Convergence Rate for Predictive Densities of Bayesian Generalized Linear Models with a Scalar Response and Symmetric Tensor Predictor

Rajarshi Guhaniyogi^{a,*}, Sharmistha Guha^b

^a1156 High Street, Santa Cruz, CA 95064 ^b206 Old Chem. Bldg., Durham, NC 27708

Abstract

This article investigates statistical convergence rates for predictive densities of a novel Bayesian generalized linear model (GLM) framework with a scalar response and a symmetric tensor predictor with labeled "nodes." GLM frameworks involving a symmetric tensor predictor and a scalar response may appear in a variety of real life applications, including diffusion weighted magnetic resonance imaging (DWI) and functional magnetic resonance imaging (fMRI), among others. This article specifically focuses on a class of such models where the over-arching goal is to identify nodes and cells of the symmetric tensor influential in predicting the response. We establish a near optimal convergence rate for the posterior predictive density from the proposed model to the true density, depending on how the number of tensor nodes grows with the sample size. Moreover, we show that the method has adaptivity to the unknown rank of the true tensor, i.e., the near optimal rate is achieved even if the rank of the true tensor coefficient is not known a priori.

Keywords: Low-rank tensor decomposition, Posterior convergence rate, Symmetric tensor, Spike-and-slab prior

Preprint submitted to Journal of Multivariate Analysis

^{*}Corresponding author

Email addresses: rguhaniy@ucsc.edu (Rajarshi Guhaniyogi), sharmistha.guha@duke.edu (Sharmistha Guha)

1. Introduction

Of late, scientific applications often involve predictors having a multidimensional array or tensor structure, which are higher order analogues to vectors and matrices. Analogous to the rows and columns of a matrix, various axes of a tensor are known as tensor modes and the indices of a tensor mode are often re-

- ferred to as "tensor nodes." Entries in a tensor are known as "tensor cells." This article considers *symmetric* tensors, which are invariant upon interchanging the modes. We specifically focus on developing a regression relationship between a scalar response and a symmetric tensor predictor, with the ability of identifying
- tensor nodes influential in predicting the response. One major application of such modeling framework appears in brain connectome data, where the goal is to predict a brain related phenotype from the brain connectome network of subjects, with an emphasis of drawing inference on brain regions of interests(ROIs) related to the phenotype [1].
- In developing a modeling approach to address our problem of interest, one can possibly proceed to vectorize the symmetric tensor response and regress it on the predictors, leading to a high dimensional vector regression problem [2, 3]. This approach is able to make use of the expanding literature on Bayesian high dimensional regression [4, 5] but appears to be less than adequate to achieve all
- ²⁰ of our inferential goals simultaneously for a few reasons. First, the ordinary high dimensional regression framework assumes the coefficients corresponding to the tensor cells to be exchangeable, although, intuitively, the coefficients related to the same tensor node should be correlated a priori. Second, the strategy of reshaping a symmetric tensor into a vector leads to a massive dimensional vector
- ²⁵ predictor with applications often involving a limited number of samples. From an inferential point of view, Bayesian high-dimensional regression frameworks may be statistically inefficient when the number of predictors far exceeds the sample size [6]. More importantly, identification of important tensor nodes is not one of the inferential objectives of these approaches. Recent developments ³⁰ on tensor regression [7, 8] provide a solution to the problem by exploiting the

tensor structure of the predictor in the model and prior development. However, these approaches do not generally take into account the symmetry constraint in the tensor predictor, tend to focus mainly on prediction, and are not specifically designed to detect important nodes impacting the response.

- A recent approach to address all the inferential objectives mentioned above is outlined in [1] in the context of symmetric predictor matrices. More specifically, [1] develop a novel shrinkage prior on the symmetric matrix coefficients by combining ideas from low-rank matrix factorization and the Bayesian shrinkage prior literature. The structure offers parsimony by allowing identification of impor-
- tant tensor node specific coefficient vectors using a spike-and-slab prior on them. The framework exhibits good empirical performance with precise predictive inference as well as accurate identification of important tensor nodes. Moreover, the proposed prior allows auto tuning of all the hyperparameters with Markov chain Monte Carlo chains showing reasonably rapid mixing. While [1] provide
- the methodological and empirical motivations regarding the prior construction, rigorous theoretical understanding of Bayesian symmetric tensor regressions is yet to be established. Furthermore, the modeling framework is introduced and tested under a linear regression framework with normally distributed response variables.
- The primary focus of this article is to extend the network regression idea of [1] to a generalized linear modeling framework with a scalar response and a symmetric tensor predictor, and develop optimal posterior contraction rate for the proposed framework. Specifically, we adopt a low-rank structure for the symmetric tensor predictor coefficient and assign a spike-and-slab prior on node
- specific latent vectors within the low-rank structure to determine the tensor nodes significantly related to the scalar response a posteriori. Our main contribution is in developing conditions on the ranks and magnitudes of the true tensor coefficients and the number of tensor nodes for the near optimal learning of the proposed GLM. Note that several influential articles have emerged
- ⁶⁰ in the last few years detailing conditions for posterior contraction in ordinary high dimensional regression models, both with various point-mass priors in the

many normal-means models [9, 10, 11], and with classes of continuous shrinkage priors [12, 13]. In contrast, there is a dearth of papers studying posterior contraction properties for generalized linear models with tensor predictors in the

- ⁶⁵ Bayesian paradigm. A few recent articles [14, 8] offer conditions for consistency or optimal rates for posterior contraction with tensor predictors without the symmetry constraint, and with a different class of multiway shrinkage priors [8]. As a result, the theoretical construction in 8 does not find ready extension to our framework. Additionally, we relax the key assumption in 8 that both
- the tensor predictor coefficient generating the data (also referred to as the *true* tensor coefficient) and the fitted tensor coefficient have the same low-rank decompositions. In practice, the rank of the true tensor coefficient is never known. The current article is based upon a more realistic assumption that the rank of the fitted tensor coefficient is greater than or equal to the rank of the true tensor coefficient.

75 coefficient.

The rest of the article proceeds as follows. Section 2 develops the notations, defines the GLM framework for the fitted model and the true data generating model, and details out the prior distributions on the parameters. Section 3 describes the posterior contraction rate results for the predictive distribution.

Finally Section 4 concludes the article with an eye towards future work. Proofs of all theoretical results can be found in Appendix A and B.

2. Problem Setting

2.1. Notations

A D-way tensor $\mathbf{\Gamma} \in \bigotimes_{l=1}^{D} \mathbb{R}^{V_l}$ is a multidimensional array whose $(k_1, ..., k_D)$ th cell is denoted by $\Gamma_{(k_1,...,k_D)}$, $1 \leq k_1 \leq V_1,..., 1 \leq k_D \leq V_D$. When D = 2, a tensor corresponds to a matrix. This article mainly focuses on symmetric tensors with dummy diagonal entries (set at 0 for definiteness) ensuring $V_1 = \cdots = V_D = V$ and $\Gamma_{(k_1,...,k_D)} = \Gamma_{(P(k_1),...,P(k_D))}$, for any permutation $P(\cdot)$ of $\{k_1,...,k_D\}$ and $\Gamma_{(k_1,...,k_D)} = 0$, if any two of the indices k_l and $k_{l'}$ are equal. Similar to row and column indices of a matrix, the indices $\mathcal{N} = \{1, 2, ..., V\} \text{ for symmetric tensors are referred to as tensor nodes. Let } \mathcal{K} = \{(k_1, ..., k_D) : 1 \leq k_1 < \cdots < k_D \leq V\} \text{ be a set of indices with cardinality } q = \frac{V(V-1)\cdots(V-D+1)}{D!}. \text{ While expressing a symmetric tensor } \boldsymbol{\Gamma} \text{ with 0 diagonal entries, it is enough to specify } \boldsymbol{\Gamma}_{\boldsymbol{k}} \text{ for } \boldsymbol{k} \in \mathcal{K}. \text{ This holds since for any } \boldsymbol{k} \notin \mathcal{K}$ ⁹⁵ \exists a permutation $P(\cdot)$ s.t. $(P(k_1), ..., P(k_D)) \in \mathcal{K}.$ Then, by the property of the symmetric tensor, $\boldsymbol{\Gamma}_{P(\boldsymbol{k})} = \boldsymbol{\Gamma}_{\boldsymbol{k}}.$ A symmetric tensor with 0 diagonal entries $\boldsymbol{\Gamma}$ assumes a rank-1 PARAFAC decomposition if $\boldsymbol{\Gamma}_{\boldsymbol{k}}$ for $\boldsymbol{k} \in \mathcal{K}$ can be expressed as $\boldsymbol{\Gamma}_{\boldsymbol{k}} = \gamma_{k_1} \cdots \gamma_{k_D}$, for $\boldsymbol{\gamma} = (\gamma_1, ..., \gamma_V)' \in \mathbb{R}^V$. A rank R symmetric PARAFAC decomposition expresses $\boldsymbol{\Gamma}_{\boldsymbol{k}}$ as $\boldsymbol{\Gamma}_{\boldsymbol{k}} = \sum_{r=1}^{R} \gamma_{k_1}^{(r)} \cdots \gamma_{k_D}^{(r)}$, where $\boldsymbol{\gamma}^{(r)} = (\gamma_1^{(r)}, ..., \gamma_V^{(r)})' \in \mathbb{R}^V$. Importantly, for two symmetric tensors \boldsymbol{A} and \boldsymbol{B} with zero diagonal entries, the Frobenius inner product between \boldsymbol{A} and \boldsymbol{B} are given by $\langle \boldsymbol{A}, \boldsymbol{B} \rangle = D! \sum_{\boldsymbol{k} \in \mathcal{K}} a_{\boldsymbol{k}} b_{\boldsymbol{k}}.$ Finally, $||\boldsymbol{\Gamma}|| = \sqrt{\sum_{k_1=1}^{V_1} \cdots \sum_{k_D=1}^{V_D} \boldsymbol{\Gamma}_{(k_1,...,k_D)}^2$ and $||\boldsymbol{\Gamma}||_{\infty} = \max_{(k_1,...,k_D)} |\boldsymbol{\Gamma}_{(k_1,...,k_D)}|$ denote the l_2 and l_{∞} norms, respectively, for a tensor $\boldsymbol{\Gamma}$. The l_2 and l_{∞} norms of vectors are defined analogously.

105 2.2. Modeling Framework

For i = 1, ..., n, let y_i be the scalar response and $\mathbf{X}_i = ((x_{i,(k_1,...,k_D)}))_{k_1,...,k_D=1}^V \in \mathbb{R}^{V \times \cdots \times V}$ denote the symmetric tensor predictor with 0 diagonal entries. We assume that the data are generated from the generalized linear model given by the following density function

$$g_0(y_i|\boldsymbol{X}_i) = \exp(a(\alpha_0)y_i + b(\alpha_0) + c(y_i)), \ \alpha_0 = \sum_{\boldsymbol{k}\in\mathcal{K}} x_{i,\boldsymbol{k}}\Gamma_{0,\boldsymbol{k}},$$
(1)

where $\Gamma_{0,\mathbf{k}}$ corresponds to the $\mathbf{k} = (k_1, ..., k_D)$ -th entry of a symmetric tensor (with 0 diagonal entries) Γ_0 , a(h) and b(h) are continuously differentiable functions, with a(h) having a nonzero derivative. This parameterization includes some popular classes of densities, including binary logit and probit regressions of y on \mathbf{X} , Poisson regression of y on \mathbf{X} with count valued response, and normal regression with known error variance for continuous response y [15]. The conditional density of y_i given \mathbf{X}_i fitted to the data is also assumed to belong to the same class of generalized linear models, and is given by

$$g(y_i|\boldsymbol{X}_i) = \exp(a(\alpha)y_i + b(\alpha) + c(y_i)), \ \alpha = \sum_{\boldsymbol{k}\in\mathcal{K}} x_{i,\boldsymbol{k}}\Gamma_{\boldsymbol{k}},$$
(2)

where $\Gamma_{\mathbf{k}}$ is the \mathbf{k} -th entry of Γ , which is a symmetric tensor with 0 diagonal entries.

Suppose Γ_0 and Γ assume symmetric rank- R_0 and rank-R PARAFAC decompositions, respectively, for $R \geq R_0$, so that

$$\Gamma_{0,\boldsymbol{k}} = \sum_{r=1}^{R_0} \gamma_{0,k_1}^{(r)} \cdots \gamma_{0,k_D}^{(r)}, \ \Gamma_{\boldsymbol{k}} = \sum_{r=1}^R \lambda_r \gamma_{k_1}^{(r)} \cdots \gamma_{k_D}^{(r)}, \tag{3}$$

where $\boldsymbol{\gamma}^{(r)} = (\gamma_1^{(r)}, ..., \gamma_V^{(r)})$ and $\boldsymbol{\gamma}_0^{(r)} = (\gamma_{0,1}^{(r)}, ..., \gamma_{0,V}^{(r)}) \in \mathbb{R}^V$ for all r = 1, ..., R. Since the rank of the fitted symmetric tensor coefficient $\boldsymbol{\Gamma}$ is assumed to be higher than the rank of the true tensor coefficient $\boldsymbol{\Gamma}_0$, rank specific binary inclusion variables $\lambda_r \in \{0, 1\}$ are added in order to *switch-off* the contribution of unnecessary summands. The assumed low-rank decomposition offers parsimony by reducing the number of estimable parameters from $V(V-1)\cdots(V-D+1)/D!$ to RV, typically with $R \ll V$. When D = 2, the formulation assumes further simplification. To see this, denote $\tilde{\boldsymbol{\gamma}}_h = (\boldsymbol{\gamma}_h^{(1)}, ..., \boldsymbol{\gamma}_h^{(R)})'$, h = 1, ..., Vand $\boldsymbol{\Lambda} = diag(\lambda_1, ..., \lambda_R)$. The $\boldsymbol{k} = (k_1, k_2)$ th entry of $\boldsymbol{\Gamma}$ then simplifies as $\Gamma_{\boldsymbol{k}} = \tilde{\boldsymbol{\gamma}}'_{k_1} \boldsymbol{\Lambda} \tilde{\boldsymbol{\gamma}}_{k_2}, \, \boldsymbol{k} \in \mathcal{K}$, which represents a *bilinear* [16] interaction between $\tilde{\boldsymbol{\gamma}}_{k_1}$ and $\tilde{\boldsymbol{\gamma}}_{k_2}$. Accordingly, the significance of the \boldsymbol{k} th tensor cell of $\boldsymbol{\Gamma}$ in explaining the response increases with the similarity in the positions of $\tilde{\boldsymbol{\gamma}}_{k_1}$ and $\tilde{\boldsymbol{\gamma}}_{k_2}$.

¹²⁰ the similarity being measure variables in the latent space.

From (3), the *h*th tensor node of the symmetric tensor predictor X is deemed to have no impact on the response if $\tilde{\gamma}_h = \mathbf{0}$, $h \in \mathcal{N}$. The *k*th cell is considered unrelated to the response if $\Gamma_{\mathbf{k}} = 0$. Since $\Gamma_{\mathbf{k}} = 0$ if $\tilde{\gamma}_{k_l} = 0$ for some k_l , the proposed formulation assumes that the contribution of the *k*th cell of the tensor predictor to the response is insignificant if k_l th node is unrelated to the response, for some k_l . Our modeling framework is pertinent to a variety of applications, a few of them are presented below.

- Example 1 (Brain Connectome Data): In many neuroscientific applications, it is of interest to build a predictive model of a brain related phenotype (e.g., presence of a neuronal disease) on the connectivity network in a human brain (referred to as the *brain connectome*) (for e.g., see 17). To quantify brain connectivity, important regions of interest (ROI) in the brain are identified and
- the number of neurons connecting different ROIs is measured from the brain white matter using a brain imaging technique known as *diffusion tensor imaging* (DTI). Alternatively, the brain connectome tensor can also be constructed by computing the correlation of *functional magnetic resonance imaging* (fMRI) signals for different pairs of regions after suitably thresholding them to zero below
- ¹⁴⁰ a certain pre-specified cut-off. The inferential interest here lies in predicting the phenotypic response from the brain connectome matrix, as well as identifying ROIs significantly related to the response. This appears to be a direct application of (1), with y and X as the phenotype and the symmetric brain connectome matrix, respectively, and nodes in the matrix representing the ROIs. [1] ana-
- ¹⁴⁵ lyze this dataset with a regression framework similar to ours involving a scalar response and an undirected network predictor, assuming normally distributed errors in the regression.
- Example 2 (International Trade Data): Developing a regression relationship between world gross domestic product (GDP) and multilateral trade between countries is an informative exercise in international trade theory. Analysis of datasets with such information is important to statistically identify countries which are major economic drivers of the world, and also to direct significant world economic policies by international financial institutions [18, 19]. In the
- context of (1), the response and predictors would be the world GDP and multilateral trade relationships (which constitute a symmetric higher order tensor), respectively. The countries are the tensor nodes to draw inference on. In this context, it is generally believed that free trade agreements between countries could benefit the overall economic health of the world. For example, one can

consider the trilateral free trade agreement between China-Japan-South Korea [20], or between the U.S.-Canada-Mexico (referred to as the North Atlantic Free Trade Agreement or NAFTA) [21]. It is instructive to statistically analyze important economic outcomes like GDP in relation to such multi-lateral free trade agreements.

165 2.3. Prior Structure

To assess if the *h*th tensor node is active in predicting the response, we assign a spike-and-slab mixture prior distribution on $\tilde{\gamma}_h$ as

$$\tilde{\boldsymbol{\gamma}}_h \sim \zeta_h N(\boldsymbol{0}, \boldsymbol{I}) + (1 - \zeta_h) \delta_{\boldsymbol{0}}, \ \zeta_h \sim Ber(\Delta), \ \Delta \sim U(0, 1),$$
 (4)

where $\delta_{\mathbf{0}}$ is the Dirac-delta function at $\mathbf{0}$, Δ corresponds to the probability of the nonzero mixture component and ζ_h is a binary indicator set to 0 if $\tilde{\boldsymbol{\gamma}}_h = \mathbf{0}$. Thus, the posterior distributions of the ζ_h 's are analyzed to ascertain which nodes are influential in predicting the response. Notably, $(\tilde{\boldsymbol{\gamma}}_h, \zeta_h)$ are i.i.d. over

¹⁷⁰ h given Δ . Finally, to infer on how many ranks are necessary to express Γ , the rank specific binary inclusion variables, the λ_r 's, are assigned a hierarchical prior, $\lambda_r | \nu_r \stackrel{ind.}{\sim} Ber(\nu_r), \nu_r \stackrel{ind.}{\sim} Beta(1, r^{\eta})$, over r. Choosing $\eta > 1$ ensures increasing shrinkage on λ_r as r grows. Thus a low-rank solution to Γ is favored a priori, which helps avoid over-fitting.

Analysis of datasets using the model (2) involving a continuous scalar response and symmetric tensor predictors are available in some recent work [22], though a rigorous theoretical treatment of such models is missing in the literature. The overarching goal of this article is to develop theoretical conditions to draw optimal predictive inference from such models. It will be shown in due course that the posterior predictive loss (defined in Section 3) of our model de-

cays at the "near" optimal rate to 0 under fairly mild assumptions. Moreover, such theoretical results will be obtained for an easily computable posterior with standard Markov chain Monte Carlo updates for all the parameters.

3. Convergence Rate Analysis

This article assesses the predictive accuracy of the proposed model $g(y|\mathbf{X})$ in estimating the true model $g_0(y|\mathbf{X})$, following the notion of convergence described in 15. Define the Hellinger distance between g and g_0 as

$$d_H(g,g_0) = \sqrt{\int \int (\sqrt{g(y|\boldsymbol{X})} - \sqrt{g_0(y|\boldsymbol{X})})^2 \nu_y(dy) \nu_{\boldsymbol{X}}(d\boldsymbol{X})},$$

where $\nu_{\mathbf{X}}$ is the unknown probability measure for \mathbf{X} , and ν_y is the dominating measure for g and g_0 . We focus on showing $E_{g_0} \Pi[d_H(g, g_0) > \epsilon_n | \{y_i, \mathbf{X}_i\}_{i=1}^n] < \xi_n$, for large n, for some sequences ϵ_n, ξ_n converging to 0 as $n \to \infty$, where $\Pi(\mathcal{S}|\{y_i, \mathbf{X}_i\}_{i=1}^n)$ is the posterior probability of the set \mathcal{S} . The result implies that the posterior probability outside of a shrinking neighborhood around the

true predictive density g_0 converges to 0 as $n \to \infty$. Specifically, we focus on identifying conditions that lead to convergence rate ϵ_n of the order of $n^{-1/2}$ up to a $\log(n)$ factor.

3.1. Framework and Main Results

Without loss of generality, the predictor X_i satisfies $|x_{i,k}| < 1$ for all i and $k \in \mathcal{K}$. In what follows, we add the subscript n to the number of tensor nodes V_n , the rank R_n of Γ and rank $R_{0,n}$ of the true symmetric tensor coefficient Γ_0 . We assume V_n , R_n and $R_{0,n}$ are all non-decreasing functions of n, with $R_n < V_n$ and $R_n > R_{0,n}$ for all large n. Hence, the number of elements in \mathcal{K} , given by $q_n = V_n(V_n - 1)...(V_n - D + 1)/D!$, is a function of n. This paradigm attempts to capture the fact that q_n grows much faster than n, and a higher rank CP decomposition of Γ can be estimated more precisely in the presence of a larger sample size n.

One of the key quantities in proving posterior convergence rate results is the concentration of the prior distribution. The prior concentration can be quantified by $\mathcal{E}_n(\kappa)$, defined, for each $\kappa > 0$ by

$$\mathcal{E}_n(\kappa) = -\log\left\{\Pi(||\mathbf{\Gamma} - \mathbf{\Gamma}_0||_{\infty} \le \kappa)\right\}.$$
(5)

In order to achieve an optimal rate of convergence for the posterior, one expects the prior to put considerable mass around Γ_0 . Since Γ_0 is not known, it is not desirable to have a lot of prior mass around one point or a few points. Rather, the prior mass should be spread judiciously, taking into account the wide range of possibilities for Γ_0 . Prior concentration provides such a quantification of prior mass around the truth. Instead of characterizing the prior concentration function $\mathcal{E}_n(\kappa)$, we evaluate the prior concentration conditional on a set \mathcal{C} given by $\mathcal{C} = \{\lambda_1 = 1, ..., \lambda_{R_{0,n}} = 1, \lambda_{R_{0,n}+1} = 0, ..., \lambda_{R_n} = 0\}$, with Lemma 5.2 in the Appendix A quantifying a lower bound on $P(\mathcal{C})$. The prior concentration conditional on the set \mathcal{C} is given by

$$\mathcal{E}_n(\kappa|\mathcal{C}) = -\log\left\{\Pi(||\mathbf{\Gamma} - \mathbf{\Gamma}_0||_{\infty} \le \kappa|\mathcal{C})\right\}$$
(6)

Lemma 5.3 in Appendix A presents an upper bound on the conditional prior concentration corresponding to our proposed prior distribution in Section 2.3. We now state the main theorem involving the contraction of the fitted predictive density to the true predictive density.

Theorem 3.1. Define the function $H(\kappa) = 1 + \kappa \sup_{\substack{|w| \le \kappa}} |a'(w)| \sup_{\substack{|w| \le \kappa}} |b'(w)/a'(w)|$, where a'(w) and b'(w) are derivatives of the functions a(w) and b(w) in (1) and (2), respectively. For a sequence ϵ_n satisfying $0 < \epsilon_n < 1$, $n\epsilon_n^2 \to \infty$, and another sequence C_n , let the following conditions hold

(a) $R_n V_n \log(V_n) = o(n\epsilon_n^2)$

205

210

- (b) $R_n V_n \log(1/\epsilon_n^2) = o(n\epsilon_n^2)$
- (c) $R_n V_n \log(H(R_n C_n^D V_n^D)) = o(n\epsilon_n^2),$
- (d) $(1 \Phi(C_n)) \le e^{-4n\epsilon_n^2}$, for all large n
- 215 (e) $\limsup_{n \to \infty} \sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}|| < \infty$, where $\tilde{\gamma}_{0,h} = (\gamma_{0,h}^{(1)}, ..., \gamma_{0,h}^{(R_n)})'$.

Then,
$$E_{g_0} \prod\{d_H(g, g_0) > 4\epsilon_n | \{Y_i, x_i\}_{i=1}^n\} < 4e^{-n\epsilon_n^2/2}$$
, for all large n

The following remarks characterize $H(\kappa)$ and its implications for various regression settings under GLM.

Remark 1: For ordinary linear regression with normal errors, $H(\kappa)$ grows at most at the order of $|\kappa|^2$. Thus, assumption (c) becomes equivalent to $R_n V_n \log(C_n) = o(n\epsilon_n^2)$, considering assumption (a).

Remark 2: For binary regression with logit or probit links, $H(\kappa)$ grows at most linearly with $|\kappa|$. Thus, assumption (c) becomes equivalent to $R_n V_n \log(C_n) = o(n\epsilon_n^2)$, considering assumption (a).

For our theoretical exposition, we will focus on continuous and binary regression only. Theorem 3.1, together with the functional properties of $H(\kappa)$ mentioned, leads to the following result on the convergence rate ϵ_n of the proposed model.

Corollary 3.2. Let, $\limsup_{n\to\infty} \sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}|| < \infty$, where $\tilde{\gamma}_{0,h} = (\gamma_{0,h}^{(1)}, ..., \gamma_{0,h}^{(R_n)})'$. Assume that for some $0 < \xi < 1$, $V_n \leq M_1 n^{\xi}$ (for some constant $M_1 > 0$) and the tensor rank R_n grows at a much slower rate of $(\log n)^{z_1}$ for some z_1 , i.e., $R_n \leq M_2 (\log n)^{z_1}$, for some constant M_2 . Choose C_n such that $n^{\phi_1} \leq C_n \leq n^{\phi_2}$, satisfying $0 < \xi/2 < \phi_1 < \phi_2$. Then the convergence rate ϵ_n can be expressed as $\epsilon_n \sim n^{-(1-\xi)/2} (\log n)^{z_1/2+1}$ for the linear regression model, as well as the binary regression model with logistic or probit link functions.

Remark 3: Note that whatever be the value of z_1 , $(\log n)^{z_1/2+1} \le n^{\xi/2}$ for all large n, so that one can achieve a convergence rate of $n^{-(1-2\xi)/2}$. Depending on V_n , ξ can be made very small to achieve a rate close to the "finite-dimensional" rate of $n^{-1/2}$.

Remark 4: Note that the condition $\limsup_{n\to\infty} \sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}|| < \infty$ includes as a special case the scenario in which only a fixed and finite number of $||\tilde{\gamma}_{0,h}||$'s are nonzero, while also allowing a more realistic setup with many small $||\tilde{\gamma}_{0,h}||$, none of which are exactly zero. The convergence rate also depends on how V_n and R_n grows with n. In fact, the convergence rate deteriorates as ξ becomes higher, i.e., the number of tensor nodes grows faster as a function of n.

4. Conclusion

This article investigates the convergence rate of the predictive distribution for generalized linear models involving a scalar response and a symmetric tensor predictor. Under mild assumptions, we provide a "near optimal" convergence ²⁵⁰ rate for the predictive distribution of the proposed model. The theoretical results proved here allow the number of tensor cells to grow much faster than the sample size. The near optimal rate is rank adaptive, i.e., it holds even if the rank of the symmetric tensor coefficient for the true data generating regression model is unknown. Most importantly, the bound on the predictive accuracy is achieved for a prior that leads to an easily computable posterior, as observed in a few recent articles [1, 23].

Several future directions of research emerge from this article. For example, it might be of interest to relax assumption (e) in Theorem 3.1 and investigate convergence rate by allowing $\sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}||$ to vary slowly as an increasing function of *n*. Another interesting future direction constitutes extending this theoretical set up to prove the tensor node selection consistency for the proposed model.

5. Acknowledgement

The work of the first author is partially supported from Office of Naval Research award no. N00014-18-1-2741 and National Science Foundation DMS-1854662.

Appendix A

We begin by stating a series of lemmas. Lemma 5.1 provides a bound on the root of a monic polynomial. Lemma 5.2 quantifies the a lower bound on $P(\mathcal{C})$ where \mathcal{C} is defined in Section 3.1. Finally, Lemma 5.3 presents an upper

²⁷⁰ bound on the conditional prior concentration corresponding to our proposed prior distribution in Section 2.3. All these three lemmas will be crucial to prove Theorem 3.1. **Lemma 5.1.** Let J(x) be a monic polynomial given by $J(x) = x^D + b_{D-1}x^{D-1} + \cdots + b_1x - b_0, \ b_0, \dots, b_{D-1} \ge 0$. If x_0 is a real positive root of the equation J(x) = 0, then $1/x_0 \le 1 + (b_1/b_0)$.

Proof Let z = 1/x. Then J(x) = 0 implies J(1/z) = 0, i.e., $z^D - (b_1/b_0)z^{D-1} - \cdots - (b_{D-1}/b_0)z - (1/b_0) = 0$. Since this is a monic polynomial with $1/x_0$ as one of its positive real roots, by the Lagrange-Maclaurin theorem, $1/x_0 \le 1 + (b_1/b_0)$.

Lemma 5.2. For $\lambda_r | \nu_r \stackrel{ind.}{\sim} Ber(\nu_r)$ and $\nu_r \stackrel{ind.}{\sim} Beta(1, r^{\eta}), r = 1, ..., R_n$, and $\eta > 0$,

$$P(\mathcal{C}) = P(\lambda_1 = 1, ..., \lambda_{R_{0,n}} = 1, \lambda_{R_{0,n}+1} = 0, ..., \lambda_{R_n} = 0) \ge \frac{R_{0,n}^{\eta(R_n - R_{0,n})}}{(1 + R_{0,n}^{\eta})^{R_n}}.$$

Proof $P(\lambda_r = 1) = E(\nu_r) = \frac{1}{1+r^{\eta}}$ for $r = 1, ..., R_n$. Then,

$$P(\lambda_{1} = 1, ..., \lambda_{R_{0,n}} = 1, \lambda_{R_{0,n+1}} = 0, ..., \lambda_{R_{n}} = 0) = \prod_{r=1}^{R_{0,n}} \frac{1}{(1+r^{\eta})} \prod_{r=R_{0,n+1}}^{R_{n}} \frac{r^{\eta}}{(1+r^{\eta})}$$
$$\geq \frac{1}{(1+R_{0,n}^{\eta})^{R_{0,n}}} \left\{ \frac{R_{0,n}^{\eta}}{(1+R_{0,n}^{\eta})} \right\}^{R_{n}-R_{0,n}} = \frac{R_{0,n}^{\eta(R_{n}-R_{0,n})}}{(1+R_{0,n}^{\eta})^{R_{n}}}.$$

The first inequality follows due to the fact that $r^{\eta}/(1+r^{\eta})$ is a monotone increasing function of r and $1/(1+r^{\eta})$ is a monotone decreasing function of r.

Lemma 5.3. Let $\tilde{\gamma}_{0,h} = (\gamma_{0,h}^{(1)}, ..., \gamma_{0,h}^{(R_{0,n})})'$ and for $\mathbf{k} \in \mathcal{K}$, let $u_{\mathbf{k},n}$ be the only positive root of the equation

$$x\prod_{s=2}^{D}(x+||\tilde{\gamma}_{0,k_s}||)+||\tilde{\gamma}_{0,k_1}||x\prod_{s=3}^{D}(x+||\tilde{\gamma}_{0,k_s}||)+\dots+x\prod_{s=1}^{D-1}||\tilde{\gamma}_{0,k_s}||=v_n.$$
(7)

Assume $u_n = \min_{\mathbf{k} \in \mathcal{K}} u_{\mathbf{k},n}$. Then,

$$\begin{aligned} \mathcal{E}_n(v_n|\mathcal{C}) &\leq \sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}||^2 / 2 + \frac{R_{0,n}V_n}{2}\log\left(2\pi\right) + \log\left(\frac{R_{0,n}V_n + 1}{R_{0,n}V_n}\right) + R_{0,n}V_n\log\left(R_{0,n}/(2u_n)\right) \\ &+ V_n u_n^2 / R_{0,n} \end{aligned}$$

 $\begin{array}{l} \mathbf{Proof \ Let \ } \mathcal{J} = \{ \mathbf{\Gamma} : ||\mathbf{\Gamma} - \mathbf{\Gamma}_{0}||_{\infty} \leq \upsilon_{n} \}. \ \text{Under } \mathcal{C}, \ \text{for } \mathbf{k} \in \mathcal{K}, \\ |\Gamma_{\mathbf{k}} - \Gamma_{0,\mathbf{k}}| = |\sum_{r=1}^{R_{0,n}} \gamma_{k_{1}}^{(r)} \cdots \gamma_{k_{D}}^{(r)} - \sum_{r=1}^{R_{0,n}} \gamma_{0,k_{1}}^{(r)} \cdots \gamma_{0,k_{D}}^{(r)}| = |\sum_{r=1}^{R_{0,n}} (\gamma_{k_{1}}^{(r)} - \gamma_{0,k_{1}}^{(r)}) \prod_{s=2}^{D} \gamma_{k_{s}}^{(r)}| + \\ \cdots + |\sum_{r=1}^{R_{0,n}} (\gamma_{k_{D}}^{(r)} - \gamma_{0,k_{D}}^{(r)}) \prod_{s=1}^{D-1} \gamma_{0,k_{s}}^{(r)}| \leq ||\tilde{\boldsymbol{\gamma}}_{k_{1}} - \tilde{\boldsymbol{\gamma}}_{0,k_{1}}|| \prod_{s=2}^{D} ||\tilde{\boldsymbol{\gamma}}_{k_{s}}|| + \cdots + ||\tilde{\boldsymbol{\gamma}}_{k_{D}} - \\ \tilde{\boldsymbol{\gamma}}_{0,k_{D}}||\prod_{s=1}^{D-1} ||\tilde{\boldsymbol{\gamma}}_{0,k_{s}}|| \leq ||\tilde{\boldsymbol{\gamma}}_{k_{1}} - \tilde{\boldsymbol{\gamma}}_{0,k_{1}}|| \prod_{s=2}^{D} (||\tilde{\boldsymbol{\gamma}}_{k_{s}} - \tilde{\boldsymbol{\gamma}}_{0,k_{s}}|| + ||\tilde{\boldsymbol{\gamma}}_{0,k_{s}}||) + \cdots + \\ |\tilde{\boldsymbol{\gamma}}_{k_{D}} - \tilde{\boldsymbol{\gamma}}_{0,k_{D}}|| \prod_{s=1}^{D-1} ||\tilde{\boldsymbol{\gamma}}_{0,k_{s}}||. \end{array}$

If $||\tilde{\boldsymbol{\gamma}}_{h} - \tilde{\boldsymbol{\gamma}}_{0,h}|| \leq u_{n}, h = 1, ..., V_{n}$, the above inequality implies that $|\Gamma_{\boldsymbol{k}} - \Gamma_{0,\boldsymbol{k}}| \leq u_{n} \prod_{s=2}^{D} (u_{n} + ||\tilde{\boldsymbol{\gamma}}_{0,k_{s}}||) + \cdots + u_{n} \prod_{s=1}^{D-1} ||\tilde{\boldsymbol{\gamma}}_{0,k_{s}}|| \leq v_{n}$. Thus $\Pi(||\mathbf{\Gamma} - \mathbf{\Gamma}_{0}||_{\infty} \leq v_{n}) \geq \Pi(||\tilde{\boldsymbol{\gamma}}_{h} - \tilde{\boldsymbol{\gamma}}_{0,h}|| \leq u_{n}, h = 1, ..., V_{n})$. Therefore,

$$\Pi(\mathcal{J}|\mathcal{C}) \ge \Pi(||\tilde{\gamma}_{h} - \tilde{\gamma}_{0,h}|| \le u_{n}, h = 1, ..., V_{n})$$

$$\ge E\left[\Pi(||\tilde{\gamma}_{h} - \tilde{\gamma}_{0,h}|| \le u_{n}, h = 1, ..., V_{n}|\boldsymbol{\zeta})\right] \ge E\left[\prod_{h=1}^{V_{n}} \left\{\exp\left\{-||\tilde{\gamma}_{0,h}||^{2}/2\right\} \Pi(||\tilde{\gamma}_{h}|| \le u_{n}|\boldsymbol{\zeta})\right\}\right]$$

$$= \exp\left\{-\sum_{h=1}^{V_{n}} ||\tilde{\gamma}_{0,h}||^{2}/2\right\} E\left[\prod_{h=1}^{V_{n}} \Pi(||\tilde{\gamma}_{h,n}|| \le u_{n}|\boldsymbol{\zeta})\right], \qquad (8)$$

where the second inequality follows from Anderson Lemma. We will now make use of the fact that $\int_{-a}^{a} e^{-x^2/2} dx \ge e^{-a^2} 2a$ to conclude

$$\begin{split} \Pi(||\tilde{\boldsymbol{\gamma}}_{h}|| \leq u_{n}|\Delta) &\geq \prod_{r=1}^{R_{0,n}} \Pi(|\boldsymbol{\gamma}_{h}^{(r)}| \leq u_{n}/R_{0,n}|\Delta) = \prod_{r=1}^{R_{0,n}} \left((1-\Delta) + \frac{\Delta}{\sqrt{2\pi}} \int_{-u_{n}/R_{0,n}}^{u_{n}/R_{0,n}} \exp(-x^{2}/2) \right) \\ &\geq \prod_{r=1}^{R_{0,n}} \left((1-\Delta) + \frac{\Delta}{\sqrt{2\pi}} \exp\left(-\frac{u_{n}^{2}}{R_{0,n}^{2}}\right) \left(\frac{2u_{n}}{R_{0,n}}\right) \right) \\ &\geq \left[(1-\Delta) + \frac{\Delta}{\sqrt{2\pi}} \exp\left(-\frac{u_{n}^{2}}{R_{0,n}^{2}}\right) \left(\frac{2u_{n}}{R_{0,n}}\right) \right]^{R_{0,n}}. \end{split}$$

$$\begin{split} &\prod_{h=1}^{V_n} \Pi(||\tilde{\gamma}_h|| \le u_n) \ge E\left[\left(1 - \Delta\right) + \frac{\Delta}{\sqrt{2\pi}} \exp\left(-\frac{u_n^2}{R_{0,n}^2}\right) \left(\frac{2u_n}{R_{0,n}}\right) \right]^{R_{0,n}V_n} \\ &= E\left[\sum_{l=0}^{R_{0,n}V_n} \binom{R_{0,n}V_n}{l} \left(1 - \Delta\right)^l \left(\frac{\Delta}{\sqrt{2\pi}}\right)^{R_{0,n}V_n - l} \left(\frac{2u_n}{R_{0,n}}\right)^{R_{0,n}V_n - l} \exp\left(-(R_{0,n}V_n - l)\frac{u_n^2}{R_{0,n}^2}\right) \right] \\ &\ge \left(\frac{1}{\sqrt{2\pi}}\right)^{R_{0,n}V_n} \sum_{l=0}^{R_{0,n}V_n} \binom{R_{0,n}V_n}{l} Beta(R_{0,n}V_n - l + 1, l + 1) \left(\frac{2u_n}{R_{0,n}}\right)^{R_{0,n}V_n - l} \exp\left(-(R_{0,n}V_n - l)\frac{u_n^2}{R_{0,n}^2}\right) \\ &\ge \left(\frac{1}{\sqrt{2\pi}}\right)^{R_{0,n}V_n} \sum_{l=0}^{R_{0,n}V_n} \frac{(R_{0,n}V_n)!}{l!(R_{0,n}V_n - l)!} \frac{l!(R_{0,n}V_n - l)!}{(R_{0,n}V_n + 1)!} \left(\frac{2u_n}{R_{0,n}}\right)^{R_{0,n}V_n - l} \exp\left(-(R_{0,n}V_n - l)\frac{u_n^2}{R_{0,n}^2}\right) \\ &\ge \left(\frac{1}{\sqrt{2\pi}}\right)^{R_{0,n}V_n} \frac{R_{0,n}V_n}{R_{0,n}V_n + 1} \left(\frac{2u_n}{R_{0,n}}\right)^{R_{0,n}V_n} \exp\left(-V_n\frac{u_n^2}{R_{0,n}}\right). \end{split}$$

Aggregating all pieces together

$$\begin{split} \Pi(||\mathbf{\Gamma} - \mathbf{\Gamma}_{0}||_{\infty} &\leq v_{n}|\mathcal{C}) \geq \exp\left(-\frac{\sum_{h=1}^{V_{n}} ||\tilde{\boldsymbol{\gamma}}_{0,h}||^{2}}{2}\right) \left(\frac{1}{\sqrt{2\pi}}\right)^{R_{0,n}V_{n}} \frac{R_{0,n}V_{n}}{R_{0,n}V_{n}+1} \left(\frac{2u_{n}}{R_{0,n}}\right)^{R_{0,n}V_{n}} \\ &\exp\left(-V_{n}\frac{u_{n}^{2}}{R_{0,n}}\right). \end{split}$$

Appendix B

Proof of Theorem 3.1

To begin, we define a few metrics of discrepancy between g and g_0 as below:

$$d_0(g,g_0) = \int \int g_0(y|\mathbf{X}) \log\left(\frac{g_0(y|\mathbf{X})}{g(y|\mathbf{X})}\right) \nu_{\mathbf{X}}(d\mathbf{X}) \nu_y(dy),$$

$$d_t(g,g_0) = (1/t) \left\{ \int \int g_0(y|\mathbf{X}) \left\{ \frac{g_0(y|\mathbf{X})}{g(y|\mathbf{X})} \right\}^t \nu_y(dy) \nu_{\mathbf{X}}(d\mathbf{X}) - 1 \right\}.$$

For every n, define a set of probability densities given by \mathcal{P}_n . Let the minimum number of Hellinger balls of radius ϵ_n required to cover \mathcal{P}_n be given by $\mathcal{N}_{\epsilon_n}(\mathcal{P}_n)$. To prove the theorem, it suffices to show that conditions (i)-(iii) hold for all large n:

(i)
$$\log \mathcal{N}_{\epsilon_n}(\mathcal{P}_n) \le n\epsilon_n^2$$

290

- (ii) $\Pi(\mathcal{P}_n^c) \le \exp(-2n\epsilon_n^2)$
- (iii) For t = 1, $\Pi[g : d_t(g, g_0) \le \epsilon_n^2/4] \ge e^{-n\epsilon_n^2/4}$,

using Proposition 1 of 15. Below we show (i)-(iii) for the proposed model.

- Proof of condition (i): Define \mathcal{P}_n as the set of all densities s.t. at most m_n among $\tilde{\gamma}_1, ..., \tilde{\gamma}_{V_n}$ are nonzero and each element in a nonzero $\tilde{\gamma}_h$ satisfies $|\gamma_h^{(r)}| \leq C_n$, for $h = 1, ..., V_n$. Let $g_{\boldsymbol{\zeta}}$ denote a density in \mathcal{P}_n expressed with the binary variables $\boldsymbol{\zeta} = (\zeta_1, ..., \zeta_{V_n})'$. With $|\boldsymbol{\zeta}| = \sum_{h=1}^{V_n} \zeta_h$, \mathcal{P}_n contains densities $g_{\boldsymbol{\zeta}}$ s.t. $|\boldsymbol{\zeta}| \leq m_n$. Note that, each $g_{\boldsymbol{\zeta}} \in \mathcal{P}_n$ is represented by $|\boldsymbol{\zeta}|$ nonzero $\tilde{\gamma}_h$'s with each component $\gamma_h^{(r)}, r = 1, ..., R_n$ of a nonzero $\tilde{\gamma}_h$ is bounded between $[-C_n, C_n]$. It takes at most $(1 + \frac{C_n}{\kappa})^{R_n |\boldsymbol{\zeta}|}$ balls of the form $[\boldsymbol{\xi}_h^{(r)} - \kappa, \boldsymbol{\xi}_h^{(r)} + \kappa]$ (with their
- It takes at most $(1 + \frac{C_n}{\kappa})^{n_n |\boldsymbol{\zeta}|}$ balls of the form $[\xi_h^{(r)} \kappa, \xi_h^{(r)} + \kappa]$ (with their centers $\xi_h^{(r)}$'s satisfying $|\xi_h^{(r)}| \leq C_n$) to cover the parameter space of $g_{\boldsymbol{\zeta}}$. There are at most V_n^l models satisfying $|\boldsymbol{\zeta}| = l$. Hence, the total number of balls to cover the parameter space of regression functions in \mathcal{P}_n is given by $N(\kappa) =$ $\sum_{l \leq m_n} V_n^l (1 + \frac{C_n}{\kappa})^{R_n l} \leq (m_n + 1) \left[V_n (1 + \frac{C_n}{\kappa}) \right]^{R_n m_n}$.

Let $p_{\boldsymbol{\zeta}}$ be any density in \mathcal{P}_n , with $p_{\boldsymbol{\zeta}}(y|\boldsymbol{X}) = \exp(a(\mu)y + b(\mu) + c(y))$, $\mu = \sum_{\boldsymbol{k} \in \mathcal{K}} x_{\boldsymbol{k}} F_{\boldsymbol{k}}$, where $|\boldsymbol{\zeta}| \leq m_n$ and $F_{\boldsymbol{k}} = \sum_{r=1}^{R_n} \lambda_r f_{k_1}^{(r)} \dots f_{k_D}^{(r)}$, with $|f_h^{(r)}| \leq C_n$ for all $h \in \mathcal{A}$, $r = 1, \dots, R_n$. There exists a density $g_{\boldsymbol{\zeta}} \in \mathcal{P}_n$ given by $g_{\boldsymbol{\zeta}}(y|\boldsymbol{X}) = \exp(a(\alpha)y + b(\alpha) + c(y))$, with $\alpha = \sum_{\boldsymbol{k} \in \mathcal{K}} x_{\boldsymbol{k}} \Gamma_{\boldsymbol{k}}$. $\Gamma_{\boldsymbol{k}} = \sum_{r=1}^{R_n} \lambda_r \gamma_{k_1}^{(r)} \dots \gamma_{k_D}^{(r)}$, where $\gamma_h^{(r)}$'s are such that $f_h^{(r)} \in (\gamma_h^{(r)} - \kappa, \gamma_h^{(r)} + \kappa)$ for every r and h.

Applying Taylor expansion on $d_0(p_{\boldsymbol{\zeta}}, g_{\boldsymbol{\zeta}})$ to show that $d_0(p_{\boldsymbol{\zeta}}, g_{\boldsymbol{\zeta}}) = E_{\boldsymbol{X}} \left[\left\{ a'(\alpha_{\mu}) \left(-\frac{b'(\alpha)}{a'(\alpha)} \right) + b'(\alpha_{\mu}) \right\} (\alpha - \mu) \right],$ where α_{μ} is an intermediate point between α and μ . Let $\mathcal{B} = \{ \boldsymbol{k} \in \mathcal{K} : \zeta_{k_1} = 1, .., \zeta_{k_D} = 1 \}.$ Now note that,

$$|\alpha - \mu| = |\sum_{\boldsymbol{k} \in \mathcal{B}} x_{i,\boldsymbol{k}} \Gamma_{\boldsymbol{k}} - \sum_{\boldsymbol{k} \in \mathcal{B}} x_{i,\boldsymbol{k}} F_{\boldsymbol{k}}| \le \sum_{\boldsymbol{k} \in \mathcal{B}} |\Gamma_{\boldsymbol{k}} - F_{\boldsymbol{k}}| \le m_n^D \max_{\boldsymbol{k} \in \mathcal{B}} |\Gamma_{\boldsymbol{k}} - F_{\boldsymbol{k}}|.$$

It follows from the above that,

$$\begin{aligned} |\Gamma_{\boldsymbol{k}} - F_{\boldsymbol{k}}| &= |\sum_{r=1}^{R_n} \lambda_r \gamma_{k_1}^{(r)} \dots \gamma_{k_D}^{(r)} - \sum_{r=1}^{R_n} \lambda_r f_{k_1}^{(r)} \dots f_{k_D}^{(r)}| \le |\sum_{r=1}^{R_n} \gamma_{k_1}^{(r)} \dots \gamma_{k_D}^{(r)} - \sum_{r=1}^{R_n} f_{k_1}^{(r)} \dots f_{k_D}^{(r)}| \\ &\le \sum_{r=1}^{R_n} \left\{ |\gamma_{k_1}^{(r)} - f_{k_1}^{(r)}| \prod_{l=2}^{D} |\gamma_{k_l}^{(r)}| + |f_{k_1}^{(r)}| |\gamma_{k_2}^{(r)} - f_{k_2}^{(r)}| \prod_{l=3}^{D} |\gamma_{k_l}^{(r)}| + \dots + \prod_{l=1}^{D-1} |f_{k_l}^{(r)}| |\gamma_{k_D}^{(r)} - f_{k_D}^{(r)}| \right\} \\ &\le R_n \kappa C_n^{D-1}. \end{aligned}$$

Thus, $|\alpha - \mu| \leq m_n^D R_n \kappa C_n^{D-1}$. Similarly, $|\alpha|$, $|\mu|$ (and therefore $|\alpha_{\mu}|$) being

bounded by $R_n C_n^D m_n^D$. Hence,

$$d_{H}(p_{\boldsymbol{\zeta}}, g_{\boldsymbol{\zeta}}) \leq \left\{ d_{0}(p_{\boldsymbol{\zeta}}, g_{\boldsymbol{\zeta}}) \right\}^{1/2} \leq \left\{ 2 \sup_{|w| \leq R_{n} C_{n}^{D} m_{n}^{D}} |a'(w)| \sup_{|w| \leq R_{n} C_{n}^{D} m_{n}^{D}} \left| \frac{b'(w)}{a'(w)} \right| \kappa R_{n} m_{n}^{D} C_{n}^{D-1} \right\}^{1/2}.$$
Choosing $\kappa = \frac{\epsilon_{n}^{2}}{2 \sup_{|w| \leq R_{n} m_{n}^{D} C_{n}^{D}} |a'(w)|} \sup_{|w| \leq R_{n} m_{n}^{D} C_{n}^{D}} \left| \frac{b'(w)}{a'(w)} \right| R_{n} m_{n}^{D} C_{n}^{D-1}}, \text{ we obtain } d_{H}(g_{\boldsymbol{\zeta}}, p_{\boldsymbol{\zeta}}) \leq \epsilon_{n}.$ Hence

$$\log \mathcal{N}_{\epsilon_{n}}(\mathcal{P}_{n}) \leq \log N(\kappa)$$

$$\leq \log(m_{n}+1) + R_{n}m_{n}\log(V_{n}) + R_{n}m_{n}\log\left(1 + \frac{2\sup_{|w| \leq R_{n}m_{n}^{D}C_{n}^{D}}|u'(w)|\sup_{|w| \leq R_{n}m_{n}^{D}C_{n}^{D}}\left|\frac{b'(w)}{a'(w)}\right|R_{n}m_{n}^{D}C_{n}^{D}}{\epsilon_{n}^{2}}\right)$$

$$\leq \log(m_{n}+1) + R_{n}m_{n}\log(V_{n}) + R_{n}m_{n}\log(2/\epsilon_{n}^{2}) + R_{n}m_{n}\log(H(R_{n}m_{n}^{D}C_{n}^{D}))$$

 $\leq n\epsilon_n^2$, for large *n*, by assumptions (a)-(c).

Proof of condition (ii): Define, $\mathcal{A} = \{h \in \mathcal{N} : \zeta_h = 1\}$. Then for all large n,

where the last inequality follows from assumptions (a) and (d).

Proof of Condition (iii): Using the mean value theorem, there exists v such that $d_t(g, g_0) = E_{\mathbf{X}} \{g'(v)(\alpha - \alpha_0)\}$, where $g'(\cdot)$ represents the continuous derivative function of g in the neighborhood of g_0 . Let $\tau_n = \frac{\epsilon_n^2}{8q_n}$. If for each $\mathbf{k} \in \mathcal{K}, \Gamma_{\mathbf{k}} \in (\Gamma_{0,\mathbf{k}} - \tau_n, \Gamma_{0,\mathbf{k}} + \tau_n)$, then

$$|\alpha - \alpha_0| = |\sum_{\boldsymbol{k} \in \mathcal{K}} x_{\boldsymbol{k}} \Gamma_{\boldsymbol{k}} - \sum_{\boldsymbol{k} \in \mathcal{K}} x_{\boldsymbol{k}} \Gamma_{0,\boldsymbol{k}}| \le \sum_{\boldsymbol{k} \in \mathcal{B}} |\Gamma_{\boldsymbol{k}} - \Gamma_{0,\boldsymbol{k}}| \le q_n \tau_n \le \epsilon_n^2/8,$$

for large *n*. Again, $|v| \leq |\alpha - \alpha_0| + |\alpha_0| \leq q_n \tau_n + \omega_n = \epsilon_n^2/8 + \omega_n$, where $\omega_n = |\alpha_0| = |\sum_{\boldsymbol{k} \in \mathcal{K}} x_{\boldsymbol{k}} \Gamma_{0,\boldsymbol{k}}| \leq \sum_{\boldsymbol{k} \in \mathcal{K}} |\Gamma_{0,\boldsymbol{k}}| \leq \sum_{\boldsymbol{k} \in \mathcal{K}} ||\tilde{\gamma}_{0,k_1}|| \cdots ||\tilde{\gamma}_{0,k_D}|| \leq \varepsilon_n ||\tilde{\gamma}_{0,k_1}||$

³¹⁵ $(\sum_{h=1}^{V_n} || \tilde{\gamma}_{0,h} ||)^D$, which is bounded by assumption (e), for sufficiently large n. Hence ||g'(v)|| is bounded for sufficiently large n. Thus, $d_t(g, g_0) = E_{\mathbf{X}} \{g(v)(\alpha - \alpha_0)\} \leq C_0 q_n \tau_n \leq \epsilon_n^2/4$ for large n, for some constant C_0 .

Let $C_1 = \{ \boldsymbol{\Gamma} : \Gamma_{\boldsymbol{k}} \in (\Gamma_{0,\boldsymbol{k}} - \tau_n, \Gamma_{0,\boldsymbol{k}} + \tau_n), \forall \boldsymbol{k} \in \mathcal{K} \}$ and $C_2 = \{ \lambda_1 = 1, ..., \lambda_{R_{0,n}} = 1, \lambda_{R_{0,n}+1} = 0, ..., \lambda_{R_n} = 0 \}$. This implies that

$$\Pi(\{g: d_t(g, g_0) \le \epsilon_n^2/4\}) \ge \Pi(\mathcal{C}_1 \cap \mathcal{C}_2) = \Pi(\mathcal{C}_2)\Pi(\mathcal{C}_1|\mathcal{C}_2).$$

By Lemma 5.2, $\Pi(\mathcal{C}_2) \geq \frac{1}{(1+R_{0,n}^{\eta})^{R_n}} R_{0,n}^{\eta(R_n-R_{0,n})}$. By Lemma 5.3, $-\log \Pi(\mathcal{C}_1|\mathcal{C}_2) = -\log \Pi(||\mathbf{\Gamma} - \mathbf{\Gamma}_0||_{\infty} \leq \tau_n |\mathcal{C}_2) \leq \sum_{h=1}^{V_n} ||\tilde{\boldsymbol{\gamma}}_{0,h}||^2 / 2 + (R_{0,n}V_n/2)\log(2\pi) + \log(1 + (1/(R_{0,n}V_n))) + R_{0,n}V_n\log(R_{0,n}) + R_{0,n}V_n\log(1/(2u_n)) + V_n u_n^2 / R_{0,n}$. Here u_n is the minimum of the root of the equation (7) with v_n replaced by τ_n .

Since $||\tilde{\gamma}_{0,h}|| \geq 0$, $\sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}||^2 \leq (\sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}||)^2$ is bounded for large n, by assumption (e). By assumption (a), $R_{0,n}V_n\log(R_{0,n}) = o(n\epsilon_n^2)$ (hence $R_{0,n}V_n = o(n\epsilon_n^2)$). Using the Lagrange-Maclaurin bound on the positive root of a monic polynomial of degree D, we have $u_n \leq 1 + \tau_n^{1/D}$, implying $V_n u_n^2/R_{0,n} =$

a monic polynomial of degree D, we have $u_n \leq 1 + \tau_n^{1/2}$, implying $V_n u_n^2 / R_{0,n} = o(n\epsilon_n^2)$, for all large n, by assumption (a). Using Lemma 5.1, $1/u_n \leq (\sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}||)^D / \tau_n + 1$. If $G_0 = \limsup_{n \to \infty} \sum_{h=1}^{V_n} ||\tilde{\gamma}_{0,h}||$, then $R_n V_n \log(1/u_n) \leq R_n V_n \log(G_0^D / \tau_n + 1) = DR_n V_n \log(G_0) + R_n V_n \log(8q_n) + R_n V_n \log(1/\epsilon_n^2) = o(n\epsilon_n^2)$, where the last line follows from assumptions (a) and (b).

All the aforementioned calculations yield $-\log \Pi(\mathcal{C}_1 \cap \mathcal{C}_2) \leq n\epsilon_n^2/4$, for all large n, which implies $\Pi(\{g : d_t(g, g_0) \leq \epsilon_n^2/4\}) \geq \exp(-n\epsilon_n^2/4)$ for all large n. This concludes the proof.

References

References

- [1] S. Guha, A. Rodriguez, Bayesian regression with undirected network predictors with an application to brain connectome data, Journal of the American Statistical Association (2020) 1–13.
 - [2] R. C. Craddock, P. E. Holtzheimer III, X. P. Hu, H. S. Mayberg, Disease state prediction from resting state functional connectivity, Magnetic Res-

- onance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine 62 (6) (2009) 1619–1628.
 - [3] J. Richiardi, H. Eryilmaz, S. Schwartz, P. Vuilleumier, D. Van De Ville, Decoding brain states from fmri connectivity graphs, Neuroimage 56 (2) (2011) 616–626.
- [4] T. Park, G. Casella, The Bayesian lasso, Journal of the American Statistical Association 103 (482) (2008) 681–686.
 - [5] C. M. Carvalho, N. G. Polson, J. G. Scott, The horseshoe estimator for sparse signals, Biometrika 97 (2) (2010) 465–480.
 - [6] A. Armagan, D. B. Dunson, J. Lee, W. U. Bajwa, N. Strawn, Posterior consistency in linear models under shrinkage priors, Biometrika 100 (4) (2013) 1011–1018.
 - [7] H. Zhou, L. Li, H. Zhu, Tensor regression with applications in neuroimaging data analysis, Journal of the American Statistical Association 108 (502) (2013) 540–552.
- [8] R. Guhaniyogi, S. Qamar, D. B. Dunson, Bayesian tensor regression, Journal of Machine Learning Research 18 (79) (2017) 1–31.
 - [9] I. Castillo, A. van der Vaart, et al., Needles and straw in a haystack: Posterior concentration for possibly sparse sequences, The Annals of Statistics 40 (4) (2012) 2069–2101.
- 360 [10] E. Belitser, N. Nurushev, Needles and straw in a haystack: robust confidence for possibly sparse sequences, arXiv preprint arXiv:1511.01803.
 - [11] R. Martin, R. Mess, S. G. Walker, et al., Empirical bayes posterior concentration in sparse high-dimensional linear models, Bernoulli 23 (3) (2017) 1822–1847.
- ³⁶⁵ [12] Q. Song, F. Liang, Nearly optimal bayesian shrinkage for high dimensional regression, arXiv preprint arXiv:1712.08964.

340

350

- [13] R. Wei, S. Ghosal, Contraction properties of shrinkage priors in logistic regression, Preprint at http://www4. stat. ncsu. edu/~ ghoshal/papers.
- [14] R. Guhaniyogi, Convergence rate of bayesian supervised tensor modeling with multiway shrinkage priors, Journal of Multivariate Analysis 160 (2017) 157–168.

370

- [15] W. Jiang, Bayesian variable selection for high dimensional generalized linear models: convergence rates of the fitted densities, The Annals of Statistics 35 (4) (2007) 1487–1511.
- 375 [16] P. D. Hoff, Bilinear mixed-effects models for dyadic data, Journal of the American Statistical Association 100 (469) (2005) 286–295.
 - [17] J. D. A. Relión, D. Kessler, E. Levina, S. F. Taylor, et al., Network classification with applications to brain connectomics, The Annals of Applied Statistics 13 (3) (2019) 1648–1677.
- ³⁸⁰ [18] S. Chan, C.-C. Kuo, Trilateral trade relations among china, japan and south korea: Challenges and prospects of regional economic integration, East Asia 22 (1) (2005) 33–50.
 - [19] G. C. Hufbauer, NAFTA revisited: Achievements and challenges, Peterson Institute, 2005.
- ³⁸⁵ [20] M.-H. Chiang, The potential of china-japan-south korea free trade agreement, East Asia 30 (3) (2013) 199–216.
 - [21] D. K. Brown, A. V. Deardorff, R. M. Stern, Estimates of a north american free trade agreement, in: Modeling North American Economic Integration, Springer, 1995, pp. 59–74.
- ³⁹⁰ [22] S. Guha, A. Rodriguez, Bayesian regression with undirected network predictors with an application to brain connectome data, arXiv preprint arXiv:1803.10655.

[23] S. Guha, R. Guhaniyogi, Bayesian generalized sparse symmetric tensor-onvector regression, Technometrics (2020) 1–11.