

# Paper Sketch for AISTATS

Tentative title: “Binary tensor regression with multi-mode features”

first draft on 08/13, updated on 08/22

## 1 Preliminaries

We use lower-case letters  $(a, b, \dots)$  for scalars and vectors, upper-case boldface letters  $(\mathbf{A}, \mathbf{B}, \dots)$  for matrices, and calligraphy letter  $(\mathcal{A}, \mathcal{B}, \dots)$  for tensors of order 3 or greater. Let  $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  denote an order- $K$   $(d_1, \dots, d_K)$ -dimensional tensor. We say that an event  $A$  occurs “with very high probability” if  $\mathbb{P}(A)$  tends to 1 faster than any polynomial of  $d_{\min} = \min\{d_1, \dots, d_K\}$ . We use  $\mathbf{S}^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$  to denote the Euclidean unit sphere in dimension  $d$ .

**Property 1.** Let  $\mathbf{X} \in \mathbb{R}^{d \times p}$  be a rank- $r$  matrix with  $p \leq d$ . The SVD of  $\mathbf{X}$  can be expressed as  $\mathbf{X} = \mathbf{P}\Delta\mathbf{Q}^T$ , where  $\mathbf{P} \in \mathbb{R}^{d \times r}$  and  $\mathbf{Q} \in \mathbb{R}^{p \times r}$  consist of, respectively, the left and right singular vectors, and  $\Delta \in \mathbb{R}^{r \times r}$  is the diagonal matrix consisting of non-zero singular values. The following properties hold:

1.  $(\mathbf{X}^T \mathbf{X})^{-1/2} = \mathbf{Q}\Delta^{-1}\mathbf{Q}^T$ .
2. Let  $\tilde{\mathbf{X}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1/2}$ . Then  $\tilde{\mathbf{X}} = \mathbf{P}\mathbf{Q}^T$ .
3.  $\tilde{\mathbf{X}}^T \mathbf{X} = \mathbf{Q}\Delta\mathbf{Q}^T$ .

## 2 Results

Suppose we observe an order- $K$  binary tensor  $\mathcal{Y} \in \{0, 1\}^{d_1 \times \dots \times d_K}$ , along with a set of covariate matrices  $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$  for  $k = 1, \dots, K$ . Consider a tensor regression model:

$$\text{logit}(\mathbb{E}(\mathcal{Y})) = \mathcal{B} \times_1 \mathbf{X}_1 \times_2 \dots \times_K \mathbf{X}_K, \quad (1)$$

where  $\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K}$  is a coefficient tensor of interest. Furthermore, the tensor  $\mathcal{B}$  is assumed to (i) be entrywise bounded, and (ii) admit a low-rank Tucker decomposition; that is,  $\text{rank}(\mathcal{B}) = \mathbf{r} \equiv (r_1, \dots, r_K)^T$ , where  $r_k \leq p_k \leq d_k$ . The parameter space we consider is

$$\mathcal{P} = \mathcal{P}(\mathbf{r}, \alpha) = \{\mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_K} : \text{rank}(\mathcal{B}) \leq \mathbf{r}, \text{ and } \|\mathcal{B}\|_\infty \leq \alpha\}.$$

In the following analysis, we assume both the multilinear rank  $\mathbf{r}$  and entrywise bound  $\alpha$  are known. The adaptation of unknown rank will be addressed in the next note.

**Remark 1.** Model (1) incorporates the following examples as special cases:

(1) **Binary tensor decomposition.** In the absence of side information, set  $\mathbf{X} = \mathbf{I}_k$  to be identity matrix and  $p_k = d_k$  for  $k = 1, \dots, K$ . Then the model (1) reduces to unsupervised binary tensor decomposition.

(2) **Network link prediction model.** Suppose  $K = 2$  and  $\mathbf{X}_1 = \mathbf{X}_2$ . Then the model (1) reduced to the matrix logistic model [Baldin and Berthet, 2018] that is commonly used in the network analysis:

$$\text{logit}(\mathbb{E}(\mathbf{Y})) = \mathbf{X}^T \mathbf{B} \mathbf{X}, \quad \text{where} \quad \text{rank}(\mathbf{B}) \leq r.$$

(3) **Semi-supervised decomposition.** Suppose the covariate information is available only for a subset of modes. Without loss of generality, suppose the covariates  $\mathbf{X}_k \neq \mathbf{I}$  are available in modes  $1, \dots, L$ , where  $L < K$ . Then the model (1) reduces to a semi-supervised decomposition model:

$$\text{logit}(\mathbb{E}(\mathcal{Y})) = \underbrace{\mathcal{B}}_{\in \mathbb{R}^{p_1 \times \dots \times p_L \times d_{L+1} \times \dots \times d_K}} \times_1 \underbrace{\mathbf{X}_1}_{\in \mathbb{R}^{d_1 \times p_1}} \times_2 \dots \times_L \underbrace{\mathbf{X}_L}_{\in \mathbb{R}^{d_L \times p_L}}.$$

For parsimony, we do not distinguish modes with available side information from those without side information. We focus on the general tensor regression model (1) with mild assumption on  $\{\mathbf{X}_k\}$ . Specifically, the covariates  $\{\mathbf{X}_k\}$  are assumed to satisfy the following restricted isometry property (RIP) assumption.

**Assumption 1** (Restricted Isometry Property). *Let  $d = \prod_k d_k$ . The covariates  $\{\mathbf{X}_k\}$  are called to satisfy the RIP condition if there exists a positive constant  $\delta_{\mathbf{r}, \alpha} \in (0, 1)$  such that*

$$d(1 - \delta_{\mathbf{r}, \alpha}) \|\mathcal{B}\|_F^2 \leq \|\mathcal{B} \times_1 \mathbf{X}_1 \times_2 \dots \times_K \mathbf{X}_K\|_F^2 \leq d(1 + \delta_{\mathbf{r}, \alpha}) \|\mathcal{B}\|_F^2,$$

*holds for all tensors  $\mathcal{B} \in \mathcal{P}(\mathbf{r}, \alpha)$  in the parameter space.*

**Remark 2.** The RIP assumption requires the covariates at each of the modes are nearly orthonormal when restricted to the desired parameter space. The RIP condition is a relatively mild assumption on the covariates. In particular, both fixed design and Gaussian random design satisfy the RIP condition.

**Example 1** (Gaussian Random design). Suppose  $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$  are random design matrices with i.i.d. standard Gaussian entries, where  $\frac{p_k}{d_k} = \lambda_k \in (0, 1)$  for all  $k = 1, \dots, K$ . Then, with very high probability,  $\{\mathbf{X}_k\}$  satisfy the RIP condition. In particular, the RIP constant can be chosen as

$$\delta_{\mathbf{r}, \alpha} = \max \left\{ -1 + \prod_k \sqrt{1 + \lambda_k}, 1 - \prod_k \sqrt{1 - \lambda_k} \right\} \in (0, 1).$$

**Example 2** (Fixed design). Suppose  $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$  are deterministic design matrices with full rank. Then,  $\{\mathbf{X}_k\}$  (upon proper rescaling) satisfy the RIP condition.

**Theorem 1** (Main Results). *Consider a tensor regression model (1) with  $\mathcal{Y} \in \{0, 1\}^{d_1 \times \dots \times d_K}$  the response and  $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$  the mode- $k$  covariates. Let  $\hat{\mathcal{B}}_{MLE}$  be the restricted rank- $\mathbf{r}$  maximum*

likelihood estimate of the coefficient tensor, where  $\mathbf{r} = (r_1, \dots, r_K)$ ,

$$\hat{\mathcal{B}}_{MLE} = \arg \min_{\mathcal{B}: \text{rank}(\mathcal{B})=\mathbf{r}, \|\mathcal{B}\|_\infty \leq \alpha} \text{Log-lik}(\mathcal{B}; \mathcal{Y}, \{\mathbf{X}_k\}).$$

Suppose the covariates  $\mathbf{X}_k$  satisfy the RIP condition with RIP constant  $\delta \in (0, 1)$ . Then, with very high probability,

$$\left\| \hat{\mathcal{B}}_{MLE} - \mathcal{B}_{true} \right\|_F \leq \frac{C_\alpha}{\sqrt{\prod_k d_k}} \sqrt{\frac{(1 + \delta_{2\mathbf{r}, 2\alpha})}{(1 - \delta_{2\mathbf{r}, 2\alpha})^2} \frac{\prod_{k=1}^K r_k}{r_{\max}} \sum_{k=1}^K p_k},$$

where  $C_\alpha > 0$  is a constant independent of the tensor dimension or rank.

**Theorem 2** (KL-Divergence and Hellinger Loss). See Zhuoyan’s note “Evidence theory on prediction error” (08/09) and Jiaxin’s note “Boundaries for different prediction error metrics” (08/09).

### 3 Experiments

**Algorithm sketch.** Please refer to earlier notes for algorithm.

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

$$\hat{\mathbf{r}} = \arg \min_{\mathbf{r}=(r_1, \dots, r_K)} BIC(\mathbf{r}) = \arg \min_{\mathbf{r}=(r_1, \dots, r_K)} \left[ -2\text{Log-lik}(\hat{\Theta}) + p_e(r_1, \dots, r_K) \log \left( \prod_k d_k \right) \right],$$

where  $p_e(r_1, \dots, r_K) \stackrel{\text{def}}{=} \sum_k (d_k - 1)r_k + \prod_k r_k$  is the effective number of parameters in the model.

**Performance.** There are several issues summarized in the note “Summary\_of\_last\_semester\_and\_summer\_plan.pdf”. The main issues are

- The choice of distribution to generate the core tensor. Bad estimation when the core is generated i.i.d.  $N(0, 1)$ ,  $N(10, 1)$ ,  $\text{Unif}[0, 1]$  or  $\text{Unif}[0, 10]$ .
- When to stop? “It is reported that when we run more iterations, the result sometime seems not better or even worse.”
- Unbounded estimation. “Sometimes, when we check the real scale of  $U$ , there may some entries with extremely large values”.

We have proposed to impose an infinity-norm constraint to the optimization. Which of the following leads to the best result?

1. Update the core tensor to the feasible region via conjugate gradient solver. (good MSE, but haven not checked its performance on “bad” distributions. )

2. Update the core by global downscaling.
3. Update the factor matrixes by backward linear search (computationally slow).

MSE exhibits good agreement between simulation and theory. More figures to come.

To-do list:

- Assess the sensitivity of algorithm on the “distribution” on core tensor.
- Understand the bad convergence that arises in simulations.
- Evaluate the BIC rank selection accuracy via simulation.
- Evaluate the MSE for supervised regression.
- Apply the method to a real dataset. Possible dataset: multi-relational social networks with user attributes; air-flight connection map with flight attributes; or nation data with multilayer political networks (already tried; see earlier note).

## 4 Proofs

*Proof of Theorem 1.* Following the similar argument as in [Wang and Li, 2019], we have  $\text{Log-lik}(\mathcal{B}_{\text{true}}) \leq \text{Log-lik}(\hat{\mathcal{B}}_{\text{MLE}})$ . By Taylor expansion,

$$\|(\hat{\mathcal{B}}_{\text{MLE}} - \mathcal{B}_{\text{true}}) \times_1 \mathbf{X}_1 \times_2 \cdots \times_K \mathbf{X}_K\|_F^2 \leq C_\alpha \langle \mathcal{S}, (\hat{\mathcal{B}}_{\text{MLE}} - \mathcal{B}_{\text{true}}) \times_1 \mathbf{X}_1 \times_2 \cdots \times_K \mathbf{X}_K \rangle, \quad (2)$$

where  $\mathcal{S} \in \mathbb{R}^{d_1 \times \cdots \times d_K}$  is a random tensor consisting of i.i.d. bounded random entries. Applying the RIP condition to  $(\hat{\mathcal{B}}_{\text{MLE}} - \mathcal{B}_{\text{true}}) \in \mathcal{P}(2\mathbf{r}, 2\alpha)$  in the inequality (2) yields

$$\begin{aligned} & d(1 - \delta_{2\mathbf{r}, 2\alpha}) \|(\hat{\mathcal{B}}_{\text{MLE}} - \mathcal{B}_{\text{true}})\|_F^2 \\ & \leq \|(\hat{\mathcal{B}}_{\text{MLE}} - \mathcal{B}_{\text{true}}) \times_1 \mathbf{X}_1 \times_2 \cdots \times_K \mathbf{X}_K\|_F^2 \\ & \leq C_\alpha \times \|\hat{\mathcal{B}}_{\text{MLE}} - \mathcal{B}_{\text{true}}\|_F \times \sqrt{d(1 + \delta_{2\mathbf{r}, 2\alpha}) \frac{\prod_k r_k}{r_{\max}} \sum_k p_k}, \end{aligned}$$

where  $d = \prod_k d_k$  and the last line uses the Lemma 2. Therefore,

$$\|\hat{\mathcal{B}}_{\text{MLE}} - \mathcal{B}_{\text{true}}\|_F \leq \frac{C_\alpha}{\sqrt{\prod_k d_k}} \sqrt{\frac{(1 + \delta_{2\mathbf{r}, 2\alpha})}{(1 - \delta_{2\mathbf{r}, 2\alpha})^2} \frac{\prod_k r_k}{r_{\max}} \sum_k p_k}.$$

□

**Lemma 1.** Suppose the matrices  $\{\mathbf{X}_k\}$  satisfy the RIP condition with constant  $\delta_{\mathbf{r}, \alpha} \in (0, 1)$ . Then the matrices  $\{\tilde{\mathbf{X}}_k^T \mathbf{X}_k\}$  also satisfy the RIP condition with the same RIP constant.

*Proof.* Let  $\mathbf{X}_k = \mathbf{P}_k \Delta_k \mathbf{Q}_k^T$  be the SVD of  $\mathbf{X}_k$ , and by Property 1,  $\tilde{\mathbf{X}}_k^T \mathbf{X}_k = \mathbf{Q} \Delta_k \mathbf{Q}^T \in \mathbb{R}^{p_k \times p_k}$ . Note that the F-norm is invariant under orthonormal transformation. Hence,

$$\begin{aligned} \|\mathcal{B} \times_1 \mathbf{X}_1 \times_2 \cdots \times_K \mathbf{X}_K\|_F &= \|\mathcal{B} \times_1 (\mathbf{P}_1 \Delta_1 \mathbf{Q}_1^T) \times_2 \cdots \times_K (\mathbf{P}_K \Delta_K \mathbf{Q}_K^T)\|_F \\ &= \|\mathcal{B} \times_1 (\mathbf{Q} \Delta_1 \mathbf{Q}^T) \times_2 \cdots \times_K (\mathbf{Q} \Delta_K \mathbf{Q}^T)\|_F \\ &= \|\mathcal{B} \times_1 (\tilde{\mathbf{X}}_1 \mathbf{X}_1^T)^{1/2} \times_2 \cdots \times_K (\tilde{\mathbf{X}}_K \mathbf{X}_K^T)^{1/2}\|_F. \end{aligned}$$

The proof is complete by invoking the Assumption 1.  $\square$

**Lemma 2.** Let  $\mathcal{B} \in \mathcal{P}(\mathbf{r}, \alpha)$  be a fixed tensor in the parameter space  $\mathcal{P}(\mathbf{r}, \alpha)$  and  $\mathcal{S} \in \mathbb{R}^{d_1 \times \cdots \times d_K}$  be a random tensor with i.i.d. bounded random entries. Denote  $d = \prod_k d_k$ . Suppose  $\{\mathbf{X}_k\}$  satisfy the RIP condition with RIP constant  $\delta_{\mathbf{r}, \alpha}$ . Then, with very high probability,

$$\langle \mathcal{S}, \mathcal{B} \times_1 \mathbf{X}_1 \times_2 \cdots \times_K \mathbf{X}_K \rangle \leq \|\mathcal{B}\|_F \times \sqrt{d(1 + \delta_{\mathbf{r}, \alpha}) \frac{\prod_{k=1}^K r_k}{r_{\max}} \sum_{k=1}^K p_k}.$$

*Proof.* Let  $\tilde{\mathbf{X}}_k = \mathbf{X}_k (\mathbf{X}_k^T \mathbf{X}_k)^{-1/2} = \mathbf{P}_k$ , where  $\mathbf{P}_k$  consists of left singular vectors of  $\mathbf{X}_k$ . By the definition of inner product,

$$\begin{aligned} &\langle \mathcal{S}, \mathcal{B} \times_1 \mathbf{X}_1 \times_2 \cdots \times_K \mathbf{X}_K \rangle \\ &= \left\langle \underbrace{\mathcal{S} \times_1 \tilde{\mathbf{X}}_1^T \times_2 \cdots \times_K \tilde{\mathbf{X}}_K^T}_{:= \mathcal{E} \in \mathbb{R}^{p_1 \times \cdots \times p_K} \text{ is a sub-Gaussian(1) tensor by Lemma 3}}, \mathcal{B} \times_1 (\tilde{\mathbf{X}}_1^T \mathbf{X}_1) \times_2 \cdots \times_K (\tilde{\mathbf{X}}_K^T \mathbf{X}_K) \right\rangle \\ &\leq \|\mathcal{E}\|_\sigma \times \left\| \mathcal{B} \times_1 (\tilde{\mathbf{X}}_1^T \mathbf{X}_1) \times_2 \cdots \times_K (\tilde{\mathbf{X}}_K^T \mathbf{X}_K) \right\|_* \\ &\leq \|\mathcal{E}\|_\sigma \times \sqrt{\frac{\prod_k r_k}{r_{\max}}} \times \left\| \mathcal{B} \times_1 (\tilde{\mathbf{X}}_1^T \mathbf{X}_1) \times_2 \cdots \times_K (\tilde{\mathbf{X}}_K^T \mathbf{X}_K) \right\|_F \\ &\leq \sqrt{\frac{\prod_k r_k}{r_{\max}}} \times \|\mathcal{E}\|_\sigma \times \sqrt{d(1 + \delta_{\mathbf{r}, \alpha})} \|\mathcal{B}\|_F, \end{aligned}$$

where the last line comes from the RIP condition of  $\{\tilde{\mathbf{X}}_k^T \mathbf{X}_k\}$  by Lemma 1. Combining with the fact that  $\|\mathcal{E}\|_\sigma \asymp \mathcal{O}(\sqrt{\sum_k p_k})$  (c.f. Theorem 1 in Tommioka and Suzuki, 2014), we have

$$\langle \mathcal{S}, \mathcal{B} \times_1 \mathbf{X}_1 \times_2 \cdots \times_K \mathbf{X}_K \rangle \leq \|\mathcal{B}\|_F \times \sqrt{d(1 + \delta_{\mathbf{r}, \alpha}) \frac{\prod_k r_k}{r_{\max}} \sum_k p_k}.$$

$\square$

**Lemma 3.** Let  $\mathcal{S}$  be an  $sG(\sigma)$  tensor of dimension  $(d_1, \dots, d_K)$  and  $\tilde{\mathbf{X}}_k \in \mathbb{R}^{d_k \times p_k}$  be column-wise orthogonal matrices. Then  $\mathcal{E} = \mathcal{S} \times_1 \tilde{\mathbf{X}}_1^T \times_2 \cdots \times_K \tilde{\mathbf{X}}_K^T$  is an  $sG(\sigma)$  tensor of dimension  $(p_1, \dots, p_K)$ .

*Proof.* (Extended from Zhuoyan's note version 4.0) To show  $\mathcal{E}$  is an sG tensor, it suffices to show that the  $\mathcal{E}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K) \stackrel{\text{def}}{=} \langle \mathcal{E}, \mathbf{u}_1 \otimes \cdots \otimes \mathbf{u}_K \rangle$  is a sub-Gaussian random variable with parameter

$\sigma$ , where  $\mathbf{u}_k \in \mathbf{S}^{p_k-1}$  for all  $k = 1, \dots, K$ .

Note that,

$$\mathcal{E}(\mathbf{u}_1, \dots, \mathbf{u}_K) = \mathcal{S}(\tilde{\mathbf{X}}_1 \mathbf{u}_1, \dots, \tilde{\mathbf{X}}_K \mathbf{u}_K).$$

Because  $\tilde{\mathbf{X}}_k \in \mathbb{R}^{d \times p}$  are column-wise orthogonal matrices, so  $\|\tilde{\mathbf{X}}_k \mathbf{u}_k\|_2 = \|\mathbf{u}_k\|_2 = 1$ . By definition of sub-Gaussian tensor,  $\mathcal{S}(\tilde{\mathbf{X}}_1 \mathbf{u}_1, \dots, \tilde{\mathbf{X}}_K \mathbf{u}_K)$  is a sub-Gaussian random variable with parameter  $\sigma$ , so is the  $\mathcal{E}(\mathbf{u}_1, \dots, \mathbf{u}_K)$ .  $\square$