

Canonical Polyadic Decomposition

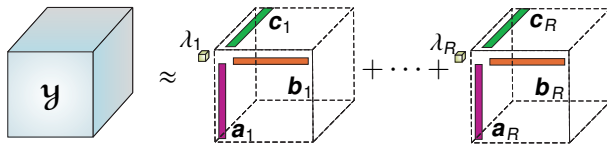
Center for Computational and Data-Intensive Science and Engineering (CDISE)
Skolkovo Institute of Science and Technology (SKOLTECH), Moscow, Russia

February 2, 2022

Part I

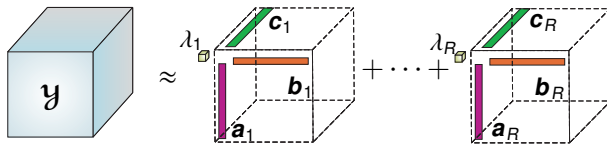
Research Summary: Challenging Problems in CANDECOP/PARAFAC

CANDECOMP/PARAFAC (CPD)



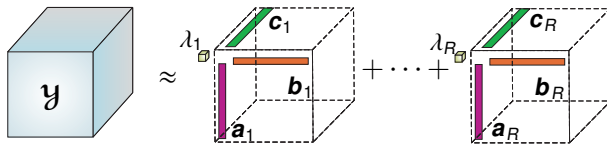
- ▶ Application of Tensors in a raw format may be intractable, as the required storage memory and a number of operations grow exponentially with the tensor order.
- ▶ To tackle this issue, represent higher-order tensors through multi-way operations over their latent components.
- ▶ Canonical Polyadic Decomposition (CPD): an extension of matrix factorization to extract rank-1 tensor patterns from multiway data.

CANDECOMP/PARAFAC (CPD)



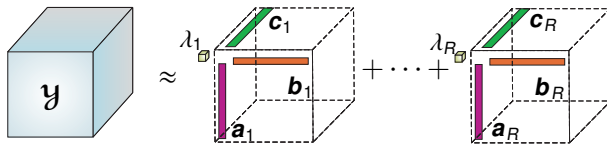
- ▶ 1927 first studied Hitchcock (1927)
- ▶ 1969 determined the complexity of matrix multiplication Strassen (1969).
- ▶ 1972 later known as parallel factor analysis (PARAFAC), a tool for chemometric analysis popularized by Harshman (1972), Carroll and Chang (1970) and Kruskal (1977)
- ▶ Since the 1990's, the CPD has attracted attention in signal processing Sidiropoulos et al. (2000); Yeredor (2000); Comon and Rajih (2006).

CANDECOMP/PARAFAC (CPD)



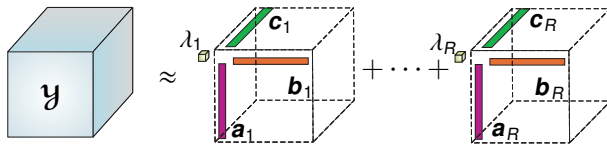
- ▶ 2000s New applications of CPD in feature extraction, data reconstruction, image completion and various tracking scenarios Nion and Sidiropoulos (2009); Phan et al. (2015).
- ▶ 2014-CPD of kernel tensors in convolutional can accelerate the inference of convolutional neural networks Jaderberg et al. (2014).
- ▶ In cyber security, the tensor can represent network traffic volume: time \times senderIP \times receiverIP \times port.

CANDECOMP/PARAFAC (CPD)



- ▶ Over its long history, the model has been extended with various constraints with many efficient algorithms
- ▶ Despite the great successes, there are still many challenging problems in the CP tensor representation/decomposition.

CANDECOMP/PARAFAC (CPD)



- ▶ Computation for large volume and high order tensors, e.g., those of order $N = 10$ or 20
- ▶ Degeneracy

Challenging Problems in CPD

CPD can be achieved by minimizing the Frobenius norm of the error

$$\min \quad \|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2$$

where $\hat{\mathcal{Y}}$ is a low-rank approximation of \mathcal{Y} in the CP format

$$\hat{y}_{i_1, i_2, \dots, i_N} = \sum_{r=1}^R \lambda_r a_{i_1, r}^{(1)} a_{i_2, r}^{(2)} \cdots a_{i_N, r}^{(N)}$$

The tensor rank R can exceed the tensor dimension $R \geq I_n$.

Challenging Problems in CPD I

Prob.1: Nonlinear optimization

The optimisation is nonlinear with respect to all the parameters, but linear in parameters in each $\mathbf{A}^{(n)}$, or parameters in non-overlapping partitions of different factor matrices

$$D(\mathbf{y} \parallel \hat{\mathbf{y}}) = \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|_F^2 = \frac{1}{2} \|\mathbf{Y}_{(n)} - \mathbf{A}^{(n)} \mathbf{Z}_n^T\|_F^2$$

where $\mathbf{Z}_n = \mathbf{A}^{(N)} \odot \dots \odot \mathbf{A}^{(n+1)} \odot \mathbf{A}^{(n-1)} \odot \dots \odot \mathbf{A}^{(1)}$, \odot Khatri-Rao product of $\mathbf{A}^{(m)}$.

Prob.2: High computational cost $O(N^N R^3)$ due to the gradient

Computational costs of algorithms grow exponentially with the tensor order and tensor rank, due to computation of $\mathbf{Y}_{(n)} \mathbf{Z}_n$.

→ Most algorithms are applicable to only tensors of low order, which never exceeds three.

Challenging Problems in CPD II

Prob.3: Degeneracy

Often occurs in some difficult scenarios, e.g.,

- ▶ when the rank exceeds the tensor dimension ($R > I_n$)
- ▶ when the loading components, $\mathbf{a}_r^{(n)}$ and $\mathbf{a}_s^{(n)}$, are highly collinear
- ▶ or when CPD does not provide an optimal solution

Prob.4: Second order optimisation is inadequate

Because inverse of large Hessian matrix of size $NIR \times NIR$.

Prob.5: CPD with additional constraints

nonnegativity, orthogonality, sparsity, smoothness constraints

Prob.6: A measure of performance

A scalar index to assess performance of a decomposition or of an algorithm

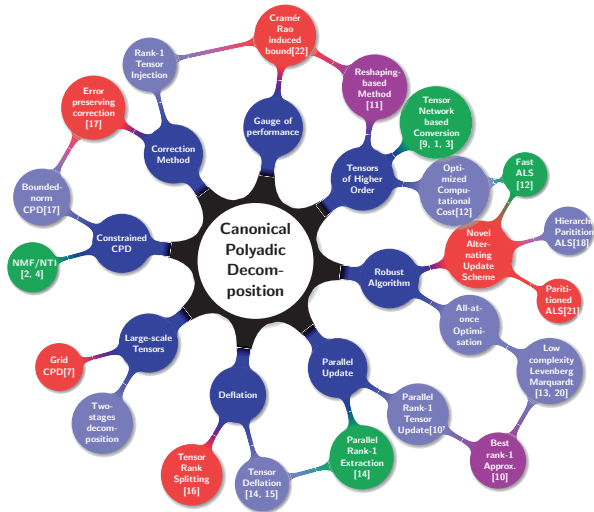
Prob.7: Decomposition of massive tensor data

Large-scale tensor cannot be handled in a computer.

Prob.8: Deflation (the most challenging problem)

Sequential extraction of rank-1 tensors of \mathcal{Y} ?

- ▶ In matrix factorisations, e.g., SVD, rank-1 matrix components can be estimated via a deflation process.
- ▶ Unfortunately, this property does not hold for tensors because subtracting the best rank-1 tensor may increase the tensor rank (Stegeman and Comon 2010).



We deal with the main challenging problems for CPD and learn state-of-the-art algorithms for the problems.

Part II

Algorithms for CANDECOMP/
PARAFAC

Trilinear and Direct Trilinear Decompositions I

Generalized Rank Annihilation Method (GRAM) Booksh and Kowalski (1994); Booksh et al. (1994); Faber et al. (1994).

- ▶ Given \mathcal{Y} a tensor of size $J \times J \times 2$ with only two frontal slices: \mathbf{Y}_1 and \mathbf{Y}_2 .

CPD of \mathcal{Y} by $\mathbf{A} \in \mathbb{R}^{J \times J}$, $\mathbf{B} \in \mathbb{R}^{J \times J}$ and $\mathbf{C} = [\underline{\mathbf{c}}_1^T, \underline{\mathbf{c}}_2^T]^T \in \mathbb{R}^{2 \times J}$

$$\mathbf{Y}_1 = \mathbf{A} \mathbf{D}_1 \mathbf{B}^T, \quad \mathbf{Y}_2 = \mathbf{A} \mathbf{D}_2 \mathbf{B}^T, \quad (1)$$

where $\mathbf{D}_1 = \text{diag}\{\underline{\mathbf{c}}_1\}$, $\mathbf{D}_2 = \text{diag}\{\underline{\mathbf{c}}_2\}$.

Trilinear and Direct Trilinear Decompositions II

- ▶ Generalized Eigenvalue problem of two square matrices \mathbf{Y}_1 and \mathbf{Y}_2

$$\begin{aligned}\mathbf{Y}_1 (\mathbf{B}^T)^{-1} &= \mathbf{A} \mathbf{D}_1 \mathbf{B}^T (\mathbf{B}^T)^{-1} = \mathbf{A} \mathbf{D}_2 \mathbf{D}_2^{-1} \mathbf{D}_1 = \mathbf{A} \mathbf{D}_2 \mathbf{B}^T (\mathbf{B}^T)^{-1} \mathbf{D} \\ &= \mathbf{Y}_2 (\mathbf{B}^T)^{-1} \mathbf{D},\end{aligned}\tag{2}$$

where $\mathbf{D} = \mathbf{D}_2^{-1} \mathbf{D}_1 = \text{diag}\{\underline{\mathbf{c}}_1 \oslash \underline{\mathbf{c}}_2\}$ is a diagonal matrix of generalized eigenvalues, and columns of $(\mathbf{B}^T)^{-1}$ are corresponding eigenvectors.

- ▶ Due to the scaling ambiguity, we can set the second row vector $\underline{\mathbf{c}}_2 = \mathbf{1}^T$, hence, the first row vector $\underline{\mathbf{c}}_1$ represents generalized eigenvalues:

$$\underline{\mathbf{c}}_1 = \text{diag}\{\mathbf{D}\}^T.\tag{3}$$

Trilinear and Direct Trilinear Decompositions III

Algorithm 1: GRAM

Input: $\underline{\mathbf{Y}}$: input data of size $J \times J \times 2$, J : number of basis components

Output: $\mathbf{A} \in \mathbb{R}^{J \times J}$, $\mathbf{B} \in \mathbb{R}^{J \times J}$ and $\mathbf{C} \in \mathbb{R}^{2 \times J}$ such that
 $\underline{\mathbf{Y}} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$

```
1 begin
2    $[\Phi, \mathbf{D}] = \text{eig}(\mathbf{Y}_1, \mathbf{Y}_2);$ 
3    $\mathbf{B} = (\Phi^T)^{-1};$ 
4    $\mathbf{A} = \mathbf{Y}_2 \Phi;$ 
5    $\mathbf{C} = [\text{diag}\{\mathbf{D}\}, \mathbf{1}]^T;$ 
6 end
```

Trilinear and Direct Trilinear Decompositions IV

Direct Tri-Linear Decomposition (DTLD) (Sanchez and Kowalski)

Nonsquare matrices, \mathbf{Y}_1 and \mathbf{Y}_2 , can be transformed to matrices \mathbf{T}_q of size $J \times J$ by projecting on the two orthonormal matrices $\mathbf{U} \in \mathbb{R}^{I \times J}$ and $\mathbf{V} \in \mathbb{R}^{T \times J}$, usually via a singular value decomposition

$$\begin{aligned} [\mathbf{U}, \mathbf{S}, \mathbf{V}] &= \text{svd}(\mathbf{Y}_1 + \mathbf{Y}_2, J) \\ \mathbf{T}_q &= \mathbf{U}^T \mathbf{Y}_q \mathbf{V}, \quad (q = 1, 2). \end{aligned}$$

Applying the GRAM algorithm to the tensor \mathcal{T} , factors \mathbf{A} , \mathbf{B} and \mathbf{C} can be found as

$$[\Phi, \mathbf{D}] = \text{eig}(\mathbf{T}_1, \mathbf{T}_2) \tag{4}$$

$$\mathbf{A} = \mathbf{U} \mathbf{T}_2 \Phi \tag{5}$$

$$\mathbf{B} = \mathbf{V} (\Phi^T)^{-1} \tag{6}$$

$$\mathbf{C} = [\text{diag}\{\mathbf{D}\}, \mathbf{1}]^T. \tag{7}$$

Trilinear and Direct Trilinear Decompositions V

The GRAM method is restricted to the tensor having only two slices.

For tensor with more than 2 slides, $\underline{\mathbf{Y}} \in \mathbb{R}^{I \times T \times Q}$ using the Tucker model with a core tensor $\underline{\mathbf{G}} \in \mathbb{R}^{J \times J \times 2}$

$$\underline{\mathbf{Y}} = \llbracket \underline{\mathbf{G}}; \mathbf{A}_t, \mathbf{B}_t, \mathbf{C}_t \rrbracket, \quad (8)$$

where $\mathbf{A}_t \in \mathbb{R}^{I \times J}$, $\mathbf{B}_t \in \mathbb{R}^{T \times J}$, and $\mathbf{C}_t \in \mathbb{R}^{Q \times 2}$ are factors obtained by using the ALS Tucker algorithm.

The core tensor $\underline{\mathbf{G}}$ is then factorized using the GRAM algorithm

$$\underline{\mathbf{G}} = \llbracket \mathbf{A}_f, \mathbf{B}_f, \mathbf{C}_f \rrbracket. \quad (9)$$

The final model is given as $\underline{\mathbf{Y}} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ where

$$\mathbf{A} = \mathbf{A}_t \mathbf{A}_f \quad (10)$$

$$\mathbf{B} = \mathbf{B}_t \mathbf{B}_f \quad (11)$$

$$\mathbf{C} = \mathbf{C}_t \mathbf{C}_f. \quad (12)$$

Trilinear and Direct Trilinear Decompositions VI

Algorithm 2: DTLD

Input: $\underline{\mathbf{Y}}$: input data of size $I \times T \times Q$, J : number of basis components

Output: $\mathbf{A} \in \mathbb{R}^{I \times J}$, $\mathbf{B} \in \mathbb{R}^{T \times J}$ and $\mathbf{C} \in \mathbb{R}^{Q \times J}$ such that
 $\underline{\mathbf{Y}} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$

```
1 begin
2    $\underline{\mathbf{Y}} = \llbracket \underline{\mathbf{G}}; \mathbf{A}_t, \mathbf{B}_t, \mathbf{C}_t \rrbracket$            /* Tucker ALS */;
3    $[\Phi, \mathbf{D}] = \text{eig}(\mathbf{G}_1, \mathbf{G}_2)$ ;
4    $\mathbf{A} = \mathbf{A}_t \mathbf{G}_2 \Phi$ ;
5    $\mathbf{B} = \mathbf{B}_t (\Phi^T)^{-1}$ ;
6    $\mathbf{C} = \mathbf{C}_t [\text{diag}\{\mathbf{D}\}, \mathbf{1}]^T$            /* or  $\mathbf{C} = \mathbf{Y}_{(3)} (\mathbf{B} \odot \mathbf{A})$  */;
7 end
```

DTLD

- ▶ Efficient initialization for inexact case or constrained CPD
- ▶ Work only for order-3 tensor
- ▶ rank- R must not exceed tensor dimensions

Generalized rank annihilation method for Higher Order CPD

Is there a closed form solution for CPD of tensors of higher order?

How???

The higher order tensor is first unfolded to be an order-3 tensor of size $I \times J \times K$ so that the two largest dimensions are greater than or equal to rank R , then factorized using DTLD.

Generalized rank annihilation method for Higher Order CPD

Is there a closed form solution for CPD of tensors of higher order?

How???

The higher order tensor is first unfolded to be an order-3 tensor of size $I \times J \times K$ so that the two largest dimensions are greater than or equal to rank R , then factorized using DTLD.

Generalized DTLD method for Higher Order CPD

Example

Given a rank-3 tensor, \mathcal{Y} , of size $4 \times 4 \times 5 \times 6$. Find its CPD

$\mathcal{Y} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D} \rrbracket$ using DTLD.

- ▶ Reshaping \mathcal{Y} to an order-3 tensor, $\mathcal{Y}_{(3,4)}$, of size $4 \times 4 \times 30$ which is still of rank-3.
- ▶ Apply DTLD to $\mathcal{Y}_{(3,4)}$ and obtain a rank-3 CPD
 $\mathcal{Y}_{(3,4)} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{F} \rrbracket$
- ▶ Since $\mathbf{F} = \mathbf{D} \odot \mathbf{C}$ has Khatri-Rao product,
 $\mathbf{c}_r \mathbf{d}_r^T$ is the best rank-1 of the folding matrix, \mathbf{F}_r of size 5×6
from $\mathbf{f}_r = \text{vec}(\mathbf{F}_r)$

$$\mathbf{F}_r \approx \mathbf{c}_r \mathbf{d}_r^T$$

using truncated SVD

Generalized DTLD method for Higher Order CPD

Example

Given a rank-3 tensor, \mathcal{Y} , of size $4 \times 4 \times 5 \times 6 \times 7$. Find its CPD $\mathcal{Y} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{E} \rrbracket$ using DTLD.

- ▶ Reshaping \mathcal{Y} to an order-3 tensor, $\mathcal{Y}_{(2,3),(4,5)}$, of size $4 \times 20 \times 42$ which is still of rank-3.
- ▶ Apply DTLD to $\mathcal{Y}_{(3,4)}$ and obtain a rank-3 CPD $\mathcal{Y}_{(2,3),(4,5)} = \llbracket \mathbf{A}, \mathbf{U}, \mathbf{V} \rrbracket$
- ▶ It can be shown that $\mathbf{U} = \mathbf{B} \odot \mathbf{C}$, $\mathbf{V} = \mathbf{E} \odot \mathbf{D}$, hence $\mathbf{b}_r \mathbf{c}_r^T = \text{reshape}(\mathbf{u}_r, [45])$ is the best rank-1 of the matrix of size 4×5 folded from \mathbf{u}_r
 $\mathbf{d}_r \mathbf{e}_r^T = \text{reshape}(\mathbf{v}_r, [67])$ is the best rank-1 of the matrix of size 6×7 folded from \mathbf{v}_r

Alternating Least Squares Algorithms

Cost function

$$D(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2 = \frac{1}{2}\|\mathbf{Y}_{(n)} - \mathbf{A}^{(n)}\{\mathbf{A}\}^{\odot_{-n}^T}\|_F^2.$$

1. Initialize all $\mathbf{A}^{(n)}$ randomly or by using SVD Berry et al. (2007), or by selected fibers from the data \mathcal{Y} .
2. Estimate $\mathbf{A}^{(n)}$ from the approximation by solving

$$\min_{\mathbf{A}^{(n)}} D_F(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2, \text{ with fixed } \mathbf{A}^{(m)}, m \neq n.$$

3. For nonnegative components, set all negative elements of $\mathbf{A}^{(n)}$ to zero or a small positive value ε .
4. Estimate other factors until convergence.

Alternating Least Squares Algorithms

Cost function

$$D(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2 = \frac{1}{2}\|\mathbf{Y}_{(n)} - \mathbf{A}^{(n)}\{\mathbf{A}\}^{\odot_{-n}^T}\|_F^2.$$

1. Initialize all $\mathbf{A}^{(n)}$ randomly or by using SVD Berry et al. (2007), or by selected fibers from the data \mathcal{Y} .
2. Estimate $\mathbf{A}^{(n)}$ from the approximation by solving

$$\min_{\mathbf{A}^{(n)}} D_F(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2, \text{ with fixed } \mathbf{A}^{(m)}, m \neq n.$$

3. For nonnegative components, set all negative elements of $\mathbf{A}^{(n)}$ to zero or a small positive value ε .
4. Estimate other factors until convergence.

Alternating Least Squares Algorithms

Cost function

$$D(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2 = \frac{1}{2}\|\mathbf{Y}_{(n)} - \mathbf{A}^{(n)}\{\mathbf{A}\}^{\odot_{-n}^T}\|_F^2.$$

1. Initialize all $\mathbf{A}^{(n)}$ randomly or by using SVD Berry et al. (2007), or by selected fibers from the data \mathcal{Y} .
2. Estimate $\mathbf{A}^{(n)}$ from the approximation by solving

$$\min_{\mathbf{A}^{(n)}} D_F(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2, \text{ with fixed } \mathbf{A}^{(m)}, m \neq n.$$

3. For nonnegative components, set all negative elements of $\mathbf{A}^{(n)}$ to zero or a small positive value ε .
4. Estimate other factors until convergence.

Alternating Least Squares Algorithms

Cost function

$$D(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2 = \frac{1}{2}\|\mathbf{Y}_{(n)} - \mathbf{A}^{(n)}\{\mathbf{A}\}^{\odot_{-n}^T}\|_F^2.$$

1. Initialize all $\mathbf{A}^{(n)}$ randomly or by using SVD Berry et al. (2007), or by selected fibers from the data \mathcal{Y} .
2. Estimate $\mathbf{A}^{(n)}$ from the approximation by solving

$$\min_{\mathbf{A}^{(n)}} D_F(\mathcal{Y}||\hat{\mathcal{Y}}) = \frac{1}{2}\|\mathcal{Y} - \hat{\mathcal{Y}}\|_F^2, \text{ with fixed } \mathbf{A}^{(m)}, m \neq n.$$

3. For nonnegative components, set all negative elements of $\mathbf{A}^{(n)}$ to zero or a small positive value ε .
4. Estimate other factors until convergence.

Gradient of the cost function w.r.t $\mathbf{A}^{(n)}$

$$\mathbf{G}^{(n)} = \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) - \mathbf{A}^{(n)} \left(\bigotimes_{k \neq n} \mathbf{A}^{(k)T} \mathbf{A}^{(k)} \right) \in \mathbb{R}^{I_n \times R},$$

ALS Algorithm for CP

$$\mathbf{A}^{(n)} \leftarrow \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) \left(\left\{ \mathbf{A}^T \mathbf{A} \right\}^{\otimes -n} \right)^{\dagger}$$

ε is a small constant (typically, 10^{-16}). \mathbf{A}^{\dagger} is the Moore-Penrose inverse of \mathbf{A} ,

Gradient of the cost function w.r.t $\mathbf{A}^{(n)}$

$$\mathbf{G}^{(n)} = \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) - \mathbf{A}^{(n)} \left(\bigotimes_{k \neq n} \mathbf{A}^{(k)T} \mathbf{A}^{(k)} \right) \in \mathbb{R}^{I_n \times R},$$

ALS Algorithm for CP

$$\mathbf{A}^{(n)} \leftarrow \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) \left(\left\{ \mathbf{A}^T \mathbf{A} \right\}^{\otimes -n} \right)^{\dagger}$$

ε is a small constant (typically, 10^{-16}). \mathbf{A}^{\dagger} is the Moore-Penrose inverse of \mathbf{A} ,

Pros: simple, work well for general data.

Cons: high cost due to computing $\{\mathbf{A}\}^{\otimes -n}$, $\mathbf{Y}_{(n)} \{\mathbf{A}\}^{\otimes -n}$
not guaranteed to converge to a global minimum or even a stationary point, but only to a solution where the cost functions cease to decrease Kolda and Bader (2009); Berry et al. (2007).
relatively slow for nearly collinear data Langville et al. (2006).

Gradient of the cost function w.r.t $\mathbf{A}^{(n)}$

$$\mathbf{G}^{(n)} = \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) - \mathbf{A}^{(n)} \left(\bigotimes_{k \neq n} \mathbf{A}^{(k)T} \mathbf{A}^{(k)} \right) \in \mathbb{R}^{I_n \times R},$$

ALS Algorithm for CP

$$\mathbf{A}^{(n)} \leftarrow \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) \left(\left\{ \mathbf{A}^{(k)T} \mathbf{A}^{(k)} \right\}^{\otimes -n} \right)^{\dagger}$$

ε is a small constant (typically, 10^{-16}). \mathbf{A}^{\dagger} is the Moore-Penrose inverse of \mathbf{A} ,

Pros: simple, work well for general data.

Cons: high cost due to computing $\{\mathbf{A}^{(k)T} \mathbf{A}^{(k)}\}^{\otimes -n}$, $\mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right)$
not guaranteed to converge to a global minimum or even a stationary point, but only to a solution where the cost functions cease to decrease Kolda and Bader (2009); Berry et al. (2007).
relatively slow for nearly collinear data Langville et al. (2006).

Gradient of the cost function w.r.t $\mathbf{A}^{(n)}$

$$\mathbf{G}^{(n)} = \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) - \mathbf{A}^{(n)} \left(\bigotimes_{k \neq n} \mathbf{A}^{(k)T} \mathbf{A}^{(k)} \right) \in \mathbb{R}^{I_n \times R},$$

ALS Algorithm for CP

$$\mathbf{A}^{(n)} \leftarrow \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right) \left(\left\{ \mathbf{A}^T \mathbf{A} \right\}^{\otimes -n} \right)^{\dagger}$$

ε is a small constant (typically, 10^{-16}). \mathbf{A}^{\dagger} is the Moore-Penrose inverse of \mathbf{A} ,

Pros: simple, work well for general data.

Cons: high cost due to computing $\{\mathbf{A}\}^{\odot -n}$, $\mathbf{Y}_{(n)} \{\mathbf{A}\}^{\odot -n}$
not guaranteed to converge to a global minimum or even a stationary point, but only to a solution where the cost functions cease to decrease Kolda and Bader (2009); Berry et al. (2007).
relatively slow for nearly collinear data Langville et al. (2006).

Fast Computation of CP gradient and FastALS I

CP gradients - the most expensive steps in CP algorithms.

Computation of CP gradients $\mathbf{G}^{(n)} = \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right)$ costs $O(NRl_1 l_2 \dots l_N)$.

Remark: Order of dimensions

For the least cost to compute $\mathbf{G}^{(n)}$, the dimensions of \mathcal{Y} should be in the ascending or descending order.

Without loss of generality, we assume $l_1 \leq l_2 \leq \dots \leq l_N$.

Fast Computation of CP gradient and FastALS II

Each column of $\mathbf{G}^{(n)}$ is efficiently computed from “left-to-right”

$$\mathbf{g}_r^{(n)} = \mathbf{Y}_{(n)} \left(\bigotimes_{k \neq n} \mathbf{a}_r^{(k)} \right) = \left(\mathbf{y} \bar{\times}_{k=1}^{n-1} \mathbf{a}_r^{(k)} \right) \bar{\times}_{l=2}^{N-n+1} \mathbf{a}_r^{(l+n-1)},$$

or from “right-to-left”

$$\mathbf{Y}_{(n)} \left(\bigotimes_{k \neq n} \mathbf{a}_r^{(k)} \right) = \left(\mathbf{y} \bar{\times}_{l=n+1}^N \mathbf{a}_r^{(l)} \right) \bar{\times}_{k=1}^{n-1} \mathbf{a}_r^{(k)}.$$

Let $\mathcal{L}^{(r,n)} \triangleq \mathbf{y} \bar{\times}_{k=1}^{n-1} \mathbf{a}_r^{(k)}$ and $\mathcal{R}^{(r,n)} \triangleq \mathbf{y} \bar{\times}_{k=n+1}^N \mathbf{a}_r^{(k)}$. Then

Recursive computation

$$\begin{aligned} \mathcal{L}^{(r,n)} &= \mathbf{y} \bar{\times}_{k=1}^{n-1} \mathbf{a}_r^{(k)} = \mathcal{L}^{(r,n-1)} \bar{\times}_1 \mathbf{a}_r^{(n-1)}, \\ \mathcal{R}^{(r,n)} &= \mathbf{y} \bar{\times}_{k=n+1}^N \mathbf{a}_r^{(k)} = \mathcal{R}^{(r,n+1)} \bar{\times}_{n+1} \mathbf{a}_r^{(n+1)}, \end{aligned}$$

Fast Computation of CP gradient and FastALS III

Note that $\mathcal{L}^{(r,n)}$ and $\mathcal{R}^{(r,n)}$ are computed without tensor permutation

$$\begin{aligned}\mathcal{L}^{(r,n)} &= \text{reshape} \left(\mathbf{Y}_{(1:n-1)}^T \left(\bigotimes_{k=1}^{n-1} \mathbf{a}_r^{(k)} \right), [l_n, \dots, l_N] \right), \\ \mathcal{R}^{(r,n)} &= \text{reshape} \left(\mathbf{Y}_{(1:n)} \left(\bigotimes_{k=n+1}^N \mathbf{a}_r^{(k)} \right), [l_1, \dots, l_n] \right).\end{aligned}$$

Fast CP-gradients over all modes

Fast Computation of CP gradient and FastALS IV

- ▶ Compute $J_n = l_1 l_2 \cdots l_n$ and $K_n = l_{n+1} l_{n+2} \cdots l_N$, and find n^* such that

$$n^* = \max\{n : J_n \leq K_{n-1}\}$$

- ▶ First compute gradient $\mathbf{G}^{(n^*)}$ (or $\mathbf{G}^{(n^*+1)}$)
- ▶ Sequentially compute other gradients $\mathbf{G}^{(n)}$ in the following order $n = n^* - 1, n^* - 2, \dots, 1, n^* + 1, n^* + 2, \dots, N$.

Computational cost: $R(2J_N + 2J_{n^*} + 2K_{n^*})$

Fast CP gradients over all modes

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, N factor matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

Output: $\mathbf{G}^{(n)} = \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right)$: $I_n \times R$, $n = 1, 2, \dots, N$

```

1  begin
3       $n^* = \max\{n : J_n \leq K_{n-1}\}$  where  $J_n = I_1 I_2 \dots I_n, K_n = I_{n+1}$ 
4      for  $n = n^*, n^* - 1, \dots, 1, n^* + 1, n^* + 2, \dots, N$  do
5          if  $(n == n^*)$  then
7               $\mathbf{R}_{(n)}^{(n)} = \text{reshape}(\mathcal{Y}, [J_n, K_n]) \left( \bigodot_{k=n+1}^N \mathbf{A}^{(k)} \right)$ ;
9               $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{R}^{(n)}, \{\mathbf{A}\}, \text{'right'})$ 
10             else if  $(n == n^* + 1)$  then
12                  $\mathbf{L}_{(1)}^{(n)} = \left( \bigodot_{k=1}^{n-1} \mathbf{A}^{(k)} \right)^T \text{reshape}(\mathcal{Y}, [J_{n-1}, K_{n-1}])$ 
14                  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{L}^{(n)}, \{\mathbf{A}\}, \text{'left'})$ 
15             else if  $(n < n^*)$  then
17                 for  $r = 1, 2, \dots, R$  do  $\text{vec}(\mathcal{R}^{(r,n)}) \leftarrow \text{reshape}(\mathcal{R}^{(r,n+1)}, [J_n, I_{n+1}]) \mathbf{a}_r^{(n+1)}$ ;
19                  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{R}^{(n)}, \{\mathbf{A}\}, \text{'right'})$ 
20             else
22                 for  $r = 1, 2, \dots, R$  do  $\text{vec}(\mathcal{L}^{(r,n)}) \leftarrow \text{reshape}(\mathcal{L}^{(r,n-1)}, [I_{n-1}, K_{n-1}])^T \mathbf{a}_r^{(n-1)}$ ;
24                  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{L}^{(n)}, \{\mathbf{A}\}, \text{'left'})$ 
25             end
26         end
27     end

```

```

1  function  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{Z}^{(n)}, \{\mathbf{A}\}, \text{side})$ 
2  for  $r = 1, 2, \dots, R$  do
3      switch side do
4          case 'right': do  $g_r^{(n)} =$ 
5                           $\text{reshape}(\mathcal{Z}^{(r,n)}, [J_{n-1}, I_n])^T \left( \bigotimes_{k=1}^{n-1} \mathbf{a}_r^{(k)} \right)$ 
6                          end
7          case 'left': do  $g_r^{(n)} =$ 
5                           $\text{reshape}(\mathcal{Z}^{(r,n)}, [I_n, K_n]) \left( \bigotimes_{k=n+1}^N \mathbf{a}_r^{(k)} \right)$ 
6                          end
7      end

```

Fast CP gradients over all modes

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, N factor matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

Output: $\mathbf{G}^{(n)} = \mathbf{Y}_{(n)} \left(\bigodot_{k \neq n} \mathbf{A}^{(k)} \right)$: $I_n \times R$, $n = 1, 2, \dots, N$

```

1  begin
3       $n^* = \max\{n : J_n \leq K_{n-1}\}$  where  $J_n = I_1 I_2 \dots I_n, K_n = I_{n+1}$ 
4      for  $n = n^*, n^* - 1, \dots, 1, n^* + 1, n^* + 2, \dots, N$  do
5          if  $(n == n^*)$  then
7               $\mathbf{R}_{(n)}^{(n)} = \text{reshape}(\mathcal{Y}, [J_n, K_n]) \left( \bigodot_{k=n+1}^N \mathbf{A}^{(k)} \right)$ ;
9               $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{R}^{(n)}, \{\mathbf{A}\}, \text{'right'})$ 
10             else if  $(n == n^* + 1)$  then
12                  $\mathbf{L}_{(1)}^{(n)} = \left( \bigodot_{k=1}^{n-1} \mathbf{A}^{(k)} \right)^T \text{reshape}(\mathcal{Y}, [J_{n-1}, K_{n-1}])$ 
14                  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{L}^{(n)}, \{\mathbf{A}\}, \text{'left'})$ 
15             else if  $(n < n^*)$  then
17                 for  $r = 1, 2, \dots, R$  do  $\text{vec}(\mathcal{R}^{(r,n)}) \leftarrow \text{reshape}(\mathcal{R}^{(r,n+1)}, [J_n, I_{n+1}]) \mathbf{a}_r^{(n+1)}$ ;
19                  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{R}^{(n)}, \{\mathbf{A}\}, \text{'right'})$ 
20             else
22                 for  $r = 1, 2, \dots, R$  do  $\text{vec}(\mathcal{L}^{(r,n)}) \leftarrow \text{reshape}(\mathcal{L}^{(r,n-1)}, [I_{n-1}, K_{n-1}])^T \mathbf{a}_r^{(n-1)}$ ;
24                  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{L}^{(n)}, \{\mathbf{A}\}, \text{'left'})$ 
25             end
26         end
27     end

```

```

1  function  $\mathbf{G}^{(n)} = \text{cp\_gradient}(\mathcal{Z}^{(n)}, \{\mathbf{A}\}, \text{side})$ 
2  for  $r = 1, 2, \dots, R$  do
3      switch side do
4          case 'right': do  $\mathbf{g}_r^{(n)} =$ 
5                           $\text{reshape}(\mathcal{Z}^{(r,n)}, [J_{n-1}, I_n])^T \left( \bigotimes_{k=1}^{n-1} \mathbf{a}_r^{(k)} \right)$ 
6          case 'left': do  $\mathbf{g}_r^{(n)} =$ 
5                           $\text{reshape}(\mathcal{Z}^{(r,n)}, [I_n, K_n]) \left( \bigotimes_{k=n+1}^N \mathbf{a}_r^{(k)} \right)$ 
6      end
7  end

```

FastALS vs ALS for order-4 Tensors

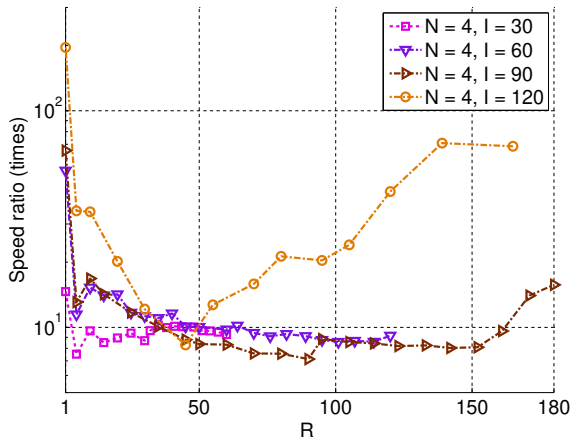
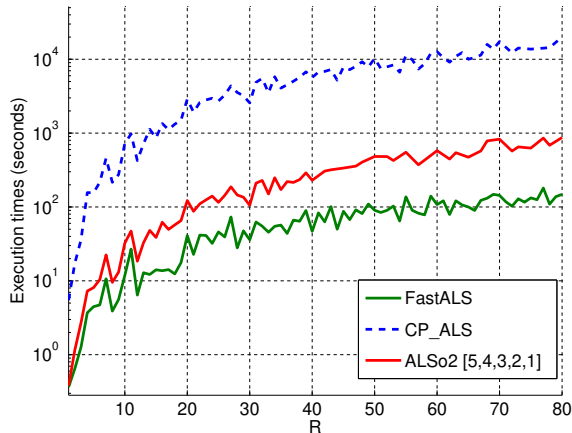


Figure: Speed-up ratios per iteration for the FastALS algorithm in comparison with the standard CP_ALS algorithm for factorizations of order-4 tensors with sizes $l_n = l$ for all n .

FastALS vs ALS for order-5 Tensors



ALS02 Eigenvector (2012)

Figure: Execution times of CP algorithms in factorization of an order-5 EEG MI tensor of size 2 dictionaries \times 23 frequency bins (8-30 Hz) \times 50 time frames \times 62 channels \times 200 trials.

Bro's Line search (LS)

- Predict the factors $\mathbf{A}^{(n)}$ from their previous estimates

$$\mathbf{A}_t^{(n)} = \mathbf{A}_{t-1}^{(n)} + \eta \left(\mathbf{A}_{t-1}^{(n)} - \mathbf{A}_{t-2}^{(n)} \right), \quad n = 1, 2, \dots, N, \quad (13)$$

where t denotes the iteration index, and η is stepsize.

- A heuristic stepsize η is iteratively searched to improve the fit

$$\eta \leftarrow \alpha \eta, \quad (14)$$

where α is a suitably chosen scalar constant.

- The LS method can find a suitable step value but not an optimal step.

Line Search for order-3 CPD

$$\begin{aligned}\hat{\mathcal{Y}} &= \llbracket \mathbf{A}_{t+1}, \mathbf{B}_{t+1}, \mathbf{C}_{t+1} \rrbracket \\ &= \llbracket \mathbf{A}_t + \eta \Delta \mathbf{A}_t, \mathbf{B}_t + \eta \Delta \mathbf{B}_t, \mathbf{C}_t + \eta \Delta \mathbf{C}_t \rrbracket \\ &= \llbracket \mathbf{A}_t, \mathbf{B}_t, \mathbf{C}_t \rrbracket \\ &\quad + \eta (\llbracket \Delta \mathbf{A}_t, \mathbf{B}_t, \mathbf{C}_t \rrbracket + \llbracket \mathbf{A}_t, \Delta \mathbf{B}_t, \mathbf{C}_t \rrbracket + \llbracket \mathbf{A}_t, \mathbf{B}_t, \Delta \mathbf{C}_t \rrbracket) \\ &\quad + \eta^2 (\llbracket \Delta \mathbf{A}_t, \Delta \mathbf{B}_t, \mathbf{C}_t \rrbracket + \llbracket \Delta \mathbf{A}_t, \mathbf{B}_t, \Delta \mathbf{C}_t \rrbracket + \llbracket \mathbf{A}_t, \Delta \mathbf{B}_t, \Delta \mathbf{C}_t \rrbracket) \\ &\quad + \eta^3 \llbracket \Delta \mathbf{A}_t, \Delta \mathbf{B}_t, \Delta \mathbf{C}_t \rrbracket \\ &= \llbracket \mathbf{A}_t, \mathbf{B}_t, \mathbf{C}_t \rrbracket \\ &\quad + \eta \llbracket [\Delta \mathbf{A}_t, \mathbf{A}_t, \mathbf{A}_t], [\mathbf{B}_t, \Delta \mathbf{B}_t, \mathbf{B}_t], [\mathbf{C}_t, \mathbf{C}_t, \Delta \mathbf{C}_t] \rrbracket \\ &\quad + \eta^2 \llbracket [\Delta \mathbf{A}_t, \Delta \mathbf{A}_t, \mathbf{A}_t], [\Delta \mathbf{B}_t, \mathbf{B}_t, \Delta \mathbf{B}_t], [\mathbf{C}_t, \Delta \mathbf{C}_t, \Delta \mathbf{C}_t] \rrbracket \\ &\quad + \eta^3 \llbracket \Delta \mathbf{A}_t, \Delta \mathbf{B}_t, \Delta \mathbf{C}_t \rrbracket\end{aligned}$$

Line search III

Exact line search

Enhanced or Exact Line Search (ELS) Rajih et al. (2008) finds optimal step values at every even iteration by minimizing the modified cost function

$$D(\mathcal{Y} \parallel \hat{\mathcal{Y}}) = \frac{1}{2} \left\| \mathcal{Y} - [\mathbf{A}_t^{(1)} + \eta_1 \Delta \mathbf{A}_t^{(1)}, \mathbf{A}_t^{(2)} + \eta_2 \Delta \mathbf{A}_t^{(2)}, \mathbf{A}_t^{(3)} + \eta_2 \Delta \mathbf{A}_t^{(3)}] \right\|_F^2,$$

where η_1, η_2, η_3 are stepsizes, and $\Delta \mathbf{A}_t^{(n)} = \mathbf{A}_{t-1}^{(n)} - \mathbf{A}_{t-2}^{(n)}$.

When $\eta_1 = \eta_2 = \eta_3 = \eta$, ELS find roots of a polynomial $\mathbf{p}(\mu)$ of degree 5. An optimal stepsize is the root of that polynomial $\mathbf{p}(\mu)$ which yields the smallest cost function value.

- ▶ ELS alleviates bottlenecks more efficiently than LS.
- ▶ But less efficient if the directions provided by the algorithm are bad Comon et al. (2009).
- ▶ High computational cost for higher order tensors, due to building up $(2N - 1)$ - order polynomials $\mathbf{p}(\eta)$.

Line search: Example

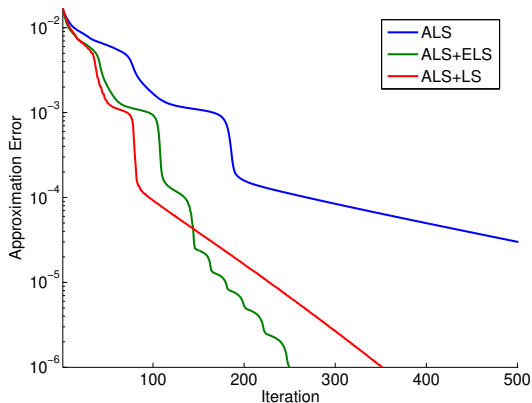


Figure: Approximation errors of the ALS with/without line searches for factorizations of tensors of size $30 \times 30 \times 30$ and rank $R = 10$ composed from highly correlated components which have $\mathbf{a}_r^{(n)T} \mathbf{a}_s^{(n)} \in [0.8, 0.999]$, $r \neq s$.

Rank-one Update Algorithm

Algorithm updates only one component during iterations.

Assume component $\mathbf{a}_r^{(n)}$ to be updated, all other components are fixed. $\hat{\mathbf{y}}$ is split into two parts as follows

A rank-one tensor $\hat{\mathbf{y}}^{(1)}$ is built up from components $\mathbf{a}_r^{(n)}$

$$\hat{\mathbf{y}}^{(1)} = \mathbf{a}_1^{(n)} \otimes \mathbf{a}_2^{(n)} \otimes \dots \otimes \mathbf{a}_r^{(n)} \quad (10)$$

A rank- $(r-1)$ tensor $\hat{\mathbf{y}}^{(2)}$ is composed by all factors

$$\hat{\mathbf{y}}^{(2)} = [\mathbf{a}_1^{(n)} \otimes \mathbf{a}_2^{(n)} \otimes \dots \otimes \mathbf{a}_{r-1}^{(n)}] \quad \text{that is,}$$

$$\hat{\mathbf{y}}^{(2)} = \sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} \dots \sum_{i_{r-1}=1}^{I_{r-1}} \mathbf{a}_{i_1}^{(n)} \otimes \mathbf{a}_{i_2}^{(n)} \otimes \dots \otimes \mathbf{a}_{i_{r-1}}^{(n)} \quad (11)$$

$$\hat{\mathbf{y}} = \hat{\mathbf{y}}^{(1)} + \hat{\mathbf{y}}^{(2)} = \mathbf{a}_1^{(n)} \otimes \mathbf{a}_2^{(n)} \otimes \dots \otimes \mathbf{a}_r^{(n)} + \hat{\mathbf{y}}^{(2)} \quad (12)$$

Rank-one Update Algorithm

Algorithm updates only one component during iterations.

Assume component $\mathbf{a}_r^{(n)}$ to be updated, all other components are fixed. $\hat{\mathbf{y}}$ is split into two parts as follows

Rank-one Update Algorithm

Algorithm updates only one component during iterations.

Assume component $\mathbf{a}_r^{(n)}$ to be updated, all other components are fixed. $\hat{\mathbf{y}}$ is split into two parts as follows

- ▶ A rank-one tensor $\hat{\mathbf{y}}^{(r)}$ is built up from components $\mathbf{a}_r^{(n)}$

$$\hat{\mathbf{y}}^{(r)} = \mathbf{a}_r^{(1)} \circ \mathbf{a}_r^{(2)} \circ \dots \circ \mathbf{a}_r^{(N)}. \quad (15)$$

- ▶ A rank- $(R-1)$ CP tensor $\hat{\mathbf{y}}^{(-r)}$ is composed by N factors $\mathbf{A}_{-r}^{(N)} = [\mathbf{a}_1^{(n)}, \dots, \mathbf{a}_{r-1}^{(n)}, \mathbf{a}_{r+1}^{(n)}, \dots, \mathbf{a}_R^{(n)}]$, that is

$$\begin{aligned} \hat{\mathbf{y}}^{(-r)} &= \sum_{k \neq r} \mathbf{a}_k^{(1)} \circ \mathbf{a}_k^{(2)} \circ \dots \circ \mathbf{a}_k^{(N)} \\ &= \mathbf{J} \times_1 \mathbf{A}_{-r}^{(1)} \times_2 \mathbf{A}_{-r}^{(2)} \dots \times_N \mathbf{A}_{-r}^{(N)}. \end{aligned} \quad (16)$$

Rank-one Update Algorithm

Algorithm updates only one component during iterations.

Assume component $\mathbf{a}_r^{(n)}$ to be updated, all other components are fixed. $\hat{\mathbf{y}}$ is split into two parts as follows

- ▶ A rank-one tensor $\hat{\mathbf{y}}^{(r)}$ is built up from components $\mathbf{a}_r^{(n)}$

$$\hat{\mathbf{y}}^{(r)} = \mathbf{a}_r^{(1)} \circ \mathbf{a}_r^{(2)} \circ \dots \circ \mathbf{a}_r^{(N)}. \quad (15)$$

- ▶ A rank- $(R-1)$ CP tensor $\hat{\mathbf{y}}^{(-r)}$ is composed by N factors $\mathbf{A}_{-r}^{(N)} = [\mathbf{a}_1^{(n)}, \dots, \mathbf{a}_{r-1}^{(n)}, \mathbf{a}_{r+1}^{(n)}, \dots, \mathbf{a}_R^{(n)}]$, that is

$$\begin{aligned} \hat{\mathbf{y}}^{(-r)} &= \sum_{k \neq r} \mathbf{a}_k^{(1)} \circ \mathbf{a}_k^{(2)} \circ \dots \circ \mathbf{a}_k^{(N)} \\ &= \mathbf{J} \times_1 \mathbf{A}_{-r}^{(1)} \times_2 \mathbf{A}_{-r}^{(2)} \dots \times_N \mathbf{A}_{-r}^{(N)}. \end{aligned} \quad (16)$$

Rank-one update

- ▶ Approximation $\mathbf{y} = \hat{\mathbf{y}}^{(-r)} + \hat{\mathbf{y}}^{(r)} + \boldsymbol{\varepsilon}$.
- ▶ Define a residual tensor $\tilde{\mathbf{y}}^{(r)} = \mathbf{y} - \hat{\mathbf{y}}^{(-r)} = \hat{\mathbf{y}}^{(r)} + \boldsymbol{\varepsilon}$.

Update rule for $\mathbf{a}_r^{(n)}$

$$\mathbf{a}_r^{(n)} \leftarrow \tilde{\mathbf{Y}}_{(n)}^{(r)} \{\mathbf{a}_r\}^{\odot -n} \left(\{\mathbf{a}_r^T \mathbf{a}_r\}^{\otimes -n} \right)^{\dagger} = \frac{\tilde{\mathbf{y}}^{(r)} \bar{\mathbf{X}}_{-n}(\mathbf{a}_r)}{\gamma_r^{(n)}}$$

$$\gamma_r^{(n)} = \{\mathbf{a}_r^T \mathbf{a}_r\}^{\otimes -n} = \begin{cases} \mathbf{a}_r^{(N)T} \mathbf{a}_r^{(N)}, & n \neq N \\ 1, & n = N. \end{cases}$$

Due to normalization $\mathbf{a}_r^{(n)} = \mathbf{a}_r^{(n)} / \|\mathbf{a}_r^{(n)}\|_2$ for $n = 1, 2, \dots, N-1$ after each iteration step, scaling factor $\gamma_r^{(n)}$ can be omitted.

Rank-one update

- ▶ Approximation $\mathbf{y} = \hat{\mathbf{y}}^{(-r)} + \hat{\mathbf{y}}^{(r)} + \boldsymbol{\varepsilon}$.
- ▶ Define a residual tensor $\tilde{\mathbf{y}}^{(r)} = \mathbf{y} - \hat{\mathbf{y}}^{(-r)} = \hat{\mathbf{y}}^{(r)} + \boldsymbol{\varepsilon}$.

Update rule for $\mathbf{a}_r^{(n)}$

$$\mathbf{a}_r^{(n)} \leftarrow \tilde{\mathbf{Y}}_{(n)}^{(r)} \{\mathbf{a}_r\}^{\odot -n} \left(\{\mathbf{a}_r^T \mathbf{a}_r\}^{\otimes -n} \right)^{\dagger} = \frac{\tilde{\mathbf{y}}^{(r)} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\}}{\gamma_r^{(n)}}$$

$$\gamma_r^{(n)} = \{\mathbf{a}_r^T \mathbf{a}_r\}^{\otimes -n} = \begin{cases} \mathbf{a}_r^{(N)T} \mathbf{a}_r^{(N)}, & n \neq N \\ 1, & n = N. \end{cases}$$

Due to normalization $\mathbf{a}_r^{(n)} = \mathbf{a}_r^{(n)} / \|\mathbf{a}_r^{(n)}\|_2$ for $n = 1, 2, \dots, N-1$ after each iteration step, scaling factor $\gamma_r^{(n)}$ can be omitted.

Rank-one update

- ▶ Approximation $\mathbf{y} = \hat{\mathbf{y}}^{(-r)} + \hat{\mathbf{y}}^{(r)} + \boldsymbol{\varepsilon}$.
- ▶ Define a residual tensor $\tilde{\mathbf{y}}^{(r)} = \mathbf{y} - \hat{\mathbf{y}}^{(-r)} = \hat{\mathbf{y}}^{(r)} + \boldsymbol{\varepsilon}$.

Update rule for $\mathbf{a}_r^{(n)}$

$$\mathbf{a}_r^{(n)} \leftarrow \tilde{\mathbf{Y}}_{(n)}^{(r)} \{\mathbf{a}_r\}^{\odot -n} \left(\{\mathbf{a}_r^T \mathbf{a}_r\}^{\otimes -n} \right)^{\dagger} = \frac{\tilde{\mathbf{y}}^{(r)} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\}}{\gamma_r^{(n)}}$$

$$\gamma_r^{(n)} = \{\mathbf{a}_r^T \mathbf{a}_r\}^{\otimes -n} = \begin{cases} \mathbf{a}_r^{(N)T} \mathbf{a}_r^{(N)}, & n \neq N \\ 1, & n = N. \end{cases}$$

Due to normalization $\mathbf{a}_r^{(n)} = \mathbf{a}_r^{(n)} / \|\mathbf{a}_r^{(n)}\|_2$ for $n = 1, 2, \dots, N-1$ after each iteration step, scaling factor $\gamma_r^{(n)}$ can be omitted.

Efficient version

$$\begin{aligned}\mathbf{a}_r^{(n)} &\leftarrow \tilde{\mathbf{y}}^{(r)} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} = \mathbf{y} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} - \mathbf{y}^{(-r)} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} \\ &= \mathbf{y} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} - \mathbf{J} \times_n \mathbf{A}_{-r}^{(n)} \bar{\mathbf{x}}_{-n}\{\mathbf{A}_{-r}^T \mathbf{a}_r\} \\ &= \mathbf{y} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} - \mathbf{A}_{-r}^{(n)} \left\{ \mathbf{A}_{-r}^T \mathbf{a}_r \right\}^{\otimes -n},\end{aligned}$$

where $\left\{ \mathbf{A}_{-r}^T \mathbf{a}_r \right\} \triangleq \left\{ \mathbf{A}_{-r}^{(1)T} \mathbf{a}_r^{(1)}, \dots, \mathbf{A}_{-r}^{(N)T} \mathbf{a}_r^{(N)} \right\}$.

Note

Efficient version

$$\begin{aligned}\mathbf{a}_r^{(n)} &\leftarrow \tilde{\mathbf{y}}^{(r)} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} = \mathbf{y} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} - \mathbf{y}^{(-r)} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} \\ &= \mathbf{y} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} - \mathbf{J} \times_n \mathbf{A}_{-r}^{(n)} \bar{\mathbf{x}}_{-n}\{\mathbf{A}_{-r}^T \mathbf{a}_r\} \\ &= \mathbf{y} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\} - \mathbf{A}_{-r}^{(n)} \{\mathbf{A}_{-r}^T \mathbf{a}_r\}^{\otimes -n},\end{aligned}$$

where $\{\mathbf{A}_{-r}^T \mathbf{a}_r\} \triangleq \{\mathbf{A}_{-r}^{(1)T} \mathbf{a}_r^{(1)}, \dots, \mathbf{A}_{-r}^{(N)T} \mathbf{a}_r^{(N)}\}$.

Note

- ▶ $\mathbf{y} \bar{\mathbf{x}}_{-n}\{\mathbf{a}_r\}$ does not demand significant storage.
- ▶ The ALS rule computes $\{\mathbf{A}\}^{\otimes -n}$ resulting a large-scale matrix $I^{N-1} \times R$.

Efficient version

$$\begin{aligned}\mathbf{a}_r^{(n)} &\leftarrow \tilde{\mathbf{y}}^{(r)} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\} = \mathbf{y} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\} - \mathbf{y}^{(-r)} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\} \\ &= \mathbf{y} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\} - \mathbf{J} \times_n \mathbf{A}_{-r}^{(n)} \bar{\mathbf{x}}_{-n} \{\mathbf{A}_{-r}^T \mathbf{a}_r\} \\ &= \mathbf{y} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\} - \mathbf{A}_{-r}^{(n)} \{\mathbf{A}_{-r}^T \mathbf{a}_r\}^{\otimes -n},\end{aligned}$$

where $\{\mathbf{A}_{-r}^T \mathbf{a}_r\} \triangleq \{\mathbf{A}_{-r}^{(1)T} \mathbf{a}_r^{(1)}, \dots, \mathbf{A}_{-r}^{(N)T} \mathbf{a}_r^{(N)}\}$.

Note

- ▶ $\mathbf{y} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\}$ does not demand significant storage.
- ▶ The ALS rule computes $\{\mathbf{A}\}^{\odot -n}$ resulting a large-scale matrix $I^{N-1} \times R$.

Algorithm 5: Rank-1 Update Algorithm for CP and NTF

```
1 begin
2   Initialize all  $\mathbf{A}^{(n)}$ 
3   Normalize all  $\mathbf{a}_{r_n}^{(n)}$  to unit length
4   repeat
5     for  $n = 1$  to  $N$  do                                     /* Update  $\mathbf{A}^{(n)}$  */
6       for  $r = 1$  to  $R$  do
7          $\mathbf{a}_r^{(n)} \leftarrow \mathcal{Y} \bar{\mathbf{x}}_{-n} \{\mathbf{a}_r\} - \mathbf{A}_{-r}^{(n)} \{\mathbf{A}_{-r}^T \mathbf{a}_r\}^{\otimes -n}$ 
8          $\mathbf{a}_r^{(n)} \leftarrow [\mathbf{a}_r^{(n)}]_+$  /* Retrieve nonnegative entries */
9         if  $n \neq N$  then /* Normalize and fix scaling */
10           $\mathbf{a}_r^{(N)} = \frac{\mathbf{a}_r^{(N)} \|\mathbf{a}_r^{(n)}\|_2}{\|\mathbf{a}_r^{(N)}\|_2^2}, \quad \mathbf{a}_r^{(n)} = \frac{\mathbf{a}_r^{(n)}}{\|\mathbf{a}_r^{(n)}\|_2}$ 
11        end
12      end
13    end
14  until a stopping criterion is met
15 end
```

Part III

Higher order CPD

High order CPD

- ▶ Most algorithms for order-3 CPD can be straightforwardly extended to higher order tensors.
- ▶ However, the algorithms become more complicated, computationally demanding, and often not applicable to relatively large scale tensors.
 - ▶ Complexity of CP gradients grows linearly with the tensor order N with a cost of $\mathcal{O}\left(NR \prod_{n=1}^N I_n\right)$.
 - ▶ More time consuming due to more tensor unfoldings $\mathbf{Y}_{(n)}$ ($n = 2, 3, \dots, N - 1$).
 - ▶ Line search extrapolation methods Harshman (1970); Andersson and Bro (2000); Tomasi (2006); Rajih et al. (2008) become more complicated, and high computationally demanding.
- ▶ In practice, CPD algorithms may take a lot of iterations, and are very slow, when some factor matrices become nearly rank deficient Comon et al. (2009).

NOTE

An approach for high dimensional large scale tensor is prerequisite.

The First Example

- ▶ Given an order-4 tensor $\mathcal{Y} = \mathcal{J} \times_1 \mathbf{A}^{(1)} \times_2 \mathbf{A}^{(2)} \times_3 \mathbf{A}^{(3)} \times_4 \mathbf{A}^{(4)}$.
- ▶ An unfolding tensor $\mathcal{Y}_{\llbracket 1,2,(3,4) \rrbracket}$ is given by

$$\mathcal{Y}_{\llbracket 1,2,(3,4) \rrbracket} = \mathcal{J} \times_1 \mathbf{A}^{(1)} \times_2 \mathbf{A}^{(2)} \times_3 \mathbf{B}$$

where $\mathbf{B} = \mathbf{A}^{(4)} \odot \mathbf{A}^{(3)}$.

- ▶ Since CPD is unique under “mild” conditions, one can perform higher order CPD through a lower order CPD.
For example, $\mathbf{a}_r^{(3)}$ and $\mathbf{a}_r^{(4)}$ are the best rank-one approximation of $\text{reshape}(\mathbf{b}_r, [l_3, l_4])$ using SVD.

Definition (Reshaping)

The reshape operator for a tensor $\mathcal{Y} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ to a size specified by a vector $\mathbf{L} = [L_1, L_2, \dots, L_M]$ with $\prod_{m=1}^M L_m = \prod_{n=1}^N I_n$ returns an order- M tensor \mathcal{X} , such that $\text{vec}(\mathcal{Y}) = \text{vec}(\mathcal{X})$, and is expressed as $\mathcal{X} = \text{reshape}(\mathcal{Y}, \mathbf{L}) \in \mathbb{R}^{L_1 \times L_2 \times \dots \times L_M}$.

Definition (Tensor transposition Ragnarsson and Loan (2012))

If $\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ and \mathbf{p} is a permutation of $[1, 2, \dots, N]$, then \mathbf{p} -transpose of \mathcal{A} is defined by

$$\mathcal{A}^{<\mathbf{p}>}(i_{p_1}, \dots, i_{p_N}) = \mathcal{A}(i_1, \dots, i_N), \quad \mathbf{1} \leq \mathbf{i} \leq \mathbf{I} = [I_1, I_2, \dots, I_N].$$

Definition (Generalized tensor unfolding)

Reshaping a \mathbf{p} -transpose $\mathcal{Y}^{<\mathbf{p}>}$ to an order- M tensor of size $\mathbf{L} = [L_1, L_2, \dots, L_M]$ with $L_m = \prod_{k \in I_m} l_k$, where $[l_1, l_2, \dots, l_M] \equiv [p_1, p_2, \dots, p_N]$, $I_m = [I_m(1), \dots, I_m(K_m)]$

$$\mathcal{Y}_{[\mathbf{I}]} \triangleq \text{reshape}(\mathcal{Y}^{<\mathbf{p}>}, \mathbf{L}), \quad \mathbf{I} = [I_1, I_2, \dots, I_M].$$

Higher order CPD through Tensor UnFolding

Simple unFolding CPD (FCP)

1

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, rank R ,
Unfolding $I = [I_1, I_2, \dots, I_M]$ where $I_m = [I_m(1), \dots, I_m(K_m)]$, $M \geq 3$

Output: $\lambda \in \mathbb{R}^N$, N matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

2 **begin**

3 /* **Stage 1: Optional compression of unfolding tensor $\mathcal{Y}_{[I]}$** */

5 $[\mathcal{G}, \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(M)}] = \text{Tucker_decomposition}(\mathcal{Y}_{[I]}, \min(I, R))$

6 /* Unfolding \mathcal{Y} to be a lower order tensor $\mathcal{Y}_{[I]}$ and */

8 /* Perform multilinear low-rank approximation to $\mathcal{Y}_{[I]}$
 using High order SVD or Higher Order Orthogonal
 Iteration. */

9 **end**

Higher order CPD through Tensor UnFolding

Simple unFolding CPD (FCP)

1

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, rank R ,
Unfolding $I = [I_1, I_2, \dots, I_M]$ where $I_m = [I_m(1), \dots, I_m(K_m)]$, $M \geq 3$

Output: $\lambda \in \mathbb{R}^N$, N matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

2 **begin**

3 /* **Stage 1:** Optional compression of unfolding tensor $\mathcal{Y}_{[I]}$ */

5 $[\mathcal{G}, \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(M)}] = \text{TD}(\mathcal{Y}_{[I]}, \min(I, R))$

6 /* **Stage 2:** CPD of \mathcal{G} and back projection----- */

9 $[\lambda; \mathbf{B}^{(1)}, \dots, \mathbf{B}^{(M)}] = \text{CPD}(\mathcal{G}, R)$

11 **for** $m = 1, 2, \dots, M$ **do** $\mathbf{B}^{(m)} \leftarrow \mathbf{U}^{(m)} \mathbf{B}^{(m)}$

12

13 **end**

Higher order CPD through Tensor UnFolding

Simple unFolding CPD (FCP)

1

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, rank R ,
Unfolding $I = [I_1, I_2, \dots, I_M]$ where $I_m = [I_m(1), \dots, I_m(K_m)]$, $M \geq 3$

Output: $\lambda \in \mathbb{R}^N$, N matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

2 **begin**

3 /* **Stage 1: Optional compression of unfolding tensor $\mathcal{Y}_{[I]}$** */

5 $[\mathcal{G}, \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(M)}] = \text{TD}(\mathcal{Y}_{[I]}, \min(I, R))$

6 /* **Stage 2: CPD of \mathcal{G} and back projection-----** */

9 $[\lambda; \mathbf{B}^{(1)}, \dots, \mathbf{B}^{(M)}] = \text{CPD}(\mathcal{G}, R)$

11 **for** $m = 1, 2, \dots, M$ **do** $\mathbf{B}^{(m)} \leftarrow \mathbf{U}^{(m)} \mathbf{B}^{(m)}$

12

13 /* **Stage 3: Estimate factors for high order tensors-** */

16 **for** $m = 1, 2, \dots, M$ **do**

18 **for** $r = 1, 2, \dots, R$ **do**

20 $\mathcal{B}_r = \text{reshape}(\mathbf{b}_r^{(m)}, [I_m(1), \dots, I_m(K_m)])$

22 $[\mathcal{g}; \mathbf{a}_r^{(I_m(1))}, \dots, \mathbf{a}_r^{(I_m(K_m))}] = \text{TD}(\mathcal{B}_r, 1)$

24 $\lambda_r \leftarrow \lambda_r \mathcal{g}$

25 **end**

Higher order CPD through Tensor UnFolding

Simple unFolding CPD (FCP)

1

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, rank R ,
Unfolding $I = [I_1, I_2, \dots, I_M]$ where $I_m = [I_m(1), \dots, I_m(K_m)]$, $M \geq 3$

Output: $\lambda \in \mathbb{R}^N$, N matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

2 **begin**

3 /* **Stage 1: Optional compression of unfolding tensor $\mathcal{Y}_{[I]}$** */

5 $[\mathcal{G}, \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(M)}] = \text{TD}(\mathcal{Y}_{[I]}, \min(I, R))$

6 /* **Stage 2: CPD of \mathcal{G} and back projection-----** */

9 $[\lambda; \mathbf{B}^{(1)}, \dots, \mathbf{B}^{(M)}] = \text{CPD}(\mathcal{G}, R)$

11 **for** $m = 1, 2, \dots, M$ **do** $\mathbf{B}^{(m)} \leftarrow \mathbf{U}^{(m)} \mathbf{B}^{(m)}$

12

13 /* **Stage 3: Estimate factors for high order tensors-** */

16 **for** $m = 1, 2, \dots, M$ **do**

18 **for** $r = 1, 2, \dots, R$ **do**

20 $\mathcal{B}_r = \text{reshape}(\mathbf{b}_r^{(m)}, [I_m(1), \dots, I_m(K_m)])$

22 $[\mathcal{g}; \mathbf{a}_r^{(I_m(1))}, \dots, \mathbf{a}_r^{(I_m(K_m))}] = \text{TD}(\mathcal{B}_r, 1)$

24 $\lambda_r \leftarrow \lambda_r \mathcal{g}$

25

end

Loop over all columns of matrices $\mathbf{B}^{(m)}$

reshape them to tensors/matrices

Higher order CPD through Tensor UnFolding

Simple unFolding CPD (FCP)

1

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, rank R ,
Unfolding $I = [I_1, I_2, \dots, I_M]$ where $I_m = [I_m(1), \dots, I_m(K_m)]$, $M \geq 3$

Output: $\lambda \in \mathbb{R}^N$, N matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

2 **begin**

3 /* **Stage 1: Optional compression of unfolding tensor $\mathcal{Y}_{\llbracket I \rrbracket}$** */

5 $\llbracket \mathcal{G}, \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(M)} \rrbracket = \text{TD}(\mathcal{Y}_{\llbracket I \rrbracket}, \min(I, R))$

6 /* **Stage 2: CPD of \mathcal{G} and back projection-----** */

9 $\llbracket \lambda; \mathbf{B}^{(1)}, \dots, \mathbf{B}^{(M)} \rrbracket = \text{CPD}(\mathcal{G}, R)$

11 **for** $m = 1, 2, \dots, M$ **do** $\mathbf{B}^{(m)} \leftarrow \mathbf{U}^{(m)} \mathbf{B}^{(m)}$

12

13 /* **Stage 3: Estimate factors for high order tensors-** */

16 **for** $m = 1, 2, \dots, M$ **do**

18 **for** $r = 1, 2, \dots, R$ **do**

20 $\mathcal{B}_r = \text{reshape}(\mathbf{b}_r^{(m)}, [I_m(1), \dots, I_m(K_m)])$

22 $\llbracket g; \mathbf{a}_r^{(I_m(1))}, \dots, \mathbf{a}_r^{(I_m(K_m))} \rrbracket = \text{TD}(\mathcal{B}_r, 1)$

24 $\lambda_r \leftarrow \lambda_r g$

25 **end**

Rank-one tensor approximation

Higher order CPD through Tensor UnFolding

Simple unFolding CPD (FCP)

1

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \dots \times I_N)$, rank R ,
Unfolding $I = [I_1, I_2, \dots, I_M]$ where $I_m = [I_m(1), \dots, I_m(K_m)]$, $M \geq 3$
Output: $\lambda \in \mathbb{R}^N$, N matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$

~ Iteration

Simple or Rank-one FCP

- ▶ Sensitive to unfolding rules.
- ▶ Worst performance if unfolding modes having low-collinearity.

1.

12

13

16

18

20

22

24

25

```
/* Stage 3: Estimate factors for high order tensors- */  
for m = 1, 2, ..., M do  
    for r = 1, 2, ..., R do  
         $\mathcal{B}_r = \text{reshape}(\mathbf{b}_r^{(m)}, [I_m(1), \dots, I_m(K_m)])$   
         $[\mathbf{g}; \mathbf{a}_r^{(I_m(1))}, \dots, \mathbf{a}_r^{(I_m(K_m))}] = \text{TD}(\mathcal{B}_r, 1)$   
         $\lambda_r \leftarrow \lambda_r \mathbf{g}$ 
```

Loss of Accuracy of the simple FCP

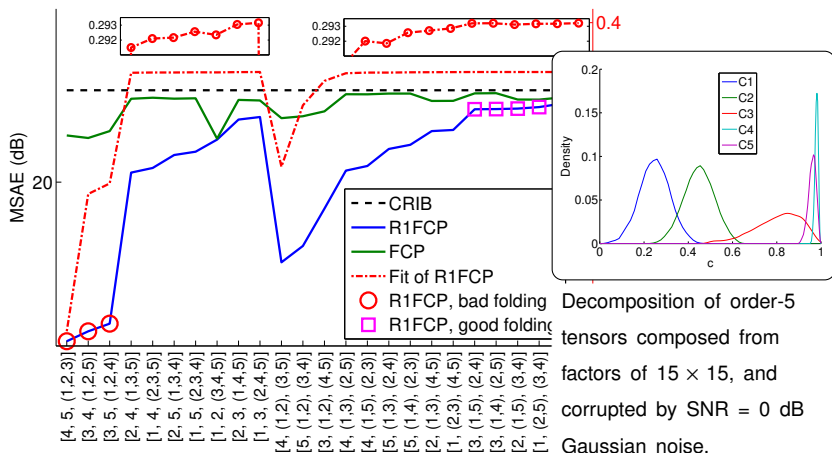


Figure: R1FCP caused high loss of accuracy when combining modes 1 and 2 with low collinearity degrees.

The Second Example I

- ▶ Consider a rank- R CPD and a simple unfolding $\mathbf{I} = [1, 2, (3, 4)]$

$$\begin{aligned}\mathcal{Y} &= \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}, \mathbf{A}^{(4)} \rrbracket + \mathcal{E}_1, \quad \boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_R] \\ \mathcal{Y}_{\llbracket \mathbf{I} \rrbracket} &= \llbracket \boldsymbol{\mu}; \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \mathbf{B}^{(3)} \rrbracket + \mathcal{E}_2, \quad \boldsymbol{\mu} = [\mu_1, \dots, \mu_R]\end{aligned}$$

- ▶ If the noise variance is low, $\boldsymbol{\lambda} \approx \boldsymbol{\mu}$, $\mathbf{A}^{(k)} \approx \mathbf{B}^{(k)}$, for $k = 1, 2$.
- ▶ Put $\mathbf{F}_r = \text{reshape}(\mathbf{b}_r^{(3)}, [l_3, l_4])$, and let $\mathbf{F}_r = \mathbf{U}_r \boldsymbol{\Sigma}_r \mathbf{V}_r^T$ be Singular Value Decomposition of \mathbf{F}_r

$$\begin{aligned}\mathbf{U}_r &= [\mathbf{u}_{r1}, \dots, \mathbf{u}_{rT}], \quad \mathbf{V}_r = [\mathbf{v}_{r1}, \dots, \mathbf{v}_{rT}] \\ \boldsymbol{\sigma}_r &= [\sigma_{r1}, \sigma_{r2}, \dots, \sigma_{rT}], \quad 1 \geq \sigma_{r1} \geq \sigma_{r2} \geq \dots \geq \sigma_{rT} \geq 0.\end{aligned}$$

- ▶ Theoretically, \mathbf{F}_r has rank one, and

$$\mathcal{Y} \approx \mathcal{Y}_R = \llbracket \boldsymbol{\mu}; \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \dots, \widetilde{\mathbf{B}}^{(N-1)}, \widetilde{\mathbf{B}}^{(N)} \rrbracket$$

where $\widetilde{\mathbf{B}}^{(N-1)} = [\mathbf{u}_{11}, \dots, \mathbf{u}_{R1}]$ and $\widetilde{\mathbf{B}}^{(N)} = [\mathbf{v}_{11}, \dots, \mathbf{v}_{R1}]$.

The Second Example II

- ▶ In general, a numerical rank of \mathbf{F}_r is $J_r \geq 1$, $r = 1, \dots, R$ defined as $\sum_{j=1}^{J_r} \sigma_{rj}^2 \geq \tau$ where $0 < \tau < 1$, e.g., $\tau = 0.98$.
- ▶ We can write a rank- J tensor approximation of \mathcal{Y} , $J = \sum_{r=1}^R J_r$

$$\begin{aligned}\mathcal{Y}_J &= \sum_{r=1}^R \mu_r \mathbf{b}_r^{(1)} \circ \mathbf{b}_r^{(2)} \circ \left(\sum_{j=1}^{J_r} \sigma_{rj} (\mathbf{u}_{rj} \circ \mathbf{v}_{rj}) \right) \\ &= [\![\tilde{\lambda}; \tilde{\mathbf{A}}^{(1)}, \tilde{\mathbf{A}}^{(2)}, \tilde{\mathbf{A}}^{(3)}, \tilde{\mathbf{A}}^{(4)}]\!]\end{aligned}$$

where

$$\tilde{\mathbf{A}}^{(k)} = \mathbf{B}^{(k)} \mathbf{M} \quad \text{for } k = 1, 2$$

$$\tilde{\mathbf{A}}^{(3)} = [\mathbf{U}_1, \dots, \mathbf{U}_R]$$

$$\tilde{\mathbf{A}}^{(4)} = [\mathbf{V}_1, \dots, \mathbf{V}_R]$$

$$\mathbf{M} = \text{bdiag}(\mathbf{1}_{1 \times J_1}, \dots, \mathbf{1}_{1 \times J_R})$$

$$\tilde{\lambda} = [\mu_1 \sigma_{11}, \dots, \mu_1 \sigma_{1J_1}, \dots, \mu_R \sigma_{R1}, \dots, \mu_R \sigma_{RJ_R}].$$

The Second Example III

- ▶ The desired rank- R approximation of \mathcal{Y} can be obtained by applying a structured CPD algorithm to \mathcal{Y}_J using \mathcal{Y}_R as an initial value.
- ▶ The method can be easily extended for higher order tensor, and multimode unfoldings.

1

Input: Data tensor \mathcal{Y} : ($I_1 \times I_2 \times \dots \times I_N$), rank R , threshold $\tau (\geq 0.98)$
 Unfolding rule $\mathbf{I} = [\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_M]$, $\mathbf{I}_m = [I_m(1), I_m(2), \dots, I_m(K_m)]$

Output: Rank- R Kruskal tensor $\mathcal{Y}_R = \llbracket \lambda; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket$

2 **begin**

```

3   /* Stage 1: Tensor unfolding and compression----- */
5    $\llbracket \mathcal{G}, \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(M)} \rrbracket = \text{TD}(\mathcal{Y}_{\llbracket \mathbf{I} \rrbracket}, \min(\mathbf{I}, R))$ 
6   /* Stage 2: CPD of  $\mathcal{G}$  and back projection ----- */
9    $\llbracket \lambda; \mathbf{B}^{(1)}, \dots, \mathbf{B}^{(M)} \rrbracket = \text{CPD}(\mathcal{G}, R)$ 
11   $\mathcal{Y}_R = \llbracket \lambda; \mathbf{U}^{(1)} \mathbf{B}^{(1)}, \dots, \mathbf{U}^{(M)} \mathbf{B}^{(M)} \rrbracket$  /* rank- $R$  K-tensor */
13  /* Stage 3: Sequential estimation ----- */
14  foreach group  $\mathbf{I}_m$  with  $K_m \geq 2$  do
16       $n = K_1 + \dots + K_{m-1}$ 
17      for  $k = 1, 2, \dots, K_m$  do
18          /* Stage 3a: Construction of rank- $J$  K-tensor----- */
19           $[\mathcal{Y}_J, \mathcal{Y}_R] = \text{twomodeFCP}(\mathcal{Y}_R, n + k, I_{m(k)})$ 
20          /* Stage 3b: Rank- $J$  to rank- $R$  K-tensor----- */
23          if  $J > R$  then  $\mathcal{Y}_R = \text{structuredCPD}(\mathcal{Y}_J, R, \mathcal{Y}_R)$ 
24      end
25  end
26  /* Stage 4: Refinement if needed ----- */
29   $\mathcal{Y}_R = \text{CPD}(\mathcal{Y}, R, \mathcal{Y}_R)$ 
30 end
```

1

Input: Data tensor \mathcal{Y} : $(I_1 \times I_2 \times \cdots \times I_N)$, rank R , threshold $\tau (\geq 0.98)$
 Unfolding rule $I = [I_1, I_2, \dots, I_M]$, $I_m = [I_m(1), I_m(2), \dots, I_m(K_m)]$
Output: Rank- R Kruskal tensor $\mathcal{Y}_R = \llbracket \lambda; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket$

2 **begin**

```

1 function  $[\mathcal{Y}_J, \mathcal{Y}_R] = \text{twomodeFCP}(\mathcal{Y}_R, n, T)$ 
   Input: Order- $N$  rank- $R$  Kruskal tensor  $\mathcal{Y}_R = \llbracket \lambda; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket$ , with
        $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$  to be reconstructed
   Output: Two order- $(N+1)$  rank- $J$  and rank- $R$  K-tensors  $\mathcal{Y}_J, \mathcal{Y}_R$ 
2 for  $r = 1, 2, \dots, R$  do
3      $\mathbf{F}_r = \text{reshape}(\mathbf{a}_r^{(n)}, [T, I_n/T])$ 
4      $\mathbf{F}_r \approx \mathbf{U}_r \text{diag}(\boldsymbol{\sigma}_r) \mathbf{V}_r^T$  /* rank- $J_r$  truncated SVD,  $\|\boldsymbol{\sigma}_r\|_2^2 \geq \tau$  */
5 end
6 /* Construct a rank- $J$  K-tensor,  $J = J_1 + J_2 + \cdots + J_R$  ----- */
7  $\mathbf{M} = \text{bdiag}(\mathbf{1}_{1 \times J_1}, \dots, \mathbf{1}_{1 \times J_R})$ 
8  $\tilde{\lambda} = [\lambda_1 \sigma_{11}, \lambda_2 \sigma_{21}, \dots, \lambda_R \sigma_{R1}]$ ,  $\mathbf{G} = [\mathbf{U}_1, \dots, \mathbf{U}_R]$ ,  $\mathbf{H} = [\mathbf{V}_1, \dots, \mathbf{V}_R]$ 
9  $\mathcal{Y}_J = \llbracket \tilde{\lambda}; \mathbf{A}^{(1)} \mathbf{M}, \dots, \mathbf{A}^{(n-1)} \mathbf{M}, \mathbf{G}, \mathbf{H}, \mathbf{A}^{(n+1)} \mathbf{M}, \dots, \mathbf{A}^{(N)} \mathbf{M} \rrbracket$ 
10 /* Construct a rank- $R$  Kruskal tensor ----- */
11  $\lambda \leftarrow [\lambda_1 \sigma_{11}, \lambda_2 \sigma_{21}, \dots, \lambda_R \sigma_{R1}]$ ,  $\mathbf{G} = [\mathbf{u}_{11}, \dots, \mathbf{u}_{R1}]$ ,  $\mathbf{H} = [\mathbf{v}_{11}, \dots, \mathbf{v}_{R1}]$ 
12  $\mathcal{Y}_R = \llbracket \lambda; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(n-1)}, \mathbf{G}, \mathbf{H}, \mathbf{A}^{(n+1)}, \dots, \mathbf{A}^{(N)} \rrbracket$ 

```

29 $\mathcal{Y}_R = \text{CPD}(\mathcal{Y}, \mathbf{H}, \mathcal{Y}_R)$ 30 **end**

Decomposition of the ORL face database

- ▶ An order-5 Gabor feature tensor of size $16 \times 16 \times 8 \times 4 \times 400$ was constructed from the ORL faces for 8 different orientations at 4 scales.
- ▶ \mathcal{Y} was unfolded to be an order-3 tensor using $I = [1, 2, (3, 4, 5)]$.
- ▶ The factor $\mathbf{A}^{(5)} \in \mathbb{R}^{400 \times R}$ comprised compressed features which were used to cluster faces using the K-means algorithm.

R	Algorithm	Fit (%)	Time (seconds)	Ratio $\frac{\text{ALS}}{\text{FCP}}$	ACC (%)	NMI (%)
30	FCP	60.59	24	39	85.00	92.91
	ALS	60.56	927		85.25	93.22
40	FCP	62.46	39	41	84.25	92.57
	ALS	62.63	1599		85.50	93.68
60	FCP	65.47	105	162	83.00	91.62
	ALS	65.64	16962		81.38	91.44

Part IV

Damped Gauss-Newton Algorithms

Damped Gauss-Newton Algorithm ?

$$D = \|\mathcal{Y} - \mathcal{J} \times_1 \mathbf{A}^{(1)} \cdots \times_N \mathbf{A}^{(N)}\|_F^2 + P_I(\{\mathbf{A}^{(n)}\}) + P_S(\{\mathbf{A}^{(n)}\}),$$

$$P_I = - \sum_{n=1}^N \alpha_n \sum_{i_n=1}^{I_n} \sum_{r=1}^R \log(a_{i_n r}^{(n)}), \quad P_S = \sum_{n=1}^N \beta_n \|\mathbf{A}^{(n)}\|_1,$$

$$\alpha_n \geq 0, \beta_n \geq 0 \quad \forall n.$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g},$$

$$\mathbf{a} = \left[\text{vec}(\mathbf{A}^{(1)})^T \text{vec}(\mathbf{A}^{(2)})^T \cdots \text{vec}(\mathbf{A}^{(N)})^T \right]^T$$

\mathbf{g} gradient $R \sum I_n \times 1$

\mathbf{H} approximate Hessian $R \sum I_n \times R \sum I_n$

$\mu > 0$: damping parameter,

Damped Gauss-Newton Algorithm ?

$$D = \|\mathcal{Y} - \mathcal{J} \times_1 \mathbf{A}^{(1)} \cdots \times_N \mathbf{A}^{(N)}\|_F^2 + P_l(\{\mathbf{A}^{(n)}\}) + P_s(\{\mathbf{A}^{(n)}\}),$$

$$P_l = - \sum_{n=1}^N \alpha_n \sum_{i_n=1}^{I_n} \sum_{r=1}^R \log(a_{i_n r}^{(n)}), \quad P_s = \sum_{n=1}^N \beta_n \|\mathbf{A}^{(n)}\|_1,$$

$$\alpha_n \geq 0, \beta_n \geq 0 \quad \forall n.$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g},$$

$$\mathbf{a} = \left[\text{vec}(\mathbf{A}^{(1)})^T \text{vec}(\mathbf{A}^{(2)})^T \cdots \text{vec}(\mathbf{A}^{(N)})^T \right]^T$$

\mathbf{g} gradient $R \sum I_n \times 1$

\mathbf{H} approximate Hessian $R \sum I_n \times R \sum I_n$

$\mu > 0$: damping parameter,

Damped Gauss-Newton Algorithm ?

$$D = \|\mathcal{Y} - \mathcal{J} \times_1 \mathbf{A}^{(1)} \cdots \times_N \mathbf{A}^{(N)}\|_F^2 + P_I(\{\mathbf{A}^{(n)}\}) + P_S(\{\mathbf{A}^{(n)}\}),$$

$$P_I = - \sum_{n=1}^N \alpha_n \sum_{i_n=1}^{I_n} \sum_{r=1}^R \log(a_{i_n r}^{(n)}), \quad P_S = \sum_{n=1}^N \beta_n \|\mathbf{A}^{(n)}\|_1,$$

$$\alpha_n \geq 0, \beta_n \geq 0 \quad \forall n.$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g},$$

$$\mathbf{a} = \left[\text{vec}(\mathbf{A}^{(1)})^T \text{vec}(\mathbf{A}^{(2)})^T \cdots \text{vec}(\mathbf{A}^{(N)})^T \right]^T$$

$$\mathbf{g} \quad \text{gradient } R \sum I_n \times 1$$

$$\mathbf{H} \quad \text{approximate Hessian } R \sum I_n \times R \sum I_n$$

$$\mu > 0 : \quad \text{damping parameter,}$$

Damped Gauss-Newton Algorithm ?

$$D = \|\mathcal{Y} - \mathcal{J} \times_1 \mathbf{A}^{(1)} \cdots \times_N \mathbf{A}^{(N)}\|_F^2 + P_l(\{\mathbf{A}^{(n)}\}) + P_s(\{\mathbf{A}^{(n)}\}),$$

$$P_l = - \sum_{n=1}^N \alpha_n \sum_{i_n=1}^{I_n} \sum_{r=1}^R \log(a_{i_n r}^{(n)}), \quad P_s = \sum_{n=1}^N \beta_n \|\mathbf{A}^{(n)}\|_1,$$

$$\alpha_n \geq 0, \beta_n \geq 0 \quad \forall n.$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g},$$

$$\mathbf{a} = \left[\text{vec}(\mathbf{A}^{(1)})^T \text{vec}(\mathbf{A}^{(2)})^T \cdots \text{vec}(\mathbf{A}^{(N)})^T \right]^T$$

\mathbf{g} gradient $R \sum I_n \times 1$

\mathbf{H} approximate Hessian $R \sum I_n \times R \sum I_n$

$\mu > 0$: damping parameter,

Damped Gauss-Newton Algorithm ?

$$D = \|\mathcal{Y} - \mathcal{J} \times_1 \mathbf{A}^{(1)} \cdots \times_N \mathbf{A}^{(N)}\|_F^2 + P_l(\{\mathbf{A}^{(n)}\}) + P_s(\{\mathbf{A}^{(n)}\}),$$

$$P_l = - \sum_{n=1}^N \alpha_n \sum_{i_n=1}^{I_n} \sum_{r=1}^R \log(a_{i_n r}^{(n)}), \quad P_s = \sum_{n=1}^N \beta_n \|\mathbf{A}^{(n)}\|_1,$$

$$\alpha_n \geq 0, \beta_n \geq 0 \quad \forall n.$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g},$$

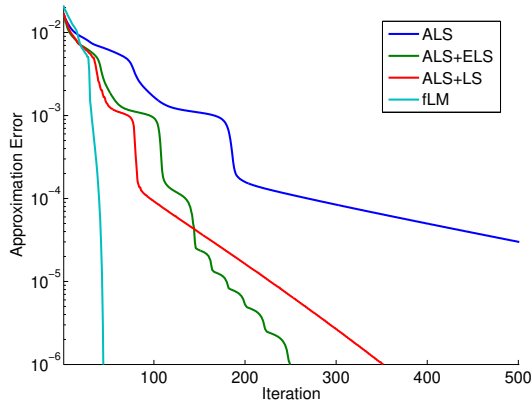
$$\mathbf{a} = \left[\text{vec}(\mathbf{A}^{(1)})^T \text{vec}(\mathbf{A}^{(2)})^T \cdots \text{vec}(\mathbf{A}^{(N)})^T \right]^T$$

\mathbf{g} gradient $R \sum I_n \times 1$

\mathbf{H} approximate Hessian $R \sum I_n \times R \sum I_n$

$\mu > 0$: damping parameter,

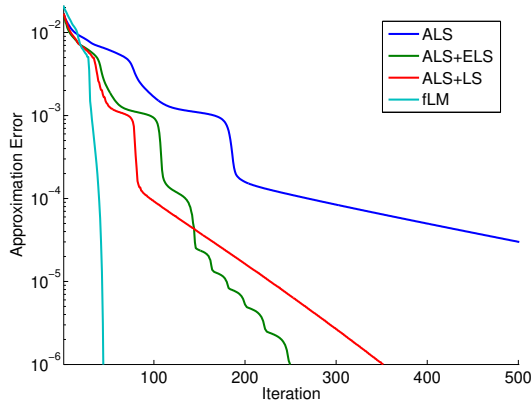
damped Gauss-Newton vs ALS



Tensor of size $30 \times 30 \times 30$ and rank- $R = 10$ composed by highly correlated components which have $\mathbf{a}_r^{(n)T} \mathbf{a}_s^{(n)} \in [0.8, 0.999]$, $r \neq s$.

- ▶ LM works well even for difficult cases such as tensors with highly collinear components, or different magnitudes of factors.
- ▶ But high computational cost due to computing Jacobian, gradient, Hessian and inverse of Hessian.

damped Gauss-Newton vs ALS



Tensor of size $30 \times 30 \times 30$ and rank- $R = 10$ composed by highly correlated components which have $\mathbf{a}_r^{(n)T} \mathbf{a}_s^{(n)} \in [0.8, 0.999]$, $r \neq s$.

- ▶ LM works well even for difficult cases such as tensors with highly collinear components, or different magnitudes of factors.
- ▶ But high computational cost due to computing Jacobian, gradient, Hessian and inverse of Hessian.

Jacobian, Hessian and gradient matrices

$$\mathbf{J} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}} = \left[\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(1)}} \quad \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(2)}} \quad \cdots \quad \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(N)}} \right] \in \mathbb{R}^{(\Pi l_n) \times R \Sigma l_n}.$$

$$\mathbf{g} = \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) + \frac{\partial P_l}{\partial \mathbf{a}} + \frac{\partial P_s}{\partial \mathbf{a}}$$

$$= \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) - \left[\left(\alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-1]} - \beta_n \right)^T \right]_{n=1}^N \in \mathbb{R}^{R \Sigma l_n}.$$

$$\mathbf{H} = \mathbf{J}^T \mathbf{J} + \frac{\partial^2 P_l}{\partial \mathbf{a}^2} + \frac{\partial^2 P_s}{\partial \mathbf{a}^2}$$

$$= \mathbf{J}^T \mathbf{J} + \bigoplus_{n=1}^N \text{diag} \left\{ \alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-2]} \right\} \in \mathbb{R}^{R \Sigma l_n \times R \Sigma l_n}.$$

Update rule is **computationally demanding** due to construction of Jacobian, and **inverse of large Hessians**.

Jacobian, Hessian and gradient matrices

$$\mathbf{J} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}} = \left[\begin{array}{ccc} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(1)}} & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(2)}} & \cdots & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(N)}} \end{array} \right] \in \mathbb{R}^{(\prod l_n) \times R \sum l_n}.$$

$$\begin{aligned} \mathbf{g} &= \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) + \frac{\partial P_l}{\partial \mathbf{a}} + \frac{\partial P_s}{\partial \mathbf{a}} \\ &= \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) - \left[\left(\alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-1]} - \beta_n \right)^T \right]_{n=1}^N \in \mathbb{R}^{R \sum l_n}. \end{aligned}$$

$$\begin{aligned} \mathbf{H} &= \mathbf{J}^T \mathbf{J} + \frac{\partial^2 P_l}{\partial \mathbf{a}^2} + \frac{\partial^2 P_s}{\partial \mathbf{a}^2} \\ &= \mathbf{J}^T \mathbf{J} + \bigoplus_{n=1}^N \text{diag} \left\{ \alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-2]} \right\} \in \mathbb{R}^{R \sum l_n \times R \sum l_n}. \end{aligned}$$

Update rule is **computationally demanding** due to construction of Jacobian, and **inverse of large Hessians**.

Jacobian, Hessian and gradient matrices

$$\mathbf{J} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}} = \left[\begin{array}{ccc} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(1)}} & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(2)}} & \cdots & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(N)}} \end{array} \right] \in \mathbb{R}^{(\prod l_n) \times R \sum l_n}.$$

$$\begin{aligned} \mathbf{g} &= \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) + \frac{\partial P_l}{\partial \mathbf{a}} + \frac{\partial P_s}{\partial \mathbf{a}} \\ &= \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) - \left[\left(\alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-1]} - \beta_n \right)^T \right]_{n=1}^N \in \mathbb{R}^{R \sum l_n}. \end{aligned}$$

$$\begin{aligned} \mathbf{H} &= \mathbf{J}^T \mathbf{J} + \frac{\partial^2 P_l}{\partial \mathbf{a}^2} + \frac{\partial^2 P_s}{\partial \mathbf{a}^2} \\ &= \mathbf{J}^T \mathbf{J} + \bigoplus_{n=1}^N \text{diag} \left\{ \alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-2]} \right\} \in \mathbb{R}^{R \sum l_n \times R \sum l_n}. \end{aligned}$$

Update rule is computationally demanding due to construction of Jacobian, and inverse of large Hessians.

Jacobian, Hessian and gradient matrices

$$\mathbf{J} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}} = \left[\begin{array}{ccc} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(1)}} & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(2)}} & \cdots & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(N)}} \end{array} \right] \in \mathbb{R}^{(\Pi I_n) \times R \Sigma I_n}.$$

$$\begin{aligned} \mathbf{g} &= \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) + \frac{\partial P_l}{\partial \mathbf{a}} + \frac{\partial P_s}{\partial \mathbf{a}} \\ &= \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) - \left[\left(\alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-1]} - \beta_n \right)^T \right]_{n=1}^N \in \mathbb{R}^{R \Sigma I_n}. \end{aligned}$$

$$\begin{aligned} \mathbf{H} &= \mathbf{J}^T \mathbf{J} + \frac{\partial^2 P_l}{\partial \mathbf{a}^2} + \frac{\partial^2 P_s}{\partial \mathbf{a}^2} \\ &= \mathbf{J}^T \mathbf{J} + \bigoplus_{n=1}^N \text{diag} \left\{ \alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-2]} \right\} \in \mathbb{R}^{R \Sigma I_n \times R \Sigma I_n}. \end{aligned}$$

Update rule is **computationally demanding** due to construction of Jacobian, and **inverse of large Hessians**.

Jacobian matrix

$$\mathbf{J} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}} = \left[\begin{array}{ccc} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(1)}} & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(2)}} & \cdots & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(N)}} \end{array} \right].$$

Jacobian matrix

$$\mathbf{J} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}} = \begin{bmatrix} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}(1)} & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}(2)} & \cdots & \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}(N)} \end{bmatrix}.$$

$$\begin{aligned} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_r^{(n)}} &= \frac{\partial \text{vec}(\hat{\mathbf{y}})}{\partial \mathbf{a}_r^{(n)}} = \mathbf{P}_n \frac{\partial \text{vec}(\hat{\mathbf{Y}}_{(n)})}{\partial \mathbf{a}_r^{(n)}} = \mathbf{P}_n \frac{\partial \sum_{r=1}^R \mathbf{a}_r^{(n)} (\mathbf{a}_r^{\otimes -n})^T}{\partial \mathbf{a}_r^{(n)}} \\ &= \mathbf{P}_n (\mathbf{a}_r^{\otimes -n} \otimes \mathbf{I}_{l_n}). \end{aligned}$$

\mathbf{P}_n is a permutation matrix: $\text{vec}(\mathbf{y}) = \mathbf{P}_n \text{vec}(\mathbf{Y}_{(n)})$

Fast computation of Gradient

Jacobian

$$\mathbf{J} = \left[\mathbf{P}_1 (\mathbf{A}^{\odot -1} \otimes \mathbf{I}_{l_1}) \quad \cdots \quad \mathbf{P}_n (\mathbf{A}^{\odot -n} \otimes \mathbf{I}_{l_n}) \quad \cdots \quad \mathbf{P}_N (\mathbf{A}^{\odot -N} \otimes \mathbf{I}_{l_N}) \right]$$

Fast computation of Gradient

Jacobian

$$\mathbf{J} = \left[\mathbf{P}_1 (\mathbf{A}^{\odot -1} \otimes \mathbf{I}_{l_1}) \quad \cdots \quad \mathbf{P}_n (\mathbf{A}^{\odot -n} \otimes \mathbf{I}_{l_n}) \quad \cdots \quad \mathbf{P}_N (\mathbf{A}^{\odot -N} \otimes \mathbf{I}_{l_N}) \right]$$

Gradient matrix

$$\mathbf{g} = \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) + \frac{\partial P_l}{\partial \mathbf{a}} + \frac{\partial P_s}{\partial \mathbf{a}}$$

Fast computation of Gradient

Jacobian

$$\mathbf{J} = \begin{bmatrix} \mathbf{P}_1 (\mathbf{A}^{\odot-1} \otimes \mathbf{I}_{l_1}) & \cdots & \mathbf{P}_n (\mathbf{A}^{\odot-n} \otimes \mathbf{I}_{l_n}) & \cdots & \mathbf{P}_N (\mathbf{A}^{\odot-N} \otimes \mathbf{I}_{l_N}) \end{bmatrix}$$

Gradient matrix

$$\mathbf{g} = \mathbf{J}^T (\hat{\mathbf{y}} - \mathbf{y}) + \frac{\partial P_l}{\partial \mathbf{a}} + \frac{\partial P_s}{\partial \mathbf{a}}$$

$$\begin{aligned} \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{A}^{(n)}} \right)^T (\mathbf{y} - \hat{\mathbf{y}}) &= (\mathbf{A}^{\odot-n T} \otimes \mathbf{I}_n) \mathbf{P}_n^T \text{vec}(\mathbf{y} - \hat{\mathbf{y}}) \\ &= (\mathbf{A}^{\odot-n T} \otimes \mathbf{I}_n) \text{vec}(\mathbf{Y}_{(n)} - \hat{\mathbf{Y}}_{(n)}) \\ &= \text{vec}(\mathbf{Y}_{(n)} \mathbf{A}^{\odot-n} - \mathbf{A}^{(n)} \Gamma^{(n,n)}). \end{aligned}$$

$$\Gamma^{(n,n)} = (\mathbf{A}^{(N)T} \mathbf{A}^{(N)}) \circledast \cdots \circledast (\mathbf{A}^{(n+1)T} \mathbf{A}^{(n+1)}) \circledast (\mathbf{A}^{(n-1)T} \mathbf{A}^{(n-1)}) \circledast \cdots \circledast (\mathbf{A}^{(1)T} \mathbf{A}^{(1)})$$

Fast computation of Gradient

Gradient

$$\mathbf{g} = \begin{bmatrix} \text{vec}\left(-\mathbf{Y}_{(1)}\mathbf{A}^{\odot-1} + \mathbf{A}^{(1)}\Gamma^{(1,1)} - \alpha_1 \left(\mathbf{A}^{(1)}\right)^{\bullet[-1]} + \beta_1\right) \\ \vdots \\ \text{vec}\left(-\mathbf{Y}_{(n)}\mathbf{A}^{\odot-n} + \mathbf{A}^{(n)}\Gamma^{(n,n)} - \alpha_n \left(\mathbf{A}^{(n)}\right)^{\bullet[-1]} + \beta_n\right) \\ \vdots \\ \text{vec}\left(-\mathbf{Y}_{(N)}\mathbf{A}^{\odot-N} + \mathbf{A}^{(N)}\Gamma^{(N,N)} - \alpha_N \left(\mathbf{A}^{(N)}\right)^{\bullet[-1]} + \beta_N\right) \end{bmatrix}.$$

$$\Gamma^{(n,n)} = \bigotimes_{k \neq n} \mathbf{A}^{(k)T} \mathbf{A}^{(k)} \in \mathbb{R}^{R \times R}$$

Fast computation of Hessian

- ▶ Inverse of \mathbf{H} is too computationally expensive and of order $O(R^3(l_1 + l_2 + \dots + l_N)^3)$.

$$\mathbf{H} + \mu \mathbf{I} = \mathbf{J}^T \mathbf{J} + \bigoplus_{n=1}^N \text{diag} \left\{ \alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-2]} + \mu \right\} \in \mathbb{R}^{R \sum l_n \times R \sum l_n}$$

- ▶ But approximate Hessian is rank-deficient. Expression of \mathbf{H} by low-rank adjustment can reduce computational cost.

$$\mathbf{H} + \mu \mathbf{I} = \begin{bmatrix} \mathbf{H}^{(1,1)} & \dots & \mathbf{H}^{(1,m)} & \dots & \mathbf{H}^{(1,N)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{H}^{(n,1)} & \dots & \mathbf{H}^{(n,m)} & \dots & \mathbf{H}^{(n,N)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{H}^{(N,1)} & \dots & \mathbf{H}^{(N,m)} & \dots & \mathbf{H}^{(N,N)} \end{bmatrix}.$$

Fast computation of Hessian

- ▶ Inverse of \mathbf{H} is too computationally expensive and of order $O(R^3(l_1 + l_2 + \dots + l_N)^3)$.

$$\mathbf{H} + \mu \mathbf{I} = \mathbf{J}^T \mathbf{J} + \bigoplus_{n=1}^N \text{diag} \left\{ \alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-2]} + \mu \right\} \in \mathbb{R}^{R \sum l_n \times R \sum l_n}$$

- ▶ But approximate Hessian is rank-deficient. Expression of \mathbf{H} by low-rank adjustment can reduce computational cost.

$$\mathbf{H} + \mu \mathbf{I} = \begin{bmatrix} \mathbf{H}^{(1,1)} & \dots & \mathbf{H}^{(1,m)} & \dots & \mathbf{H}^{(1,N)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{H}^{(n,1)} & \dots & \mathbf{H}^{(n,m)} & \dots & \mathbf{H}^{(n,N)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{H}^{(N,1)} & \dots & \mathbf{H}^{(N,m)} & \dots & \mathbf{H}^{(N,N)} \end{bmatrix}.$$

Expression of block matrix $\mathbf{H}^{(n,m)}$ |

$$\begin{aligned}\mathbf{H} &= \mathbf{J}^T \mathbf{J} \\ &= \begin{bmatrix} \mathbf{H}_{1,1}^{(1,1)} & \cdots & \mathbf{H}_{1,R}^{(1,1)} & \cdots & \mathbf{H}_{r,s}^{(1,m)} & \cdots & \mathbf{H}_{R,R}^{(1,N)} \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{H}_{1,1}^{(n,1)} & \cdots & \mathbf{H}_{1,R}^{(n,1)} & \cdots & \mathbf{H}_{r,s}^{(n,m)} & \cdots & \mathbf{H}_{R,R}^{(n,N)} \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{H}_{1,1}^{(N,1)} & \cdots & \mathbf{H}_{1,R}^{(N,1)} & \cdots & \mathbf{H}_{r,s}^{(N,m)} & \cdots & \mathbf{H}_{R,R}^{(N,N)} \end{bmatrix}, \quad (17)\end{aligned}$$

where block matrices $\mathbf{H}_{r,s}^{(n,m)} \in \mathbb{R}^{l_n \times l_m}$ are computed as

$$\mathbf{H}_{r,s}^{(n,m)} = \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_r^{(n)}} \right)^T \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_s^{(m)}} \right), \quad (18)$$

$n, m = 1, 2, \dots, N, r, s = 1, 2, \dots, R.$

Expression of block matrix $\mathbf{H}^{(n,m)}$ II

Theorem (Permutation in a Kronecker product)

Multiplication of a permutation matrix \mathbf{P}_n with a Kronecker product of N vectors $\mathbf{a}^{(n)} \in \mathbb{R}^{I_n}$, $n = 1, 2, \dots, N$ will move the n -th item $\mathbf{a}^{(n)}$ to the last term, that is

$$\begin{aligned} & (\mathbf{a}^{(N)} \otimes \dots \otimes \mathbf{a}^{(2)} \otimes \mathbf{a}^{(1)}) \\ &= \mathbf{P}_n (\mathbf{a}^{(N)} \otimes \dots \otimes \mathbf{a}^{(n+1)} \otimes \mathbf{a}^{(n-1)} \otimes \dots \otimes \mathbf{a}^{(1)} \otimes \mathbf{a}^{(n)}) . \end{aligned} \quad (19)$$

Theorem (Rank-one block matrices $\mathbf{H}_{r,s}^{(n,m)}$)

Block matrices $\mathbf{H}_{r,s}^{(n,m)}$, $n \neq m$ are rank-one matrices

$$\mathbf{H}_{r,s}^{(n,m)} = \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_r^{(n)}} \right)^T \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_s^{(m)}} \right) = \gamma_{rs}^{(n,m)} \mathbf{a}_r^{(n)} \mathbf{a}_s^{(m)T} \in \mathbb{R}^{I_n \times I_m} . \quad (20)$$

Expression of block matrix $\mathbf{H}^{(n,m)}$ III

$$\begin{aligned} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_r^{(n)}} &= \mathbf{P}_n \left[\mathbf{a}_r^{\otimes -n} \otimes \mathbf{e}_1^{(n)} \quad \mathbf{a}_r^{\otimes -n} \otimes \mathbf{e}_2^{(n)} \quad \cdots \quad \mathbf{a}_r^{\otimes -n} \otimes \mathbf{e}_{l_n}^{(n)} \right] \\ &= \left[\mathbf{P}_n \left(\mathbf{a}_r^{\otimes -n} \otimes \mathbf{e}_{i_n}^{(n)} \right) \right]_{i_n=1}^{l_n} \in \mathbb{R}^{(l_1 l_2 \dots l_N) \times l_n}, \end{aligned} \quad (21)$$

where $\mathbf{e}_{i_n}^{(n)}$ is the i_n -th column of the identity matrix $\mathbf{I}_{l_n} \in \mathbb{R}^{l_n \times l_n}$.

$$\begin{aligned} \left[\mathbf{H}_{r,s}^{(n,m)} \right]_{i_n, i_m} &= \left[\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_r^{(n)}} \right]_{i_n}^T \left[\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_s^{(m)}} \right]_{i_m} \\ &= \left(\mathbf{a}_r^{\otimes -n T} \otimes \mathbf{e}_{i_n}^{(n) T} \right) \mathbf{P}_n^T \mathbf{P}_m \left(\mathbf{a}_s^{\otimes -m} \otimes \mathbf{e}_{i_m}^{(m)} \right) \end{aligned}$$

Expression of block matrix $\mathbf{H}^{(n,m)}$ IV

$$\begin{aligned}
 \left[\mathbf{H}_{r,s}^{(n,m)} \right]_{i_n, i_m} &= \left(\mathbf{a}_r^{(N)T} \otimes \dots \otimes \mathbf{a}_r^{(n+1)T} \otimes \mathbf{e}_{i_n}^{(n)T} \otimes \mathbf{a}_r^{(n-1)T} \otimes \dots \otimes \mathbf{a}_r^{(1)T} \right) \\
 &\quad \left(\mathbf{a}_s^{(N)} \otimes \dots \otimes \mathbf{a}_s^{(m+1)} \otimes \mathbf{e}_{i_m}^{(m)} \otimes \mathbf{a}_s^{(m-1)} \otimes \dots \otimes \mathbf{a}_s^{(1)} \right) \\
 &= \left(\prod_{k \neq n, m} \left(\mathbf{a}_r^{(k)T} \mathbf{a}_s^{(k)} \right) \right) \left(\left(\mathbf{a}_r^{(m)T} \mathbf{e}_{i_m}^{(m)} \right) \otimes \left(\mathbf{e}_{i_n}^{(n)T} \mathbf{a}_s^{(n)} \right) \right) \\
 &= \gamma_{rs}^{(n,m)} \mathbf{a}_{i_m r}^{(m)} \mathbf{a}_{i_n s}^{(n)}, \tag{22}
 \end{aligned}$$

where $\gamma_{rs}^{(n,m)}$ is the (r, s) element of $\Gamma^{(n,m)} = (\mathbf{A}^T \mathbf{A})^{\otimes - (n,m)}$. This leads to a compact

$$\mathbf{H}_{r,s}^{(n,m)} = \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_r^{(n)}} \right)^T \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_s^{(m)}} \right) = \gamma_{rs}^{(n,m)} \mathbf{a}_r^{(n)} \mathbf{a}_s^{(m)T} \in \mathbb{R}^{I_n \times I_m}. \tag{23}$$

Expression of block matrix $\mathbf{H}^{(n,m)}$

Theorem (Diagonal block matrices $\mathbf{H}_{r,s}^{(n,n)}$)

Block matrices $\mathbf{H}_{r,s}^{(n,n)}$, $n \neq m$ are diagonal matrices.

$$\begin{aligned}\mathbf{H}_{r,s}^{(n,n)} &= \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_r^{(n)}} \right)^T \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_s^{(n)}} \right) = \left(\mathbf{a}_r^{\otimes -n T} \otimes \mathbf{I}_n \right) \mathbf{P}_n^T \mathbf{P}_n \left(\mathbf{a}_s^{\otimes -n T} \otimes \mathbf{I}_n \right) \\ &= \left(\prod_{k \neq n} \mathbf{a}_r^{(k) T} \mathbf{a}_s^{(k)} \right) \mathbf{I}_n = \gamma_{rs}^{(n,n)} \mathbf{I}_n,\end{aligned}\tag{24}$$

where $\gamma_{rs}^{(n,n)}$ are the (r, s) element of $\Gamma^{(n,n)} = (\mathbf{A}^T \mathbf{A})^{\otimes -(n)}$.

Expression of block matrix $\mathbf{H}^{(n,m)}$ VI

Theorem (Diagonal block $\mathbf{H}^{(n,n)}$)

A block matrix $\mathbf{H}^{(n,n)}$ for $n = 1, 2, \dots, N$ on the diagonal of the Hessian \mathbf{H}

$$\mathbf{H}^{(n,n)} = \Gamma^{(n,n)} \otimes \mathbf{I}_n.$$

Proof.

Expression of block matrix $\mathbf{H}^{(n,m)}$ VII

$$\begin{aligned}\mathbf{H}^{(n,n)} &= \begin{bmatrix} \gamma_{1,1}^{(n,n)} \mathbf{I}_n & \cdots & \gamma_{r,1}^{(n,n)} \mathbf{I}_n & \cdots & \gamma_{R,1}^{(n,n)} \mathbf{I}_n \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \gamma_{1,s}^{(n,n)} \mathbf{I}_n & \cdots & \gamma_{r,s}^{(n,n)} \mathbf{I}_n & \cdots & \gamma_{R,s}^{(n,n)} \mathbf{I}_n \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \gamma_{1,R}^{(n,n)} \mathbf{I}_n & \cdots & \gamma_{r,R}^{(n,n)} \mathbf{I}_n & \cdots & \gamma_{R,R}^{(n,n)} \mathbf{I}_n \end{bmatrix} \\ &= \Gamma^{(n,n)} \otimes \mathbf{I}_n.\end{aligned}\tag{25}$$

□

Expression of block matrix $\mathbf{H}^{(n,m)}$ VIII

Theorem (Block matrix $\mathbf{H}^{(n,m)}$)

$\mathbf{H}^{(n,m)}$ for $\forall n \neq m$ can be decomposed as

$$\mathbf{H}^{(n,m)} = (\mathbf{I}_R \otimes \mathbf{A}^{(n)}) \mathbf{P}_{R,R} \text{diag}(\text{vec}(\Gamma^{(n,m)})) (\mathbf{I}_R \otimes \mathbf{A}^{(m)T}) .$$

Expression of block matrix $\mathbf{H}^{(n,m)}$ IX

Definition

Vectorization of an $I \times J$ matrix \mathbf{X} can be established via its transpose matrix \mathbf{X}^T by an $IJ \times IJ$ permutation matrix $\mathbf{P}_{I,J}$

$$\text{vec}(\mathbf{X}) = \mathbf{P}_{I,J} \text{vec}(\mathbf{X}^T) . \quad (26)$$

Theorem (Interchange in Kronecker product)

Consider the permutation matrices $\mathbf{P}_{I,P}$ and $\mathbf{P}_{J,Q}$. Then

$$\mathbf{B} \otimes \mathbf{A} = \mathbf{P}_{I,P}^T (\mathbf{A} \otimes \mathbf{B}) \mathbf{P}_{J,Q} , \quad (27)$$

for all matrices $\mathbf{A} \in \mathbb{R}^{I \times J}$ and $\mathbf{B} \in \mathbb{R}^{P \times Q}$.

Proof.

This theorem is replicated from the problem 4.3.10 in Horn and Johnson (1991). See proof in Horn and Johnson (1991). □

Expression of block matrix $\mathbf{H}^{(n,m)}$

Theorem (Concatenation of Kronecker products)

Concatenation of the Kronecker products forms a Kronecker product of the concatenation matrix

$$[\mathbf{A} \otimes \mathbf{C} \quad \mathbf{D} \otimes \mathbf{C}] = [\mathbf{A} \quad \mathbf{D}] \otimes \mathbf{C}, \quad (28)$$

and

$$[\mathbf{A} \otimes \mathbf{B} \quad \mathbf{A} \otimes \mathbf{C}] = (\mathbf{A} \otimes [\mathbf{B} \quad \mathbf{C}]) \mathbf{P}, \quad (29)$$

where \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} are matrices of size $I \times J$, $P \times Q$, $P \times R$ and $I \times K$, respectively, and \mathbf{P} denotes a permutation matrix.

Expression of block matrix $\mathbf{H}^{(n,m)} \mathbf{X}$

$$\begin{aligned} [\mathbf{A} \otimes \mathbf{C} \quad \mathbf{D} \otimes \mathbf{C}] &= \begin{bmatrix} \mathbf{a}_1 \otimes \mathbf{C} & \cdots & \mathbf{a}_M \otimes \mathbf{C} & \mathbf{d}_1 \otimes \mathbf{C} & \cdots & \mathbf{d}_K \otimes \mathbf{C} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_M & \mathbf{d}_1 & \cdots & \mathbf{d}_K \end{bmatrix} \otimes \mathbf{C} \\ &= \begin{bmatrix} \mathbf{A} & \mathbf{D} \end{bmatrix} \otimes \mathbf{C}. \end{aligned} \quad (30)$$

By employing Theorem 12, the second concatenation (29) is proved as follows

$$\begin{aligned} [\mathbf{A} \otimes \mathbf{B} \quad \mathbf{A} \otimes \mathbf{C}] &= \begin{bmatrix} \mathbf{P}_{P,I}^T (\mathbf{B} \otimes \mathbf{A}) \mathbf{P}_{Q,J} & \mathbf{P}_{P,I}^T (\mathbf{C} \otimes \mathbf{A}) \mathbf{P}_{R,J} \end{bmatrix} \\ &= \mathbf{P}_{P,I}^T ([\mathbf{B} \quad \mathbf{C}] \otimes \mathbf{A}) \begin{bmatrix} \mathbf{P}_{Q,J} \\ \mathbf{P}_{R,J} \end{bmatrix} \\ &= \mathbf{P}_{P,I}^T \mathbf{P}_{P,I} (\mathbf{A} \otimes [\mathbf{B} \quad \mathbf{C}]) \mathbf{P}_{Q+R,J}^T \begin{bmatrix} \mathbf{P}_{Q,J} \\ \mathbf{P}_{R,J} \end{bmatrix} \\ &= (\mathbf{A} \otimes [\mathbf{B} \quad \mathbf{C}]) \mathbf{P}. \end{aligned} \quad (31)$$

Expression of block matrix $\mathbf{H}^{(n,m)}$ XII

$$\begin{bmatrix} \mathbf{A} \otimes \mathbf{b}_1 & \cdots & \mathbf{A} \otimes \mathbf{b}_q & \cdots & \mathbf{A} \otimes \mathbf{b}_Q \end{bmatrix} = (\mathbf{A} \otimes \mathbf{B}) \mathbf{P}_{J,Q}.$$

$$\begin{aligned} \mathbf{H}^{(n,m)} &= \begin{bmatrix} \gamma_{1,1}^{(n,m)} \mathbf{a}_1^{(n)} \mathbf{a}_1^{(m)T} & \cdots & \gamma_{r,1}^{(n,m)} \mathbf{a}_r^{(n)} \mathbf{a}_1^{(m)T} & \cdots & \gamma_{R,s}^{(n,m)} \mathbf{a}_R^{(n)} \mathbf{a}_1^{(m)T} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \gamma_{1,s}^{(n,m)} \mathbf{a}_1^{(n)} \mathbf{a}_s^{(m)T} & \cdots & \gamma_{r,s}^{(n,m)} \mathbf{a}_r^{(n)} \mathbf{a}_s^{(m)T} & \cdots & \gamma_{R,s}^{(n,m)} \mathbf{a}_R^{(n)} \mathbf{a}_s^{(m)T} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \gamma_{1,R}^{(n,m)} \mathbf{a}_1^{(n)} \mathbf{a}_R^{(m)T} & \cdots & \gamma_{r,R}^{(n,m)} \mathbf{a}_r^{(n)} \mathbf{a}_R^{(m)T} & \cdots & \gamma_{R,R}^{(n,m)} \mathbf{a}_R^{(n)} \mathbf{a}_R^{(m)T} \end{bmatrix} \\ &= \begin{bmatrix} (\mathbf{I}_R \otimes \mathbf{a}_1^{(n)}) \mathbf{D}_1^{(n,m)} \mathbf{A}^{(n)T} & \cdots & (\mathbf{I}_R \otimes \mathbf{a}_r^{(n)}) \mathbf{D}_r^{(n,m)} \mathbf{A}^{(n)T} & \cdots & (\mathbf{I}_R \otimes \mathbf{a}_R^{(n)}) \mathbf{D}_R^{(n,m)} \mathbf{A}^{(n)T} \end{bmatrix} \\ &= \begin{bmatrix} (\mathbf{I}_R \otimes \mathbf{a}_1^{(n)}) & \cdots & (\mathbf{I}_R \otimes \mathbf{a}_r^{(n)}) & \cdots & (\mathbf{I}_R \otimes \mathbf{a}_R^{(n)}) \end{bmatrix} \mathbf{D}^{(n,m)} (\mathbf{I}_R \otimes \mathbf{A}^{(m)T}) \\ &= (\mathbf{I}_R \otimes \mathbf{A}^{(n)}) \mathbf{P}_{R,R} \mathbf{D}^{(n,m)} (\mathbf{I}_R \otimes \mathbf{A}^{(m)T}), \end{aligned}$$

where $\mathbf{D}_r^{(n,m)} = \text{diag}(\gamma_r^{(n,m)})$ is the diagonal matrix whose diagonal entries are the r -th row of the matrix $\Gamma^{(n,m)}$, $\mathbf{D}^{(n,m)} = \text{diag}(\text{vec}(\Gamma^{(n,m)}))$ is the diagonal matrix whose the main diagonal is $\text{vec}(\Gamma^{(n,m)})$, $\mathbf{P}_{R,R}$ is the permutation matrix of an $R \times R$ matrix.

Expression of block matrix $\mathbf{H}^{(n,m)}$ XIII

Theorem (Expression of block matrix $\mathbf{H}^{(n,m)}$)

A block matrix $\mathbf{H}^{(n,m)}$ can be expressed by

$$\mathbf{H}^{(n,m)} = \begin{cases} \Gamma^{(n,n)} \otimes \mathbf{I}_n + \text{diag} \left\{ \alpha_n \text{vec}(\mathbf{A}^{(n)})^{\bullet[-2]} + \mu \right\}, & n = m \\ (\mathbf{I}_R \otimes \mathbf{A}^{(n)}) \Pi_{R,R} \text{diag}(\text{vec}(\Gamma^{(n,m)})) (\mathbf{I}_R \otimes \mathbf{A}^{(m)T}), & n \neq m \end{cases}$$

where $\Gamma^{(n,m)} = \bigcirc_{k \neq m,n} \mathbf{A}^{(k)T} \mathbf{A}^{(k)}$,

$\Pi_{R,R}$ is the commutation matrix: $\text{vec}(\mathbf{X}_{R \times R}) = \Pi_{R,R} \text{vec}(\mathbf{X}^T)$.

Theorem (Low rank Adjustment for the Hessian \mathbf{H})

The Hessian \mathbf{H} can be decomposed into low rank matrices under the form as

$$\mathbf{H} = \mathbf{G} + \mathbf{Z}\mathbf{K}\mathbf{Z}^T. \quad (33)$$

where \mathbf{G} is an invertible matrix, matrix \mathbf{Z} consists of $N R^2$ columns, and kernel matrix \mathbf{K} has size of $N R^2 \times N R^2$.

Theorem (Low rank Adjustment for the approximate Hessian \mathbf{H})

The approximate Hessian \mathbf{H} can be decomposed as

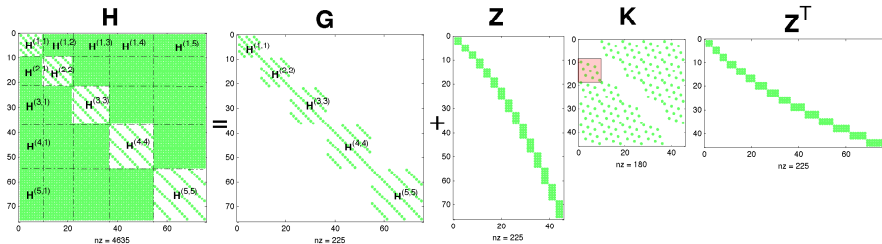
$$\mathbf{H} = \mathbf{G} + \mathbf{Z}\mathbf{K}\mathbf{Z}^T.$$

$$\mathbf{G} = \mathbf{P} \left(\bigoplus_{n=1}^N \bigoplus_{i_n=1}^{I_n} \left(\Gamma^{(n,n)} + \text{diag}\{\alpha_n \mathbf{a}_{i_n}^{(n)\bullet[-2]} + \mu\} \right) \right) \mathbf{P}^T,$$

$$\mathbf{Z} = \bigoplus_{n=1}^N \left(\mathbf{I}_R \otimes \mathbf{A}^{(n)} \right) \in \mathbb{R}^{\prod I_n \times N R^2},$$

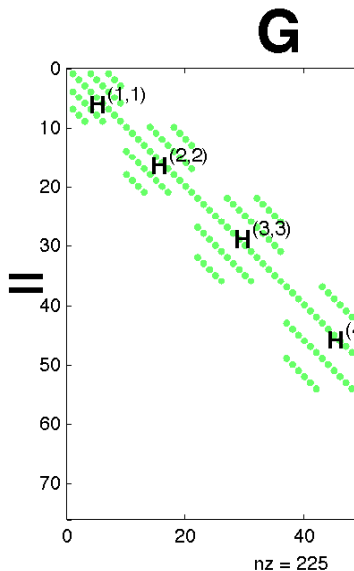
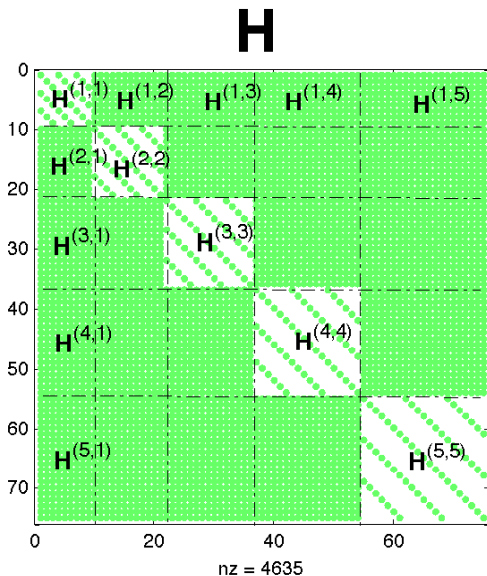
$$\mathbf{K} = \left[\mathbf{K}^{(n,m)} \right]_{n,m} \in \mathbb{R}^{N R^2 \times N R^2}, \quad \mathbf{K}^{(n,m)} = \begin{cases} \Pi_{R,R} \text{diag}(\text{vec}(\Gamma^{(n,m)})) & m \neq n, \\ \mathbf{0}, & m = n. \end{cases}$$

Low Rank Adjustment for \mathbf{H}

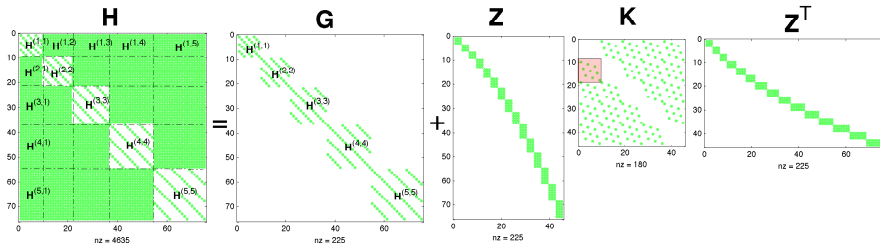


Approximate Hessian for a 5-D tensor. Green dots indicate nonzero elements.

Low Rank Adjustment for \mathbf{H}



Low Rank Adjustment for \mathbf{H}

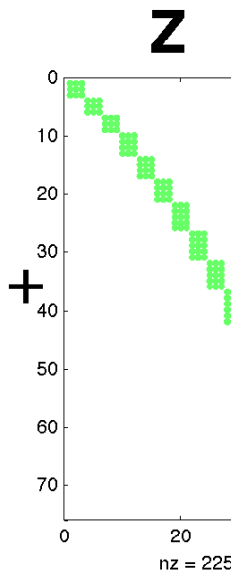
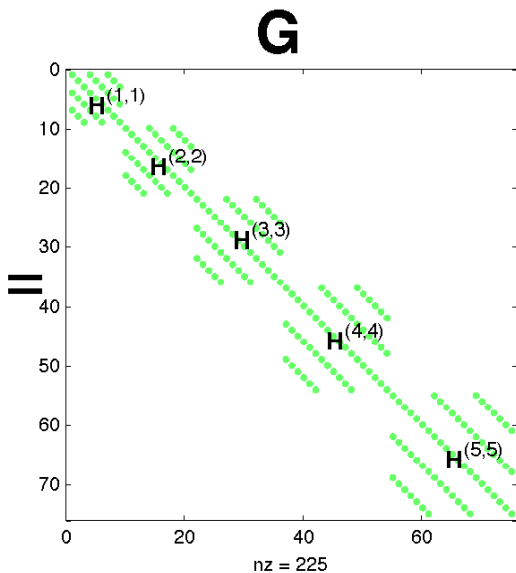
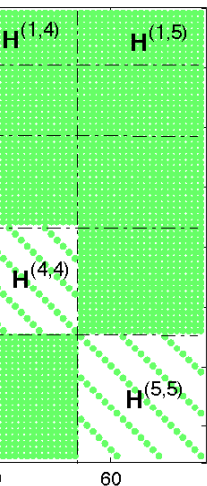


Approximate Hessian for a 5-D tensor. Green dots indicate nonzero elements.

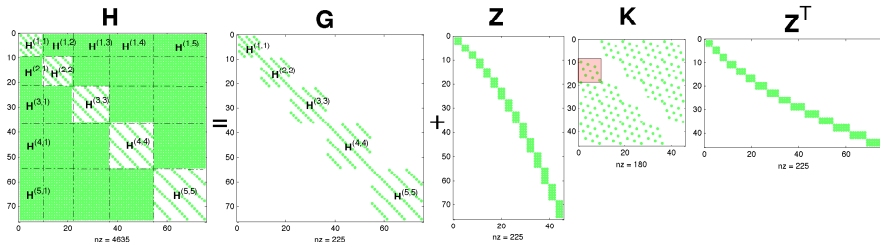
- ▶ According to the binomial inverse theorem Horn and Johnson (2012)

$$\mathbf{H}^{-1} = \mathbf{G}^{-1} - \mathbf{G}^{-1} \mathbf{Z} (\mathbf{K}^{-1} + \mathbf{Z}^T \mathbf{G}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{G}^{-1}.$$

Low Rank Adjustment for \mathbf{H}



Low Rank Adjustment for \mathbf{H}



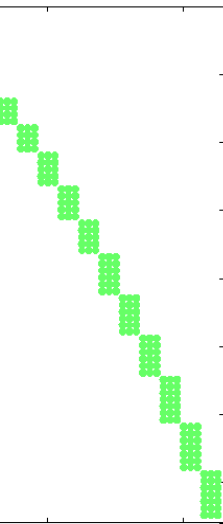
Approximate Hessian for a 5-D tensor. Green dots indicate nonzero elements.

$$\blacktriangleright \mathbf{G}^{-1} = \mathbf{P} \left(\bigoplus_{n=1}^N \bigoplus_{i_n=1}^{I_n} \Theta_{i_n}^{(n)} \right) \mathbf{P}^T,$$

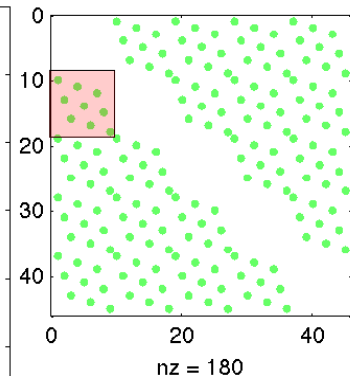
$$\text{where } \Theta_{i_n}^{(n)} = \left(\Gamma^{(n,n)} + \text{diag} \left\{ \alpha_n \mathbf{a}_{i_n:}^{(n) \bullet [-2]} + \mu \right\} \right)^{-1} \in \mathbb{R}^{R \times R}.$$

Low Rank Adjustment for \mathbf{H}

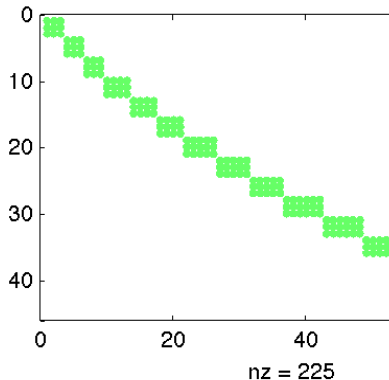
\mathbf{Z}



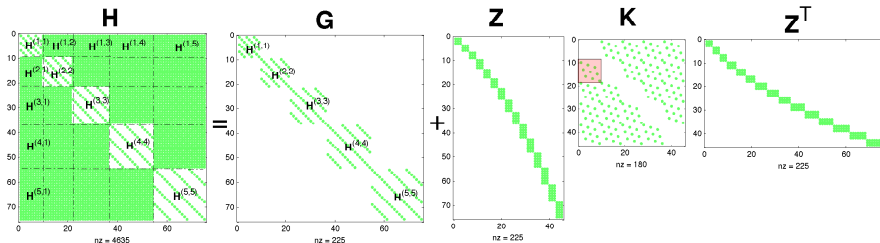
\mathbf{K}



\mathbf{Z}^T



Low Rank Adjustment for \mathbf{H}



Approximate Hessian for a 5-D tensor. Green dots indicate nonzero elements.

► $\mathcal{K} = \mathbf{K}^{-1}$ is a partitioned matrix

$$\mathcal{K}^{(n,m)} = \begin{cases} -\frac{N-2}{N-1} \Pi_{R,R} \text{diag} \left(\text{vec}(\mathbf{A}^{(n)T} \mathbf{A}^{(n)}) \oslash \text{vec}(\Gamma^{(n,n)}) \right), & n = m, \\ \frac{1}{N-1} \Pi_{R,R} \text{diag} \left(\mathbf{1} \oslash \text{vec}(\Gamma^{(n,m)}) \right), & n \neq m. \end{cases}$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}, \quad (34)$$

Fast dGN for NTF

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{G}^{-1} \mathbf{g} + \mathbf{L} \Phi^{-1} (\mathbf{L}^T \mathbf{g}), \quad (35)$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}, \quad (34)$$

Fast dGN for NTF

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{G}^{-1} \mathbf{g} + \mathbf{L} \Phi^{-1} (\mathbf{L}^T \mathbf{g}), \quad (35)$$

$$\Phi = \mathfrak{K} + (\mathbf{I}_N \otimes \Pi_{R,R}) \bigoplus_{n=1}^N \mathbf{W}^{(n)} (\mathbf{I}_N \otimes \Pi_{R,R}) \in \mathbb{R}^{NR^2 \times NR^2}$$

$$\mathbf{W}^{(n)} = \begin{bmatrix} \mathbf{w}_{r,s}^{(n)} \end{bmatrix}, \quad \mathbf{w}_{r,s}^{(n)} = \sum_{i_n=1}^{I_n} \Theta_{i_n}^{(n)} a_{i_n,r}^{(n)} a_{i_n,s}^{(n)} \in \mathbb{R}^{R \times R}$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}, \quad (34)$$

Fast dGN for NTF

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{G}^{-1} \mathbf{g} + \mathbf{L} \Phi^{-1} (\mathbf{L}^T \mathbf{g}), \quad (35)$$

$$\Phi = \mathfrak{K} + (\mathbf{I}_N \otimes \Pi_{R,R}) \left(\bigoplus_{n=1}^N \mathbf{W}^{(n)} \right) (\mathbf{I}_N \otimes \Pi_{R,R}) \in \mathbb{R}^{NR^2 \times NR^2}$$

$$\mathbf{L} = \mathbf{G}^{-1} \mathbf{Z}^T = \mathbf{P} \left(\bigoplus_{n=1}^N \Upsilon^{(n)} \right) (\mathbf{I}_N \otimes \Pi_{R,R}) \in \mathbb{R}^{R \sum l_n \times NR^2}$$

$$\Upsilon^{(n)} = \left[\mathbf{a}_{1:}^{(n)T} \otimes \Theta_1^{(n)} \cdots \mathbf{a}_{l_n:}^{(n)T} \otimes \Theta_{l_n}^{(n)} \right]^T \in \mathbb{R}^{R l_n \times R^2}$$

Damped Gauss-Newton update rule

$$\mathbf{a} \leftarrow \mathbf{a} - (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}, \quad (34)$$

Fast dGN for NTF

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{G}^{-1} \mathbf{g} + \mathbf{L} \Phi^{-1} (\mathbf{L}^T \mathbf{g}), \quad (35)$$

- ▶ For $R \ll I_n, \forall n$, Φ is much smaller than $\mathbf{H} \in \mathbb{R}^{RT \times RT}$.
- ▶ Φ^{-1} requires computational complexity of order $O(N^3 R^6)$, while \mathbf{H}^{-1} is of order $O(N^3 R^3 I^3)$ ($I_1 = \dots = I_N = I$).
- ▶ We don't need to construct the approximate Hessian \mathbf{H} and the Jacobian \mathbf{J} .

Fast LM with cost $O(NR^6)$

Fast inverse of the approximate Hessian \mathbf{H} ($O(N^3R^6)$)

$$\mathbf{H}^{-1} = \mathbf{G}^{-1} - \mathbf{G}^{-1} \mathbf{Z} (\mathbf{K}^{-1} + \mathbf{Z}^T \mathbf{G}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{G}^{-1}.$$

Fast LM with cost $O(NR^6)$

Fast inverse of the approximate Hessian \mathbf{H} ($O(N^3R^6)$)

$$\mathbf{H}^{-1} = \mathbf{G}^{-1} - \mathbf{G}^{-1} \mathbf{Z} (\mathbf{K}^{-1} + \mathbf{Z}^T \mathbf{G}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{G}^{-1}.$$

Low-rank adjustment form of \mathbf{K}

$$\mathbf{K} = \mathbf{K}_0 + \mathbf{D} \mathbf{F} \mathbf{D}^T \in \mathbb{R}^{NR^2 \times NR^2}$$

$$\mathbf{K}_0 = \text{bdiag}(\mathbf{K}_n)_{n=1}^N, \quad \mathbf{K}_n = -\mathbf{P}_R \text{dvec}(\Gamma_{nn} \oslash \mathbf{C}_n)$$

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}_1 \\ \vdots \\ \mathbf{D}_N \end{bmatrix}, \quad \mathbf{D}_n = \text{dvec}(1 \oslash \mathbf{C}_n)$$

$$\mathbf{F} = \mathbf{P}_R \text{dvec}(\Gamma) \in \mathbb{R}^{R^2 \times R^2}, \quad \Gamma = \bigotimes_{n=1}^N \mathbf{C}_n.$$

Fast LM with cost $O(NR^6)$

Fast inverse of the approximate Hessian \mathbf{H} ($O(N^3R^6)$)

$$\mathbf{H}^{-1} = \mathbf{G}^{-1} - \mathbf{G}^{-1} \mathbf{Z} (\mathbf{K}^{-1} + \mathbf{Z}^T \mathbf{G}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{G}^{-1}.$$

Fast inverse of \mathbf{K} ($O(NR^6)$)

$$\mathbf{K}^{-1} = \mathbf{K}_0^{-1} - \mathbf{K}_0^{-1} \mathbf{D} (\mathbf{F}^{-1} + \mathbf{D}^T \mathbf{K}_0^{-1} \mathbf{D})^{-1} \mathbf{D}^T \mathbf{K}_0^{-1}$$

Fast LM with cost $O(NR^6)$

Fast inverse of the approximate Hessian \mathbf{H} ($O(N^3R^6)$)

$$\mathbf{H}^{-1} = \mathbf{G}^{-1} - \mathbf{G}^{-1} \mathbf{Z} (\mathbf{K}^{-1} + \mathbf{Z}^T \mathbf{G}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{G}^{-1}.$$

Fast inverse of \mathbf{H} ($O(NR^6)$)

$$\mathbf{H}^{-1} = \mathbf{G}^{-1} - \mathbf{G}^{-1} \mathbf{Z} (\mathbf{J}_\mu^{-1} + \mathbf{L}_\mu \mathbf{S}_\mu^{-1} \mathbf{L}_\mu^T) \mathbf{Z}^T \mathbf{G}^{-1}$$

where

$$\mathbf{S}_\mu = \mathbf{Q} - \sum_{n=1}^N \mathbf{D}_n \mathbf{K}_n^{-1} \mathbf{J}_n^{-1} \mathbf{K}_n^{-1} \mathbf{D}_n$$

$$\mathbf{L}_\mu = \mathbf{J}_\mu^{-1} \mathbf{K}_0^{-1} \mathbf{D}.$$

$$\mathbf{Q} = -(N-1) \mathbf{P}_R \text{dvec}(1 \otimes \Gamma)$$

$$\mathbf{J}_\mu = \text{bdiag}(\mathbf{J}_n)_{n=1}^N,$$

$$\mathbf{J}_n = -\mathbf{P}_R \text{dvec}(\mathbf{C}_n \otimes \Gamma_{nn}) + (\Gamma_{nn} + \mu \mathbf{I}_R)^{-1} \otimes \mathbf{C}_n.$$

Selection of barrier and sparse parameters

- ▶ An heuristic approach: initialize $\alpha = \alpha_1 = \dots = \alpha_N$ by a large enough value, then slowly descend down to near-zero after each 10 iterations ?.
- ▶ An alternative approach is based on the Karush-Kuhn-Tucker condition Rojas and Steihaug (2002); Ding et al. (2006)
maximize α subject to $\alpha \mathbf{A}^{(n)} \leq \mathbf{F}^{(n)}$
where $\mathbf{A}^{(n)} = [\alpha \mathbf{a}_1^{(n)} \dots \alpha \mathbf{a}_N^{(n)}]$ and $\mathbf{F}^{(n)} = [\mathbf{f}_1^{(n)} \dots \mathbf{f}_N^{(n)}]$
- ▶ For NTF without P_s ,
$$\alpha_n = \max\left(0, \min\left(\text{vec}\left(-\mathbf{A}^{(n)} \oslash \mathbf{F}^{(n)}\right)\right)\right), \quad \forall n.$$

Selection of barrier and sparse parameters

- ▶ An heuristic approach: initialize $\alpha = \alpha_1 = \dots = \alpha_N$ by a large enough value, then slowly descend down to near-zero after each 10 iterations ?.
- ▶ An alternative approach is based on the Karush-Kuhn-Tucker condition Rojas and Steihaug (2002); Ding et al. (2006)

$$\mathbf{a} \geq 0, \quad \mathbf{a} \circledast \mathbf{g} = 0, \quad \mathbf{g} \geq 0.$$

$$\mathbf{A}^{(n)} \circledast \left(\mathbf{F}^{(n)} + \alpha_n \left(\mathbf{A}^{(n)} \right)^{\bullet[-1]} - \beta_n \right) = \mathbf{0}, \quad \forall n, \quad (36)$$

$$\mathbf{F}^{(n)} + \alpha_n \left(\mathbf{A}^{(n)} \right)^{\bullet[-1]} - \beta_n \leq \mathbf{0}, \quad \forall n. \quad (37)$$

$$\min_{\mathbf{x}} \|\mathbf{B}^{(n)} \mathbf{x} - \mathbf{u}^{(n)}\|_2^2 \quad \text{such that} \quad \begin{cases} \mathbf{B}^{(n)} \mathbf{x} \leq \mathbf{u}^{(n)}, \\ \mathbf{x} = [\alpha_n \beta_n]^T \geq \mathbf{0}, \end{cases} \quad \forall n,$$

where $\mathbf{B}^{(n)} = [\mathbf{1} \quad \text{vec}(-\mathbf{A}^{(n)})]$, and $\mathbf{u}^{(n)} = \text{vec}(-\mathbf{A}^{(n)} \circledast \mathbf{F}^{(n)})$.

- ▶ For NTF without P_s ,

$$\alpha_n = \max\left(0, \min\left(\text{vec}(-\mathbf{A}^{(n)} \circledast \mathbf{F}^{(n)})\right)\right), \quad \forall n,$$

Selection of barrier and sparse parameters

- ▶ An heuristic approach: initialize $\alpha = \alpha_1 = \dots = \alpha_N$ by a large enough value, then slowly descend down to near-zero after each 10 iterations ?.
- ▶ An alternative approach is based on the Karush-Kuhn-Tucker condition Rojas and Steihaug (2002); Ding et al. (2006)

$$\min_{\mathbf{x}} \|\mathbf{B}^{(n)} \mathbf{x} - \mathbf{u}^{(n)}\|_2^2 \quad \text{such that} \quad \begin{cases} \mathbf{B}^{(n)} \mathbf{x} \leq \mathbf{u}^{(n)}, \\ \mathbf{x} = [\alpha_n \beta_n]^T \geq \mathbf{0}, \end{cases} \quad \forall n,$$

where $\mathbf{B}^{(n)} = [\mathbf{1} \quad \text{vec}(-\mathbf{A}^{(n)})]$, and $\mathbf{u}^{(n)} = \text{vec}(-\mathbf{A}^{(n)} \circledast \mathbf{F}^{(n)})$.

- ▶ For NTF without P_s ,
 $\alpha_n = \max(0, \min(\text{vec}(-\mathbf{A}^{(n)} \circledast \mathbf{F}^{(n)})))$, $\forall n$.

Selection of barrier and sparse parameters

- ▶ An heuristic approach: initialize $\alpha = \alpha_1 = \dots = \alpha_N$ by a large enough value, then slowly descend down to near-zero after each 10 iterations ?.
- ▶ An alternative approach is based on the Karush-Kuhn-Tucker condition Rojas and Steihaug (2002); Ding et al. (2006)

$$\min_{\mathbf{x}} \|\mathbf{B}^{(n)} \mathbf{x} - \mathbf{u}^{(n)}\|_2^2 \quad \text{such that} \quad \begin{cases} \mathbf{B}^{(n)} \mathbf{x} \leq \mathbf{u}^{(n)}, \\ \mathbf{x} = [\alpha_n \beta_n]^T \geq \mathbf{0}, \end{cases} \quad \forall n,$$

where $\mathbf{B}^{(n)} = [\mathbf{1} \quad \text{vec}(-\mathbf{A}^{(n)})]$, and $\mathbf{u}^{(n)} = \text{vec}(-\mathbf{A}^{(n)} \circledast \mathbf{F}^{(n)})$.

- ▶ For NTF without P_s ,
 $\alpha_n = \max\left(0, \min\left(\text{vec}(-\mathbf{A}^{(n)} \circledast \mathbf{F}^{(n)})\right)\right), \quad \forall n.$

fast LM vs damped Gauss-Newton (dGN)

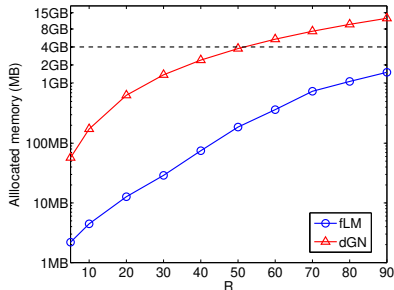
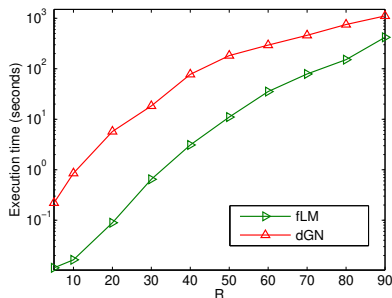


Figure: Memory requirements and execution time per iteration of fLM algorithms in approximation of $100 \times 100 \times 100$ dimensional tensors by rank- R tensors where $R = 5, 10, 20, \dots, 90$.

dGN: damped Gauss-Newton algorithm in the package PARAFAC3W.Tomasi (2003)

References I

- Andersson, C. and Bro, R. (2000). The N-way toolbox for MATLAB. *Chemometrics and Intelligent Laboratory Systems*, 52(1):1–4.
- Berry, M., Browne, M., Langville, A., Pauca, P., and Plemmons, R. (2007). Algorithms and applications for approximate nonnegative matrix factorization. *Computational Statistics and Data Analysis*, 52(1):155–173.
- Booksh, K. and Kowalski, B. (1994). Error analysis of the generalized rank annihilation method. *Journal of Chemometrics*, 8.
- Booksh, K., Lin, Z., Wang, Z., and Kowalski, B. (1994). Extension of trilinear decomposition method with an application to the flow probe sensor. *Analytical Chemistry*, 66.
- Carroll, J. and Chang, J. (1970). Analysis of individual differences in multidimensional scaling via an n -way generalization of Eckart–Young decomposition. *Psychometrika*, 35(3):283–319.

References II

- Comon, P., Luciani, X., and de Almeida, A. L. F. (2009). Tensor decompositions, alternating least squares and other tales. *Journal of Chemometrics*, 23.
- Comon, P. and Rajih, M. (2006). Blind identification of under-determined mixtures based on the characteristic function. *Signal Processing*, 86(9):2271 – 2281. Special Section: Signal Processing in {UWB} Communications.
- Ding, C., Li, T., Peng, W., and Park, H. (2006). Orthogonal nonnegative matrix tri-factorizations for clustering. In *KDD06*, pages 126–135, New York, NY, USA. ACM Press.
- Eigenvector (2012). PLS Toolbox.
- Faber, N., Buydens, L., and Kateman, G. (1994). Generalized rank annihilation. iii: Practical implementation. *Journal of Chemometrics*, 8.

References III

- Harshman, R. (1970). Foundations of the PARAFAC procedure: Models and conditions for an explanatory multimodal factor analysis. *UCLA Working Papers in Phonetics*, 16:1–84.
- Harshman, R. A. (1972). Determination and proof of minimum uniqueness conditions for PARAFAC1. *UCLA Working Papers in Phonetics*, 22.
- Hitchcock, F. (1927). Multiple invariants and generalized rank of a p -way matrix or tensor. *Journal of Mathematics and Physics*, 7:39–79.
- Horn, R. A. and Johnson, C. R. (1991). *Topics in matrix analysis*. Cambridge University Press, Cambridge.
- Horn, R. A. and Johnson, C. R. (2012). *Matrix Analysis*. Cambridge University Press.

References IV

- Jaderberg, M., Vedaldi, A., and Zisserman, A. (2014). Speeding up convolutional neural networks with low rank expansions. *CoRR*, abs/1405.3866.
- Kolda, T. and Bader, B. (2009). Tensor decompositions and applications. *SIAM Review*, 51(3):455–500.
- Kruskal, J. (1977). Three-way arrays: Rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear Algebra Appl.*, 18:95–138.
- Langville, A. N., Meyer, C. D., and Albright, R. (2006). Initializations for the nonnegative matrix factorization. In *Proc. of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Philadelphia, USA.
- Nion, D. and Sidiropoulos, N. (2009). Adaptive Algorithms to Track the PARAFAC Decomposition of a Third-Order Tensor. *IEEE Transactions on Signal Processing*, 57(6):2299–2310.

- Phan, A.-H., Tichavský, P., and Cichocki, A. (2015). Tensor deflation for CANDECOMP/PARAFAC. Part 1: Alternating Subspace Update Algorithm. *IEEE Transaction on Signal Processing*, 63(12):5924–5938.
- Ragnarsson, S. and Loan, C. F. V. (2012). Block tensor unfoldings. *SIAM J. Matrix Analysis Applications*, 33(1):149–169.
- Rajih, M., Comon, P., and Harshman, R. A. (2008). Enhanced line search: A novel method to accelerate PARAFAC. *SIAM Journal of Matrix Analysis and Applications*, 30(3):1128–1147.
- Rojas, M. and Steihaug, T. (2002). An interior-point trust-region-based method for large-scale non-negative regularization. *Inverse Problems*, 18:1291–1307.
- Sidiropoulos, N., Bro, R., and Giannakis, G. (2000). Parallel factor analysis in sensor array processing. *IEEE Transactions on Signal Processing*, 48(8):2377–2388.

References VI

Strassen, V. (1969). Gaussian elimination is not optimal. *Numer. Math.*, 13(4):354–356.

Tomasi, G. (2003). INDAFAC and PARAFAC3W.
<http://www.models.kvl.dk/source/indafac/index.asp>.

Tomasi, G. (2006). *Practical and Computational Aspects in Chemometric Data Analysis*. PhD thesis, Royal Veterinary and Agricultural University, Frederiksberg, Denmark.

Yeredor, A. (2000). Blind source separation via the second characteristic function. *Signal Processing*, 80(5):897–902.