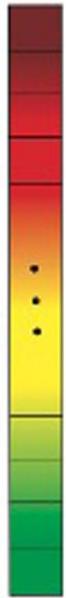


# Efficient Nonconvex Regularized Tensor Completion with Structure-aware Proximal Iterations

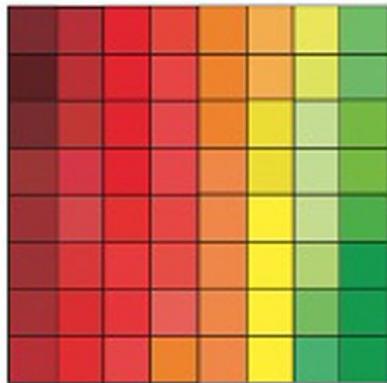
Presenter: *Quanming Yao* (4Paradigm)

Joint work with *James T. Kwok* (HKUST) and *Bo Han* (RIKEN)

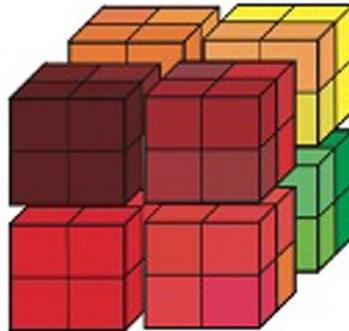
# What is Tensor?



Vector



Matrix



Tensor

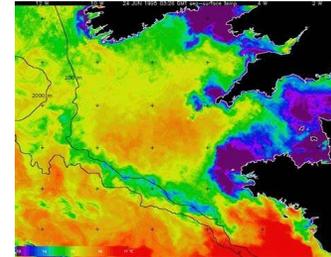


Capture higher order interactions inside the data

Color images



Remote sensing data



Spatiotemporal recommendation  
*where to eat: (user, location, action)*

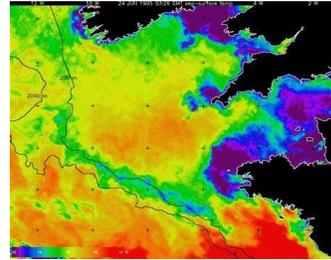
Examples of 3-order tensors

# Why needs Tensor Completion?

Color images



Image inpainting



Remote sensing data



Missing super-pixel / bands

Spatiotemporal recommendation  
*where to eat: (user, location, action)*



Predict unknown triplet

Tensor completion: predict missing entries in the tensor

On 2-order tensor: reduce to matrix completion

# How? - Overlapped nuclear norm

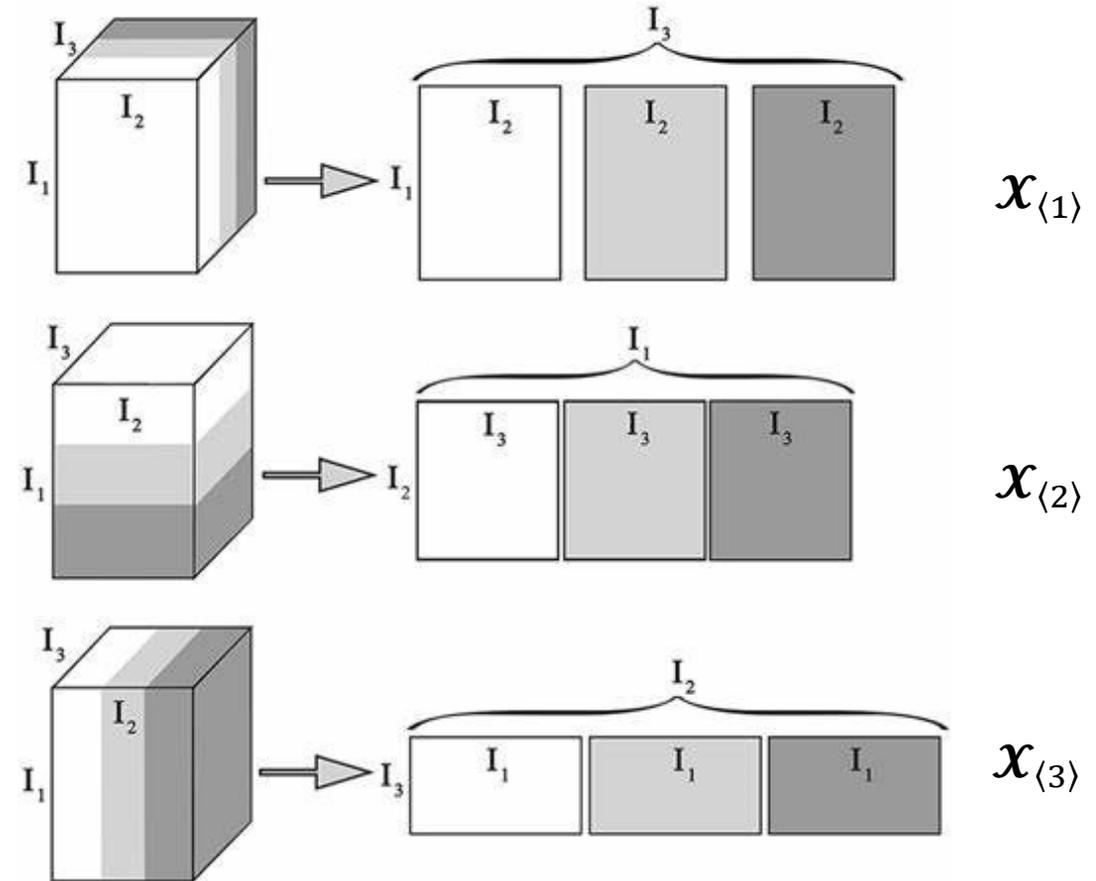
Nuclear norm  $\|\mathbf{X}\|_*$  [Candes & Recht, 2009]

- Summation of all **singular values** of a matrix
- convex envelope of the matrix rank function

Tensor: **overlapped nuclear norm** [Tomioka et al., 2010]

**Definition 1.** For a  $M$ -order tensor  $\mathcal{X}$ , the overlapped nuclear norm is  $\|\mathcal{X}\|_{\text{overlap}} = \sum_{m=1}^M \lambda_m \|\mathcal{X}_{\langle m} \rangle\|_*$ , where  $\{\lambda_m \geq 0\}$  are hyperparameters.

- $\mathcal{X}_{\langle m} \rangle$  unfold tensor along with  $m$ th mode
- encourage **all unfolded matrix** to be low-rank



unfolding operations

folding is the inverse of unfolding

# Tensor Completion with Overlapped Nuclear Norm [Tomioka et al., 2010]

- Redundancy and correlations → low-rank approach is a power method in tensor completion
- Overlapped nuclear norm is a sound approach with **statistical and convergence guarantee** (compared with other tensor low-rank approaches [Tomioka et al., 2011; Liu et al., 2013; Guo et al., 2017])

$$\min_{\mathbf{x}} \underbrace{\frac{1}{2} \|P_{\Omega}(\mathbf{x} - \mathbf{O})\|_F^2}_{\text{Squared loss on observed entries}} + \underbrace{\sum_{d=1}^D \lambda_m \|\mathbf{x}_{\langle m} \|\ast}_{\text{Overlapped nuclear norm}}$$

However, two critical limitations

1. **Inferior** empirical performance
  - nuclear norm over penalize singular values
2. **Expensive** optimization
  - full tensor needs to be maintained due to folding / unfolding operations

# Proposed NORT: Nonconvex regularized tensor completion

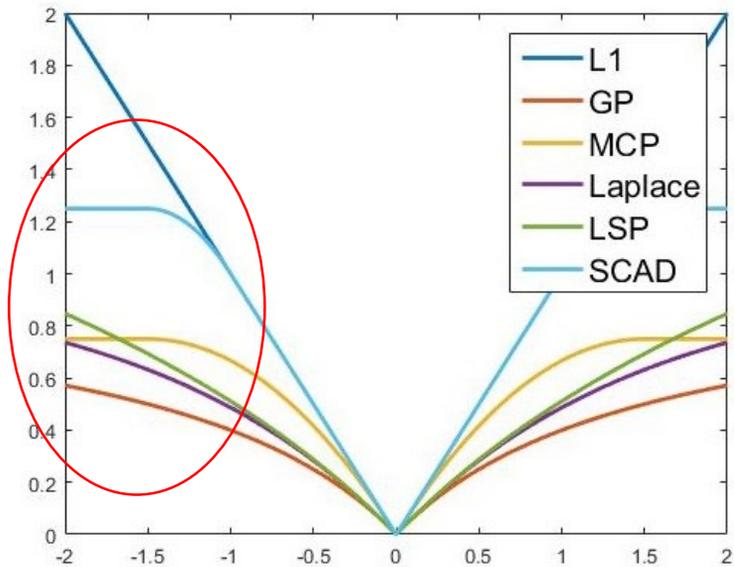
Our contributions, propose NORT algorithm

1. **Improve the performance** of overlapped nuclear norm
  - Extend **nonconvex regularization** with overlapped nuclear norm
2. **Speedup optimization** with structure aware proximal iterations
  - **Cheap iteration**: propose a special way to do matrix multiplication without tensor folding/unfolding
  - **Fast convergence**: enhance proximal average with adaptive momentum

# Improve Performance: nonconvex regularization

Nonconvex regularization

Objective: 
$$\min_{\mathbf{x}} F(\mathbf{X}) \equiv \frac{1}{2} \|P_{\Omega}(\mathbf{X} - \Theta)\|_F^2 + \sum_{d=1}^D \frac{\lambda_d}{D} \phi(\mathbf{x}_{\langle d \rangle}).$$
 where 
$$\phi(\mathbf{X}) = \sum_{i=1}^n \kappa(\sigma_i(\mathbf{X})),$$



Common examples of  $\kappa(\sigma_i(\mathbf{X}))$ . Here,  $\theta$  is a constant. For capped- $\ell_1$ , LSP and MCP,  $\theta > 0$ ; for SCAD,  $\theta > 2$ ; and for TNN,  $\theta$  is a positive integer.

	$\kappa(\sigma_i(\mathbf{X}))$
nuclear norm	$\sigma_i(\mathbf{X})$
capped- $\ell_1$	$\min(\sigma_i(\mathbf{X}), \theta)$
LSP	$\log(\sigma_i(\mathbf{X})/\theta + 1)$
TNN [27]	$\begin{cases} \sigma_i(\mathbf{X}) & \text{if } i > \theta \\ 0 & \text{otherwise} \end{cases}$
SCAD	$\begin{cases} \sigma_i(\mathbf{X}) & \text{if } \sigma_i(\mathbf{X}) \leq 1 \\ \frac{2\theta\sigma_i(\mathbf{X}) - \sigma_i(\mathbf{X})^2 - 1}{2(\theta - 1)} & \text{if } 1 < \sigma_i(\mathbf{X}) \leq \theta \\ (\theta + 1)^2/2 & \text{otherwise} \end{cases}$
MCP	$\begin{cases} \sigma_i(\mathbf{X}) - \alpha^2/2\theta & \text{if } \sigma_i(\mathbf{X}) \leq \theta \\ \theta^2/2 & \text{otherwise} \end{cases}$

**Less penalize** large singular values, which are more informative

# Speedup optimization: Structure-aware proximal iterations

Proximal average algorithm [Bauschke et al., 2008; Yu, 2013]

$$\mathbf{x}_t = \frac{1}{K} \sum_{i=1}^K \mathbf{y}_t^i,$$

*maintain low-rank factorization*

$$\mathbf{z}_t = \mathbf{x}_t - \frac{1}{\tau} \nabla f(\mathbf{x}_t),$$

*sparse plus low-rank structure*

$$\mathbf{y}_{t+1}^i = \text{prox}_{\frac{\lambda_i}{\tau} g_i}(\mathbf{z}_t), \quad i = 1, \dots, K.$$

proximal step with nonconvex regularization

$$\mathbf{y}_{t+1}^i = \left[ \text{prox}_{\frac{\lambda_i}{\tau} \phi}([\mathbf{z}_t]_{\langle i \rangle}) \right]_{\langle i \rangle}$$

*matrix multiplications*

$$\mathbf{x}_t = \frac{1}{D} \sum_{i=1}^D (\mathbf{U}_t^i (\mathbf{V}_t^i)^\top)_{\langle i \rangle}$$

$$\mathbf{z}_t = \frac{1}{D} \sum_{i=1}^D (\mathbf{U}_t^i (\mathbf{V}_t^i)^\top)_{\langle i \rangle} - \frac{1}{\tau} P_\Omega(\mathbf{x}_t - \Theta).$$

*utilize sparse plus low-rank structure to efficient compute proximal step (Proposition 3.2)*

$$[\mathbf{z}_t]_{\langle i \rangle} \mathbf{b} = \frac{1}{D} \mathbf{U}_t^i [(\mathbf{V}_t^i)^\top \mathbf{b}] + \frac{1}{D} \sum_{j \neq i} [(\mathbf{U}_t^j (\mathbf{V}_t^j)^\top)_{\langle j \rangle}]_{\langle i \rangle} \mathbf{b} - \frac{1}{\tau} [P_\Omega(\mathbf{x}_t - \Theta)]_{\langle i \rangle} \mathbf{b},$$

**Needs** folding/unfolding: full tensor computation

**No** folding/unfolding: fast and need less memory

Table 1. Comparison of the proposed NORT (Algorithm 1) and direct implementations of the PA algorithm.

	per-iteration time complexity	space	convergence
direct	$O(I_{\times} \sum_{i=1}^D I_i)$	$O(I_{\times})$	slow
NORT	$O(\sum_{i=1}^D \sum_{j \neq i} (\frac{1}{I_i} + \frac{1}{I_j}) k_t^i k_{t+1}^i I_{\times} + \ \Omega\ _1 (k_t^i + k_{t+1}^i))$	$O(\sum_{i=1}^D \sum_{j \neq i} (\frac{1}{I_i} + \frac{1}{I_j}) k_t^i I_{\times} + \ \Omega\ _1)$	fast

**Algorithm 1** NOnconvex REgularized Tensor (NORT).

```

1: initialize  $\mathbf{X}_0 = \mathbf{X}_1 = 0$ ,  $\tau > \rho + DL$  and  $\gamma_1, p \in (0, 1)$ ;
2: for  $t = 1, \dots, T$  do
3:    $\mathbf{X}_{t+1} = \frac{1}{D} \sum_{i=1}^D (\mathbf{U}_{t+1}^i (\mathbf{V}_{t+1}^i)^{\top})^{(i)}$ ;
4:    $\bar{\mathbf{X}}_t = \mathbf{X}_t + \gamma_t (\mathbf{X}_t - \mathbf{X}_{t-1})$ ;
5:   if  $F_{\tau}(\bar{\mathbf{X}}_t) \leq F_{\tau}(\mathbf{X}_t)$  then
6:      $\mathbf{V}_t = \bar{\mathbf{X}}_t$ ,  $\gamma_{t+1} = \min(\frac{\gamma_t}{p}, 1)$ ;
7:   else
8:      $\mathbf{V}_t = \mathbf{X}_t$ ,  $\gamma_{t+1} = p\gamma_t$ ;
9:   end if
10:   $\mathbf{Z}_t = \mathbf{V}_t - \frac{1}{\tau} P_{\Omega}(\mathbf{V}_t - \mathbf{O})$ ;
    // compute  $P_{\Omega}(\mathbf{V}_t - \mathbf{O})$  using sparse tensor format;
11:  for  $i = 1, \dots, D$  do
12:     $\mathbf{X}_{t+1}^i = \text{prox}_{\frac{\lambda_i}{\tau} \phi}((\mathbf{Z}_t)^{(i)})$ ; // keep as  $\mathbf{U}_t^i (\mathbf{V}_t^i)^{\top}$ ;
13:  end for
14: end for
output  $\mathbf{X}_{T+1}$ .

```

adaptive  
momentum

**Theorem 3.5.** *The sequence  $\{\mathbf{X}_t\}$  generated from Algorithm 1 has at least one limit point, and all limits points are critical points of  $F_{\tau}(\mathbf{X})$ .*

**Theorem 3.7.** *Let  $r_t = F_{\tau}(\mathbf{X}_t) - F_{\tau}^{\min}$ . If  $F_{\tau}$  has the uniformized KL property, for a sufficiently large  $t_0$ , we have*

1. *If  $\beta = 1$ ,  $r_t$  reduces to zero in finite steps;*
2. *If  $\beta \in [\frac{1}{2}, 1)$ ,  $r_t \leq (\frac{d_1 C^2}{1+d_1 C^2})^{t-t_0} r_{t_0}$  where  $d_1 = \frac{2(\tau+\rho)^2}{\eta}$ ;*
3. *If  $\beta \in (0, \frac{1}{2})$ ,  $r_t \leq (\frac{C}{(t-t_0)d_2(1-2\beta)})^{1/(1-2\beta)} r_{t_0}$  where  $d_2 = \min\{\frac{1}{2d_1 C}, \frac{C}{1-2\beta} (2^{\frac{2\beta-1}{2\beta-2}} - 1)\}$ .*

- tensor size:  $I_1 \times I_2 \times I_3$
- the **speedup** can be more the **100x** on large tensors

# Experiments: synthetic data

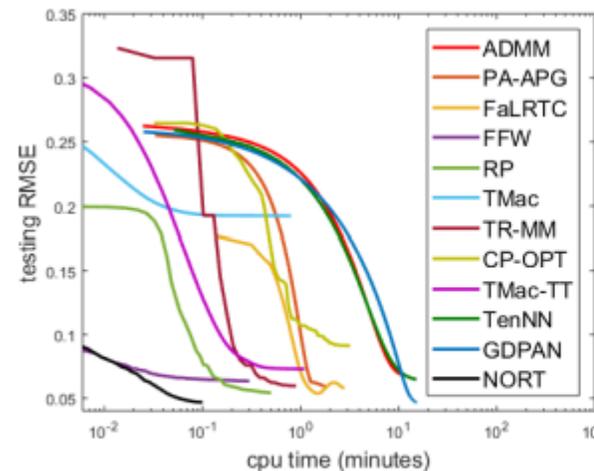
		small $I_3$ : $\bar{c} = 100$ , sparsity: 3.09%			large $I_3$ : $\hat{c} = 40$ , sparsity: 2.70%		
		RMSE	space (MB)	time (sec)	RMSE	space (MB)	time (sec)
convex	PA-APG	0.0149±0.0011	302.4±0.1	2131.7±419.9	0.0098±0.0001	4804.5±598.2	6196.4±2033.4
(nonconvex) capped- $\ell_1$	GDPAN	<b>0.0103±0.0001</b>	171.5±2.2	665.4±99.8	<b>0.0006±0.0001</b>	3243.3±489.6	3670.4±225.8
	sNORT	<b>0.0103±0.0001</b>	<b>14.0±0.8</b>	27.9±5.1	<b>0.0006±0.0001</b>	<b>44.6±0.3</b>	575.9±70.9
	NORT	<b>0.0103±0.0001</b>	14.9±0.9	<b>5.9±1.6</b>	<b>0.0006±0.0001</b>	66.3±0.6	89.4±13.4
(nonconvex) LSP	GDPAN	0.0104±0.0001	172.2±1.5	654.1±214.7	<b>0.0006±0.0001</b>	3009.3±376.2	3794.0±419.5
	sNORT	0.0104±0.0001	14.4±0.1	27.9±5.7	<b>0.0006±0.0001</b>	<b>44.6±0.2</b>	544.2±75.5
	NORT	0.0104±0.0001	15.1±0.1	<b>5.8±2.8</b>	<b>0.0006±0.0001</b>	62.1±0.5	<b>81.3±24.9</b>
(nonconvex) TNN	GDPAN	0.0104±0.0001	172.1±1.6	615.0±140.9	<b>0.0006±0.0001</b>	3009.2±412.2	3922.9±280.1
	sNORT	0.0104±0.0001	14.4±0.1	26.2±4.0	<b>0.0006±0.0001</b>	<b>44.7±0.2</b>	554.7±44.1
	NORT	<b>0.0103±0.0001</b>	15.1±0.1	<b>5.3±1.5</b>	<b>0.0006±0.0001</b>	63.1±0.6	<b>78.0±9.4</b>

- GDPAN is the direct proximal average algorithm
- Nonconvex regularization offers much lower testing RMSEs
- NORT is much **faster**, needs much **less memory** and achieves much **lower testing RMSEs**

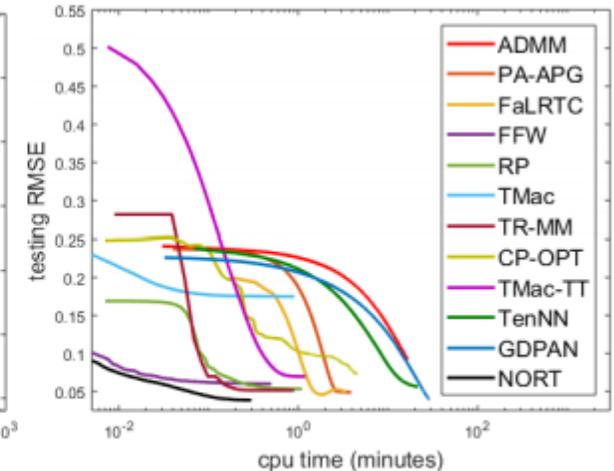
# Experiments: real data sets

Table 4. Testing RMSEs ( $\times 10^{-1}$ ) on color images.

		<i>rice</i>	<i>tree</i>	<i>windows</i>
convex	ADMM	0.680±0.003	0.915±0.005	0.709±0.004
	PA-APG	0.583±0.016	0.488±0.007	0.585±0.002
	FaLRTC	0.576±0.004	0.494±0.011	0.567±0.005
	FFW	0.634±0.003	0.599±0.005	0.772±0.004
	TR-MM	0.596±0.005	0.515±0.011	0.634±0.002
	TenNN	0.647±0.004	0.562±0.004	0.586±0.003
factorization	RP	0.541±0.011	0.524±0.010	0.388±0.026
	TMac	1.923±0.005	1.750±0.006	1.313±0.005
	CP-OPT	0.912±0.086	0.733±0.060	0.964±0.102
	TMac-TT	0.729±0.022	0.697±0.147	1.045±0.107
noncvx	GDPAN	<b>0.467±0.002</b>	<b>0.388±0.012</b>	<b>0.296±0.007</b>
	NORT	<b>0.468±0.001</b>	<b>0.386±0.009</b>	<b>0.297±0.007</b>



(a) *rice*.



(b) *tree*.

- NORT is **fast** and **achieves lower testing RMSEs** compared with other tensor low-rank approaches
- Same observations are on experiments with remote sensing data and multi-relational data (see our paper)

# Thanks.

- Questions: [yaoquanming@4paradigm.com](mailto:yaoquanming@4paradigm.com)
- Codes: available on my Github