
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 arise in applications such as neuroimaging, network analysis, and ... 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-word 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. ... (add two pictures; one for estimating network population; another for estimating link prediction)

In this article, we provide a general treatment to these seemingly different problems.

Comparison with other models

Our model is related to, but fundamentally different from, several lines of existing work.

Unsupervised tensor. supervised learning.

Tensor-predictor regression multilinear in coefficients.

Tensor-response regression multilinear in coefficient vs. multilinear in covariates.

Generalized linear model. In the high-dimensions both p and n increase while $p \leq d$. This is the case we consider. Classical GLM fixes p . Compared to Gaussian model, the log-likelihood is not strictly convex in the linear predictor. We allow various types of dependent variable.

2 Preliminaries

We begin by reviewing a few basic factors 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 $\mathbf{M}_k = \llbracket m_{i_k, j_k}^{(k)} \rrbracket \in \mathbb{R}^{s_k \times d_k}$ 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 bold-face 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.

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 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,

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

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 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,

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

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

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 $\mathbf{x}_i \in \mathbb{R}^p$ 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}

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

where $r_k(\mathcal{B})$ is the Tucker rank of the tensor at mode k . We assume that r_k is known and $r_k \leq p_k$. Other low-rankness such as CP rank is also possible. We choose the Tucker decomposition due to the following observation.

The low-rank structure in (3) implies that the coefficient tensor can be expressed as $\mathcal{B} = \mathcal{C} \times_1 \mathbf{M}_1 \times \dots \times \mathbf{M}_K$. 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.

(Question: the columns of \mathbf{X} should be normalized??)

Question: any connection to tensor completion?? If \mathbf{X}_K is a random design matrix?

4 Rank-constrained likelihood-based estimation

Note that 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 $\mathbf{X}_k = \mathbf{I}_{d_k}$ when the mode- k has no (informative) covariate. Our genral tensor regression model can be written as

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

where $\text{rank}(\mathcal{B}) = (r_1, \dots, r_K)$, (4)

where $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$ is a low-rank coefficient tensor of interest. In the following theoretical analysis, we assume the Tucker rank $\mathbf{r} = (r_1, \dots, r_K)$ is known. The adaptation of unknown \mathbf{r} will be addressed in Section 5.2.

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

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

where $c(\cdot)$ and $b(\cdot)$ are known functions and θ is the linear predictor. Note that the canonical link function f is chosen to be $f(\cdot) = b'(\cdot)$. Table 1 summarizes the canonical link function for common distributions.

In our context, the log-likelihood of (4) is the (?) divergence between the conditional distribution of $\mathcal{Y}|\Theta$ and the exponential family

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

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

Assume that we have an additional information on an upper bound $a > 0$ such that $\|\Theta\|_{\infty} \leq \alpha$. (more comments? ill-condition in the bernoulli model?) We propose the following constrained maximum likelihood estimation for the tensor coefficient

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

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

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

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 (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 $1 - \exp(-C_1 \sum_k p_k)$,

$$\text{Loss}(\mathcal{B}_{\text{true}}, \hat{\mathcal{B}}) \leq \frac{C_2 \sum_k p_k}{\prod_k d_k},$$

where $C_2 = C_2(\alpha, K, r_1, \dots, r_K) > 0$ 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 $d_1 = d_2 = \dots = d_K = d$. 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 $\mathcal{O}\left(\frac{p+r+s}{d^{mn}}\right)$. 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 $\mathbb{P}_{\mathcal{Y}_{\text{true}}}$ the distribution of \mathcal{Y} given the true $\mathcal{B}_{\text{true}}$ and $\mathbb{E}(\mathcal{Y}|\mathcal{X})$ the true mean. Let $\mathbb{P}_{\hat{\mathcal{Y}}}$ denote the distribution given the estimated $\hat{\mathcal{B}}$

Algorithm 1 Unsupervised binary tensor decomposition

- 1: **Input:** Y , the shape of core tensor r_1, r_2, r_3 ; the maximum iteration time N ; the link function f ; significant increment criterion ϵ
 - 2: **Output:** The estimate of factor matrices $\hat{M}_1, \hat{M}_2, \hat{M}_3$ and core tensor $\hat{\mathcal{C}}$.
 - 3: Use Tucker decomposition to get the initial $\mathcal{C}^{(0)}, M_1^{(0)}, M_2^{(0)}, M_3^{(0)}$; Calculate the initial log-likelihood value $l^{(0)}$; Set the iteration time $t = 0$
 - 4: **while** The increment of log-likelihood $l^{(t)} - l^{(t-1)} \geq \epsilon$ and $t \leq N$ **or** $t = 0$ **do**
 - 5: (a) Upgrade the iteration index $t = t + 1$;
 - 6: (b) Upgrade $\mathcal{C}^{(t-1)} \leftarrow \mathcal{C}^{(t)}$ by solveing one *glm* model with $r_1 \times r_2 \times r_3$ coefficients and link function f . The responses and predictors are formed by Y and $M_1^{(t-1)}, M_2^{(t-1)}, M_3^{(t-1)}$.
 - 7: **for** $i = 1, 2, 3$ **do**
 - 8: Upgrade $M_i^{(t-1)} \leftarrow M_i^{(t)}$ by solving d_i *glm* models. The responses and predictors are formed by Y and $M_j^{(t)}, M_l^{(t-1)}, j < i, l > i$.
 - 9: (c) Calculate the log-likelihood value $l^{(t)}$ with $\mathcal{C}^{(t)}, M_1^{(t)}, M_2^{(t)}, M_3^{(t)}$
 - 10: (d) After each upgrading step, orthogonalize M_i s through *tucker* decomposition,
-

and $\mathbb{E}(\widehat{\mathcal{Y}}|\mathcal{X})$ the predicted mean. We have, with probability at least $1 - \exp(-C_1 \sum_k d_k)$,

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

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

4.1.1 Comparison between Gaussian and non-Gaussian models

5 Numerical implementation

5.1 Alternating optimization

5.2 Tuning parameter selection

We propose to use Bayesian information criterion (BIC) and choose the rank that minimizes BIC; i.e.

$$\hat{\mathbf{r}} = \arg \min_{\mathbf{r}=(r_1, \dots, r_K)} \text{BIC}(\mathbf{r})$$

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

where $p_e(\mathbf{r}) \stackrel{\text{def}}{=} \sum_k (d_k - 1)r_k + \prod_k r_k$ 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.

6 Simulation

7 Data analysis

8 SUPPLEMENTARY MATERIAL

If you need to include additional appendices during submission, you can include them in the supplementary material file.

9 INSTRUCTIONS FOR CAMERA-READY PAPERS

For the camera-ready paper, if you are using L^AT_EX, please make sure that you follow these instructions. (If you are not using L^AT_EX, please make sure to achieve the same effect using your chosen typesetting package.)

1. Download `fancyhdr.sty` – the `aistats2020.sty` file will make use of it.
2. Begin your document with

```
\documentclass[twoside]{article}
\usepackage[accepted]{aistats2020}
```

The `twoside` option for the class `article` allows the package `fancyhdr.sty` to include headings for even and odd numbered pages. The option `accepted` for the package `aistats2020.sty` will write a copyright notice at the end of the first

column of the first page. This option will also print headings for the paper. For the *even* pages, the title of the paper will be used as heading and for *odd* pages the author names will be used as heading. If the title of the paper is too long or the number of authors is too large, the style will print a warning message as heading. If this happens additional commands can be used to place as headings shorter versions of the title and the author names. This is explained in the next point.

3. If you get warning messages as described above, then immediately after `\begin{document}`, write

```
\runningtitle{Provide here an  
alternative shorter version of the title  
of your paper}  
\runningauthor{Provide here the surnames  
of the authors of your paper, all  
separated by commas}
```

Note that the text that appears as argument in `\runningtitle` will be printed as a heading in the *even* pages. The text that appears as argument in `\runningauthor` will be printed as a heading in the *odd* pages. If even the author surnames do not fit, it is acceptable to give a subset of author names followed by “et al.”

4. Use the file `sample_paper.tex` as an example.
5. The camera-ready versions of the accepted papers are 8 pages, plus any additional pages needed for references.
6. If you need to include additional appendices, you can include them in the supplementary material file.
7. Please, don't change the layout given by the above instructions and by the style file.

Acknowledgements

Use the unnumbered third level heading for the acknowledgements. All acknowledgements go at the end of the paper.

References

References follow the acknowledgements. Use an unnumbered third level heading for the references section. Any choice of citation style is acceptable as long as you are consistent. Please use the same font size for references as for the body of the paper—remember that references do not count against your page length total.

References

- J. Alspector, B. Gupta, and R. B. Allen (1989). Performance of a stochastic learning microchip. In D. S. Touretzky (ed.), *Advances in Neural Information Processing Systems 1*, 748–760. San Mateo, Calif.: Morgan Kaufmann.
- F. Rosenblatt (1962). *Principles of Neurodynamics*. Washington, D.C.: Spartan Books.
- G. Tesauro (1989). Neurogammon wins computer Olympiad. *Neural Computation* **1**(3):321–323.