

ISLET: Fast and Optimal Low-Rank Tensor Regression via Importance Sketching

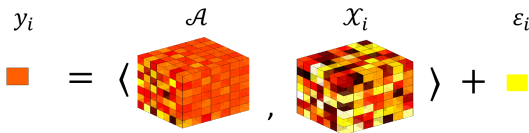
Anru R. Zhang Yuetian Luo Garvesh Raskutti Ming
Yuan

SIAM J. Mathematics of Data Science (SIMODS), 2020
DOI: 10.1137/19M126476X

Limitations of Existing Tensor Regression

Tensor regression model:

$$y_i = \langle \mathcal{X}_i, \mathcal{B} \rangle + \varepsilon_i, \quad \mathcal{X}_i \in \mathbb{R}^{p_1 \times \dots \times p_d}, \mathcal{B} \in \mathbb{R}^{p_1 \times \dots \times p_d}.$$



The diagram illustrates the tensor regression model equation $y_i = \langle \mathcal{X}_i, \mathcal{A} \rangle + \varepsilon_i$. It shows a small orange square representing y_i , followed by an equals sign, then a large orange cube representing \mathcal{A} inside angle brackets, followed by a comma, then a large yellow and black cube representing \mathcal{X}_i , followed by a plus sign and a small yellow square representing ε_i .

Existing regularization methods:

- ▶ **Convex surrogates (nuclear norms):** statistically accurate, but need repeated SVDs on large unfoldings \Rightarrow extremely slow.
- ▶ **Nonconvex factorizations:** computationally cheaper, but sensitive to initialization and with weaker guarantees.

Model

- For convenience, we focus on **order-3 low-rank tensor regression**:

$$y_i = \langle \mathcal{A}, \mathcal{X}_i \rangle + \varepsilon_i, \quad i = 1, \dots, n.$$

$$y_i = \langle \overset{\mathcal{A}}{\text{orange cube}}, \overset{\mathcal{X}_i}{\text{yellow/red cube}} \rangle + \overset{\varepsilon_i}{\text{yellow square}}$$

- Here, \mathcal{A} is **Tucker low-rank**: $\mathcal{A} = \mathcal{S} \times_1 U_1 \times_2 U_2 \times_3 U_3$, where $\mathcal{S} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$, $U_k \in \mathbb{R}^{p_k \times r_k}$, $k = 1, 2, 3$.

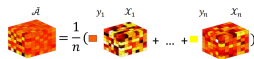
$$\overset{\mathcal{A}}{\text{orange cube}} = \overset{\mathcal{S} \times_1 U_1 \times_2 U_2 \times_3 U_3}{\text{blue, white, green, red slices}}$$

- **Goal:** estimate \mathcal{A} based on $\{y_i, \mathcal{X}_i\}_{i=1}^n$.

Step 1. Probing Importance Sketching Direction

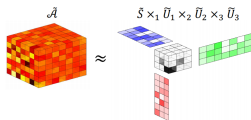
- ▶ **(Step 1.1)** Evaluate the sample covariance tensor:

$$\tilde{\mathcal{A}} = \frac{1}{n} \sum_{i=1}^n y_i \mathcal{X}_i$$



- ▶ **(Step 1.2)** Apply high-order orthogonal iteration (HOOI) to obtain a low-rank factorization:

$$\tilde{\mathcal{A}} \approx \tilde{\mathcal{S}} \times_1 \tilde{U}_1 \times_2 \tilde{U}_2 \times_3 \tilde{U}_3$$



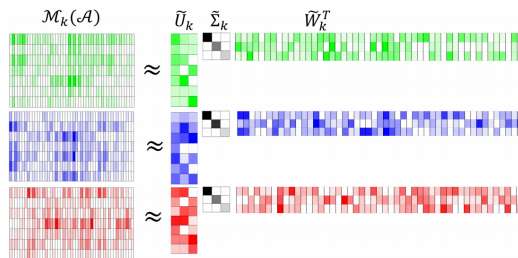
- ▶ **(Step 1.3)** Perform QR orthogonalization:

$$\tilde{V}_k = \text{QR}(\mathcal{M}_k^\top(\tilde{\mathcal{S}})).$$

- ▶ **Outcome:** $\{\tilde{U}_k, \tilde{V}_k\}_{k=1}^3$.

Interpretations of Step 1

$$\mathcal{M}_k(\mathcal{A}) \approx \tilde{U}_k \tilde{\Sigma}_k \tilde{W}_k^\top, \quad \tilde{W}_k = (\tilde{U}_{k+2} \otimes \tilde{U}_{k+1}) \tilde{V}_k.$$



- ▶ $\{\tilde{U}_k, \tilde{W}_k\}$ are **importance sketching directions**.
- ▶ They are initial sample approximations of $\{U_k, W_k\}$, i.e. left/right singular subspaces of $\mathcal{M}_k(\mathcal{A})$.
- ▶ They best align with the true tensor \mathcal{A} .

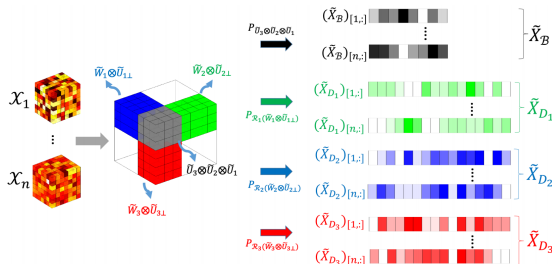
Step 2. Importance Sketching

- Construct **dimension-reduced covariates**:

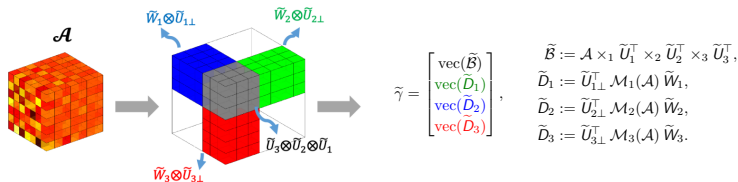
$$\hat{X}_B[i, :] = \text{vec}\left(\mathcal{X}_i \times_1 \tilde{U}_1^\top \times_2 \tilde{U}_2^\top \times_3 \tilde{U}_3^\top\right),$$

$$\hat{X}_{D_k}[i, :] = \text{vec}\left(\tilde{U}_{k\perp}^\top \mathcal{M}_k(\mathcal{X}_i) \tilde{W}_k\right), \quad k = 1, 2, 3.$$

- These sketches reduce both sample and feature dimensions.



Interpretation of Step 2



- Rewrite the regression model:

$$y_i = \langle \mathcal{X}_i, \mathcal{A} \rangle + \varepsilon_i = \tilde{X}[i, :]^\top \tilde{\gamma} + \tilde{\varepsilon}_i,$$

where

$$\tilde{X} = [\hat{X}_{\mathcal{B}}, \hat{X}_{D_1}, \hat{X}_{D_2}, \hat{X}_{D_3}], \quad \tilde{\gamma} = [\text{vec}(\tilde{\mathcal{B}}), \text{vec}(\tilde{D}_1), \text{vec}(\tilde{D}_2), \text{vec}(\tilde{D}_3)].$$

- \tilde{X} are sketching covariates.
- $\tilde{\gamma}$ is the sketch of \mathcal{A} .

Step 3. Dimension-Reduced Regression

- Perform regression in reduced space:

$$\hat{\gamma} = \arg \min_{\gamma} \|y - \tilde{X}\gamma\|_2^2.$$

- Dimension of parameter reduces from $p_1 p_2 p_3$ to

$$m = r_1 r_2 r_3 + \sum_{k=1}^3 (p_k - r_k) r_k.$$

$$\hat{\gamma} = \arg \min_{\gamma} \|y - \tilde{X}\gamma\|_2^2, \quad \tilde{X} = [\tilde{X}_B \quad \tilde{X}_{D_1} \quad \tilde{X}_{D_2} \quad \tilde{X}_{D_3}]$$

$$y \approx X_B \hat{\gamma}_B + X_{D_1} \hat{\gamma}_{D_1} + X_{D_2} \hat{\gamma}_{D_2} + X_{D_3} \hat{\gamma}_{D_3}$$
$$\approx \begin{bmatrix} X_B & X_{D_1} & X_{D_2} & X_{D_3} \end{bmatrix} \times \begin{bmatrix} \hat{\gamma}_B \\ \hat{\gamma}_{D_1} \\ \hat{\gamma}_{D_2} \\ \hat{\gamma}_{D_3} \end{bmatrix}$$

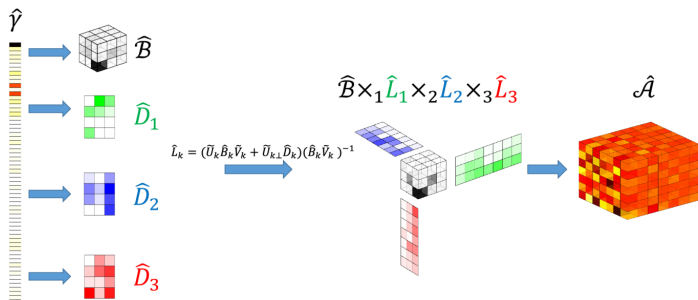
Step 4. Assembling the Final Estimate

- Reconstruct $\hat{\mathcal{A}}$ via the **Cross scheme** (Z. AoS, 2018):

$$\hat{\mathcal{A}} = \hat{\mathcal{B}} \times_1 \hat{\mathcal{L}}_1 \times_2 \hat{\mathcal{L}}_2 \times_3 \hat{\mathcal{L}}_3,$$

where each

$$\hat{\mathcal{L}}_k = (\tilde{U}_k \mathcal{M}_k(\hat{\mathcal{B}}) \tilde{V}_k + \tilde{U}_{k\perp} \hat{D}_k) (\mathcal{M}_k(\hat{\mathcal{B}}) \tilde{V}_k)^{-1}, \quad k = 1, 2, 3.$$



Algorithm: ISLET (Summary)

ISLET (Importance Sketching for Tensor Regression)

Input: samples $\{(\mathcal{X}_i, y_i)\}_{i=1}^n$, target ranks (r_1, r_2, r_3)

Output: estimate $\hat{\mathcal{A}}$

Step 1: Probing directions

$$\begin{aligned}\tilde{\mathcal{A}} &\leftarrow \frac{1}{n} \sum_{i=1}^n y_i \mathcal{X}_i \\ (\tilde{\mathcal{S}}, \tilde{U}_1, \tilde{U}_2, \tilde{U}_3) &\leftarrow \text{HOOI}(\tilde{\mathcal{A}}; r_1, r_2, r_3) \\ \tilde{V}_k &\leftarrow \text{QR}(\mathcal{M}_k(\tilde{\mathcal{S}})^\top) \text{ for } k = 1, 2, 3\end{aligned}$$

Step 2: Importance sketching (build reduced covariates)

$$\begin{aligned}\hat{X}_B[i, :] &\leftarrow \text{vec}(\mathcal{X}_i \times_1 \tilde{U}_1^\top \times_2 \tilde{U}_2^\top \times_3 \tilde{U}_3^\top) \\ \tilde{W}_k &\leftarrow (\tilde{U}_{k+2} \otimes \tilde{U}_{k+1}) \tilde{V}_k, \quad k = 1, 2, 3 \\ \hat{X}_{D_k}[i, :] &\leftarrow \text{vec}(\tilde{U}_{k\perp}^\top \mathcal{M}_k(\mathcal{X}_i) \tilde{W}_k), \quad k = 1, 2, 3 \\ \hat{X} &\leftarrow [\hat{X}_B, \hat{X}_{D_1}, \hat{X}_{D_2}, \hat{X}_{D_3}]\end{aligned}$$

Step 3: Dimension-reduced regression

$$\hat{\gamma} \leftarrow \arg \min_{\gamma} \|y - \hat{X}\gamma\|_2^2 \quad (\text{use Group Lasso if sparse})$$

Step 4: Assemble estimate

Reconstruct $\hat{\mathcal{A}}$ from $\hat{\gamma}$ and $(\tilde{U}_k, \tilde{V}_k)$ via the Cross scheme.

Oracle Inequalities (General Design)

Theorem (Oracle Inequality)

Consider order-3 tensor regression with low Tucker rank (r_1, r_2, r_3) . Under mild conditions (angle error $\theta < \frac{1}{2}$, nonsingular sketches, GRIP in the sparse case), the ISLET estimator satisfies

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{HS}}^2 \leq C \left(\frac{\sigma^2 m}{n} + \text{bias}(\theta) \right),$$

where $m = r_1 r_2 r_3 + \sum_{k=1}^3 (p_k - r_k) r_k$.

Takeaway. Error decomposes as variance $O(m/n)$ plus controlled sketching bias.

Optimal Risk under Gaussian Design

Theorem (Minimax Risk under Gaussian Design)

Suppose the observed variables is i.i.d. Gaussian and the noise is $\mathcal{N}(0, \sigma^2)$. Then

$$\mathbb{E} \|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{HS}}^2 = (1 + o(1)) \frac{m\sigma^2}{n}, \quad m = r_1 r_2 r_3 + \sum_{k=1}^3 (p_k - r_k) r_k.$$

Takeaway. Achieves minimax-optimal rate with sharp constant; m equals the degrees of freedom of the Tucker rank class.

Simulation Study: Experimental Setup

Goal. Evaluate the performance of ISLET under synthetic low-rank tensor regression.

Data Generation.

- ▶ Covariates: $\mathcal{X}_j \in \mathbb{R}^{p \times p \times p}$ with i.i.d. $\mathcal{N}(0, 1)$ entries.
- ▶ Coefficient tensor:

$$\mathcal{A} = \llbracket \mathcal{S}; E_1, E_2, E_3 \rrbracket$$

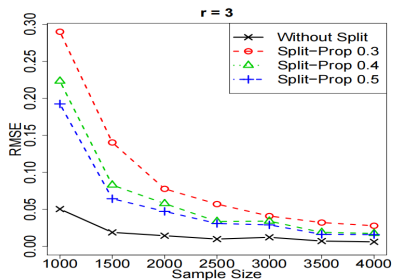
- *Nonsparse setting*: \mathcal{S} and E_k Gaussian random. - *Sparse setting*: rows of E_k randomly zeroed out, sparsity level s_k .

- ▶ Responses: $y_j = \langle \mathcal{X}_j, \mathcal{A} \rangle + \varepsilon_j$, $\varepsilon_j \sim \mathcal{N}(0, \sigma^2)$.

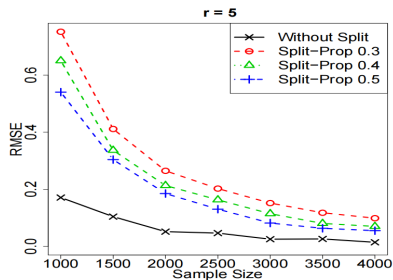
Evaluation.

- ▶ Normalized RMSE: $\|\hat{\mathcal{A}} - \mathcal{A}\|_{\text{HS}} / \|\mathcal{A}\|_{\text{HS}}$.
- ▶ Results averaged over 100 repetitions.

Simulation Study: Results

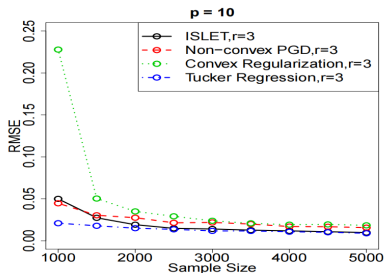


(a) $r = 3$

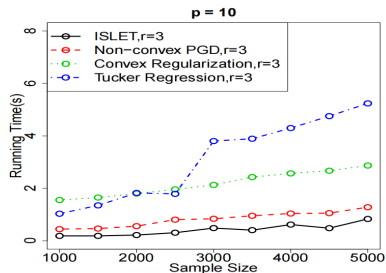


(b) $r = 5$

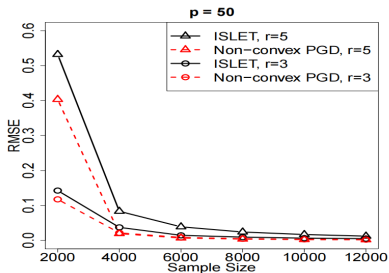
Simulation Study: Results



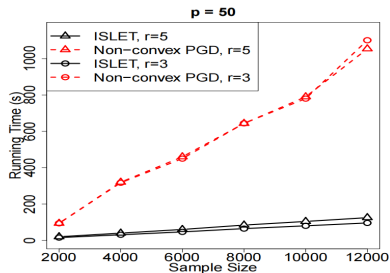
(a) RMSE



(b) Run Time



(c) RMSE



(d) Run Time