
Learning Multiple Networks via Supervised Tensor Decomposition

Anonymous Author(s)

Affiliation

Address

email

Abstract

We develop a tensor decomposition method that incorporates consider the problem of tensor decomposition with multiple side information available as interactive features. Unlike classical tensor decomposition, our supervised decomposition captures the effective dimension reduction of the data tensor confined to feature space on each mode Such problems are common in neuroimaging, network modeling, and spatial-temporal analysis. We develop a new family of exponential tensor decomposition models and establish the theoretical accuracy guarantees. An efficient alternating optimization algorithm is further developed. Our Unlike earlier methods, our proposal handles a broad range of data types, including continuous, count, and binary observations, along with available features. We apply the method to diffusion tensor imaging data from human connectome project. We and identify the key global brain connectivity pattern and pinpoint the local regions that are brain connectivity patterns associated with available features. Our method will help the practitioners efficiently analyze tensor datasets in various areas. Toward this end, all data and code has been made available to the public.

1 Introduction

Higher-order tensors have received increased attention across science and engineering. While most tensor decomposition methods are developed for a single tensor observation, scientific studies often collect side information, in the form of node features and interactions thereof, together with the tensor data. Such data problems are common in neuroimaging, network analysis, and spatial-temporal modeling. A popular example is in neuroimaging (Zhou et al., 2013). The brain connectivity networks are collected from a sample of individuals, accompanied by individual characteristics such as age, gender, and diseases status (see Figure 1a). Another example is in network analysis (Berthet and Baldin, 2020; Hoff, 2005). A typical social network consists of nodes that represent people and edges that represent friendships. (Berthet and Baldin, 2020). Side information such as people's demographic information and friendship types are often available. In both examples, it is of keen scientific interest to identify scientists are interested in identifying the variation in the data tensor (e.g., brain connectivities, social community patterns) that is affected by available features. These seemingly different scenarios pose a common yet challenging problem for tensor data modeling.

In addition to the challenge of incorporating side information, many tensor datasets consist of non-Gaussian measurements. Classical tensor decomposition methods are based on minimizing the Frobenius norm of deviation, leading to suboptimal predictions for binary- or count-valued response variables. A number of supervised tensor methods have been proposed (Narita et al., 2012; Zhao et al., 2012; Yu and Liu, 2016; Lock and Li, 2018). These methods often assume Gaussian distribution for the tensor entries, or impose random designs for the feature matrices, both of which are less suitable

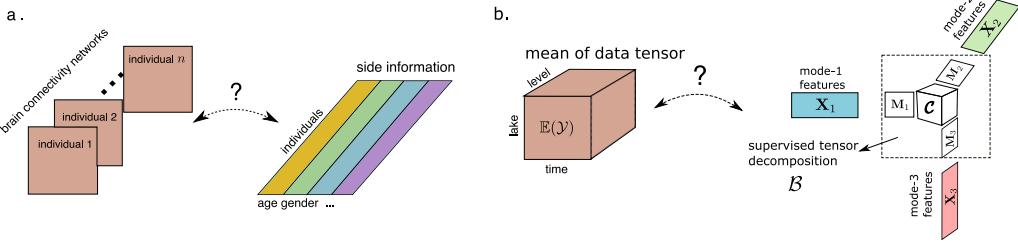


Figure 1: Examples of supervised tensor decomposition with interactive side information. (a) Network population model. (b) Spatio-temporal growth model.

for applications of our interest. The gap between theory and practice means a great opportunity to modeling paradigms and better capture the complexity in tensor data.

We present **Our contribution**. This paper presents a general model and associated method for decomposing a data tensor whose entries are from exponential family with interactive side information. We formulate the learning task as a **structured low-rank tensor** regression problem, with tensor observation serving as the response, and the multiple side information as interactive features. We leverage blend the modeling power of generalized linear model (GLM) to allow heteroscedacity due to the mean-variance relationship in the non-Gaussian data. The low-rank structure on the conditional mean tensor effectively mitigates the curse of high dimensionality. Our proposal blends the modeling power of GLM and the exploratory capability of tensor dimension reduction in order to take the best out of both worlds. Our methods greatly improves the classical tensor decomposition, and we quantify the gain in prediction through numerical experiments and data applications.

2 Method

1.1 Preliminary

We introduce the basic tensor properties used in the paper. We use lower-case letters (e.g., a, b, c) for scalars and 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. Let $\mathcal{Y} = [\![y_{i_1, \dots, i_K}]\!] \in \mathbb{R}^{d_1 \times \dots \times d_K}$ denote an order- K (d_1, \dots, d_K)-dimensional tensor. Notation. We follow the tensor notation as in Kolda and Bader (2009). The multilinear multiplication of a tensor $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$ by matrices $\mathbf{X}_k = [\![x_{i_k, j_k}]\!] \in \mathbb{R}^{p_k \times d_k}$ is defined as $\mathcal{Y} \times \{\mathbf{X}_1, \dots, \mathbf{X}_K\} = [\![\sum_{i_1, \dots, i_K} y_{i_1, \dots, i_K} x_{j_1, i_1}^{(1)} \dots x_{j_K, i_K}^{(K)}]\!]$, which results in an order- K (p_1, \dots, p_K)-dimensional tensor. For any two tensors $\mathcal{Y} = [\![y_{i_1, \dots, i_K}]\!]$, $\mathcal{Y}' = [\![y'_{i_1, \dots, i_K}]\!]$ of identical order and dimensions, their inner product The inner product between two tensors of equal size 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 tensor Frobenius norm is defined as $\|\mathcal{Y}\|_F = \langle \mathcal{Y}, \mathcal{Y} \rangle^{1/2}$, and the maximum norm is defined $\|\mathcal{Y}\|_\infty = \max_{i_1, \dots, i_K} y_{i_1, \dots, i_K}$. We let I_d denote the $d \times d$ identity matrix and $[d]$ denote the d -set $\{1, \dots, d\}$. $\langle \mathcal{Y}, \mathcal{Y}' \rangle = \sum_{i_1, \dots, i_K} y_{i_1, \dots, i_K} y'_{i_1, \dots, i_K}$. For ease of notation, we allow basic arithmetic operators (e.g., $+$, $-$) and univariate functions $f: \mathbb{R} \rightarrow \mathbb{R}$ to be applied to tensors in an element-wise manner.

1.1 General Model

2 Proposed models and motivating examples

Let $\mathcal{Y} = [\![y_{i_1, \dots, i_K}]\!] \in \mathbb{R}^{d_1 \times \dots \times d_K}$ denote an order- K data tensor. Suppose the side information is available on each of the K modes. Let $\mathbf{X}_k = [\![x_{i,j}]\!] \in \mathbb{R}^{d_k \times p_k}$ denote the feature matrix on the mode $k \in [K]$, where $x_{i,j}$ denotes the j -th feature value for the i -th tensor entity, for $(i, j) \in [d_k] \times [p_k]$, $p_k \leq d_k$. Assume We assume that, conditional on the features \mathbf{X}_k , the entries of tensor \mathcal{Y} are independent realizations from an exponential family distribution, and the conditional mean tensor

71 admits the form

$$\begin{aligned} \mathbb{E}(\mathcal{Y}|\mathbf{X}_1, \dots, \mathbf{X}_K) &= f(\Theta), \text{with } \Theta = \left(\mathcal{C} \times \{\mathbf{X}_1 \mathbf{M}_1, \dots, \mathbf{X}_K \mathbf{M}_K\} \right), \\ \text{with } \mathbf{M}_k^T \mathbf{M}_k &= \mathbf{I}_{r_k}, \mathbf{M}_k \in \mathbb{R}^{p_k \times r_k} \quad \text{for all } k = 1, \dots, K. \end{aligned} \quad (1)$$

72 where $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$ is the multilinear predictor, $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$ is the unknown parameter
73 tensor, $f(\cdot)$ is a known link function whose form depending on the data type of \mathcal{Y} . The choice
74 of link function is based on the assumed distribution family of tensor entries.

75 In classical tensor decomposition, tensor factorization is performed on either data tensor \mathcal{Y} or mean
76 tensor $\mathbb{E}(\mathcal{Y})$. In the context of supervised tensor decomposition, we propose to factorize the latent
77 parameter tensor \mathcal{B} ,

$$\mathcal{B} = \mathcal{C} \times \{\mathbf{M}_1, \dots, \mathbf{M}_K\},$$

78 where $\mathcal{C} \in \mathbb{R}^{r_1 \times \dots \times r_K}$ is a full-rank core tensor, and $\mathbf{M}_k \in \mathbb{R}^{p_k \times r_k}$ are factor matrices consisting of
79 orthonormal columns, where $r_k \leq p_k$ for all $k \in [K]$. By the definition of multilinear rank, model
80 equations and imply the low-rankness $\mathbf{r} = (r_1, \dots, r_K)$ of the conditional mean tensor under the
81 link function. We now reach our final model for supervised tensor decomposition,

$$\begin{aligned} \mathbb{E}(\mathcal{Y}|\mathbf{X}_1, \dots, \mathbf{X}_K) &= f(\mathcal{C} \times \{\mathbf{X}_1 \mathbf{M}_1, \dots, \mathbf{X}_K \mathbf{M}_K\}), \\ \text{with } \mathbf{M}_k^T \mathbf{M}_k &= \mathbf{I}_{r_k}, \mathbf{M}_k \in \mathbb{R}^{p_k \times r_k} \quad \text{for all } k = 1, \dots, K, \end{aligned}$$

82 where the parameters of interest are \mathbf{M}_k and \mathcal{C} . Note that model assumes a fixed, known rank
83 $\mathbf{r} = (r_1, \dots, r_K)$, $f(\cdot)$ is a known link function whose form depending on the data type of \mathcal{Y} .
84 Common choices of link functions include identity link for Gaussian distribution, logistic link for
85 Bernoulli distribution, and $\exp(\cdot)$ link for Poisson distribution.

86 Figure 1b provides a schematic illustration of our model. The features X_k affect the distribution of
87 tensor entries in \mathcal{Y} through the form $\mathbf{X}_k \mathbf{M}_k$, which are r_k linear combinations of features on mode
88 k . We call $\mathbf{X}_k \mathbf{M}_k$ the “supervised tensor factors” or “sufficient features” (Adragni and Cook, 2009)
89 , and call \mathbf{M}_k the “dimension reduction matrix.” The core tensor \mathcal{C} collects the interaction effects
90 between sufficient features across K modes, which links the conditional mean to the feature spaces,
91 and thereby allows the identification of variations in the tensor data attributable to the side information.
92 Our goal is to find \mathbf{M}_k and the corresponding \mathcal{C} , thereby allowing us to reveal the relationship between
93 side information \mathbf{X}_k and the observed tensor \mathcal{Y} . Note that \mathbf{M}_k and \mathcal{C} are identifiable only up to
94 orthonormal transformations.

95 3 Estimation

96 2.1 Rank-constrained M-estimator

97 We adopt the exponential family as a flexible framework for different datatypes. In a classical
98 generalized linear model with a scalar response y and feature x , the density is expressed as

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

99 where $b(\cdot)$ is a known function depending on the data types, θ is We give two examples of supervised
100 tensor decomposition models (1) that arise in practice.

101 **Example 1** (Spatio-temporal growth model). The growth curve model (Srivastava et al., 2008) was
102 originally proposed as an example of bilinear model for matrix data, and we extend it to higher-order
103 cases. Let $\mathcal{Y} = [\mathcal{Y}_{ijk}] \in \mathbb{R}^{d \times m \times n}$ denote the pH measurements of d lakes at m levels of depth
104 and for n time points. Suppose the sampled lakes belong to q types, with p lakes in each type.
105 Let $\{\ell_j\}_{j \in [m]}$ denote the sampled depth levels and $\{t_k\}_{k \in [n]}$ the time points. Assume that the
106 expected pH trend in depth is a polynomial of order at most r and that the expected trend in time
107 is a polynomial of order s . Then, the conditional mean model for the spatio-temporal growth
108 is a special case of our model (1), where $\mathbf{X}_1 = \text{blockdiag}\{\mathbf{1}_q, \dots, \mathbf{1}_q\} \in \{0, 1\}^{d \times p}$ is the linear

109 predictor, $\phi > 0$ is the dispersion parameter, and $c(\cdot)$ is a known normalizing function. In our
 110 context, we model the tensor entries y_{i_1, \dots, i_K} , conditional on θ_{i_1, \dots, i_K} , as independent draws from an
 111 exponential family. design matrix for lake types.

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

112 are the design matrices for spatial and temporal effects, respectively. The spatial-temporal mode has
 113 covariates available on each of the three modes.

114 **Example 2** (Network population model). Network response model is recently developed
 115 in the context of neuroimaging analysis. The goal is to study the relationship between
 116 network-valued response and the individual covariates. Suppose we observe n i.i.d. observations
 117 $\{(\mathbf{Y}_i, \mathbf{x}_i) : i = 1, \dots, n\}$, where $\mathbf{Y}_i \in \{0, 1\}^{d \times d}$ is the brain connectivity network on the i -th
 118 individual, and $\mathbf{x}_i \in \mathbb{R}^p$ is the individual covariate such as age, gender, cognition, etc. The
 119 network-response model (Rabusseau and Kadri, 2016) has the form

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

120 where $\mathcal{B} \in \mathbb{R}^{d \times d \times p}$ is the coefficient tensor of interest. The model (2) is also a special case of our
 121 tensor-response model, with covariates on the last mode of the tensor.

122 3 Estimation algorithms

123 We develop a likelihood-based procedure to estimate \mathcal{C} and \mathbf{M}_k in (1). Ignoring constants that do
 124 not depend on Θ , the quasi log-likelihood of (1) is equal to Bregman distance between \mathcal{Y} and $b'(\Theta)$:

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

125 where $b(\theta) = \theta^2/2$ for Gaussian response, $b(\theta) = \exp(\theta)$ for Poisson response, and
 126 $b(\theta) = \log(1 + \exp(\theta))$ for Bernoulli response. We propose a constrained maximum quasi-
 127 likelihood estimator (M-estimator),

$$(\hat{\mathcal{C}}, \hat{\mathbf{M}}_1, \dots, \hat{\mathbf{M}}_k) = \arg \max_{(\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K) \in \mathcal{P}} \mathcal{L}_{\mathcal{Y}}(\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K), \quad (3)$$

128 where parameter space $\mathcal{P} = \left\{ (\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K) \mid \mathbf{M}_k^T \mathbf{M}_k = \mathbf{I}_{r_k}, \|\Theta(\mathcal{C}, \mathbf{M}_1, \dots, \mathbf{M}_K)\|_{\infty} \leq \alpha \right\}.$
 129 The maximum norm constraint on the linear predictor Θ is a technical condition to avoid the
 130 divergence in the non-Gaussian variance.

131 3.1 Alternating optimization

132 The decision variables in the objective function (3) consist of $K + 1$ blocks of variables, one for the
 133 core tensor \mathcal{C} and K for the factor matrices \mathbf{M}_k . We notice that, if any K out of the $K + 1$ blocks of
 134 variables are known, then the optimization reduces to a simple GLM with respect to the last block
 135 of variables. This observation leads to an iterative updating scheme for one block at a time while
 136 keeping others fixed. A simplified version of the algorithm is described in Algorithm 1.

137 3.1 Statistical properties

138 In modern applications, the tensor data and features are often large-scale. We are particularly
 139 interested in the high-dimensional regime in which both d_k and p_k diverge; i.e. $d_k \rightarrow \infty$ and
 140 $p_k \rightarrow \infty$, while $p_k/d_k \rightarrow \gamma_k \in [0, 1]$. As the size of problem grows, and so does the number of
 141 unknown parameters. The classical MLE theory does not directly apply. We leverage the We provide
 142 the accuracy guarantee for the proposed M-estimator (3) by leveraging recent development in random
 143 tensor theory and high-dimensional statistics to establish the error bounds.

144 Consider a data tensor generated from model .

Algorithm 1 Supervised Tensor Decomposition with Interactive Side Information (Simplified)

Input: Response tensor $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$, feature 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 , maximum norm bound α

Output: Estimated core tensor $\hat{\mathcal{C}} \in \mathbb{R}^{r_1 \times \dots \times r_K}$ and factor matrices $\hat{\mathbf{M}}_k \in \mathbb{R}^{p_k \times r_k}$.

- 1: Random initialization of the core tensor \mathcal{C} and factor matrices \mathbf{M}_k .
- 2: **while** Do until convergence **do**
- 3: Obtain $\tilde{\mathbf{M}}_k \in \mathbb{R}^{p_k \times r_k}$ by a GLM. Orthogonalize $\tilde{\mathbf{M}}_k$ by QR factorization, for $k \in [K]$.
- 4: Update the core tensor \mathcal{C} by solving a GLM. Rescale the core tensor \mathcal{C} such that $\|\mathcal{C}\|_{\max} \leq \alpha$.
- 5: **end while**

145 **Theorem 3.1** (Convergence). Let $(\hat{\mathcal{C}}, \hat{\mathbf{M}}_1, \dots, \hat{\mathbf{M}}_K)$ be the M-estimator in (3) and $\hat{\mathcal{B}} = \hat{\mathcal{C}} \times \hat{\mathbf{M}}_1 \times \dots \times \hat{\mathbf{M}}_K$. Define $r_{\text{total}} = \prod_k r_k$ and $r_{\max} = \max_k r_k$. Under mild technical assumptions, there
146 exist two positive constants $C_1 = C_1(\alpha, K), C_2 = C_2(\alpha, K) > 0$ independent of dimensions $\{d_k\}$
147 and $\{p_k\}$ $C_1, C_2 \geq 0$, such that, with probability at least $1 - \exp(-C_1 \sum_k p_k)$,

$$\|\mathcal{B}_{\text{true}} - \hat{\mathcal{B}}\|_F^2 \leq \frac{C_2 r_{\text{total}}}{r_{\max}} \frac{\sum_k p_k}{\prod_k d_k}, \quad \text{and} \quad \sin^2 \Theta(\mathbf{M}_{k,\text{true}}, \hat{\mathbf{M}}_k) \leq \frac{C_2 r_{\text{total}}}{r_{\max} \sigma_{\min}^2(\text{Unfold}_k(\mathcal{C}_{\text{true}}))} \frac{\sum_k p_k}{\prod_k d_k},$$

149 Furthermore, if the unfolded core tensor has non-degenerate singular values at mode $k \in [K]$ where
150 $\sin \Theta(\mathbf{M}_{k,\text{true}}, \hat{\mathbf{M}}_k) = \|\mathbf{M}_{k,\text{true}}^T \hat{\mathbf{M}}_k^\perp\|_\sigma$ is the angle distance between column spaces.

151 Theorem 3.1 implies that the estimation has a convergence rate $\mathcal{O}(d^{-(K-1)})$ in the special case
152 when tensor dimensions are equal on each of the modes, i.e., $\sigma_{\min}(\text{Unfold}_k(\mathcal{C}_{\text{true}})) \geq c > 0$ for
153 some constant c , then

$$\sin^2 \Theta(\mathbf{M}_{k,\text{true}}, \hat{\mathbf{M}}_k) \leq \frac{C_2 r_{\text{total}}}{r_{\max} \sigma_{\min}^2(\text{Unfold}_k(\mathcal{C}_{\text{true}}))} \frac{\sum_k p_k}{\prod_k d_k},$$

154 where $\sin \Theta(\mathbf{M}_{k,\text{true}}, \hat{\mathbf{M}}_k) = \|\mathbf{M}_{k,\text{true}}^T \hat{\mathbf{M}}_k^\perp\|_\sigma = \max \left\{ \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} : \mathbf{x} \in \text{Span}(\mathbf{M}_{k,\text{true}}), \mathbf{y} \in \text{Span}(\hat{\mathbf{M}}_k^\perp) \right\}$
155 is the angle distance to assess the accuracy in estimating the column space, $\text{Span}(\mathbf{M}_k)$. $d_k = d$
156 for all $k \in [K]$, and feature dimension grows with tensor dimension, $p_k = \gamma d$, $\gamma \in [0, 1]$, for
157 $k \in [K]$. The convergence of our estimation becomes especially favorable as the order of tensor
158 data increases.

159 **4 Real-data-analysis**

160 **4 Numerical experiments**

161 We compare our supervised tensor decomposition (STD) with three other supervised tensor
162 methods: Higher-order low-rank regression (**HOLRR** Rabusseau and Kadri (2016)), Higher-order
163 partial least square (**HOPLS** Zhao et al. (2012)) and Subsampled tensor projected gradient
164 (**TPG** Yu and Liu (2016)). Figure 2 shows that STD outperforms others, especially in the low-signal,
165 high-rank setting. As the number of informative modes (i.e., modes with available features)
166 increases, the STD exhibits a substantial reduction in error whereas others remain unchanged
167 (Figure 2b). This showcases the benefit of incorporation of multiple features. The accuracy gain
168 in Figure 2 demonstrates the benefit of alternating algorithm – incorporation of informative modes
169 also improves the estimation in the non-informative modes.

170 The

171 We then apply our method to brain structural connectivity networks from Human Connectome
172 Project (HCP) aims to build a network map that characterizes the anatomical and functional
173 connectivity within healthy human brains (Geddes, 2016). We follow the preprocessing procedure
174 as in Zhang et al. (2018) and parcellate the brain into 68 regions of interest (Desikan et al., 2006).
175 The dataset consists of 136 brain structural networks, one for each individual. Each brain network is
176 represented as a 68-by-68 binary matrix, where the entries encode the presence or absence of fiber

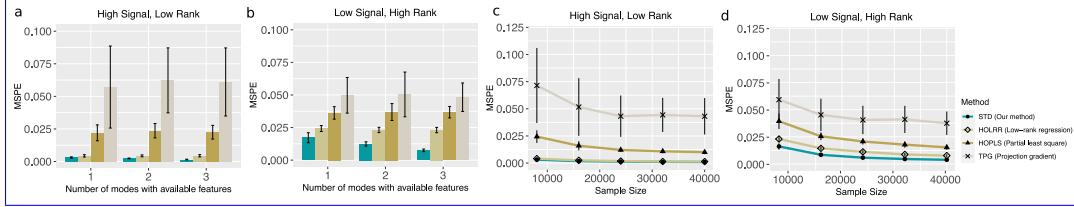


Figure 2: Comparison between different tensor methods. Panels (a) and (b) plot mean squared prediction error (MSPE) versus the number of modes with available features. Panels (c) and (d) plot MSPE versus the effective sample size d^2 . We consider rank $r = (3, 3, 3)$ (low) vs $(4, 5, 6)$ (high), and signal $\alpha = 3$ (low) vs. 6 (high).

177 connections between the 68 brain regions. We consider four individual features: gender (65 females
 178 vs. 71 males), age 22-25 ($n = 35$), age 26-30 ($n = 58$), and age 31+ ($n = 43$). The goal is to identify
 179 the connection edges that are affected by individual features.

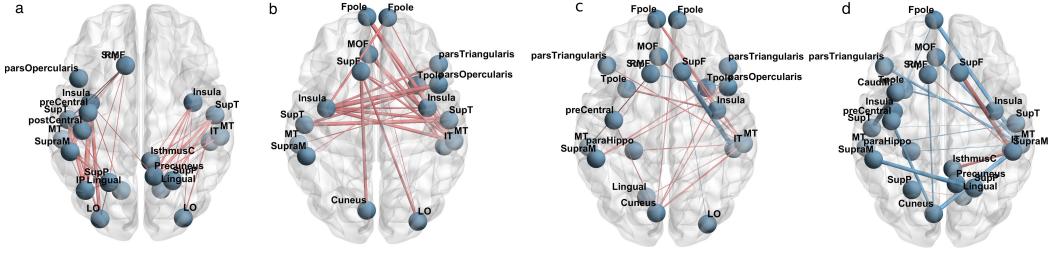


Figure 3: Top edges with large effects. (a) Global effect; (b) Female effect; (c) Age 22-25; (d) Age 31+. Red edges represent positive effects and (blue) edges represent positive (negative) effects. The edge-width is Edge-widths are proportional to the magnitudes of the effect sizesizes.

180 We perform the supervised tensor decomposition to the HCP data. The BIC selection suggests
 181 a rank $r = (10, 10, 4)$ with quasi log-likelihood $\mathcal{L}_Y = -174654.7$. We utilize the sum-to-zero
 182 contrasts in coding the feature effects, and depict only the top 3% edges whose connections are
 183 non-constant across the sample. Figure 3 shows the top edges with high effect size, overlaid on the
 184 Desikan atlas brain template (Desikan et al., 2006). We find that the global connection exhibits
 185 clear spatial separation, and that the nodes within each hemisphere are more densely connected
 186 with each other (Figure 3a). In particular, the superior-temporal (*SupT*), middle-temporal (*MT*) and
 187 Insula are the top three popular nodes in the network. Interestingly, female brains display higher
 188 inter-hemispheric connectivity, especially in the frontal, parietal and temporal lobes (Figure 3b). This
 189 is in agreement with a recent study showing that female brains are optimized for inter-hemispheric
 190 communication (Ingahalikar et al., 2014). We find several edges with declined connection in the group
 191 Age 31+. Those edges involve Frontal-pole (*Fpole*), superior-frontal (*SupF*) and Cuneus nodes. The
 192 Our results highlight the importance of Frontal-pole regionis known for its importance in memory
 193 and cognition, and the detected decline with age further highlights its biological importancefurther
 194 suggests the age effects to brain connections.

195 5 Conclusion

196 We have developed a supervised tensor decomposition method with side information on multiple
 197 modes. The empirical results demonstrate the improved interpretability and accuracy over
 198 previous approaches. Applications to the brain connection data yield conclusions with sensible
 199 interpretations, suggesting the practical utility of the proposed approach.

200 **Broader Impact**

201 Our supervised tensor decomposition method is widely applicable to network analysis, dyadic data
202 analysis, spatial-temporal model, and recommendation systems. We have shown the improved
203 predictive power and enhanced interpretability by incorporating the interactive side information in
204 tensor decomposition method. The application to the brain connection dataset yields conclusions with
205 sensible interpretations, suggesting the practical utility of the proposed approach. Tensor learning is a
206 clear challenge for further research. We believe that our model enriches the research of tensor-based
207 learning and is a powerful tool to boost scientific discoveries in various fields. We hope the work
208 opens up new inquiry that allows more machine learning researchers to contribute to this field.

209 **References**

- 210 Adragni, K. P. and Cook, R. D. (2009). Sufficient dimension reduction and prediction in regression.
211 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering*
212 *Sciences*, 367(1906):4385–4405.
- 213 Berthet, Q. and Baldin, N. (2020). Statistical and computational rates in graph logistic regression.
214 *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*,
215 108:2719–2730.
- 216 Desikan, R. S., Ségonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner,
217 R. L., Dale, A. M., Maguire, R. P., Hyman, B. T., et al. (2006). An automated labeling system
218 for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest.
219 *Neuroimage*, 31(3):968–980.
- 220 Geddes, L. (2016). Human brain mapped in unprecedented detail. *Nature*.
- 221 Hoff, P. D. (2005). Bilinear mixed-effects models for dyadic data. *Journal of the American Statistical*
222 *Association*, 100(469):286–295.
- 223 Ingallalikar, M., Smith, A., Parker, D., Satterthwaite, T. D., Elliott, M. A., Ruparel, K., Hakonarson,
224 H., Gur, R. E., Gur, R. C., and Verma, R. (2014). Sex differences in the structural connectome of
225 the human brain. *Proceedings of the National Academy of Sciences*, 111(2):823–828.
- 226 Kolda, T. G. and Bader, B. W. (2009). Tensor decompositions and applications. *SIAM review*,
227 51(3):455–500.
- 228 Lock, E. F. and Li, G. (2018). Supervised multiway factorization. *Electronic journal of statistics*,
229 12(1):1150.
- 230 Narita, A., Hayashi, K., Tomioka, R., and Kashima, H. (2012). Tensor factorization using auxiliary
231 information. *Data Mining and Knowledge Discovery*, 25(2):298–324.
- 232 Rabusseau, G. and Kadri, H. (2016). Low-rank regression with tensor responses. In *Advances in*
233 *Neural Information Processing Systems*, pages 1867–1875.
- 234 Srivastava, M. S., von Rosen, T., and Von Rosen, D. (2008). Models with a kronecker product
235 covariance structure: estimation and testing. *Mathematical Methods of Statistics*, 17(4):357–370.
- 236 Yu, R. and Liu, Y. (2016). Learning from multiway data: Simple and efficient tensor regression. In
237 *International Conference on Machine Learning*, pages 373–381.
- 238 Zhang, Z., Descoteaux, M., Zhang, J., Girard, G., Chamberland, M., Dunson, D., Srivastava, A., and
239 Zhu, H. (2018). Mapping population-based structural connectomes. *NeuroImage*, 172:130–145.
- 240 Zhao, Q., Caiafa, C. F., Mandic, D. P., Chao, Z. C., Nagasaka, Y., Fujii, N., Zhang, L., and Cichocki,
241 A. (2012). Higher order partial least squares (HOPLS): a generalized multilinear regression method.
242 *IEEE transactions on pattern analysis and machine intelligence*, 35(7):1660–1673.
- 243 Zhou, H., Li, L., and Zhu, H. (2013). Tensor regression with applications in neuroimaging data
244 analysis. *Journal of the American Statistical Association*, 108(502):540–552.