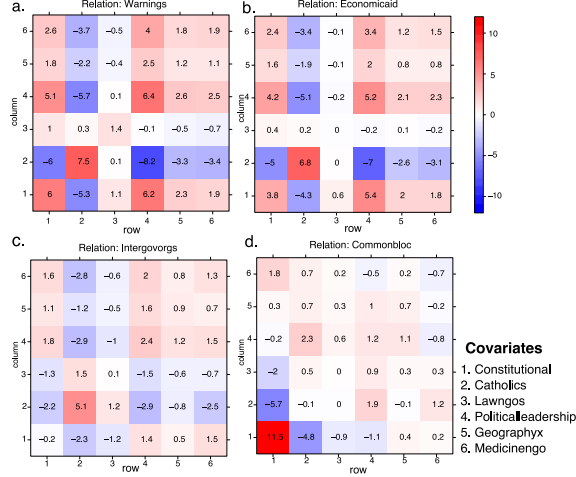


Figure 5: Estimated coefficient of country attributes towards the connection probability of relation type k .

$\hat{B} = [\hat{b}_{ijk}]$ for a few relation types (Figure 5). Note that entries \hat{b}_{ijk} can be interpreted as the contribution, at the logit scale, of covariate pair (i, j) (i th covariate for the “sender” country and j th covariate for the “receiver” country) towards the connection of relation k . Several interesting findings emerge from the estimation. We found that relations belonging to a same cluster tend to have similar covariate effects. For example, the relations *warnings* and *economicaid* were classified into Cluster II, and both exhibit similar covariate pattern (Figure 5a-b). Moreover, the diagonal entries $\hat{B}(i, i, k)$ tend to positively contribute to the connection. This is probably explained by the fact that countries with coherent attributes tend to interact more often than others. We also found that the *constitutional* attributes are greatly associated with the *commonbloc* relation, whereas such association is weaker for other relations (Figure 5d). This is not surprising, as the common block partition during Cold War is determined by capitalism vs. communism, which confounds with the *constitutional* attributes.

8 Conclusions

References

- [1] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [2] Evrim Acar, Daniel M Dunlavy, Tamara G Kolda, and Morten Mørup. Scalable tensor factorizations with missing data. In *Proceedings of the 2010 SIAM international conference on data mining*, pages 701–712. SIAM, 2010.
- [3] Miaoyan Wang and Lexin Li. Learning from binary multiway data: Probabilistic tensor decomposition and its statistical optimality. *arXiv preprint arXiv:1811.05076*, 2018.
- [4] Miaoyan Wang and Yun Song. Tensor decompositions via two-mode higher-order SVD (HOSVD). In *Artificial Intelligence and Statistics*, pages 614–622, 2017.
- [5] Emmanuel J Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717, 2009.
- [6] Navid Ghadermarzy, Yaniv Plan, and Ozgur Yilmaz. Learning tensors from partial binary measurements. *IEEE Transactions on Signal Processing*, 67(1):29–40, 2018.
- [7] Lucy F Robinson, Lauren Y Atlas, and Tor D Wager. Dynamic functional connectivity using state-based dynamic community structure: Method and application to opioid analgesia. *NeuroImage*, 108:274–291, 2015.
- [8] Miaoyan Wang and Yuchen Zeng. Multiway clustering via tensor block models. *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, to appear, page arXiv:1906.03807, 2019.
- [9] Rahul S Desikan, Florent Ségonne, Bruce Fischl, Brian T Quinn, Bradford C Dickerson, Deborah Blacker, Randy L Buckner, Anders M Dale, R Paul Maguire, Bradley T Hyman, et al. An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *Neuroimage*, 31(3):968–980, 2006.
- [10] Mingrui Xia, Jinhui Wang, and Yong He. Brainnet viewer: a network visualization tool for human brain connectomics. *PloS one*, 8(7):e68910, 2013.
- [11] Madhura Ingalkhalikar, Alex Smith, Drew Parker, Theodore D Satterthwaite, Mark A Elliott, Kosha Ruparel, Hakon Hakonarson, Raquel E Gur, Ruben C Gur, and Ragini Verma. Sex differences in the structural connectome of the human brain. *Proceedings of the National Academy of Sciences*, 111(2):823–828, 2014.
- [12] Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *International Conference on Machine Learning*, volume 11, pages 809–816, 2011.