

TenIPS: Inverse Propensity Sampling for Tensor Completion

Chengrun Yang, Lijun Ding, Ziyang Wu, Madeleine Udell

Cornell University

Tensors

On an *order-3* tensor \mathcal{B} , for each of the modes $n \in [3] := \{1, 2, 3\}$:

- size of the n -th mode: I_n
- mode- n fibers: fixing every index but the n -th. e.g., mode-1 fiber: $\mathcal{B}_{:jk}$
- mode- n unfolding: matrix $\mathcal{B}^{(n)}$, whose columns are mode- n fibers

tensor decomposition: CP, Tucker, tensor-train, ...

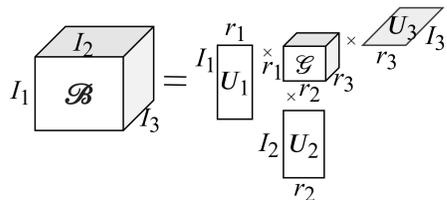


Figure 1: Tucker decomposition with multilinear rank (r_1, r_2, r_3) : $\mathcal{B} = \mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3$.

tensor completion

Given a partially observed $\mathcal{B}_{\text{obs}} \in \mathbb{R}^{I_1 \times \dots \times I_N}$, we have

- observation pattern $\Omega \in \mathbb{R}^{I_1 \times \dots \times I_N}$: $\Omega_{i_1 \dots i_N} = 1$ if $\mathcal{B}_{i_1 \dots i_N}$ is observed, and 0 otherwise
- observation probability $\mathcal{P} \in \mathbb{R}^{I_1 \times \dots \times I_N}$: $\mathcal{P}_{i_1 \dots i_N} = \mathbb{P}(\Omega_{i_1 \dots i_N} = 1) = \mathbb{P}(\mathcal{B}_{i_1 \dots i_N} \text{ is observed})$

missingness types	$\{\mathcal{P}_{i_1 \dots i_N}\}$
missing-completely-at-random (MCAR)	uniform
missing-not-at-random (MNAR)	non-uniform

1-bit matrix completion

Given a binary matrix $Y \in \{0, 1\}^{m \times n}$, predict the parameter matrix $M \in \mathbb{R}^{m \times n}$

Assumptions:

- M is approximately low rank.
- There exists a link function $\sigma: \mathbb{R} \rightarrow [0, 1]$, such that $\mathbb{P}(Y_{ij} = 1) = \sigma(M_{ij})$ for $(i, j) \in [m] \times [n]$.

Low rank surrogates for M : low nuclear norm, low max norm, ...

Our problem formulation: MNAR tensor completion

Input: MNAR data tensor $\mathcal{B}_{\text{obs}} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$

Assumptions:

- true data tensor $\mathcal{B} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is **approximately low multilinear rank**
- noiseless observation:** $(\mathcal{B}_{\text{obs}})_{i_1 \dots i_N} = \mathcal{B}_{i_1 \dots i_N}$ if $\mathcal{B}_{i_1 \dots i_N}$ is observed, and 0 otherwise
- unknown parameter tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ has the **same rank structure** as \mathcal{B}
- 1-bit observation:** With the observation propensity tensor $\mathcal{P} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$, $\mathbb{P}(\mathcal{B}_{i_1 i_2 \dots i_N} \text{ is observed}) = \mathcal{P}_{i_1 i_2 \dots i_N} = \sigma(\mathcal{A}_{i_1 i_2 \dots i_N})$, in which $\sigma: \mathbb{R} \rightarrow [0, 1]$ is a non-decreasing link function.

Algorithm Step 1: propensity recovery

Given a mask tensor Ω , get a predicted propensity tensor $\hat{\mathcal{P}}$.

	algorithm	hyperparameters
Choice 1: convex	proximal-proximal-gradient	τ and γ
Choice 2: nonconvex	gradient descent	target rank and step size

Choice 1: convex and provable

- get the **square set** and **square unfolding** [5] of Ω :
 - $c_S: [N] \rightarrow [N]$: a permutation map of the N orders that satisfies $\{c_S^{-1}(1), c_S^{-1}(2), \dots, c_S^{-1}(|S|)\} = S$
 - square set of $\Omega \in \mathbb{R}^{I_1 \times \dots \times I_N}$: $S_{\square} := \arg \min_{S \subset [N]} |\prod_{n \in S} I_n - \prod_{n \in [N] \setminus S} I_n|$
 - and the square unfolding $\Omega_{\square} := \text{reshape}(c_S \square \Omega^{(1)}, \prod_{n \in S_{\square}} I_n, \prod_{n \in [N] \setminus S_{\square}} I_n)$
- predict parameter \mathcal{A} by **logistic loss minimization** (by proximal-proximal-gradient [6])

$$\hat{\mathcal{A}}_{\square} \leftarrow \underset{\Gamma \in \mathcal{S}_{\tau, \gamma}}{\text{argmin}} \sum_{i=1}^{I_{\square}} \sum_{j=1}^{I_{\square}} -(\Omega_{\square})_{i,j} \log \sigma(\Gamma_{i,j}) - [1 - (\Omega_{\square})_{i,j}] \log [1 - \sigma(\Gamma_{i,j})],$$

$$\text{where } \mathcal{S}_{\tau, \gamma} = \left\{ \Gamma \in \mathbb{R}^{I_{\square} \times I_{\square}} : \|\Gamma\|_{\star} \leq \tau \sqrt{I_{[N]}}, \|\Gamma\|_{\max} \leq \gamma \right\}.$$

- predict propensity: $\hat{\mathcal{P}} = \sigma(\hat{\mathcal{A}})$

Choice 2: nonconvex, gradient descent

- initialization: $\mathcal{G}^A, U_1^A, \dots, U_N^A \leftarrow \mathcal{G}_0^A, (U_1)_0^A, \dots, (U_N)_0^A$

- objective function:

$$f(\mathcal{G}^A, \{U_n^A\}_{n \in [N]}) = \sum_{i_1 \dots i_N} -\Omega_{i_1 \dots i_N} \log \sigma(\hat{\mathcal{A}}_{i_1 \dots i_N}) - (1 - \Omega_{i_1 \dots i_N}) \log [1 - \sigma(\hat{\mathcal{A}}_{i_1 \dots i_N})],$$

$$\text{in which } \hat{\mathcal{A}} = \mathcal{G}^A \times_1 U_1^A \times_2 \dots \times_N U_N^A.$$

- gradient descent updates

- predict propensity: $\hat{\mathcal{P}} = \sigma(\mathcal{G}^A \times_1 U_1^A \times_2 \dots \times_N U_N^A)$

Algorithm Step 2: tensor completion

Given $\hat{\mathcal{P}}$ and MNAR observations \mathcal{B}_{obs} , get $\hat{\mathcal{B}}$

- Form an **entrywise inverse propensity estimator** for data tensor \mathcal{B} as $\tilde{\mathcal{X}}(\hat{\mathcal{P}}) = \sum_{(i_1, i_2, \dots, i_N) \in \Omega} \frac{1}{\mathcal{P}_{i_1 \dots i_N}} \mathcal{B}_{\text{obs}} \odot \mathcal{E}(i_1, \dots, i_N)$, in which
 - $\Omega := \{(i_1, \dots, i_N) | \mathcal{B}_{i_1 \dots i_N} \text{ is observed}\}$
 - $\mathcal{E}(i_1, \dots, i_N)$ is a binary tensor with the same shape as \mathcal{B} , with value 1 at the (i_1, i_2, \dots, i_N) -th entry and 0 elsewhere.
- Do **Tucker decomposition** on $\tilde{\mathcal{X}}(\hat{\mathcal{P}})$, get core tensor $\mathcal{W}(\hat{\mathcal{P}})$ and factor matrices $\{Q_n(\hat{\mathcal{P}})\}_{n \in [N]}$.
- Estimate \mathcal{B} by $\hat{\mathcal{B}}(\hat{\mathcal{P}}) = \mathcal{W}(\hat{\mathcal{P}}) \times_1 Q_1(\hat{\mathcal{P}}) \times_2 \dots \times_N Q_N(\hat{\mathcal{P}})$.

Theoretical guarantees

- Upper bound for propensity recovery error [1, 3]**
Assume that $\mathcal{P} = \sigma(\mathcal{A})$. Given a set $S \subset [N]$, together with the following assumptions:
A1. \mathcal{A}_S has bounded nuclear norm: there exists a constant $\theta > 0$ such that $\|\mathcal{A}_S\|_{\star} \leq \theta \sqrt{I_{[N]}}$.
A2. Entries of \mathcal{A} have bounded absolute value: there exists a constant $\alpha > 0$ such that $\|\mathcal{A}\|_{\max} \leq \alpha$.
 Suppose we run the convex propensity recovery algorithm with thresholds satisfying $\tau \geq \theta$ and $\gamma \geq \alpha$ to obtain an estimate $\hat{\mathcal{P}}$ of \mathcal{P} . With $L_{\gamma} := \sup_{x \in [-\gamma, \gamma]} \frac{|\sigma'(x)|}{\sigma(x)(1-\sigma(x))}$, there exists a universal constant $C > 0$ such that if $I_S + I_{S^c} \geq C$, with probability at least $1 - \frac{C}{I_S + I_{S^c}}$, the propensity estimation error $\frac{1}{I_{[N]}} \|\hat{\mathcal{P}} - \mathcal{P}\|_{\text{F}}^2 \leq 4eL_{\gamma}\tau \left(\frac{1}{\sqrt{I_S}} + \frac{1}{\sqrt{I_{S^c}}} \right)$.
- Optimality of the square unfolding for propensity recovery:** Instate the same conditions as the previous lemma on propensity recovery error, and further assume that there exists a constant $c > 0$ such that $r_n^{\text{true}} \leq cI_n$ for every $n \in [N]$. Then $S = S_{\square}$ gives the tightest upper bound on the propensity estimation error $\|\hat{\mathcal{P}} - \mathcal{P}\|_{\text{F}}$ among all unfolding sets $S \subset [N]$.

- Tensor completion error on cubical tensors** (same size in every mode):
Consider an order- N cubical tensor \mathcal{B} with size $I_1 = \dots = I_N = I$ and multilinear rank $r_1^{\text{true}} = \dots = r_N^{\text{true}} = r < I$, and two order- N cubical tensors \mathcal{P} and \mathcal{A} with the same shape as \mathcal{B} . Each entry of \mathcal{B} is observed with probability from the corresponding entry of \mathcal{P} . Assume $I \geq rN \log I$, and there exist constants $\psi, \alpha \in (0, \infty)$ such that $\|\mathcal{A}\|_{\max} \leq \alpha, \|\mathcal{B}\|_{\max} = \psi$. Further assume that for each $n \in [N]$, the condition number $\frac{\sigma_1(\mathcal{B}^{(n)})}{\sigma_r(\mathcal{B}^{(n)})} \leq \kappa$ is a constant independent of tensor sizes and dimensions. Then under the conditions of the lemma on convex propensity recovery error, with probability at least $1 - I^{-1}$, the fixed multilinear rank (r, r, \dots, r) approximation $\hat{\mathcal{B}}(\hat{\mathcal{P}})$ computed from the convex propensity recovery and tensor completion algorithms with thresholds $\tau \geq \theta$ and $\gamma \geq \alpha$ satisfies

$$\frac{\|\hat{\mathcal{B}}(\hat{\mathcal{P}}) - \mathcal{B}\|_{\text{F}}}{\|\mathcal{B}\|_{\text{F}}} \leq CN \sqrt{\frac{r \log I}{I}},$$

in which C depends on κ .

Numerics

Convex propensity recovery on a size-8 cubical tensor:

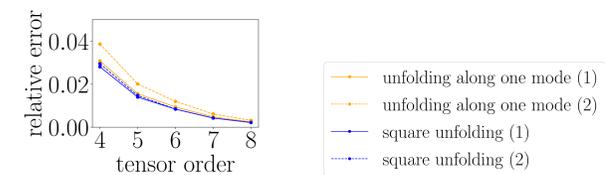


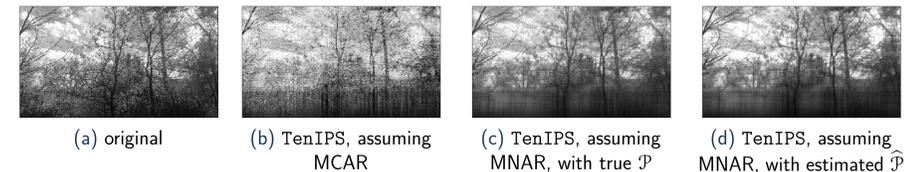
Figure 2: "(1)": setting $\tau = \theta, \gamma = \alpha$; "(2)": setting $\tau = 2\theta, \gamma = 2\alpha$

MNAR tensor completion on synthetic data:

Algorithm	time (s)	relative error $\frac{\ \hat{\mathcal{B}}(\hat{\mathcal{P}}) - \mathcal{B}\ _{\text{F}}}{\ \mathcal{B}\ _{\text{F}}}$		
		with \mathcal{P}	with $\hat{\mathcal{P}}_1$	with $\hat{\mathcal{P}}_2$
TenIPS	26	0.110	0.110	0.109
HOSVD_w [2]	35	0.129	0.116	0.110
SqUnfold	29	0.141	0.138	0.139
RectUnfold	8	0.259	0.256	0.256
LstSq	>600	-	-	-
SO-HOSVD [7]	>600	-	-	-

MNAR tensor completion on semi-synthetic data:

- real video tensor from [4]: $\mathcal{B} \in [0, 255]^{2200 \times 1080 \times 1920}$
- synthetic parameter tensor $\mathcal{A} = (\mathcal{B} - 128)/64$



Thanks!

- Chengrun Yang: cy438@cornell.edu
- Madeleine Udell: udell@cornell.edu

Bibliography

- Mark A Davenport, Yaniv Plan, Ewout Van Den Berg, and Mary Woorters. 1-bit matrix completion. *Information and Inference: A Journal of the IMA*, 3(3):189–223, 2014.
- Longxiu Huang and Deanna Needell. HOSVD-based algorithm for weighted tensor completion. *arXiv preprint arXiv:2003.08537*, 2020.
- Wei Ma and George H Chen. Missing not at random in matrix completion: The effectiveness of estimating missingness probabilities under a low nuclear norm assumption. In *Advances in Neural Information Processing Systems*, pages 14871–14880, 2019.
- Osman Asif Malik and Stephen Becker. Low-rank Tucker decomposition of large tensors using tensorsketch. In *Advances in Neural Information Processing Systems*, pages 10096–10106, 2018.
- Cun Mu, Bo Huang, John Wright, and Donald Goldfarb. Square deal: Lower bounds and improved relaxations for tensor recovery. In *International conference on machine learning*, pages 73–81, 2014.
- Ernest K Ryu and Wotao Yin. Proximal-proximal-gradient method. *arXiv preprint arXiv:1708.06908*, 2017.
- Dong Xia, Ming Yuan, and Cun-Hui Zhang. Statistically optimal and computationally efficient low rank tensor completion from noisy entries. *arXiv preprint arXiv:1711.04934*, 2017.