

Parallel Algorithms for Tensor Completion in the CP Format

Daniel Kressner

Chair of Numerical Algorithms and HPC
MATHICSE / SMA / SB / EPF Lausanne

daniel.kressner@epfl.ch <http://anchp.epfl.ch>

Joint work with:

Lars Karlsson (Umeå University)

André Uschmajew (University of Bonn)

SIAM PP 2016


14.4.2016



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

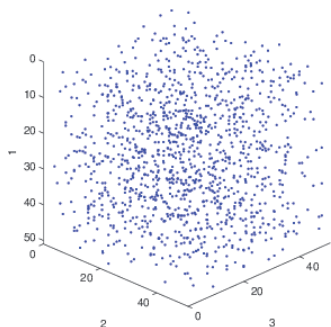


cadmos

center for advanced modeling science 

Tensor completion

Goal: Complete multivariate data.



Applications:

- ▶ Completion of corrupted hyperspectral images, CT Scans, ...
- ▶ Learning of multivariate functions
- ▶ Non-intrusive methods for stochastic/parametric PDEs
- ▶ Context-aware recommender systems
- ▶ ...

Tensor completion

Goal: Complete multivariate data.

Mathematical setting:

- ▶ $I_1 \times I_2 \times \cdots \times I_N$ tensor \mathcal{X} with very few entries known.
- ▶ $\Omega \subset [1, I_1] \cdots \times \cdots [1, I_N]$ contains (multi)indices of known entries.
- ▶ $P_\Omega : \mathbb{R}^{I_1 \times \cdots \times I_N} \rightarrow \mathbb{R}^{|\Omega|}$ orthogonal projection onto known entries.

Tensor completion:

$$\min_{\mathcal{X}} \frac{1}{2} \|\text{known entries} - P_\Omega \mathcal{X}\|^2$$

- ▶ **Ill-posed problem.**
- ▶ Regularize with (multilinear) low-dimensional model for \mathcal{X} .

Low-rank tensor completion

Goal: Complete multivariate data.

Mathematical setting:

- ▶ $I_1 \times I_2 \times \cdots \times I_N$ tensor \mathcal{X} with very few entries known.
- ▶ $\Omega \subset [1, I_1] \cdots \times \cdots [1, I_N]$ contains (multi)indices of known entries.
- ▶ $P_\Omega : \mathbb{R}^{I_1 \times \cdots \times I_N} \rightarrow \mathbb{R}^{|\Omega|}$ orthogonal projection onto known entries.

Low-rank tensor completion:

$$\begin{aligned} \min_{\mathcal{X}} \quad & \frac{1}{2} \|P_\Omega \mathcal{A} - P_\Omega \mathcal{X}\|^2 \\ \text{subject to} \quad & \mathcal{X} \text{ has tensor rank } L \end{aligned}$$

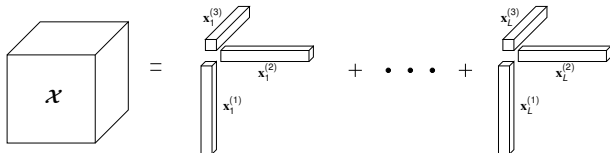
Alternatives to tensor rank (CP decomposition):

- ▶ Multilinear rank (Tucker decomposition) and mixtures.
- ▶ TT ranks/decomposition, HT ranks/decomposition, general tensor networks.

CANDECOMP/PARAFAC (CP) decomposition

\mathcal{X} has **tensor rank** L if it admits **CP decomposition**

$$\mathcal{X} = \sum_{\ell=1}^L \mathbf{x}_{\ell}^{(1)} \circ \mathbf{x}_{\ell}^{(2)} \circ \dots \circ \mathbf{x}_{\ell}^{(N)}, \quad \mathbf{x}_{\ell}^{(n)} \in \mathbb{R}^{I_n}. \quad (1)$$



Properties of CP:

- ▶ Low data complexity: For constant L , **linear (instead of exponential) complexity in N** .
- ▶ Linear wrt *each* **factor matrix** $\mathbf{X}^{(n)} = [\mathbf{x}_1^{(n)}, \mathbf{x}_2^{(n)}, \dots, \mathbf{x}_L^{(n)}] \in \mathbb{R}^{I_n \times L}$.
- ▶ Tensor rank L *not* upper semi-continuous.

Low-rank tensor completion

Inserting CP into low-rank tensor completion \rightsquigarrow

$$\min \left\| P_{\Omega} \left(\mathcal{A} - \sum_{\ell=1}^L \mathbf{x}_{\ell}^{(1)} \circ \mathbf{x}_{\ell}^{(2)} \circ \dots \circ \mathbf{x}_{\ell}^{(N)} \right) \right\|_2^2$$

This will *not* work well because of

- ▶ Non-uniqueness:

$$\alpha \mathbf{x}_{\ell}^{(1)} \circ \dots \circ \frac{1}{\alpha} \mathbf{x}_{\ell}^{(N)} = \mathbf{x}_{\ell}^{(1)} \circ \dots \circ \mathbf{x}_{\ell}^{(N)}$$

for any $\alpha \neq 0$.

- ▶ Unboundedness: $\|\mathbf{x}_{\ell}^{(n)}\|$ may become arbitrarily large.

Fix: Penalize factors with large norms.

Low-rank tensor completion

Low-rank tensor completion with regularization:

$$\min \left\| P_{\Omega} \left(\mathcal{A} - \sum_{\ell=1}^L \mathbf{x}_{\ell}^{(1)} \circ \mathbf{x}_{\ell}^{(2)} \circ \dots \circ \mathbf{x}_{\ell}^{(N)} \right) \right\|_2^2 + \lambda \sum_{\ell=1}^L \sum_{n=1}^N \|\mathbf{x}_{\ell}^{(n)}\|_2^2$$

- ▶ λ regularization parameter (typically small, e.g., $\lambda = 10^{-3}$)
- ▶ nonlinear, nonconvex optimization problem
- ▶ no closed-form solution in terms of SVD for $N \geq 3$
- ▶ quadratic convex problem in *each* individual factor matrix $\mathbf{X}^{(n)} = [\mathbf{x}_1^{(n)}, \dots, \mathbf{x}_L^{(n)}] \rightsquigarrow$ **alternating optimization methods!**

Previous parallel approaches to low-rank completion

$N = 2$ (matrix completion):

- ▶ ALS (alternating least-squares)
[Teflioudi/Makari/Gemulla'2012, Zhou/Wilkinson/Schreiber/Pan'2008]
- ▶ CCD (cyclic coordinate descent)
[Pilászy/Zibriczky/Tikk'2010, Yu et al.'2012–2013]
- ▶ SGD (stochastic gradient descent)
[Gemulla et al.'2011, Recht/Ré'2013, Makari et al.'2015]

$N > 2$ (tensor completion):

- ▶ Well studied: Complete data = approximation of complete tensor
- ▶ Incomplete data usually “fixed” via weight matrices or imputation
[Hidasi/Tikk'2012], [Acar et al.2011]
- ▶ Alternating proximal gradient method under additional nonnegativity constraints
[Xu/Yin'2013]
- ▶ Parallelization of ALS based on local CP models
[Phan/Cichoki'2011]

ALS

ALS (Alternating Least-Squares)

- ▶ ALS alternately optimizes factor matrices $\mathbf{X}^{(n)} = [\mathbf{x}_1^{(n)}, \dots, \mathbf{x}_L^{(n)}]$ for $n = 1, \dots, N$.
- ▶ Optimization of $I_1 \times L$ factor matrix $\mathbf{X}^{(1)}$ (while keeping all other factor matrices fixed):

$$\min \left\| P_{\Omega} \left(\mathcal{A} - \sum_{\ell=1}^L \mathbf{x}_{\ell}^{(1)} \circ \mathbf{x}_{\ell}^{(2)} \circ \dots \circ \mathbf{x}_{\ell}^{(N)} \right) \right\|_2^2 + \lambda \dots$$

This **decouples** wrt rows of $\mathbf{X}^{(1)}$.

- ▶ I_1 decoupled optimization problems:

$$\min_{\substack{\mathbf{i} \in \Omega \\ i_1 = \hat{i}}} \left(\mathbf{a}_{\mathbf{i}} - \sum_{\ell=1}^L \mathbf{z}_{\ell} \prod_{n=2}^N [\mathbf{x}_{\ell}^{(n)}]_{i_n} \right)^2 + \lambda \sum_{\ell=1}^L \mathbf{z}_{\ell}^2, \quad \hat{i} = 1, \dots, I_1.$$

\hat{i} th row of $\mathbf{X}^{(1)}$ determined by $|\Omega_{1, \hat{i}}| \times L$ LSQ:

$$\min \|\mathbf{a}_{1, \hat{i}} - \mathbf{H}_{1, \hat{i}} \mathbf{z}\|_2^2 + \lambda \|\mathbf{z}\|_2^2.$$

Solution $\mathbf{z}^* = (\mathbf{H}_{1, \hat{i}}^T \mathbf{H}_{1, \hat{i}} + \lambda I)^{-1} \mathbf{H}_{1, \hat{i}}^T \mathbf{a}_{1, \hat{i}}$.

ALS

Algorithm 1: ALS

```
1 Initialize all vectors  $\mathbf{x}_\ell^{(n)}$ .
2 for each outer iteration do
3   for  $\hat{n} = 1, 2, \dots, N$  do
4     for  $\hat{i} = 1, 2, \dots, I_{\hat{n}}$  do in parallel
5       Set up and solve  $\mathbf{z}^* = (\mathbf{H}_{\hat{n},\hat{i}}^\top \mathbf{H}_{\hat{n},\hat{i}} + \lambda I)^{-1} \mathbf{H}_{\hat{n},\hat{i}}^\top \mathbf{a}_{\hat{n},\hat{i}}$ .
6       Update  $L$  entries  $[\mathbf{x}_\ell^{(\hat{n})}]_{\hat{i}}, \ell = 1, 2, \dots, L$ , using  $\mathbf{z}^*$ .
```

Tempting parallelization:

- ▶ Parallelize wrt variables in inner loop.
- ▶ Extension of distributed ALS [Teflioudi/Makari/Gemulla'2012] to tensors: Each node owns corresponding data $\mathbf{a}_{\hat{n},\hat{i}}$.
- ▶ **Requires N replications of \mathcal{A} !**

Parallelization of ALS

- ▶ $\Omega = \Omega^{(1)} \cup \dots \cup \Omega^{(p)}$ partition of index set Ω over p processors
- ▶ Processor q owns a_i for all $i \in \Omega^{(q)}$.
- ▶ Factor matrices $\mathbf{X}^{(n)}$ for $n = 1, \dots, N$ of CP decomposition of \mathcal{X} replicated on all p processors.

Algorithm 2: Parallel formulation of ALS.

- 1 Initialize all vectors $\mathbf{x}_\ell^{(n)}$ in parallel;
 - 2 **for** *each outer iteration* **do**
 - 3 **for** $\hat{n} = 1, 2, \dots, N$ **do**
 - 4 Construct local contributions to $\mathbf{C}_{\hat{n}} := \mathbf{H}_{\hat{n}, \hat{n}}^\top \mathbf{H}_{\hat{n}, \hat{n}}$ and $\mathbf{d}_{\hat{n}} := \mathbf{H}_{\hat{n}, \hat{n}}^\top \mathbf{a}_{\hat{n}, \hat{n}}$ for $\hat{i} = 1, \dots, l_{\hat{n}}$.
 - 5 Reduce local contributions.
 - 6 Scatter $l_{\hat{n}}$ coefficients $\mathbf{C}_{\hat{n}}, \mathbf{d}_{\hat{n}}$ across all processors.
 - 7 Solve (roughly $l_{\hat{n}}/p$) local normal equations $(\mathbf{C}_{\hat{n}} + \lambda \mathbf{I}) \mathbf{z}_{\hat{n}} = \mathbf{d}_{\hat{n}}$.
 - 8 Replicate updated factor matrix $\mathbf{X}^{(\hat{n})}$;
-

Cyclic coordinate descent

Cyclic coordinate descent

Motivation: ALS expensive for larger tensor ranks L .

Computation $\mathcal{O}(L^3)$, Memory/communication bandwidth $\mathcal{O}(L^2)$.

Idea:

- ▶ Optimize only one term in $\mathcal{X} = \sum_{\ell=1}^L \mathbf{x}_\ell^{(1)} \circ \mathbf{x}_\ell^{(2)} \circ \dots \circ \mathbf{x}_\ell^{(N)}$.
- ▶ Attempt to optimize $\hat{\ell}$ th term by one sweep of **rank-one ALS** applied to

$$\mathcal{A} - \sum_{\ell \neq \hat{\ell}} \mathbf{x}_\ell^{(1)} \circ \mathbf{x}_\ell^{(2)} \circ \dots \circ \mathbf{x}_\ell^{(N)}.$$

- ▶ \hat{i} th entry of $\mathbf{x}_{\hat{\ell}}^{(\hat{n})}$ is replaced by

$$z^* = \frac{\sum_{\substack{\mathbf{i} \in \Omega \\ i_{\hat{n}} = \hat{i}}} \left(a_{\mathbf{i}} - \sum_{\substack{\ell=1 \\ \ell \neq \hat{\ell}}}^L \prod_{\substack{n=1 \\ n \neq \hat{n}}}^N [\mathbf{x}_\ell^{(n)}]_{i_n} \right) \prod_{\substack{n=1 \\ n \neq \hat{n}}}^N [\mathbf{x}_{\hat{\ell}}^{(n)}]_{i_n}}{\lambda + \sum_{\substack{\mathbf{i} \in \Omega \\ i_{\hat{n}} = \hat{i}}} \left(\prod_{\substack{n=1 \\ n \neq \hat{n}}}^N [\mathbf{x}_{\hat{\ell}}^{(n)}]_{i_n} \right)^2}.$$

Cyclic coordinate descent

Algorithm 3: CCD++

```
1 Initialize all vectors  $\mathbf{x}_\ell^{(n)}$ .
2 for each outer iteration do
3   for  $\hat{\ell} = 1, 2, \dots, L$  do
4     for  $\hat{n} = 1, 2, \dots, N$  do
5       for  $\hat{i} = 1, 2, \dots, I_{\hat{n}}$  do in parallel
6         Update entry  $[\mathbf{x}_{\hat{\ell}}^{(\hat{n})}]_{\hat{i}}$ .
```

- ▶ Proposed by [Yu/Hsieh/Si/Dhillon'2012] for matrix completion.
- ▶ Keeping track of

$$r_{\mathbf{i}} := a_{\mathbf{i}} - \sum_{\substack{\ell=1 \\ \ell \neq \hat{\ell}}}^L \prod_{n=1}^N [\mathbf{x}_\ell^{(n)}]_{i_n}$$

reduces cost by a factor $\approx 1/N$, at the expense of storing sparse tensor \mathcal{R} .

- ▶ Convergence to critical point follows from [Xu/Yin'2013].

Parallel cyclic coordinate descent

- ▶ $\Omega = \Omega^{(1)} \cup \dots \cup \Omega^{(p)}$ partition of index set Ω over p processors
- ▶ Processor q owns a_i and r_i for all $i \in \Omega^{(q)}$.
- ▶ Factor matrices $\mathbf{X}^{(n)}$ for $n = 1, \dots, N$ of CP decomposition of \mathcal{X} replicated on all p processors.

Algorithm 4: Parallel CCD++

- 1 Initialize in parallel all vectors $\mathbf{x}_\ell^{(n)}$ and \mathcal{R} .
- 2 **for** each outer iteration **do**
- 3 **for** $\hat{\ell} = 1, 2, \dots, L$ **do**
- 4 **for** $\hat{n} = 1, 2, \dots, N$ **do**
- 5 Construct local contributions to tensor contractions
 $\alpha_{\hat{i}} = \sum_{\substack{i \in \Omega \\ i_{\hat{n}} = \hat{i}}} r_i \gamma_i$ and $\beta_{\hat{i}} = \sum_{\substack{i \in \Omega \\ i_{\hat{n}} = \hat{i}}} \prod_{\substack{n=1 \\ n \neq \hat{n}}}^N [\mathbf{x}_\ell^{(n)}]_{i_n}^2$ for $\hat{i} = 1, \dots, I_{\hat{n}}$.
- 6 Reduce local contributions by summation.
- 7 Perform $[\mathbf{x}_\ell^{(\hat{n})}]_{\hat{i}} \leftarrow \frac{\alpha_{\hat{i}}}{\lambda + \beta_{\hat{i}}}$ on master processor.
- 8 Replicate updated variables $\mathbf{x}_\ell^{(\hat{n})}$.
- 9 Update \mathcal{R} in parallel.

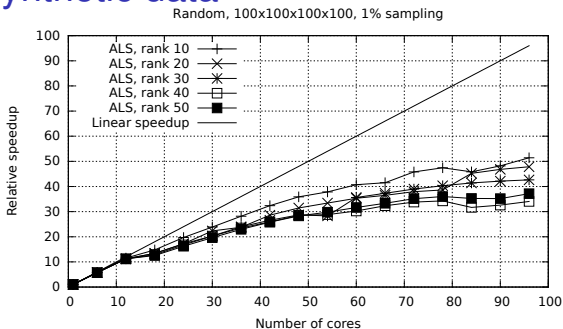
Numerical experiments

Computational environment

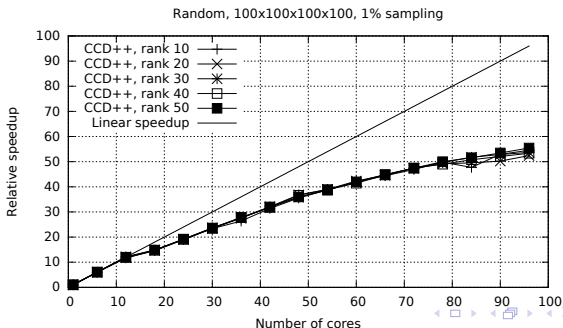
- ▶ Abisko distributed memory system at HPC2N in Umeå/Sweden.
- ▶ Each node contains four sockets with AMD Opteron 6238 processors.
- ▶ Each processor contains 12 cores partitioned into two NUMA domains.
- ▶ Nodes are interconnected with 40 Gb/s Mellanox Infiniband.
- ▶ PathScale C++ compiler with OpenMPI.

Results: Synthetic data

ALS:



CCD++



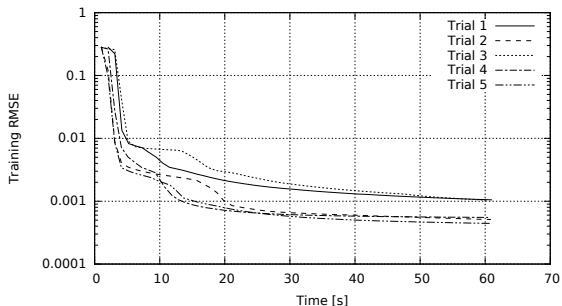
Function-related tensors

- ▶ Discretization of $\mathbf{x} \mapsto \exp(-\|\mathbf{x}\|_2)$, which has a cusp singularity at the origin, on a tensor product grid.
- ▶ $a_{i_1, i_2, \dots, i_N} = \exp(-\sqrt{\xi_{i_1}^2 + \xi_{i_2}^2 + \dots + \xi_{i_N}^2})$
- ▶ Sample $10nNL$ entries of \mathcal{A} .
- ▶ Application of tensor completion: High-dimensional integration.
- ▶ Experiments: $n = 51$, $N = 5$, $L = 100$, 48 cores.

Results: Function-related tensors

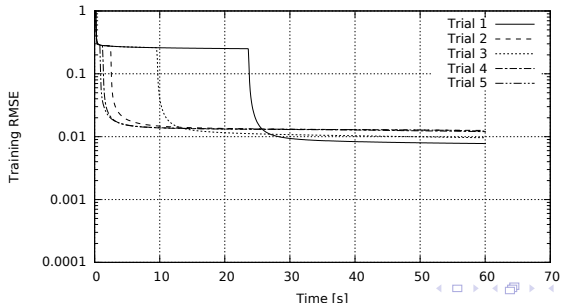
ALS:

Function tensor (N=5, n=51), ALS, rank 100, 255000 samples, 48 cores



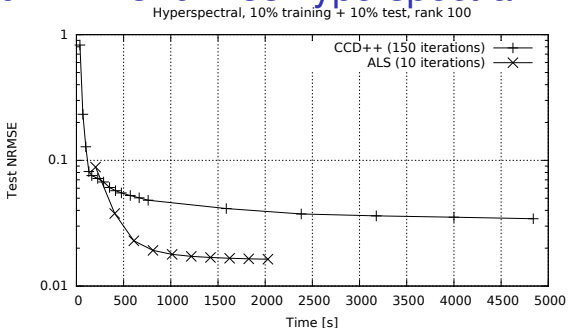
CCD++

Function tensor (N=5, n=51), CCD++, rank 100, 255000 samples, 48 cores

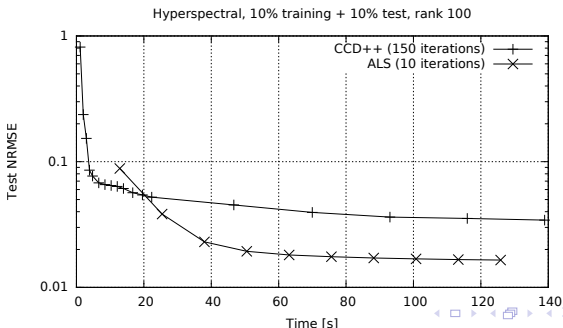


Results: $1017 \times 1340 \times 33$ hyperspectral image

1 core:

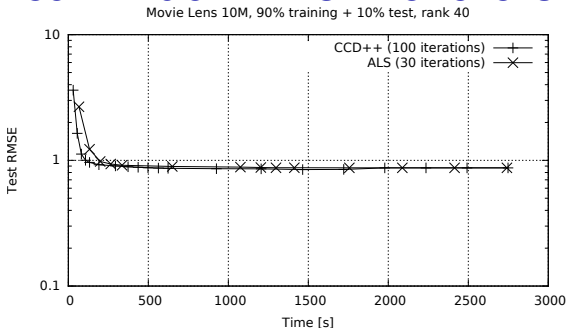


48 cores:

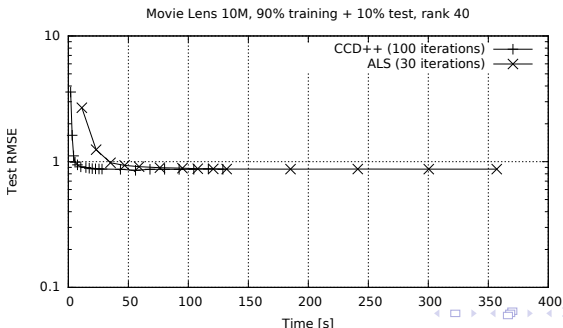


Results: $71\,567 \times 10\,677 \times 731$ MovieLens tensor

1 core:



48 cores:



Conclusions

Conclusions and ongoing work

- ▶ Parallel CCD++ method of choice for data-related applications.
- ▶ Parallel ALS sometimes better when high accuracy is required.
- ▶ Both algorithms are weakly scalable.

More details in:

- ▶ L. Karlsson, D. Kressner, and A. Uschmajew. Parallel algorithms for tensor completion in the CP format. *Parallel Computing*, 2015.

Outlook:

- ▶ (MPI+)Cuda: Ongoing joint work with Efthalia Karydi.
- ▶ Adjust distribution of Ω to increase data locality.
- ▶ Application to UQ. (Talk by Grasedyck on Tuesday)