

# A BLOCK-RANDOMIZED STOCHASTIC METHOD WITH IMPORTANCE SAMPLING FOR CP TENSOR DECOMPOSITION<sup>1</sup>

YAJIE YU

CHONGQING UNIVERSITY, CHONGQING, P.R. CHINA

*CQSIAM2023, 13<sup>RD</sup> MAY 2023*

---

<sup>1</sup>A joint work with Hanyu Li  
E-mail: zqyu@cqu.edu.cn

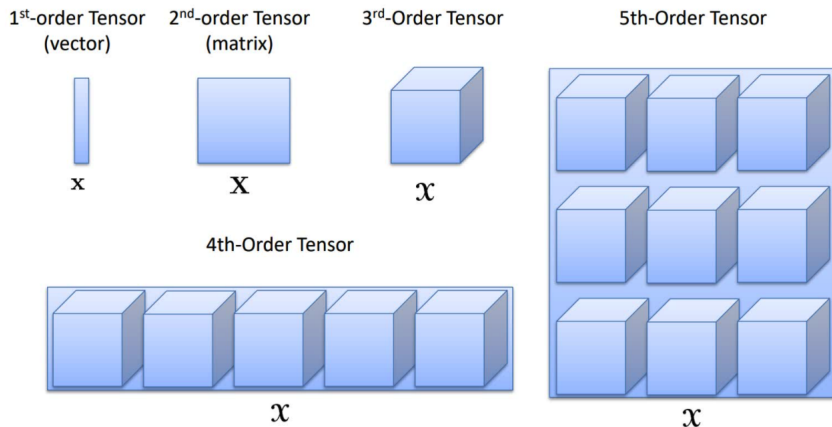


# PRESENTATION OUTLINE

- 1 Introduction
  - CANDECOMP/PARAFAC (CP) Decomposition
  - Algorithms for CP Decomposition
- 2 Proposed Method
- 3 Numerical Results
- 4 Conclusions



## A TENSOR IS AN MULTI-WAY ARRAY

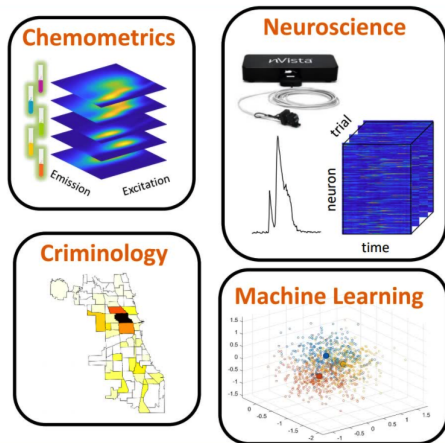


**Figure:** Graphical representation of multiway array (tensor) data.



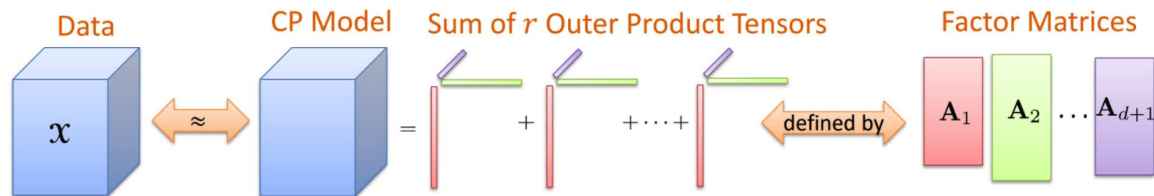
# TENSORS COME FROM MANY APPLICATIONS

TENSOR DECOMPOSITION FINDS PATTERNS IN MASSIVE DATA (UNSUPERVISED LEARNING)



- **Chemometrics:** Emission x Excitation x Samples (Fluorescence Spectroscopy)
- **Neuroscience:** Neuron x Time x Trial
- **Criminology:** Day x Hour x Location x Crime (Chicago Crime Reports)
- **Machine Learning:** Multivariate Gaussian Mixture Models Higher-Order Moments
- **Transportation:** Pickup x Dropoff x Time (Taxis)
- **Sports:** Player x Statistic x Season (Basketball)
- **Cyber-Traffic:** IP x IP x Port x Time
- **Social Network:** Person x Person x Time x Interaction-Type
- **Signal Processing:** Sensor x Frequency x Time
- **Trending Co-occurrence:** Term A x Term B x Time

# CP DECOMPOSITION APPROXIMATES AN N-TH ORDER TENSOR AS A SUM OF R RANK-ONE TENSORS



**Figure:** Graphical representation of CP decomposition.

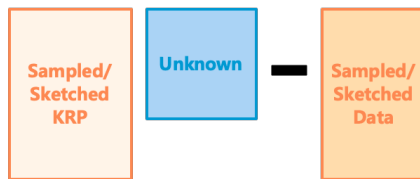
- $\mathcal{X} \approx \tilde{\mathcal{X}} = \sum_{r=1}^R \mathbf{a}_r^{(1)} \circ \mathbf{a}_r^{(2)} \circ \dots \circ \mathbf{a}_r^{(N)} = [[\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)}]]$ ;
- $\mathcal{O}(NIR)$  parameters: is linear to the tensor order  $N$ .
- **high compression ability**;
- **Uniqueness**: Under modest conditions, CP is unique up to permutation and scaling.





# RANDOMIZED CP-ALS BY MAKE FULL USE OF THE KRP STRUCTURE

EFFICIENT SAMPLING/SKETCHING WITHOUT FORMING KRP



**Figure:** CPRAND or CP-ARLS-LEV:  $\min_{\mathbf{A} \in \mathbb{R}^{I_n \times R}} \|\mathbf{SZ}^{(n)} \mathbf{A}^\top - \mathbf{SX}_{(n)}^\top\|_F^2$ .

- CPRAND<sup>3</sup>: Solve  $N$  smaller subproblems per outer iteration with **uniform sampling** and **KFJLT**;
- CP-ARLS-LEV<sup>4</sup>: Solve  $N$  smaller subproblems per outer iteration with **leverage-based sampling**.

<sup>3</sup>Casey Battaglino, Grey Ballard, and Tamara G. Kolda. "A Practical Randomized CP Tensor Decomposition". In: *SIAM J. Matrix Anal. Appl.* 39.2 (2018), pp. 876–901. doi: 10.1137/17M1112303.

<sup>4</sup>Brett W. Larsen and Tamara G. Kolda. "Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition". In: *SIAM J. Matrix Anal. Appl.* 43.3 (2022), pp. 1488–1517. doi: 10.1137/21M1441754.



## ANOTHER POPULAR APPROACH FOR COMPUTING CP DECOMPOSITION: DIRECT/ALL-AT-ONCE OPTIMIZATION

$$\min_{\{\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}\}_{n=1}^N} f(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}), \text{ where } f(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \frac{1}{2} \|\mathbf{X} - \llbracket \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)} \rrbracket\|_F^2$$

- BrasCPD<sup>5</sup>: A block-randomized stochastic method which is a judicious combination of **randomized block coordinate descent (BCD)** and **stochastic proximal gradient (SPG)**; **uniform sampling**;
- mBrasCPD<sup>6</sup>: **momentum** accelerated version;
- iBrasCPD<sup>7</sup>: **inertial** accelerated version;
- DS-MVR<sup>8</sup>: **momentum** acceleration and the **variance reduction** technique.
- Others: [AKD11; ADK11; XY13; PTC13; VD16; HF20; Van21]

<sup>5</sup>Xiao Fu et al. "Block-Randomized Stochastic Proximal Gradient for Low-Rank Tensor Factorization". In: *IEEE Trans. Signal Process.* 68 (2020), pp. 2170–2185. doi: 10.1109/TSP.2020.2982321..

<sup>6</sup>Qingsong Wang, Chunfeng Cui, and Deren Han. "A Momentum Block-Randomized Stochastic Algorithm for Low-Rank Tensor CP Decomposition". In: *Pac. J. Optim.* 17.3 (2021), pp. 433–452.

<sup>7</sup>Qingsong Wang et al. "Inertial accelerated SGD algorithms for solving large-scale lower-rank tensor CP decomposition problems". In: *J. Comput. Appl. Math.* 423 (2023), p. 114948. doi: 10.1016/j.cam.2022.114948.

<sup>8</sup>Qingsong Wang, Chunfeng Cui, and Deren Han. "Accelerated Doubly Stochastic Gradient Descent for Tensor CP Decomposition". In: *J. Optim. Theory Appl.* (2023), pp. 1–40. doi: 10.1007/s10957-023-02193-5.





# PRESENTATION OUTLINE

## 1 Introduction

## 2 Proposed Method

- Motivation
- Sampling Strategies
- Proposed Algorithms
- Theoretical analysis

## 3 Numerical Results

## 4 Conclusions

## MOTIVATIONS

The full gradient:

$$\mathbf{G}^{(n)}_{(t)} = \mathbf{A}^{(n)}_{(t)} (\mathbf{Z}^{(n)})^\top \mathbf{Z}^{(n)} - \mathbf{X}_{(n)} \mathbf{Z}^{(n)}$$

- Combine BrasCPD<sup>9</sup> with **empirical importance sampling** using **efficient sampling** technique mentioned above;
  - 1 Leverage-based sampling<sup>10</sup>;
  - 2 Euclidean-based sampling<sup>11</sup>;
- The **theoretical optimal sampling probability distribution** is given based on the idea of variance reduction.

<sup>9</sup>Xiao Fu et al. “Block-Randomized Stochastic Proximal Gradient for Low-Rank Tensor Factorization”. In: *IEEE Trans. Signal Process.* 68 (2020), pp. 2170–2185. doi: 10.1109/TSP.2020.2982321..

<sup>10</sup>Brett W. Larsen and Tamara G. Kolda. “Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition”. In: *SIAM J. Matrix Anal. Appl.* 43.3 (2022), pp. 1488–1517. doi: 10.1137/21M1441754.

<sup>11</sup>Deanna Needell, Nathan Srebro, and Rachel Ward. “Stochastic Gradient Descent, Weighted Sampling, and the Randomized Kaczmarz Algorithm”. In: *Math. Program.* 155.1 (2016), pp. 549–573. doi: 10.1007/s10107-015-0864-7.

# EXPLOIT KRP STRUCTURE TO BOUND LEVERAGE SCORES: USE LEVERAGE SCORES OF FACTOR MATRICES FOR SAMPLING PROBABILITIES

LEVERAGE SCORES KEY TO LIMITING SAMPLES(BUT TOO EXPENSIVE TO COMPUTE)

## Lemma 2.1 (Leverage Score Bounds for KRP Matrix<sup>12</sup>)

For matrices  $\mathbf{A}^{(k)} \in \mathbb{R}^{I_k \times R}$  with  $I_k > R$  for  $k = 1, \dots, K$ , let  $\ell_{i_k}$  be the leverage score of the  $i_k$ -th row of  $\mathbf{A}^{(k)}$ . Then, for the KRP matrix  $\mathbf{Z} = \mathbf{A}^{(1)} \odot \mathbf{A}^{(2)} \odot \dots \odot \mathbf{A}^{(K)}$ , the leverage score  $\ell_{i_1, \dots, i_K}$  of its  $j$ -th row corresponding to  $\{i_1, \dots, i_K\}$  satisfies

$$\ell_{i_1, \dots, i_K} \leq \prod_{k=1}^K \ell_{i_k}.$$

- The leveraged-based sampling probability for the  $j$ -th row of the above KRP matrix  $\mathbf{Z}$  can be set to be

$$p_j = \frac{1}{R^K} \prod_{k=1}^K \ell_{i_k}. \quad (2.1)$$

<sup>12</sup>Dehua Cheng et al. "SPALS: Fast Alternating Least Squares via Implicit Leverage Scores Sampling". In: *Proceedings of the 30th International Conference on Neural Information Processing Systems*. Barcelona Spain: Curran Associates Inc., 2016, pp. 721–729.



## EFFICIENT SAMPLING WITHOUT FORMING KRP

### Lemma 2.2 (Sampling KRP Efficiently<sup>13</sup>)

Let  $\mathbf{A}^{(k)} \in \mathbb{R}^{I_k \times R}$  for  $k = 1, \dots, K$ , and  $\ell(\mathbf{A}^{(k)})$  be the vector of leverage scores for  $\mathbf{A}^{(k)}$ . Let

$$i_k \sim \text{MULTINOMIAL}(\ell(\mathbf{A}^{(k)})/R) \text{ for } k = 1, \dots, K.$$

Then, the probability of selecting the multi-index  $\{i_1, \dots, i_K\}$  is  $p_j$  in (2.1).

<sup>13</sup>Brett W. Larsen and Tamara G. Kolda. "Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition". In: *SIAM J. Matrix Anal. Appl.* 43.3 (2022), pp. 1488–1517.  
doi: 10.1137/21M1441754.

## USE SQUARED EUCLIDEAN NORMS OF FACTOR MATRICES FOR SAMPLING PROBABILITIES

### Lemma 2.3

Let  $\mathbf{A}^{(k)} \in \mathbb{R}^{I_k \times R}$  with  $I_k > R$  for  $k = 1, \dots, K$ . For the KRP matrix  $\mathbf{Z} = \mathbf{A}^{(1)} \odot \mathbf{A}^{(2)} \odot \dots \odot \mathbf{A}^{(K)}$ , the squared Euclidean norms of its  $j$ -th row corresponding to  $\{i_1, \dots, i_K\}$  satisfies

$$\|\mathbf{Z}(j, :)\|_2^2 \leq \prod_{k=1}^K \|\mathbf{A}^{(k)}(i_k, :)\|_2^2.$$

- The Euclidean-based sampling probability for the  $j$ -th row of the above KRP matrix  $\mathbf{Z}$  can be set to be

$$p_j = \prod_{k=1}^K \frac{\|\mathbf{A}^{(k)}(i_k, :)\|_2^2}{\|\mathbf{A}^{(k)}\|_F^2}. \quad (2.2)$$

# EUCLIDEAN-BASED SAMPLING: EFFICIENT SAMPLING WITHOUT FORMING KRP

## Lemma 2.4

Let  $\mathbf{A}^{(k)} \in \mathbb{R}^{I_k \times R}$  for  $k = 1, \dots, K$ , and  $\mathbf{p}_k$  be the Euclidean-based probability distribution for  $\mathbf{A}^{(k)}$  with  $\beta = 1$ , i.e.,  $\mathbf{p}_k(i) = \|\mathbf{A}^{(k)}(i, :)\|_2^2 / \|\mathbf{A}^{(k)}\|_F^2$  with  $i \in [I_k]$ . Let

$$i_k \sim \text{MULTINOMIAL}(\mathbf{p}_k) \text{ for } k = 1, \dots, K.$$

Then, the probability of selecting the multi-index  $\{i_1, \dots, i_K\}$  is  $p_j$  in (2.2).

# EFFICIENT IMPORTANCE SAMPLING WITHOUT FORMING KRP

---

## Algorithm 1 SKRP-ST with importance sampling (SKRP-ST-I)<sup>14</sup>

---

```

1: function  $[\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}, \mathbf{p}_{\mathcal{F}_n}] = \text{SKRP-ST-I}(n, |\mathcal{F}_n|, \{\mathbf{p}_k\}_{k=1, k \neq n}^N, \{\mathbf{A}^{(k)}\}_{k=1, k \neq n}^N)$ 
2:   for  $k = 1, \dots, n-1, n+1, \dots, N$  do
3:      $\text{id}\mathbf{x}(:, k) = \text{RANDSAMPLE}(I_k, |\mathcal{F}_n|, \text{true}, \mathbf{p}_k)$ 
4:   end for
5:    $[\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}] = \text{SKRP-ST}^{15}(n, \text{id}\mathbf{x}, \{\mathbf{A}^{(k)}\}_{k=1, k \neq n}^N)$ 
6:    $\mathbf{p}_{\mathcal{F}_n} \leftarrow (8.2)$  with leverage/Euclidean-based probability distribution
7:   return  $\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}, \mathbf{p}_{\mathcal{F}_n}$ 
8: end function

```

---

<sup>14</sup>Brett W. Larsen and Tamara G. Kolda. “Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition”. In: *SIAM J. Matrix Anal. Appl.* 43.3 (2022), pp. 1488–1517. doi: 10.1137/21M1441754.

<sup>15</sup>SKRP-ST: Efficient uniform sampling without forming KRP in [BBK18]; see the Appendix for detailed algorithms.

# IMPORTANCE SAMPLED ROWS BY PROBABILITY OF SELECTION TO ELIMINATE BIAS

## ■ Stochastic Gradient:

$$\mathbf{G}_{(t)}^{(n)} = \frac{1}{|\mathcal{F}_n|J_n} \left( \mathbf{A}_{(t)}^{(n)} (\mathbf{D}(\mathbf{Z}_{\mathcal{F}_n}^{(n)})^\top) \mathbf{Z}_{\mathcal{F}_n}^{(n)} - \mathbf{D} \mathbf{X}_{(n)}^{\mathcal{F}_n} \mathbf{Z}_{\mathcal{F}_n}^{(n)} \right), \quad (2.3)$$

where  $\mathbf{D} = \text{diag}[\frac{1}{p_{j_1}}, \dots, \frac{1}{p_{j_{|\mathcal{F}_n|}}}]$  is from sampled leverage- or Euclidean- based probability distribution.

- Define  $\xi_{(t)} \in \{1, \dots, N\}$  and  $\zeta_{(t)} \subseteq \{1, \dots, J_{\xi_{(t)}}\}$  as the random variables (r.v.s) responsible for selecting the mode and fibers in the  $t$ -th iteration, respectively.

## Theorem 2.5 (Unbiased Gradient)

Denote  $\mathcal{B}_{(t)}$  as the filtration generated by the r.v.s

$$\{\xi_{(0)}, \zeta_{(0)}, \xi_{(1)}, \zeta_{(1)}, \dots, \xi_{(t-1)}, \zeta_{(t-1)}\}$$

such that the  $t$ -th iteration  $\boldsymbol{\theta}_{(t)}$  is determined conditioned on  $\mathcal{B}_{(t)}$ . Then the stochastic gradient in (2.3) is the unbiased estimate of the full gradient with respect to (w.r.t.)  $\mathbf{A}^{(\xi_{(t)})}$ , i.e.,

$$\mathbb{E}_{\zeta_{(t)}} \left[ \mathbf{G}_{(t)}^{(\xi_{(t)})} \mid \mathcal{B}_{(t)}, \xi_{(t)} \right] = \nabla_{\mathbf{A}^{(\xi_{(t)})}} f(\boldsymbol{\theta}_{(t)}).$$





# A DOUBLY RANDOMIZED COMPUTATIONAL FRAMEWORK FOR LARGE-SCALE CP DECOMPOSITION

- 1 Sampling a mode  $n$  from all modes  $\{1, \dots, N\}$  of the tensor;
  - 2 Forming sampled  $\mathbf{Z}^{(n)}$  of this mode by sampling latent factors;
  - 3 Sampling some fibers of this mode;
  - 4 Updating the corresponding latent factor via stochastic gradient operations.
- **Block-Randomized Weighted SGD for CPD ( BrawsCPD)**
    - Update  $\mathbf{A}_{(t+1)}^{(n)} \leftarrow \mathbf{A}_{(t)}^{(n)} - \alpha_t \mathbf{G}_{(t)}^{(n)}$ ,  $\mathbf{A}_{(t+1)}^{(n')} \leftarrow \mathbf{A}_{(t)}^{(n')}$  for  $n' \neq n$ .
  - **Adagrad version of block-randomized Weighted SGD for CPD ( AdawCPD)**
    - Update step size:  $[\boldsymbol{\eta}_{(t)}^{(n)}]_{i,r} \leftarrow \frac{\eta}{(b + \sum_{t'=1}^t [\mathbf{G}_{(t')}^{(n)}]_{i,r}^2)^{1/2}}$ ,  $i \in [I_n]$ ,  $r \in [R]$ .
    - Update  $\mathbf{A}^{(n)}$ :  $\mathbf{A}_{(t+1)}^{(n)} \leftarrow \mathbf{A}_{(t)}^{(n)} - \boldsymbol{\eta}_{(t)}^{(n)} \circledast \mathbf{G}_{(t)}^{(n)}$ ,  $\mathbf{A}_{(t+1)}^{(n')} \leftarrow \mathbf{A}_{(t)}^{(n')}$  for  $n' \neq n$



## THERE EXISTS A SUBSEQUENCE OF THE SOLUTION SEQUENCE THAT CONVERGES TO A STATIONARY POINT IN EXPECTATION

### Theorem 2.6 (Convergence Properties)

Suppose that the step size schedule follows the Robbins-Monro rule:  $\sum_{t=0}^{\infty} \alpha^t = \infty$  and  $\sum_{t=0}^{\infty} (\alpha^t)^2 < \infty$ , and the updates  $\mathbf{A}_{(t)}^{(n)}$  are bounded for all  $n$  and  $t$ . The solution sequence produced by BrowsCPD satisfies

$$\liminf_{t \rightarrow \infty} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_{(t)})\|_F^2] = 0,$$

and the solution sequence produced by AdawCPD satisfies

$$\Pr \left( \liminf_{t \rightarrow \infty} \|\nabla f(\boldsymbol{\theta}_{(t)})\|_F^2 = 0 \right) = 1.$$

# OPTIMAL SAMPLING PROBABILITY DISTRIBUTION IN THE SENSE OF MINIMIZING VARIANCE

## Theorem 2.7 (Optimal Sampling Probability Distribution)

In the setting of Theorem 2.5, suppose that  $\mathbf{p} \in \mathbb{R}^{J_n}$  is any probability distribution and  $\mathbf{R}_{(t)}^{(n)} = \mathbf{A}_{(t)}^{(n)}(\mathbf{Z}_{(t)}^{(n)})^\top - \mathbf{X}_{(n)}$ . Then if  $\mathbf{p}$  is as

$$p_i = \frac{\|\mathbf{R}_{(t)}^{(\xi(t))}(:, i)\|_2 \|\mathbf{Z}_{(t)}^{(\xi(t))}(i, :)\|_2}{\sum_{i'=1}^{J_n} \|\mathbf{R}_{(t)}^{(\xi(t))}(:, i')\|_2 \|\mathbf{Z}_{(t)}^{(\xi(t))}(i', :)\|_2}, \quad i = 1, \dots, J_n,$$

$\mathbb{E}_{\zeta(t)} \left[ \|\mathbf{G}_{(t)}^{(\xi(t))} - \nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}_{(t)})\|_F^2 \mid \mathcal{B}_{(t)}, \xi(t) \right]$  achieves its minimum as

$$\frac{1}{|\mathcal{F}_n|} \left( \sum_{j_f=1}^{J_n} \|\mathbf{R}_{(t)}^{(\xi(t))}(:, j_f)\|_2 \|\mathbf{Z}_{(t)}^{(\xi(t))}(j_f, :)\|_2 \right)^2 - \frac{1}{|\mathcal{F}_n|} \|\nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}_{(t)})\|_F^2.$$



# PRESENTATION OUTLINE

- 1 Introduction
- 2 Proposed Method
- 3 Numerical Results**
- 4 Conclusions

## SETUP

### ■ Data generation<sup>16</sup>

- Generating by factor matrices with parameters  $R_{true}$ ,  $noise$ ,  $spread$  and  $magnitude$ .
  - 1 Generating factor matrices with  $spread$  (how many non-zeros elements are added to each of these first three columns) and  $magnitude$  (those non-zero elements are chosen).
  - 2  $\mathcal{X}_{true} = \llbracket \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)} \rrbracket$ .
  - 3 Adding noise:  $\mathcal{X} = \mathcal{X}_{true} + noise \left( \frac{\|\mathcal{X}_{true}\|}{\|\mathcal{N}\|} \right) \mathcal{N}$ .

### ■ Experimental setup

- Checking convergence:  $Tol = \frac{\|\llbracket \hat{\mathbf{A}}^{(1)}, \hat{\mathbf{A}}^{(2)}, \hat{\mathbf{A}}^{(3)} \rrbracket - \mathcal{X}\|_F^2}{\|\mathcal{X}\|_F^2}$ .
- Number of trials: 10.
- Environmental: 2.3 GHz 8-Core Intel Core i9 CPU with 16 GB 2400 MHz DDR4 memory.

<sup>16</sup>Brett W. Larsen and Tamara G. Kolda. "Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition". In: *SIAM J. Matrix Anal. Appl.* 43.3 (2022), pp. 1488–1517.  
doi: 10.1137/21M1441754.



## DIFFERENT TENSOR SIZE

**Table:** Performance of the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ , the target rank  $R = 10$ ,  $noise = 0$ , and random initialization for different tensors generated by  $I \times 10$  factor matrices with different  $I$ .

Algorithms		$I = 100$	$I = 200$	$I = 300$	$I = 400$	$I = 500$
		<i>spread</i> = 15, <i>magnitude</i> = 24	<i>spread</i> = 30, <i>magnitude</i> = 30	<i>spread</i> = 45, <i>magnitude</i> = 36	<i>spread</i> = 60, <i>magnitude</i> = 42	<i>spread</i> = 75, <i>magnitude</i> = 48
AdaCPD [Fu+20]	Iterations	5962.7	4206.3	3105.8	3271.9	4229.5
	Seconds	23.997029	131.46287	472.43042	1241.3702	3224.5938
EAdawCPD	Iterations	2242.1	572.2	390.3	391.1	488.6
	Seconds	9.1444463	18.08937	58.656213	150.04822	372.96764
LAdawCPD	Iterations	2394.9	577.8	429.9	483.2	607.3
	Seconds	10.634388	18.25461	64.193687	184.35716	463.1262

■ More numerical results see the Appendix:

- 1 Different target rank;
- 2 Different noise;
- 3 Nonnegative constraint;
- 4 Real data.



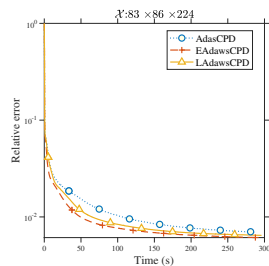
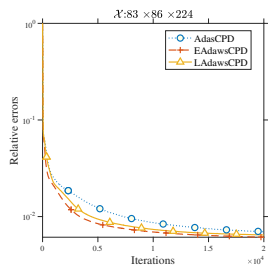
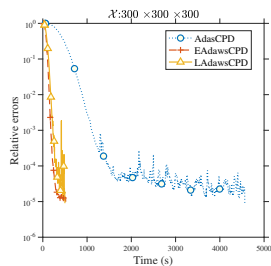
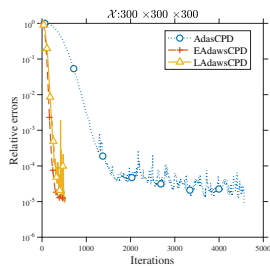
## REAL DATA: PERFORMANCE

**Table:** Performance of algorithms on real datasets.

Algorithms		SalinasA. ( $83 \times 86 \times 224$ )	Indian Pines ( $145 \times 145 \times 220$ )	Pavia Uni. ( $610 \times 340 \times 103$ )
		$R = 10,$ $ \mathcal{F}_n  = 20$	$R = 10,$ $ \mathcal{F}_n  = 20$	$R = 100,$ $ \mathcal{F}_n  = 20$
AdaCPD [Fu+20]	<i>Tol</i>	0.00697493	0.00782241	0.02831117
	Seconds	288.535306	653.276364	4387.12147
EAdawCPD	<i>Tol</i>	0.00611473	0.00725646	0.02792415
	Seconds	291.374795	659.02316	4450.21809
LAdawCPD	<i>Tol</i>	0.00644397	0.00741898	0.02796972
	Seconds	295.962095	663.648648	4559.80737



## VISUALISATION OF EXPERIMENTAL RESULTS



(a) Number of iterations v.s. Relative errors: Synthetic  $I = 300$

(b) Time v.s. Relative errors: Synthetic  $I = 300$

(c) Number of iterations v.s. Relative errors: SalinasA

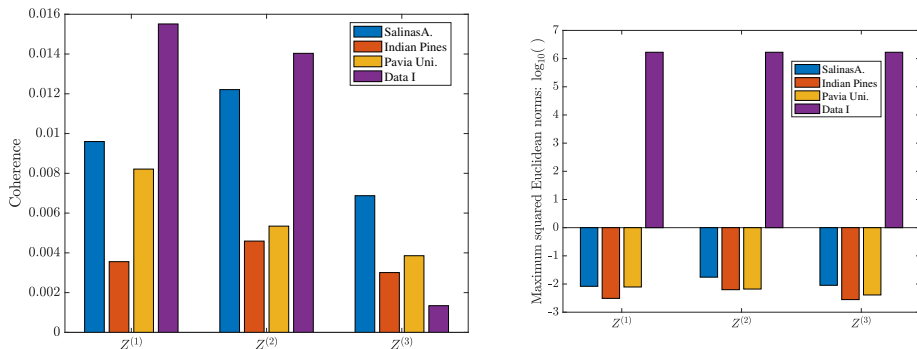
(d) Time v.s. Relative errors: SalinasA

**Figure:** (a)-(b) Output by the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ ,  $R = 10$ ,  $noise = 0$ , and random initialization for the tensor with  $I = 300$ ,  $R_{true} = 10$ ,  $spread = 45$ , and  $magnitude = 36$ ; (c)-(d) Output by the algorithms with  $|\mathcal{F}_n| = 20$  and  $R = 10$  for SalinasA.





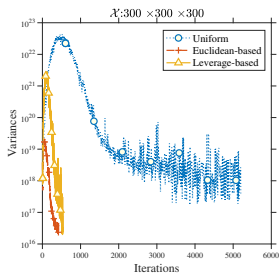
## SOME INTERPRETATIONS OF THE EXPERIMENTAL RESULTS



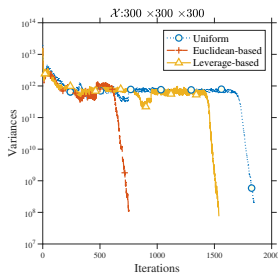
**Figure:** **Left:** Coherence for four different tensors.; **Right:** Maximum squared Euclidean norms for four different tensors.

## COMPARE VARIANCES OF DIFFERENT DATA

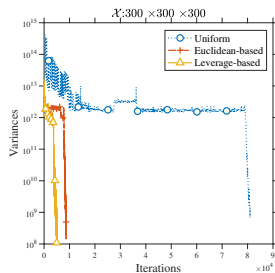
- **Data I:** Previous data.
- **Data II:** Generated by three factor matrices of size  $300 \times 10$  with independent standard Gaussian entries.
- **Data III:** The same as **Data II** except that one entry of each factor matrices is chosen uniformly and set to 20.



(a) Data I



(b) Data II



(c) Data III

**Figure:** Comparison of variances of stochastic gradients for different tensors.



# PRESENTATION OUTLINE

- 1 Introduction
- 2 Proposed Method
- 3 Numerical Results
- 4 Conclusions

# CONCLUSIONS

## ■ Conclusions

- We propose a block-randomized gradient descent method with importance sampling for CP decomposition based on two empirical sampling strategies.
- We provide a theoretical optimal sampling probability contribution.
- Numerical experiments are provided to test the proposed methods.

## ■ Future works

- Momentum version<sup>17</sup>.
- Variance reduction technique<sup>18</sup>.
- Scaled technique<sup>19</sup>.

<sup>17</sup>Qingsong Wang, Chunfeng Cui, and Deren Han. “A Momentum Block-Randomized Stochastic Algorithm for Low-Rank Tensor CP Decomposition”. In: *Pac. J. Optim.* 17.3 (2021), pp. 433–452.

<sup>18</sup>Qingsong Wang, Chunfeng Cui, and Deren Han. “Accelerated Doubly Stochastic Gradient Descent for Tensor CP Decomposition”. In: *J. Optim. Theory Appl.* (2023), pp. 1–40. doi: 10.1007/s10957-023-02193-5.

<sup>19</sup>Tian Tong, Cong Ma, and Yuejie Chi. “Accelerating Ill-Conditioned Low-Rank Matrix Estimation via Scaled Gradient Descent.”. In: *J. Mach. Learn. Res.* 22.150 (2021), pp. 1–63.



# Thanks!

Y.J. Yu and H.Y. Li. "A Block-Randomized Stochastic Method with Importance Sampling for CP Tensor Decomposition" on arXiv.

## REFERENCES I

- [ADK11] Evrim Acar, Daniel M Dunlavy, and Tamara G Kolda. “A scalable optimization approach for fitting canonical tensor decompositions”. In: *J. Chemom.* 25.2 (2011), pp. 67–86.
- [AKD11] Evrim Acar, Tamara G Kolda, and Daniel M Dunlavy. *All-at-once optimization for coupled matrix and tensor factorizations*. arXiv:1105.3422. 2011.
- [BBK18] Casey Battaglino, Grey Ballard, and Tamara G. Kolda. “A Practical Randomized CP Tensor Decomposition”. In: *SIAM J. Matrix Anal. Appl.* 39.2 (2018), pp. 876–901. doi: 10.1137/17M1112303.
- [CHE+16] Dehua Cheng et al. “SPALS: Fast Alternating Least Squares via Implicit Leverage Scores Sampling”. In: *Proceedings of the 30th International Conference on Neural Information Processing Systems*. Barcelona Spain: Curran Associates Inc., 2016, pp. 721–729.

## REFERENCES II

- [DRI+12] Petros Drineas et al. “Fast Approximation of Matrix Coherence and Statistical Leverage”. In: *J. Mach. Learn. Res.* 13.1 (2012), pp. 3475–3506.
- [Fu+20] Xiao Fu et al. “Block-Randomized Stochastic Proximal Gradient for Low-Rank Tensor Factorization”. In: *IEEE Trans. Signal Process.* 68 (2020), pp. 2170–2185. DOI: 10.1109/TSP.2020.2982321..
- [HF20] Kejun Huang and Xiao Fu. “Low-Complexity Levenberg-Marquardt Algorithm for Tensor Canonical Polyadic Decomposition”. In: *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2020, pp. 3922–3926.
- [KB09] Tamara G. Kolda and Brett W. Bader. “Tensor Decompositions and Applications”. In: *SIAM Rev.* 51.3 (2009), pp. 455–500. DOI: 10.1137/07070111X.

## REFERENCES III

- [LK22] Brett W. Larsen and Tamara G. Kolda. “Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition”. In: *SIAM J. Matrix Anal. Appl.* 43.3 (2022), pp. 1488–1517. DOI: 10.1137/21M1441754.
- [MMY15] Ping Ma, Michael Mahoney, and Bin Yu. “A statistical Perspective on algorithmic leveraging”. In: *J. Mach. Learn. Res.* 16.27 (2015), pp. 861–911. DOI: 10.1002/cem.1335.
- [NSW16] Deanna Needell, Nathan Srebro, and Rachel Ward. “Stochastic Gradient Descent, Weighted Sampling, and the Randomized Kaczmarz Algorithm”. In: *Math. Program.* 155.1 (2016), pp. 549–573. DOI: 10.1007/s10107-015-0864-7.
- [PTC13] Anh-Huy Phan, Petr Tichavský, and Andrzej Cichocki. “Low Complexity Damped Gauss–Newton Algorithms for CANDECOMP/PARAFAC”. In: *SIAM J. Matrix Anal. Appl.* 34.1 (2013), pp. 126–147. DOI: 10.1137/100808034.



## REFERENCES IV

- [TMC21] Tian Tong, Cong Ma, and Yuejie Chi. “Accelerating Ill-Conditioned Low-Rank Matrix Estimation via Scaled Gradient Descent.”. In: *J. Mach. Learn. Res.* 22.150 (2021), pp. 1–63.
- [VAN21] Michiel Vandecappelle. “Numerical Algorithms for Tensor Decompositions”. PhD thesis. Arenberg Doctoral School, 2021.
- [VD16] Nico Vervliet and Lieven De Lathauwer. “A Randomized Block Sampling Approach to Canonical Polyadic Decomposition of Large-Scale Tensors”. In: *IEEE J. Sel. Topics Signal Process.* 10.2 (2016), pp. 284–295. DOI: 10.1109/JSTSP.2015.2503260.
- [WAN+23] Qingsong Wang et al. “Inertial accelerated SGD algorithms for solving large-scale lower-rank tensor CP decomposition problems”. In: *J. Comput. Appl. Math.* 423 (2023), p. 114948. DOI: 10.1016/j.cam.2022.114948.

## REFERENCES V

- [WCH21] Qingsong Wang, Chunfeng Cui, and Deren Han. “A Momentum Block-Randomized Stochastic Algorithm for Low-Rank Tensor CP Decomposition”. In: *Pac. J. Optim.* 17.3 (2021), pp. 433–452.
- [WCH23] Qingsong Wang, Chunfeng Cui, and Deren Han. “Accelerated Doubly Stochastic Gradient Descent for Tensor CP Decomposition”. In: *J. Optim. Theory Appl.* (2023), pp. 1–40. doi: 10.1007/s10957-023-02193-5.
- [Woo14] David P. Woodruff. “Sketching as a Tool for Numerical Linear Algebra”. In: *Found. Trends Theor. Comput. Sci.* 10.1–2 (2014), pp. 1–157. doi: 10.1561/04000000060.
- [XY13] Yangyang Xu and Wotao Yin. “A Block Coordinate Descent Method for Regularized Multiconvex Optimization with Applications to Nonnegative Tensor Factorization and Completion”. In: *SIAM J. Imaging Sci.* 6.3 (2013), pp. 1758–1789. doi: 10.1137/120887795.

## PRESENTATION OUTLINE

5 Detailed Algorithms

6 More Numerical Results

7 Some Others

---

## Algorithm 2 CP-ALS<sup>20</sup>

---

1: **function**  $[\lambda, \{\mathbf{A}^{(n)}\}_{n=1}^N] = \text{CP-ALS}(\mathcal{X}, R, \{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N)$

▷  $N$ -way tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ ;

▷ rank  $R$ , initialization  $\{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N$

2:   **repeat**

3:     **for**  $n = 1, \dots, N$  **do**

4:        $\mathbf{V} \leftarrow \mathbf{A}^{(1)\top} \mathbf{A}^{(1)} * \dots * \mathbf{A}^{(n-1)\top} \mathbf{A}^{(n-1)} * \mathbf{A}^{(n+1)\top} \mathbf{A}^{(n+1)} * \dots * \mathbf{A}^{(N)\top} \mathbf{A}^{(N)}$

5:        $\mathbf{A}^{(n)} \leftarrow \mathbf{X}_{(n)} (\mathbf{A}^{(N)} \odot \dots \odot \mathbf{A}^{(n+1)} \odot \mathbf{A}^{(n-1)} \odot \dots \odot \mathbf{A}^{(1)}) \mathbf{V}^\dagger$

6:       normalize columns of  $\mathbf{A}^{(n)}$  (storing norms as  $\lambda$ )

7:     **end for**

8:   **until** some stopping criterion is reached

9:   **return**  $\lambda, \{\mathbf{A}^{(n)}\}_{n=1}^N$

10: **end function**

---

<sup>20</sup>Tamara G. Kolda and Brett W. Bader. "Tensor Decompositions and Applications". In: *SIAM Rev.* 51.3 (2009), pp. 455–500. doi: 10.1137/07070111X. < > < > < > < > < >

---

**Algorithm 3** SKRP and Sampled tensor fibers (SKRP-ST)<sup>21</sup>


---

```

1: function [ $\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}$ ] = SKRP-ST( $n, \mathbf{idx}, \{\mathbf{A}^{(k)}\}_{k=1, k \neq n}^N$ )
2:    $\mathbf{Z}_{\mathcal{F}_n}^{(n)} \leftarrow \mathbf{1}$ 
3:   for  $m = 1, \dots, n-1, n+1, \dots, N$  do
4:      $\mathbf{A}_{|\mathcal{F}_n|}^{(m)} \leftarrow \mathbf{A}^{(m)}(\mathbf{idx}(:, m), :)$ 
5:      $\mathbf{Z}_{\mathcal{F}_n}^{(n)} \leftarrow \mathbf{Z}_{\mathcal{F}_n}^{(n)} \circledast \mathbf{A}_{|\mathcal{F}_n|}^{(m)}$ 
6:   end for
7:    $\mathcal{X}^{\mathcal{F}_n} \leftarrow \mathcal{X}(\mathbf{idx}(:, 1), \dots, \mathcal{X}(\mathbf{idx}(:, n-1), :, \mathcal{X}(\mathbf{idx}(:, n+1), \dots, \mathcal{X}(\mathbf{idx}(:, N)))$ 
8:    $\mathbf{X}_{(n)}^{\mathcal{F}_n} \leftarrow \text{UNFOLDING}(\mathcal{X}^{\mathcal{F}_n}, n)$ 
9:   return  $\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}$ 
10: end function

```

▷  $\mathbf{1} \in \mathbb{R}^{|\mathcal{F}_n| \times R}$  a matrix with all elements 1

---

<sup>21</sup>Casey Battaglino, Grey Ballard, and Tamara G. Kolda. "A Practical Randomized CP Tensor Decomposition". In: *SIAM J. Matrix Anal. Appl.* 39.2 (2018), pp. 876–901. doi: 10.1137/17M1112303.

## Algorithm 4 CPRAND<sup>22</sup>

1: **function**  $[\lambda, \{\mathbf{A}^{(n)}\}_{n=1}^N] = \text{CPRAND}(\mathcal{X}, R, \{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N, |\mathcal{F}_n|)$

▷  $N$ -way tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ ;

▷ rank  $R$ , sample size  $|\mathcal{F}_n|$ , initialization  $\{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N$

2: **repeat**

3:     **for**  $n = 1, \dots, N$  **do**

4:         get  $\mathbf{id}\mathbf{x}$  using **uniform sampling**

5:          $[\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}] = \text{SKRP-ST}(n, \mathbf{id}\mathbf{x}, \{\mathbf{A}^{(k)}\}_{k=1, k \neq n}^N)$

6:          $\mathbf{A}^{(n)} \leftarrow \arg \min_{\mathbf{A}} \|\mathbf{Z}_{\mathcal{F}_n}^{(n)} \mathbf{A}^\top - \mathbf{X}_{(n)}^{\mathcal{F}_n}\|_F$

7:     **end for**

8: **until** some stopping criterion is reached

9: **return**  $\{\mathbf{A}^{(n)}\}_{n=1}^N$

10: **end function**

<sup>22</sup>Casey Battaglino, Grey Ballard, and Tamara G. Kolda. "A Practical Randomized CP Tensor Decomposition". In: *SIAM J. Matrix Anal. Appl.* 39.2 (2018), pp. 876–901. doi: 10.1137/17M1112303.

---

**Algorithm 5** CP-ALS with leveraged-based sampling<sup>23</sup>


---

```

1: function [ $\lambda, \{\mathbf{A}^{(n)}\}_{n=1}^N$ ] = CPRAND( $\mathcal{X}, R, \{\mathbf{A}^{(n)}\}_{n=1}^N, |\mathcal{F}_n|$ )
                                      $\triangleright N$ -way tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ ;
                                      $\triangleright$  rank  $R$ , sample size  $|\mathcal{F}_n|$ , initialization  $\{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N$ 
2:   Compute the leverage scores of factor matrices  $\{\mathbf{A}^{(0)}\}_{n=1}^N$  and get the probability distributions
3:   repeat
4:     for  $n = 1, \dots, N$  do
5:       get  $\mathbf{id}\mathbf{x}$  using leverage-based sampling with the probability distributions
6:        $[\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{\mathcal{F}_n}^{(n)}] = \text{SKRP-ST}(n, \mathbf{id}\mathbf{x}, \{\mathbf{A}^{(k)}\}_{k=1, k \neq n}^N)$ 
7:        $\mathbf{A}^{(n)} \leftarrow \arg \min_{\mathbf{A}} \|\mathbf{Z}_{\mathcal{F}_n}^{(n)} \mathbf{A}^\top - \mathbf{X}_{\mathcal{F}_n}^{(n)}\|_F$ 
8:       Recompute the leverage scores of  $\mathbf{A}^{(n)}$  and according probability distributions
9:     end for
10:  until some stopping criterion is reached
11:  return  $\{\mathbf{A}^{(n)}\}_{n=1}^N$ 
12: end function

```

---

<sup>23</sup>Brett W. Larsen and Tamara G. Kolda. "Practical Leverage-Based Sampling for Low-Rank Tensor Decomposition". In: *SIAM J. Matrix Anal. Appl.* 43.3 (2022), pp. 1488–1517. doi: 10.1137/21M1441754.

---

**Algorithm 6** BrasCPD<sup>24</sup>


---

1: **function**  $\{\mathbf{A}^{(n)}\}_{n=1}^N = \text{BRASCPD}(\mathcal{X}, R, |\mathcal{F}_n|, \{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N, \{\alpha^t\})$  ▷  $N$ -way tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ ;  
▷ rank  $R$ , sample size  $|\mathcal{F}_n|$ ; initialization  $\{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N$ , step size  $\{\alpha^t\}_{t=0,1,\dots}$

2:      $t \leftarrow 0$

3:     **repeat**

4:         Uniformly sample  $n$  from  $\{1, \dots, N\}$

5:         Sample  $\mathcal{F}_n$  **uniformly** from  $\{1, \dots, J_n\}$

6:         Compute  $\mathbf{G}_{(t)}^{(n)} = \frac{1}{|\mathcal{F}_n|} \left( \mathbf{A}_{(t)}^{(n)} (\mathbf{Z}_{\mathcal{F}_n}^{(n)})^\top \mathbf{Z}_{\mathcal{F}_n}^{(n)} - \mathbf{X}_{(n)}^{\mathcal{F}_n} \mathbf{Z}_{\mathcal{F}_n}^{(n)} \right)$

7:         Update  $\mathbf{A}_{(t+1)}^{(n)} \leftarrow \mathbf{A}_{(t)}^{(n)} - \alpha_t \mathbf{G}_{(t)}^{(n)}$ , and  $\mathbf{A}_{(t+1)}^{(n')} \leftarrow \mathbf{A}_{(t)}^{(n')}$  for  $n' \neq n$

8:          $t \leftarrow t + 1$

9:     **until** some stopping criterion is reached

10:    **return**  $\{\mathbf{A}^{(n)}\}_{n=1}^N$

11: **end function**

---

<sup>24</sup>Xiao Fu et al. "Block-Randomized Stochastic Proximal Gradient for Low-Rank Tensor Factorization". In: *IEEE Trans. Signal Process.* 68 (2020), pp. 2170–2185. doi: 10.1109/TSP.2020.2982321..



## Algorithm 7 BrawsCPD

---

```

1: function  $\{\mathbf{A}^{(n)}\}_{n=1}^N = \text{BRAWSCPD}(\mathcal{X}, R, |\mathcal{F}_n|, \{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N, \{\alpha_t\})$ 
2:    $t \leftarrow 0$ 
3:   repeat
4:     Uniformly sample  $n$  from  $\{1, \dots, N\}$ 
5:      $[\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}, \mathbf{p}_{\mathcal{F}_n}] = \text{SKRP-ST-I}(n, |\mathcal{F}_n|, \{\mathbf{p}_k\}_{k=1, k \neq n}^N, \{\mathbf{A}^{(N)}\}_{k=1, k \neq n}^N)$ 
6:     Form the stochastic gradient  $\mathbf{G}_{(t)}^{(n)} \leftarrow \frac{1}{|\mathcal{F}_n|J_n} \left( \mathbf{A}_{(t)}^{(n)} (\mathbf{D}(\mathbf{Z}_{\mathcal{F}_n}^{(n)})^\top) \mathbf{Z}_{\mathcal{F}_n}^{(n)} - \mathbf{D}\mathbf{X}_{(n)}^{\mathcal{F}_n} \mathbf{Z}_{\mathcal{F}_n}^{(n)} \right)$ 
7:     Update  $\mathbf{A}_{(t+1)}^{(n)} \leftarrow \mathbf{A}_{(t)}^{(n)} - \alpha_t \mathbf{G}_{(t)}^{(n)}$ ,  $\mathbf{A}_{(t+1)}^{(n')} \leftarrow \mathbf{A}_{(t)}^{(n')}$  for  $n' \neq n$ 
8:      $t \leftarrow t + 1$ 
9:   until some stopping criterion is reached
10:  return  $\{\mathbf{A}^{(n)}\}_{n=1}^N$ 
11: end function

```

---

## Algorithm 8 AdawCPD

- 1: **function**  $\{\mathbf{A}^{(n)}\}_{n=1}^N = \text{ADAWCPD}(\mathcal{X}, R, |\mathcal{F}_n|, \{\mathbf{A}_{(0)}^{(n)}\}_{n=1}^N)$
- 2:      $t \leftarrow 0$
- 3:     **repeat**
- 4:         Uniformly sample  $n$  from  $\{1, \dots, N\}$
- 5:          $[\mathbf{Z}_{\mathcal{F}_n}^{(n)}, \mathbf{X}_{(n)}^{\mathcal{F}_n}, \mathbf{p}_{\mathcal{F}_n}] = \text{SKRP-ST-I}(n, |\mathcal{F}_n|, \{\mathbf{p}_k\}_{k=1, k \neq n}^N, \{\mathbf{A}^{(N)}\}_{k=1, k \neq n}^N)$
- 6:         Form the stochastic gradient  $\mathbf{G}_{(t)}^{(n)} \leftarrow \frac{1}{|\mathcal{F}_n|J_n} \left( \mathbf{A}_{(t)}^{(n)} (\mathbf{D}(\mathbf{Z}_{\mathcal{F}_n}^{(n)})^\top) \mathbf{Z}_{\mathcal{F}_n}^{(n)} - \mathbf{D}\mathbf{X}_{(n)}^{\mathcal{F}_n} \mathbf{Z}_{\mathcal{F}_n}^{(n)} \right)$
- 7:         Determine the step size  $\eta_{(t)}^{(n)} \leftarrow \frac{\eta}{(b + \sum_{t'=1}^t [\mathbf{G}_{(t')}^{(n)}]_{i,r}^2)^{1/2}}$
- 8:         Update  $\mathbf{A}_{(t+1)}^{(n)} \leftarrow \mathbf{A}_{(t)}^{(n)} - \alpha_t \mathbf{G}_{(t)}^{(n)}$ ,  $\mathbf{A}_{(t+1)}^{(n')} \leftarrow \mathbf{A}_{(t)}^{(n')}$  for  $n' \neq n$
- 9:          $t \leftarrow t + 1$
- 10:     **until** some stopping criterion is reached
- 11:     **return**  $\{\mathbf{A}^{(n)}\}_{n=1}^N$
- 12: **end function**

# PRESENTATION OUTLINE

5 Detailed Algorithms

6 More Numerical Results

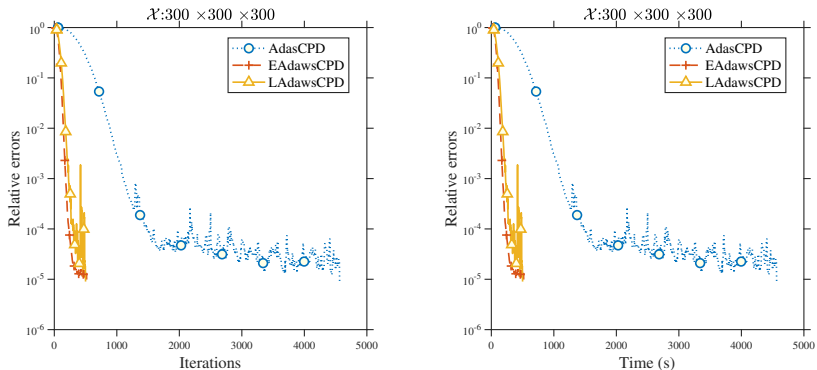
7 Some Others

## DIFFERENT TENSOR SIZE

**Table:** Performance of the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ , the target rank  $R = 10$ ,  $noise = 0$ , and random initialization for different tensors generated by  $I \times 10$  factor matrices with different  $I$ .

Algorithms		$I = 100$	$I = 200$	$I = 300$	$I = 400$	$I = 500$
		$spread = 15,$ $magnitude = 24$	$spread = 30,$ $magnitude = 30$	$spread = 45,$ $magnitude = 36$	$spread = 60,$ $magnitude = 42$	$spread = 75,$ $magnitude = 48$
AdasCPD [Fu+20]	Iterations	5962.7	4206.3	3105.8	3271.9	4229.5
	Seconds	23.997029	131.46287	472.43042	1241.3702	3224.5938
EAdawsCPD	Iterations	2242.1	572.2	390.3	391.1	488.6
	Seconds	9.1444463	18.08937	58.656213	150.04822	372.96764
LAdawsCPD	Iterations	2394.9	577.8	429.9	483.2	607.3
	Seconds	10.634388	18.25461	64.193687	184.35716	463.1262

## DIFFERENT TENSOR SIZE



**Figure:** Number of iterations v.s. Relative errors and Time v.s. Relative errors output by the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ ,  $R = 10$ ,  $noise = 0$ , and random initialization for the tensor with  $I = 300$ ,  $R_{true} = 10$ ,  $spread = 45$ , and  $magnitude = 36$ .

## DIFFERENT TARGET RANK

**Table:** Performance of the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ , different target ranks  $R$ ,  $noise = 0$ , and random initialization for the tensor generated by  $300 \times 10$  factor matrices with  $spread = 45$  and  $magnitude = 36$ .

Algorithms		$R = 5$	$R = 10$	$R = 15$	$R = 20$
AdaCPD [Fu+20]	Iterations	3059.1	3105.8	4823.7	7094.8
	Seconds	480.47423	472.43042	767.12695	1125.8787
EAdawCPD	Iterations	474	390.3	406.7	824.9
	Seconds	73.963123	58.656213	65.151505	132.01017
LAdawCPD	Iterations	535.7	429.9	567.1	827.1
	Seconds	83.791824	64.193687	90.722473	132.73588

## DIFFERENT NOISE

**Table:** Performance of the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ , the target rank  $R = 10$ , and random initialization for different tensors generated by  $300 \times 10$  factor matrices with  $spread = 45$ ,  $magnitude = 36$ , and different noises.

Algorithms		$noise = 0$	$noise = 0.01$	$noise = 0.1$	$noise = 1$
AdaCPD [Fu+20]	Iterations	3105.8	3211.7	3024.9	3511.2
	Seconds	472.430421	505.928736	481.988384	553.33599
EAdawCPD	Iterations	390.3	360.5	381.6	414.6
	Seconds	58.6562125	56.748149	60.5743992	65.3003044
LAdawCPD	Iterations	429.9	410	439	452.1
	Seconds	64.1936869	64.6965127	70.8384762	71.6683267

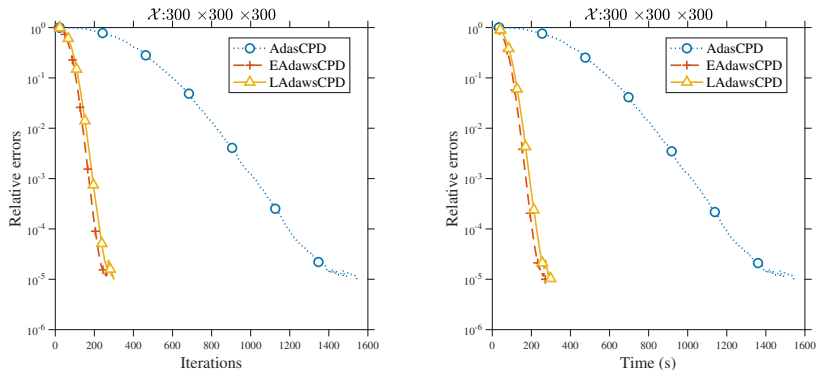
## NONNEGATIVE CONSTRAINT

**Table:** Performance of the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ , the target rank  $R = 10$ ,  $noise = 0$ , and random initialization for different tensors generated by  $I \times 10$  factor matrices with different  $I$  under nonnegative constraint.

Algorithms		$I = 100$	$I = 200$	$I = 300$	$I = 400$	$I = 500$
		<i>spread</i> = 15, <i>magnitude</i> = 24	<i>spread</i> = 30, <i>magnitude</i> = 30	<i>spread</i> = 45, <i>magnitude</i> = 36	<i>spread</i> = 60, <i>magnitude</i> = 42	<i>spread</i> = 75, <i>magnitude</i> = 48
AdaCPD [Fu+20]	Iterations	500.1	910.7	1424.5	2203.7	3372.8
	Seconds	2.0256353	28.981344	225.79166	737.83326	2550.5717
EAdawCPD	Iterations	156.6	190.6	250.6	319	394.2
	Seconds	0.6437757	6.0335162	40.041681	110.00088	299.55425
LAdawCPD	Iterations	153.7	205.9	271.5	337.5	437.6
	Seconds	0.6892614	6.6271954	43.308729	111.84547	331.55795



# NONNEGATIVE CONSTRAINT



**Figure:** Number of iterations v.s. Relative errors and Time v.s. Relative errors output by the algorithms with  $Tol = 10^{-5}$ ,  $|\mathcal{F}_n| = 18$ ,  $R = 10$ ,  $noise = 0$ , and random initialization for the tensor with  $I = 300$ ,  $R_{true} = 10$ ,  $spread = 45$ ,  $magnitude = 36$ , and  $\mathbf{A}^{(n)} \geq \mathbf{0}$ .

# REAL DATA

**Table:** Size and Type of Real Datasets.

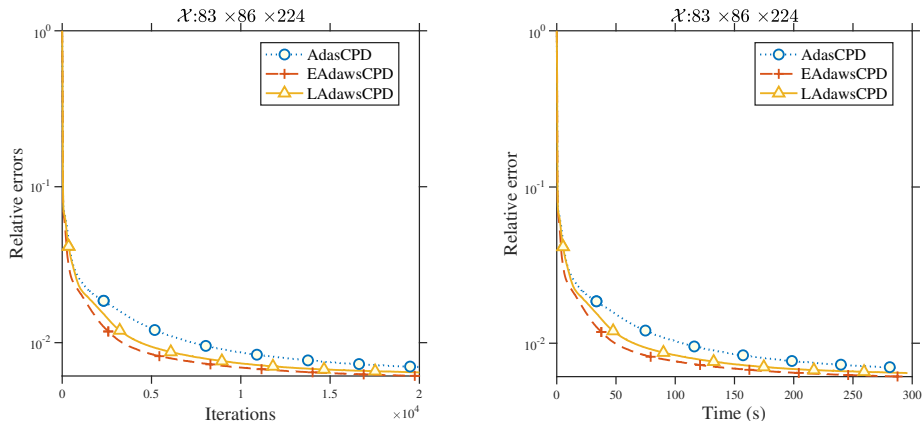
Dataset	Size	Type
SalinasA.	$83 \times 86 \times 224$	Hyperspectral
Indian Pines	$145 \times 145 \times 220$	Hyperspectral
Pavia Uni.	$610 \times 340 \times 103$	Hyperspectral

## REAL DATA: PERFORMANCE

**Table:** Performance of algorithms on real datasets.

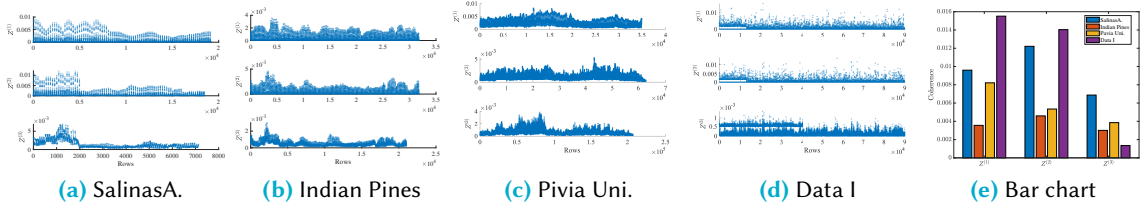
Algorithms		SalinasA.	Indian Pines	Pavia Uni.
		$R = 10,$ $ \mathcal{F}_n  = 20$	$R = 10,$ $ \mathcal{F}_n  = 20$	$R = 100,$ $ \mathcal{F}_n  = 20$
AdaCPD [Fu+20]	<i>Tol</i>	0.00697493	0.00782241	0.02831117
	Seconds	288.535306	653.276364	4387.12147
EAdawCPD	<i>Tol</i>	0.00611473	0.00725646	0.02792415
	Seconds	291.374795	659.02316	4450.21809
LAdawCPD	<i>Tol</i>	0.00644397	0.00741898	0.02796972
	Seconds	295.962095	663.648648	4559.80737

## REAL DATA: PERFORMANCE



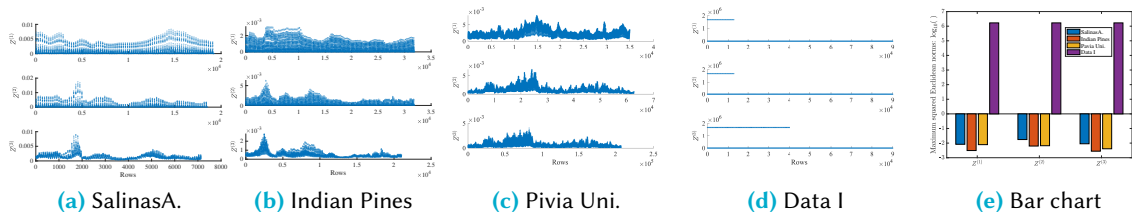
**Figure:** Number of iterations v.s. Relative errors and Time v.s. Relative errors output by the algorithms with  $|\mathcal{F}_n| = 20$  and  $R = 10$  for SalinasA..

## MORE INTERPRETATIONS OF THE EXPERIMENTAL RESULTS



**Figure:** (a)-(d) Leverage scores for four different tensors. (e) Coherence for four different tensors.

## MORE INTERPRETATIONS OF THE EXPERIMENTAL RESULTS



**Figure:** (a)-(d) Squared Euclidean norms for four different tensors. (e) Maximum squared Euclidean norms for four different tensors.

# PRESENTATION OUTLINE

5 Detailed Algorithms

6 More Numerical Results

7 Some Others

### Definition 8.1 (Leverage Scores [Dri+12])

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $m > n$ , and let  $\mathbf{Q} \in \mathbb{R}^{m \times n}$  be any orthogonal basis for the column space of  $\mathbf{A}$ . The *leverage score* of the  $i$ -th row of  $\mathbf{A}$  is given by

$$\ell_i(\mathbf{A}) = \|\mathbf{Q}(i, :)\|_2^2.$$

### Definition 8.2 (Leveraged-based Probability Distribution [Woo14])

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $m > n$ . We say a probability distribution  $\mathbf{p} = [p_1, \dots, p_m]^\top$  is a *leveraged-based probability distribution* for  $\mathbf{A}$  if  $p_i \geq \beta \frac{\ell_i(\mathbf{A})}{n}$  with  $0 < \beta \leq 1$  and  $i \in [m]$ .

### Definition 8.3 (Euclidean-based Probability Distribution)

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $m > n$ . We say a probability distribution  $\mathbf{p} = [p_1, \dots, p_m]^\top$  is an *Euclidean-based probability distribution* for  $\mathbf{A}$  if  $p_i \geq \beta \|\mathbf{A}(i, :)\|_2^2 / \|\mathbf{A}\|_F^2$  with  $0 < \beta \leq 1$  and  $i \in [m]$ .



## MINI-BATCHES AND THEIR SAMPLING PROBABILITIES

- 1 First sample  $|\mathcal{F}_n|$  rows from each  $\{\mathbf{A}^{(k)}\}_{k=1, k \neq n}^N$  using the leveraged-based probability distribution for  $\mathbf{A}^{(k)}$ , i.e.,  $\mathbf{p}_k = \frac{\ell(\mathbf{A}^{(k)})}{R}$ . Thus, we can get the  $\mathbf{id}\mathbf{x}$ :

$$\begin{bmatrix} \{i_1^{(j_1)} \quad \cdots \quad i_{n-1}^{(j_1)} \quad i_{n+1}^{(j_1)} \quad \cdots \quad i_N^{(j_1)}\} \\ \{i_1^{(j_2)} \quad \cdots \quad i_{n-1}^{(j_2)} \quad i_{n+1}^{(j_2)} \quad \cdots \quad i_N^{(j_2)}\} \\ \vdots \\ \{i_1^{(j_{|\mathcal{F}_n|})} \quad \cdots \quad i_{n-1}^{(j_{|\mathcal{F}_n|})} \quad i_{n+1}^{(j_{|\mathcal{F}_n|})} \quad \cdots \quad i_N^{(j_{|\mathcal{F}_n|})}\} \end{bmatrix}. \quad (8.1)$$

- 2 Based on (8.1), we can obtain the index set  $\mathcal{F}_n = \{j_1, j_2, \dots, j_{|\mathcal{F}_n|}\}$  of the sampled  $|\mathcal{F}_n|$  rows of  $\mathbf{Z}^{(n)}$ , and the corresponding sampling probabilities

$$\mathbf{p}_{\mathcal{F}_n} = [p_{j_1}, \dots, p_{j_{|\mathcal{F}_n|}}]^\top, \quad (8.2)$$

where  $p_{j_f} = \frac{\bar{\ell}_{j_f}(\mathbf{Z}^{(n)})}{R^{N-1}}$  and  $\bar{\ell}_{j_f}(\mathbf{Z}^{(n)}) = \prod_{k=1, k \neq n}^N \ell_{i_k^{(j_f)}}(\mathbf{A}^{(k)})$ .

- With the above  $\mathbf{id}\mathbf{x}$ , we can find  $\mathbf{Z}_{\mathcal{F}_n}^{(n)}$  and  $\mathbf{X}_{(n)}^{\mathcal{F}_n}$  using efficient sampling without forming KRP, and using (8.2) to get corresponding probabilities.

## VARIANCE OF STOCHASTIC GRADIENT

### Theorem 8.4


In the setting of Theorem 2.5, suppose that  $\mathbf{p} \in \mathbb{R}^{J_n}$  is any probability distribution proposed in previous, and  $\mathbf{R}_{(t)}^{(n)} = \mathbf{A}_{(t)}^{(n)} (\mathbf{Z}_{(t)}^{(n)})^\top - \mathbf{X}_{(n)}$ . Then

$$\begin{aligned} & \mathbb{E}_{\zeta(t)} \left[ \left\| \mathbf{G}_{(t)}^{(\xi(t))} - \nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}_{(t)}) \right\|_F^2 \mid \mathcal{B}_{(t)}, \xi(t) \right] \\ &= \frac{1}{|\mathcal{F}_n|} \sum_{j_f=1}^{J_n} \frac{1}{p_{j_f}} \left\| \mathbf{R}_{(t)}^{(\xi(t))}(:, j_f) \right\|_2^2 \left\| \mathbf{Z}_{(t)}^{(\xi(t))}(j_f, :) \right\|_2^2 - \frac{1}{|\mathcal{F}_n|} \left\| \nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}_{(t)}) \right\|_F^2. \end{aligned}$$

# VARIANCE OF STOCHASTIC GRADIENT

**Table:** Variances for different probability distributions

Probability distributions	$\mathbb{E}_{\zeta(t)} \left[ \left\  \mathbf{G}_{(t)}^{(\xi(t))} - \nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}(t)) \right\ _F^2 \mid \mathcal{B}(t), \xi(t) \right]$
Uniform <sup>25</sup>	$\frac{J_n}{ \mathcal{F}_n } \sum_{j_f=1}^{J_n} \left\  \mathbf{R}_{(t)}^{(\xi(t))}(:, j_f) \right\ _2^2 \left\  \mathbf{Z}_{(t)}^{(\xi(t))}(j_f, :) \right\ _2^2 - \frac{1}{ \mathcal{F}_n } \left\  \nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}(t)) \right\ _F^2$
Leverage-based	$\frac{R^{N-1}}{ \mathcal{F}_n } \sum_{j_f=1}^{J_n} \frac{\left\  \mathbf{R}_{(t)}^{(\xi(t))}(:, j_f) \right\ _2^2 \left\  \mathbf{Z}_{(t)}^{(\xi(t))}(j_f, :) \right\ _2^2}{\prod_{k \neq \xi(t)} \ell_{(j_f)}^{i_k}(\mathbf{A}^{(k)})} - \frac{1}{ \mathcal{F}_n } \left\  \nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}(t)) \right\ _F^2$
Euclidean-based	$\frac{\prod_{k \neq \xi(t)} \left\  \mathbf{A}^{(k)} \right\ _F^2}{ \mathcal{F}_n } \sum_{j_f=1}^{J_n} \frac{\left\  \mathbf{R}_{(t)}^{(\xi(t))}(:, j_f) \right\ _2^2 \left\  \mathbf{Z}_{(t)}^{(\xi(t))}(j_f, :) \right\ _2^2}{\prod_{k \neq \xi(t)} \left\  \mathbf{A}^{(k)}(i_k^{(j_f)}, :) \right\ _2^2} - \frac{1}{ \mathcal{F}_n } \left\  \nabla_{\mathbf{A}^{(\xi(t))}} f(\boldsymbol{\theta}(t)) \right\ _F^2$

<sup>25</sup>Ping Ma, Michael Mahoney, and Bin Yu. "A statistical Perspective on algorithmic leveraging". In: *J. Mach. Learn. Res.* 16:27 (2015), pp. 867–911. doi: 10.1002/caml.1335. 

## COMPLEXITIES

**Table:** Comparison of leading order computational complexities ( $m > R$  and allow  $|\mathcal{F}_n| \ll R$ )

Method	Complexity without initialization
CP-ALS	$\mathcal{O}(\#it \cdot NRI^N)$
CPRAND-MIX	$\mathcal{O}(I^N \log(I^N) + \#it \cdot mNIR)$
CP-ALS-LEV	$\mathcal{O}(NIR^2 + \#it \cdot mNIR)$
BrasCPD	$\mathcal{O}(\#it \cdot IR \mathcal{F}_n )$
AdaCPD	$\mathcal{O}(\#it \cdot IR \mathcal{F}_n )$
EBrawCPD	$\mathcal{O}(IR + \#it \cdot IR \mathcal{F}_n )$
LBrawCPD	$\mathcal{O}(NIR^2 + \#it \cdot IR \mathcal{F}_n )$
EAdawCPD	$\mathcal{O}(IR + \#it \cdot IR \mathcal{F}_n )$
LAdawCPD	$\mathcal{O}(NIR^2 + \#it \cdot IR \mathcal{F}_n )$