

Boosted Sparse and Low-Rank Tensor Regression

Lifang He¹, Kun Chen², Wanwan Xu², Jiayu Zhou³, Fei Wang¹

¹Department of Healthcare Policy and Research, Weill Cornell Medicine

²Department of Statistics, University of Connecticut

³Department of Computer Science and Engineering, Michigan State University

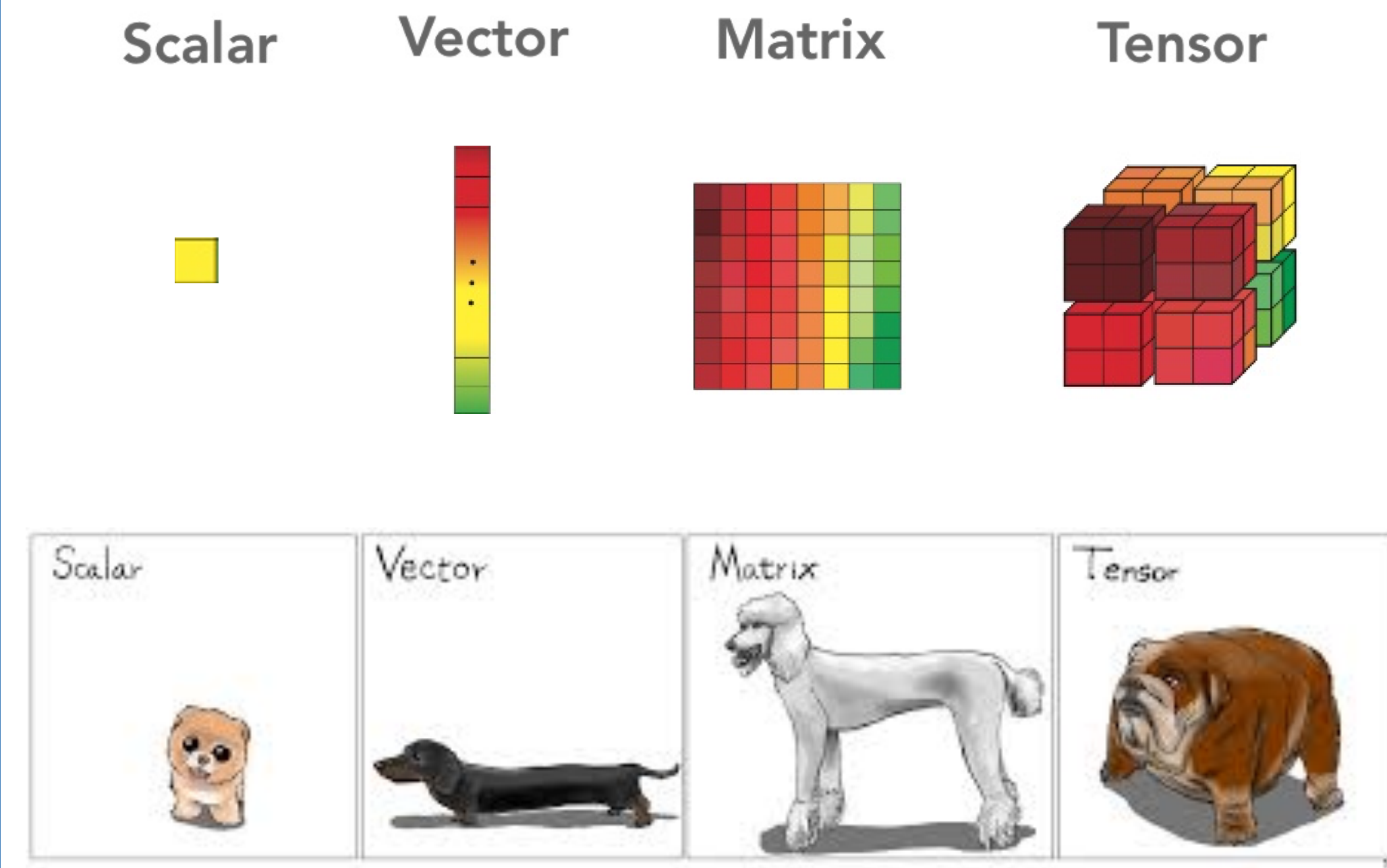


ABSTRACT

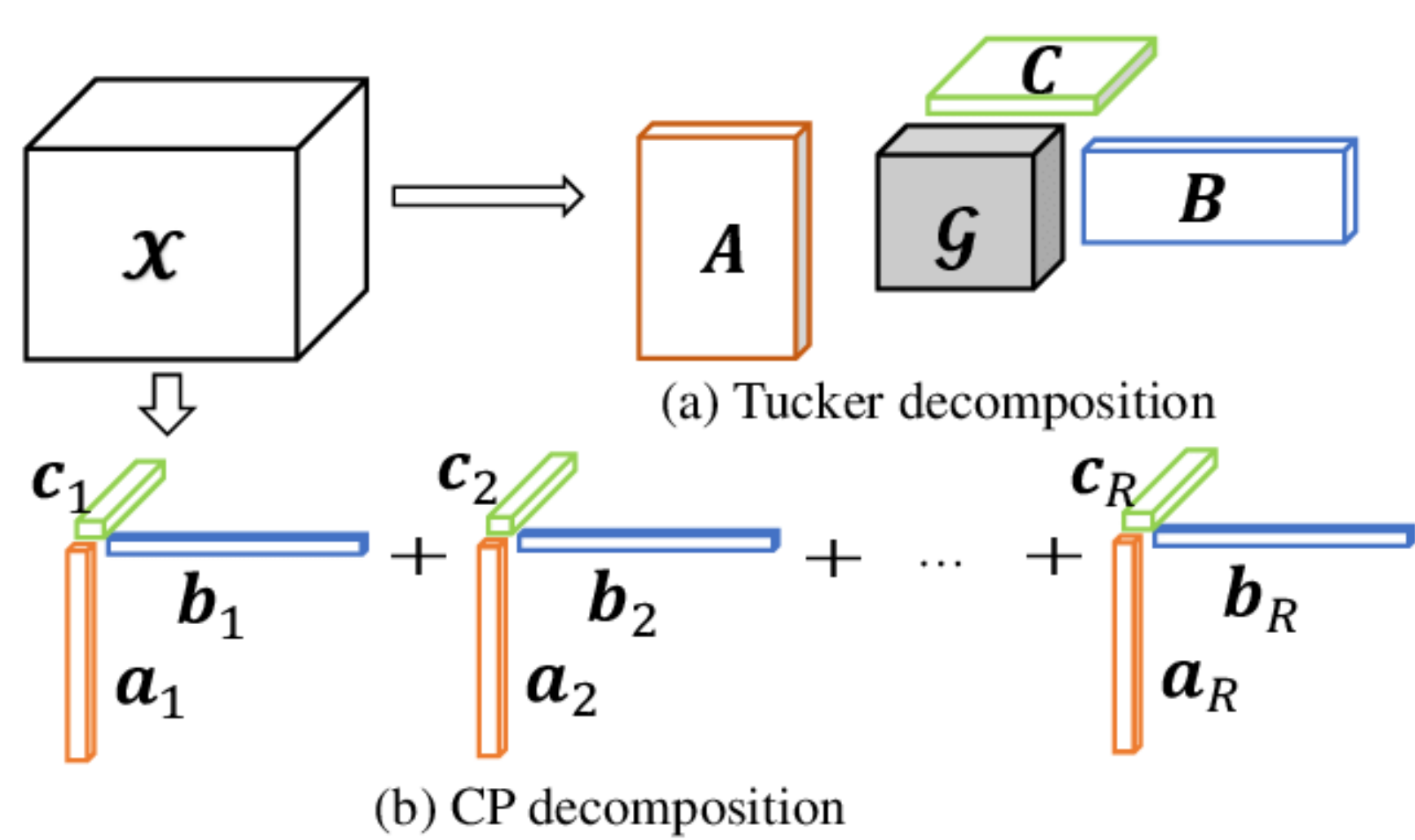
With the advanced capabilities for data acquisition, massive multiway data emerge from neuroscience: multi-channel EEG, volumetric or functional MRI, etc. Tensors and low-rank tensor decomposition are very powerful and versatile tools in machine learning for their ability to express and exploit multi-way data, where they are employed to approach a diverse number of tasks. However, tensor regression, which aims to learn a model with multilinear parameters, is especially suitable for applications with multi-directional relatedness, but has not been fully examined. In this paper, we propose a sparse and low-rank tensor regression model to relate a univariate outcome to a feature tensor, and take a divide-and-conquer strategy to simplify the task into a set of sparse unit-rank tensor factorization/regression problems (SURF). We then present a boosted estimation procedure to efficiently trace out the entire solution path. The superior performance of our approach is demonstrated on various real-world and synthetic examples.

BACKGROUND

◆ Vector to Tensor



◆ Low-Rank Tensor Decomposition



Sparse and Low-rank Tensor Regression

◆ Tensor Regression Problem

$$\min_W \frac{1}{M} \sum_{m=1}^M L(\langle \mathcal{X}^m, \mathcal{W} \rangle, y^m) + \lambda \Omega(\mathcal{W}).$$

$$y^m = \langle \mathcal{X}^m, \mathcal{W} \rangle + \varepsilon^m$$

\mathcal{X}^m : input tensor predictor; y^m : scalar response;
 \mathcal{W} : weight tensor parameter.
 $L(\cdot)$ and $\Omega(\cdot)$: loss function and regularizer.

Formulation

$$\min_{\sigma_r, \mathbf{w}_r^{(n)}} \frac{1}{M} \sum_{m=1}^M (y^m - \langle \mathcal{X}^m, \sum_{r=1}^R \sigma_r \mathbf{w}_r^{(1)} \circ \dots \circ \mathbf{w}_r^{(N)} \rangle)^2 + \sum_{r=1}^R \sum_{n=1}^N \lambda_{r,n} \|\mathbf{w}_r^{(n)}\|_1,$$

s.t. $\|\mathbf{w}_r^{(n)}\|_1 = 1, n = 1, \dots, N, r = 1, \dots, R.$

Challenge:

1. CP rank R needs to be pre-specified;
2. Parameter identifiability issues;
3. Time-consuming for many parameters need to be adjusted.

Solution:

Divide-and-Conquer: Sequential pursue for sparse tensor regression.

◆ Sparse Unit-Rank Tensor Factorization (SURF)

$$\widehat{\mathcal{W}}_r = \min_{\mathcal{W}_r} \frac{1}{M} \sum_{m=1}^M (y^m - \langle \mathcal{X}^m, \mathcal{W}_r \rangle)^2 + \lambda_r \|\mathcal{W}_r\|_1,$$

s.t. $\text{rank}(\mathcal{W}_r) \leq 1.$

Where y_r^m is the current residue of response with

$$y_r^m := \begin{cases} y^m, & \text{if } r = 1 \\ y_{r-1}^m - \langle \mathcal{X}^m, \widehat{\mathcal{W}}_{r-1} \rangle, & \text{otherwise.} \end{cases}$$

Formulation

$$\widehat{\mathcal{W}} = \min_{\mathcal{W}} \frac{1}{M} \sum_{m=1}^M (y^m - \langle \mathcal{X}^m, \mathcal{W} \rangle)^2 + \lambda \|\mathcal{W}\|_1 + \alpha \|\mathcal{W}\|_F^2,$$

s.t. $\text{rank}(\mathcal{W}) \leq 1.$

Let $\mathcal{W} = \sigma \mathbf{w}^{(1)} \circ \dots \circ \mathbf{w}^{(N)}, \mathbf{y} = [y^1, \dots, y^M]$, and $\mathcal{X} = [\mathcal{X}^1, \dots, \mathcal{X}^M]$.

$$\min_{\widehat{\mathbf{w}}^{(n)}} \frac{1}{M} \|\mathbf{y}^T - \mathbf{Z}^{(-n)T} \widehat{\mathbf{w}}^{(n)}\|_2^2 + \alpha \beta^{(-n)} \|\widehat{\mathbf{w}}^{(n)}\|_2^2 + \lambda \|\widehat{\mathbf{w}}^{(n)}\|_1,$$

where $\mathbf{Z}^{(-n)} = \mathcal{X} \times_1 \mathbf{w}^{(1)} \times_2 \dots \times_{n-1} \mathbf{w}^{(n-1)} \times_{n+1} \dots \times_N \mathbf{w}^{(N)}$,
 $\widehat{\mathbf{w}}^{(n)} = \sigma \mathbf{w}^{(n)}$ and $\beta^{(-n)} = \prod_{l \neq n} \|\mathbf{w}^{(l)}\|_2^2.$

$$\min_{\widehat{\mathbf{w}}^{(n)}} \frac{1}{M} \underbrace{\|\widehat{\mathbf{y}} - \widehat{\mathbf{Z}}^{(-n)} \widehat{\mathbf{w}}^{(n)}\|_2^2}_{J(\widehat{\mathbf{w}}^{(n)})} + \lambda \underbrace{\|\widehat{\mathbf{w}}^{(n)}\|_1}_{\Omega(\widehat{\mathbf{w}}^{(n)})}$$

where $\widehat{\mathbf{y}} = (\mathbf{y}, 0)^T$ and $\widehat{\mathbf{Z}}^{(-n)} = (\mathbf{Z}^{(-n)}, \sqrt{\alpha \beta^{(-n)}} \mathbf{I})^T$.

Fast Stagewise/Boosted Optimization

◆ The n-th Objective Function

$$\Gamma(\widehat{\mathbf{w}}^{(n)}; \lambda) = J(\widehat{\mathbf{w}}^{(n)}) + \lambda \Omega(\widehat{\mathbf{w}}^{(n)}).$$

◆ Algorithm

Algorithm 1 Fast Stagewise Unit-Rank Tensor Factorization (SURF)

Input: Training data \mathcal{D} , a small stepsize $\epsilon > 0$ and a small tolerance parameter ξ^2

Output: Solution paths of $(\sigma, \{\mathbf{w}^{(n)}\})$.

1: *Initialization:* take a forward step with $(\{\hat{i}_1, \dots, \hat{i}_N\}, \hat{s}) = \arg \min_{\{i_1, \dots, i_N\}, s=\pm \epsilon} J(s \mathbf{1}_{i_1}, \mathbf{1}_{i_2}, \dots, \mathbf{1}_{i_N})$, and

$$\sigma_0 = \epsilon, \mathbf{w}_0^{(1)} = \text{sign}(\hat{s}) \mathbf{1}_{\hat{i}_1}, \mathbf{w}_0^{(n)} = \mathbf{1}_{\hat{i}_n} (n \neq 1), \lambda_0 = (J(\{\emptyset\}) - J(\sigma_0, \{\mathbf{w}_0^{(n)}\}))/\epsilon. \quad (11)$$

Set the active index sets $I_0^{(n)} = \{\hat{i}_n\}$ for $n = 1, \dots, N; t = 0$.

2: **repeat**

3: *Backward step:*

$$(\hat{n}, \hat{i}_{\hat{n}}) := \arg \min_{n, i_n \in I_t^{(n)}} J(\widehat{\mathbf{w}}_t^{(n)} + s_{i_n} \mathbf{1}_{i_n}), \text{ where } s_{i_n} = -\text{sign}(\widehat{w}_{i_n}^{(n)}). \quad (12)$$

if $\Gamma(\widehat{\mathbf{w}}_t^{(\hat{n})} + s_{i_{\hat{n}}} \mathbf{1}_{i_{\hat{n}}}; \lambda_t) - \Gamma(\widehat{\mathbf{w}}_t^{(\hat{n})}; \lambda_t) \leq -\xi$, then

$$\sigma_{t+1} = \|\widehat{\mathbf{w}}_t^{(\hat{n})} + s_{i_{\hat{n}}} \mathbf{1}_{i_{\hat{n}}}\|_1, \mathbf{w}_{t+1}^{(\hat{n})} = (\widehat{\mathbf{w}}_t^{(\hat{n})} + s_{i_{\hat{n}}} \mathbf{1}_{i_{\hat{n}}})/\sigma_{t+1}, \mathbf{w}_{t+1}^{(-\hat{n})} = \mathbf{w}_t^{(-\hat{n})},$$

$$\lambda_{t+1} = \lambda_t, I_{t+1}^{(n)} := \begin{cases} I_t^{(n)} \setminus \{\hat{i}_{\hat{n}}\}, & \text{if } w_{(t+1)\hat{i}_{\hat{n}}}^{(\hat{n})} = 0 \\ I_t^{(n)}, & \text{otherwise.} \end{cases}$$

4: *else Forward step:*

$$(\hat{n}, \hat{i}_{\hat{n}}) := \arg \min_{n, i_n, s=\pm \epsilon} J(\widehat{\mathbf{w}}_t^{(n)} + s \mathbf{1}_{i_n}), \quad (13)$$

$$\sigma_{t+1} = \|\widehat{\mathbf{w}}_t^{(\hat{n})} + s_{i_{\hat{n}}} \mathbf{1}_{i_{\hat{n}}}\|_1, \mathbf{w}_{t+1}^{(\hat{n})} = (\widehat{\mathbf{w}}_t^{(\hat{n})} + s_{i_{\hat{n}}} \mathbf{1}_{i_{\hat{n}}})/\sigma_{t+1}, \mathbf{w}_{t+1}^{(-\hat{n})} = \mathbf{w}_t^{(-\hat{n})},$$

$$\lambda_{t+1} = \min[\lambda_t, \frac{J(\sigma_t, \{\mathbf{w}_t^{(n)}\}) - J(\sigma_{t+1}, \{\mathbf{w}_{t+1}^{(n)}\}) - \xi}{\Omega(\sigma_{t+1}, \{\mathbf{w}_{t+1}^{(n)}\}) - \Omega(\sigma_t, \{\mathbf{w}_t^{(n)}\})}], I_{t+1}^{(n)} := \begin{cases} I_t^{(n)} \cup \{\hat{i}_{\hat{n}}\}, & \text{if } n = \hat{n} \\ I_t^{(n)}, & \text{otherwise.} \end{cases}$$

5: Set $t = t + 1$.

6: **until** $\lambda_t \leq 0$

Forward Step

The solution at each forward step is

$$(\hat{n}, \hat{i}_{\hat{n}}) := \arg \max_{n, i_n} 2|\widehat{\mathbf{e}}^{(n)T} \widehat{\mathbf{Z}}^{(-n)} \mathbf{1}_{i_n}| - \epsilon \text{Diag}(\widehat{\mathbf{Z}}^{(-n)T} \widehat{\mathbf{Z}}^{(-n)})^T \mathbf{1}_{i_n},$$

$$\hat{s} = \text{sign}(\widehat{\mathbf{e}}^{(n)T} \widehat{\mathbf{Z}}^{(-n)} \mathbf{1}_{i_{\hat{n}}}) \epsilon,$$

where $\widehat{\mathbf{e}}^{(n)} = \widehat{\mathbf{y}} - \widehat{\mathbf{Z}}^{(-n)} \widehat{\mathbf{w}}^{(n)}$ is a constant at each iteration.

Backward Step

The solution at each backward step is

$$(\hat{n}, \hat{i}_{\hat{n}}) := \arg \min_{n, i_n} 2\text{sign}(\widehat{w}_{i_n}^{(n)}) \widehat{\mathbf{e}}^{T} \widehat{\mathbf{Z}}^{(-n)} \mathbf{1}_{i_n} + \epsilon \text{Diag}(\widehat{\mathbf{Z}}^{(-n)T} \widehat{\mathbf{Z}}^{(-n)})^T \mathbf{1}_{i_n}.$$

◆ Computation Analysis

At each iteration, $\mathbf{Z}^{(-n)}$ ($n \neq \hat{n}$) can be updated by

$$\mathbf{Z}_{t+1}^{(-n)} = \frac{1}{\sigma_{t+1}} (\sigma_t \mathbf{Z}_t^{(-n)} + \mathbf{Z}_t^{(-n, -\hat{n})} \times_{\hat{n}} \widehat{s}_{i_{\hat{n}}} \mathbf{1}_{i_{\hat{n}}}),$$

where $(-n, -\hat{n})$ denotes every mode except n and \hat{n} .

The computational complexity of our approach per iteration is

$$O(M \sum_{n \neq \hat{n}} (\prod_{s \neq n, \hat{n}} I_s + 5I_n) + 2MI_{\hat{n}})$$

In contrast, the ACS algorithm has to be run for each fixed λ , and within each of such problems, each iteration requires

$$O(M \prod_{n=1}^N I_n)$$

◆ Convergence Analysis

Lemma 1: For any t with $\lambda_{t+1} = \lambda_t$, we have $\Gamma(\sigma_{t+1}, \{\mathbf{w}_{t+1}^{(n)}\}; \lambda_{t+1}) \leq \Gamma(\sigma_t, \{\mathbf{w}_t^{(n)}\}; \lambda_{t+1}) - \xi$.

Lemma 2: For any t with $\lambda_{t+1} < \lambda_t$, we have $\Gamma(\widehat{\mathbf{w}}_t^{(n)} + s_{i_n} \mathbf{1}_{i_n}; \lambda_t) > \Gamma(\widehat{\mathbf{w}}_t^{(n)}; \lambda_t) - \xi$.

Lemma 1 and Lemma 2 proves the following convergence theorem.

Theorem 1: For any t such that $\lambda_{t+1} < \lambda_t$, we have $(\sigma_t, \{\mathbf{w}_t^{(n)}\}) \rightarrow (\sigma(\lambda_t), \{\mathbf{w}^{(n)}(\lambda_t)\})$ as $\epsilon, \xi \rightarrow 0$, where $(\sigma(\lambda_t), \{\mathbf{w}^{(n)}(\lambda_t)\})$ denotes a coordinate-wise minimum point of the SURF problem.

Experiments

◆ Summarization of Compared Methods

Table 1: Compared methods. α and λ are regularized parameters; R is the CP rank.							
Methods	LASSO	ENet	Remurs	orTRR	GLTRM	ACS	SURF
Input Data Type	Tensor	Tensor	Tensor	Tensor	Tensor	Tensor	Tensor
Regularization	$\ell_1(\mathbf{w})$	$\ell_1/\ell_2(\mathbf{w})$	Nuclear/ $\ell_1(W)$	$\ell_2(\mathbf{W}^{(n)})$	$\ell_1/\ell_2(\mathbf{W}^{(n)})$	$\ell_1/\ell_2(W_r)$	$\ell_1/\ell_2(W_r)$
Rank Explored	Fixed	Fixed	Fixed	Fixed	Fixed	Fixed	Fixed
Hyperparameters	λ	α, λ	λ_1, λ_2	α, R	α, λ, R	α, λ, R	α, λ, R

◆ Empirical Analysis on Synthetic Data

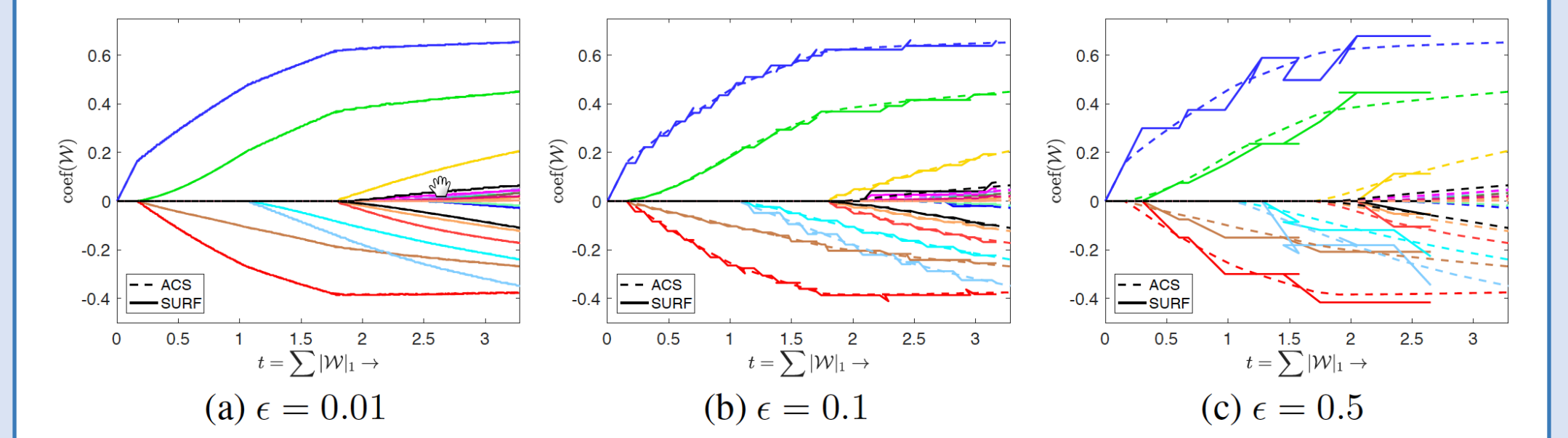


Figure 1: Comparison of solution paths of SURF (solid line) and ACS (dashed line) with different step sizes on synthetic data. The path of estimates \mathcal{W} for each λ is treated as a function of $t = \|\mathcal{W}\|_1$.

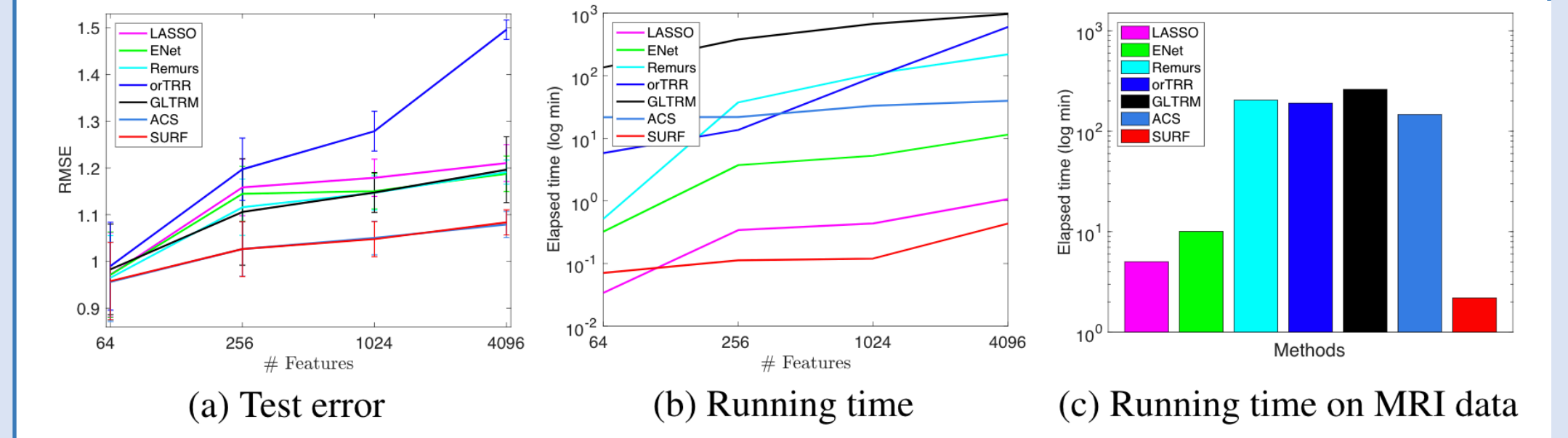


Figure 2: Results with increasing number of features on synthetic 2D data (a)-(b), and (c) real 3D MRI data of features $240 \times 175 \times 176$ with fixed hyperparameters (without cross validation).

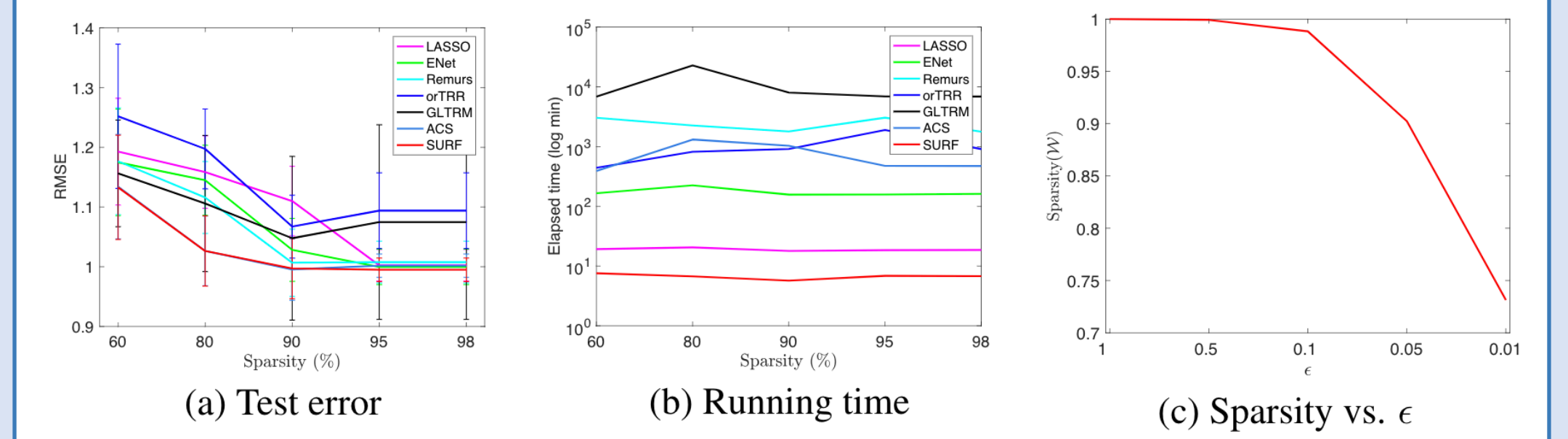


Figure 3: Results with increasing sparsity level ($S\%$) of true \mathcal{W} on synthetic 2D data (a)-(b), and (c) sparsity results of \mathcal{W} versus step size for SURF.

◆ Statistical Analysis on Real Data

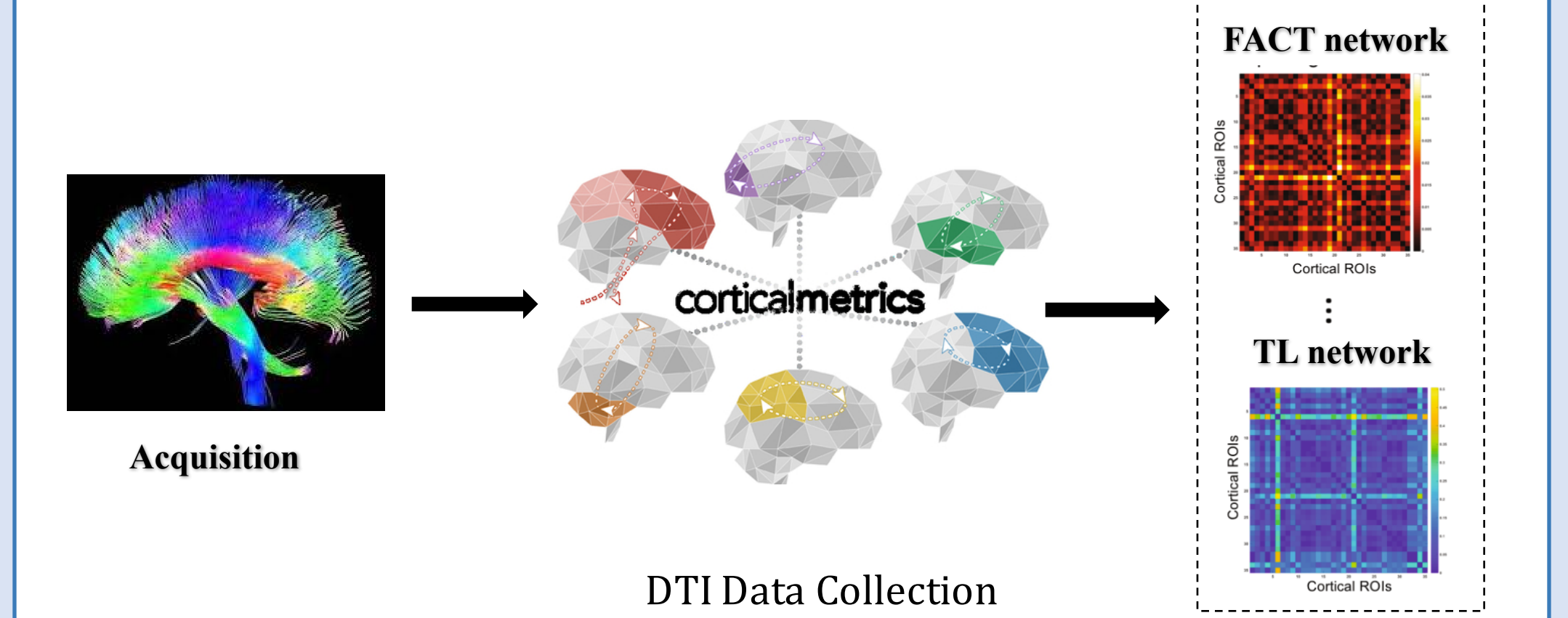


Table 2: Results on different DTI datasets (mean \pm std.). Column 2 indicates the used metrics RMSE, Sparsity of Coefficients (SC) and CPU execution time (in mins).

Datasets	Metrics	Comparative Methods						
		LASSO	ENet	Remurs	orTRR	GLTRM	ACS	SURF
DTI _{fact}	RMSE	2.94 \pm 0.34	2.92 \pm 0.32	2.91 \pm 0.32	3.48 \pm 0.21	3.09 \pm 0.35	2.81 \pm 0.24	2.81 \pm 0.23
	Sparsity	0.99 \pm 0.01	0.97 \pm 0.01	0.66 \pm 0.13	0.00 \pm 0.00	0.90 \pm 0.10	0.92 \pm 0.02	0.95 \pm 0.01
	Time	6.4 \pm 0.3	46.6 \pm 4.6	161.3 \pm 9.3	27.9 \pm 5.6	874.8 \pm 29.6	60.8 \pm 24.4	1.7 \pm 0.2
DTI _{k=2}	RMSE	3.18 \pm 0.36	3.16 \pm 0.42	2.97 \pm 0.30	3.76 \pm 0.44	3.26 \pm 0.46	2.90 \pm 0.31	2.91 \pm 0.32
	Sparsity	0.99 \pm 0.01	0.95 \pm 0.03	0.37 \pm 0.09	0.00 \pm 0.00	0.91 \pm 0.06	0.93 \pm 0.02	0.94 \pm 0.01
	Time	5.7 \pm 0.3	42.4 \pm 2.9	155.0 \pm 10.7	10.2 \pm 0.1	857.4 \pm 22.5	63.0 \pm 21.6	5.2 \pm 0.8
DTI _{el}	RMSE	3.06 \pm 0.34	2.99 \pm 0.34	2.93 \pm 0.27	3.56 \pm 0.41	3.14 \pm 0.39	2.89 \pm 0.38	2.87 \pm 0.35
	Sparsity	0.98 \pm 0.01	0.95 \pm 0.01	0.43 \pm 0.17	0.00 \pm 0.00	0.87 \pm 0.03	0.90 \pm 0.03	0.93 \pm 0.02
	Time	5.8 \pm 0.3	45.0 \pm 1.0	163.6 \pm 9.0	7.5 \pm 0.9	815.4 \pm 6.5	66.3 \pm 44.9	1.5 \pm 0.1
DTI _{el}	RMSE	3.20 \pm 0.40	3.21 \pm 0.59	2.84 \pm 0.35	3.66 \pm 0.35	3.12 \pm 0.32	2.82 \pm 0.33	2.83 \pm 0.32
	Sparsity	0.99 \pm 0.01	0.96 \pm 0.03	0.44 \pm 0.13	0.00 \pm 0.00	0.86 \pm 0.03	0.90 \pm 0.02	0.91 \pm 0.02
	Time	5.5 \pm 0.2	42.3 \pm 1.4	159.6 \pm 7.6	26.6 \pm 3.1	835.8 \pm 9.9	96.7 \pm 43.2	3.8 \pm 0.5
Combined	RMSE	3.02 \pm 0.37	2.89 \pm 0.41	2.81 \pm 0.31	3.33 \pm 0.27	3.26 \pm 0.45	2.79 \pm 0.31	2.78 \pm 0.29
	Time	0.99 \pm 0.00	0.97 \pm 0.01	0.34 \pm 0.22	0.00 \pm 0.00	0.97 \pm 0.01	0.99 \pm 0.00	0.99 \pm 0.00

NOTE: The full paper and code are available at <https://github.com/LifangHe/SURF>. Email: lifanghesur@gmail.com