


Illustration by Chris Briggman






The Canonical Polyadic Tensor Decomposition and Variants for Mining Multi-Dimensional Data

Tammy Kolda*, Danny Dunlavy#
 Sandia National Laboratories
 *Livermore, CA and #Albuquerque, NM

SIAM International Conference on Data Mining
 May 4, 2018


Sandia National Laboratories is a multi-mission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.




Special Thanks to...

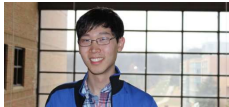
Kina Winoto (Sandia-CA) – For massive help in preparing the tutorial and helping in its presentation in other venues. All the labs that work are thanks to Kina!






Jed Duersch (Sandia-CA) – For preparing the examples for GCP, for creating additional versions of the gas data set, for preparing the social network data set.

David Hong (Univ. Michigan) – For work in preparing some of the labs for the first version of this tutorial in summer 2017.





Plus many other co-authors to be cited throughout the talk.


5/4/2018
Tensor Tutorial @ SDM18
2

Sandia National Laboratories 



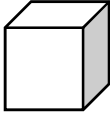
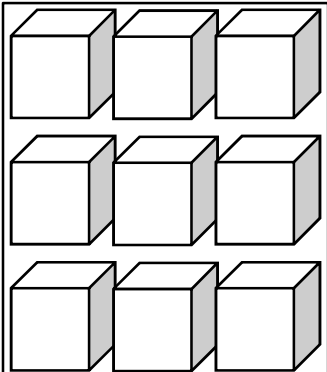
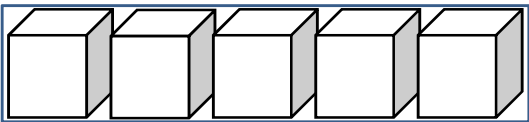
Tensors and Their Decompositions



Kolda & Bader 2009

5/4/2018 Tensor Tutorial @ SDM18 3



Sandia National Laboratories 

A Tensor is an d-Way Array

Vector $d = 1$  \mathbf{x}	Matrix $d = 2$  \mathbf{X}	3 rd -Order Tensor $d = 3$  \mathcal{X}	5 th -Order Tensor $d = 5$  \mathcal{X}
4 th -Order Tensor $d = 4$  \mathcal{X}			

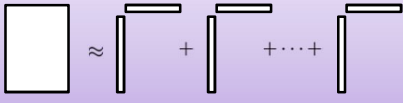
5/4/2018 Tensor Tutorial @ SDM18 4

From Matrices to Tensors: Two Points of View

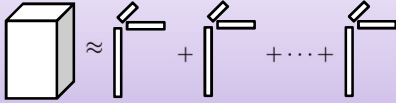



Singular value decomposition (SVD), eigendecomposition (EVD), nonnegative matrix factorization (NMF), sparse SVD, CUR, etc.

Viewpoint 1: Sum of outer products, useful for interpretation




CP Model: Sum of d -way outer products, useful for interpretation

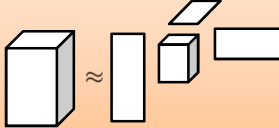


CANDECOMP, PARAFAC, Canonical Polyadic, CP

Viewpoint 2: High-variance subspaces, useful for compression





Tucker Model: Project onto high-variance subspaces to reduce dimensionality



HO-SVD, Best Rank- (R_1, R_2, \dots, R_d) decomposition

Other models for compression include hierarchical Tucker and tensor train.



5/4/2018
Tensor Tutorial @ SDM18
5

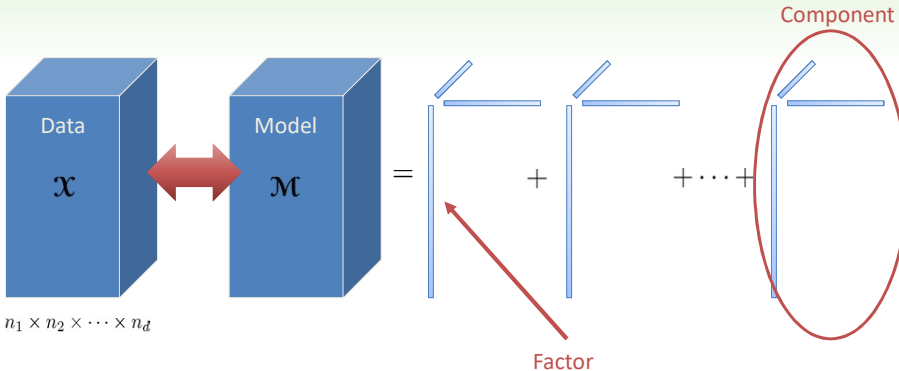



Introducing the CP Decomposition

5/4/2018
Tensor Tutorial @ SDM18
6

CP Tensor Decomposition: Sum of Outer Products






Factor

Goal: $\min \sum_i (x_i - m_i)^2$ subject to \mathcal{M} having "CP structure"


5/4/2018
Tensor Tutorial @ SDM18
7

What does CP mean?




Frank L. Hitchcock
MIT Professor
(1875–1957)




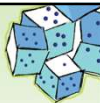
J. Douglas Carroll Bell Labs (1939-2011)	Jih-Jie Chang Bell Labs (1927-2007)
--	---

- **Polyadic form of a tensor**, Hitchcock, 1927
- **CANDECOMP** or **CAND** (Canonical Decomposition), Carroll and Chang, 1970
- **PARAFAC** (Parallel Factors), Harshman, 1970
- **CP: CANDECOMP/PARAFAC**
 - Proposed by Kiers 2000
- **CP: Canonical Polyadic**
 - Reverse-engineered circa 2010 by Comon et al.




Richard A. Harshman
Univ. Ontario
(1943-2008)

5/4/2018
Tensor Tutorial @ SDM18
8



Motivation: CP for Mouse Neural Activity



Alex Williams

A. H. Williams, T. H. Kim, F. Wang, S. Vyas, S. I. Ryu, K. V. Shenoy, M. Schnitzer, T. G. Kolda, S. Ganguli. Unsupervised Discovery of Demixed, Low-dimensional Neural Dynamics across Multiple Timescales through Tensor Components Analysis. *bioRxiv*, 2017. <https://doi.org/10.1101/211128> (accepted to *Neuron*)


5/4/2018
Tensor Tutorial @ SDM18
9

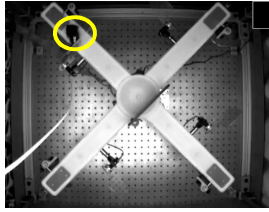
Motivating Example: Neuron Activity in Learning

Thanks to Schnitzer Group @ Stanford
Mark Schnitzer, Fori Wang, Tony Kim

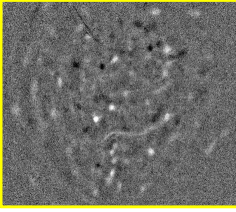
Microscope by Inscopix



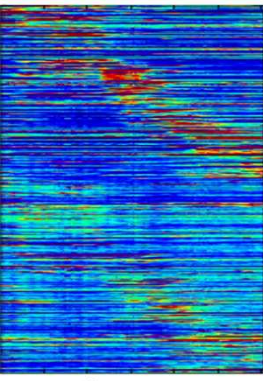
mouse in "maze"



neural activity



One Trial
300 neurons × 120 time bins



neurons



time →

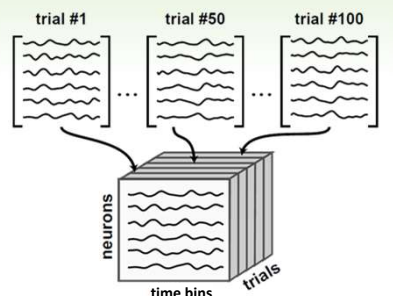
× 600 trials (over 5 days)

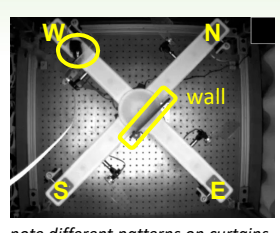
Williams, Ganguli, Kolda, et al. 2017

5/4/2018
Tensor Tutorial @ SDM18
10

Trials Vary Start Position and Strategies

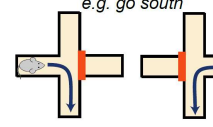




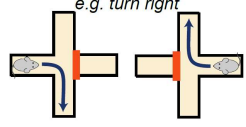
note different patterns on curtains

- 600 Trials over 5 Days
- Start West or East
- Conditions Swap Twice
 - ❖ Always Turn South
 - ❖ Always Turn Right
 - ❖ Always Turn South

Alloentric Condition
e.g. go south





Egocentric Condition
e.g. turn right

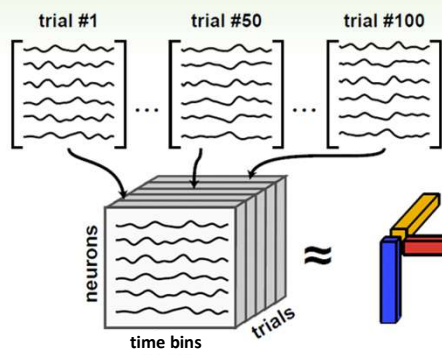


Williams, Ganguli, Kolda, et al. 2017

5/4/2018
Tensor Tutorial @ SDM18
11

CP for Simultaneous Analysis of Neurons, Time, and Trial

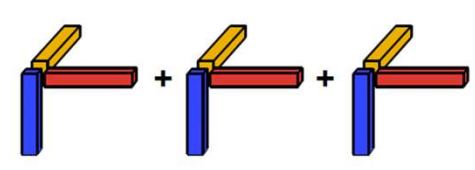





Canonical Polyadic (CP) decomposition

high-dimensional neural data

- neuron factors
- within-trial factors
- across-trial factors

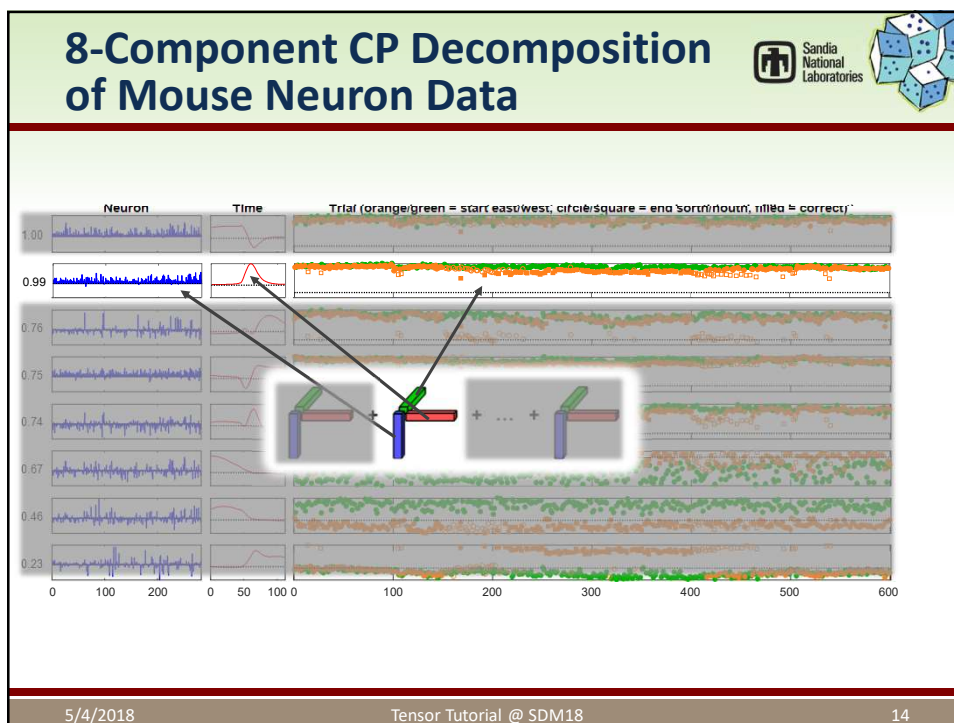
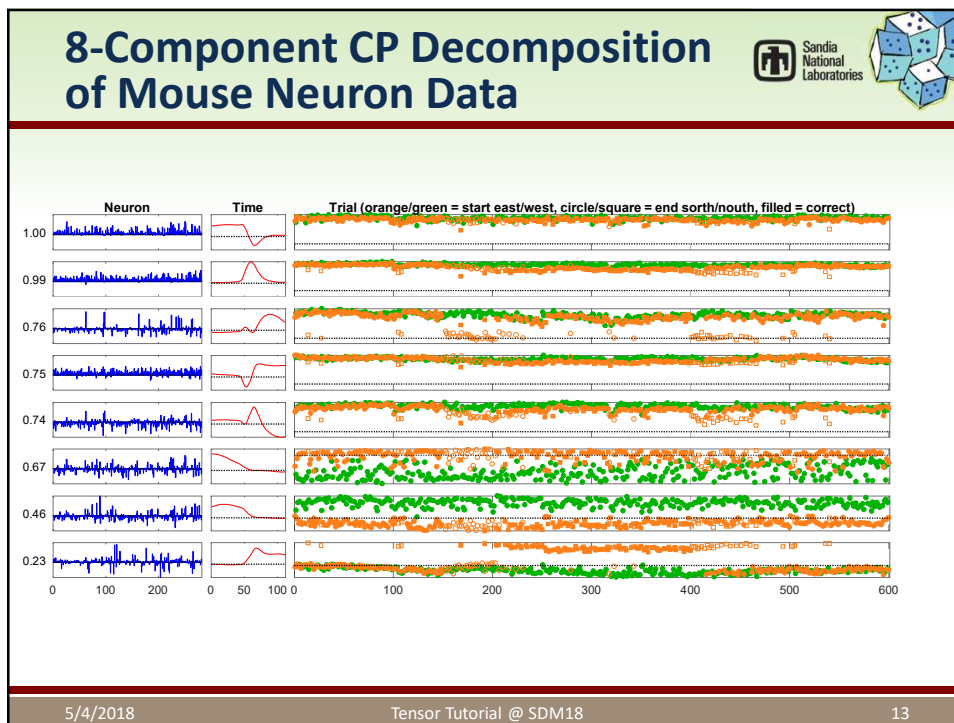


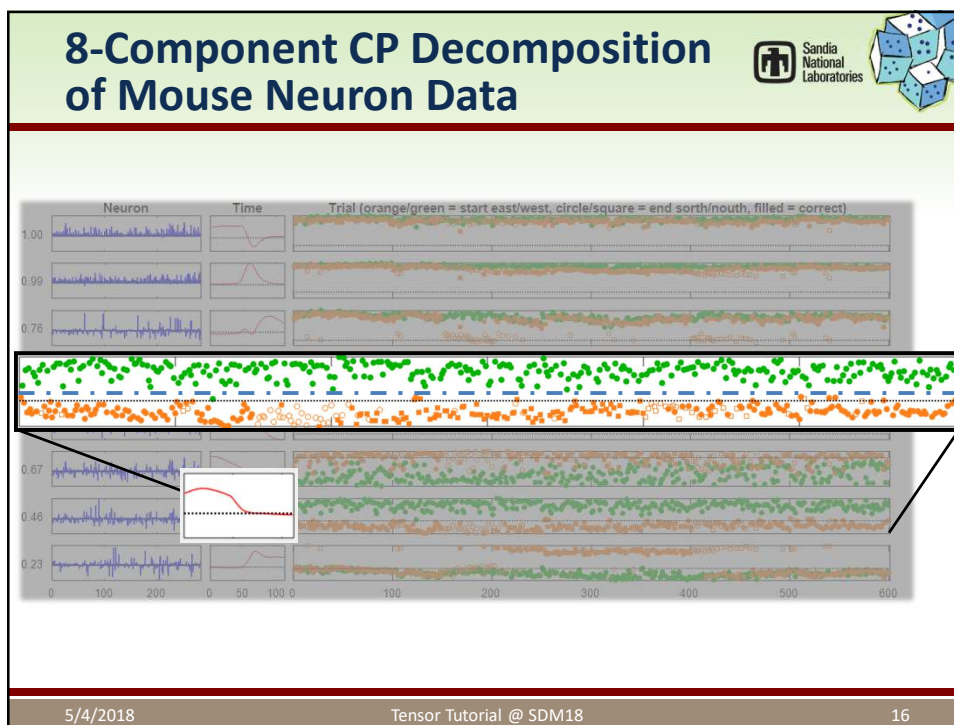
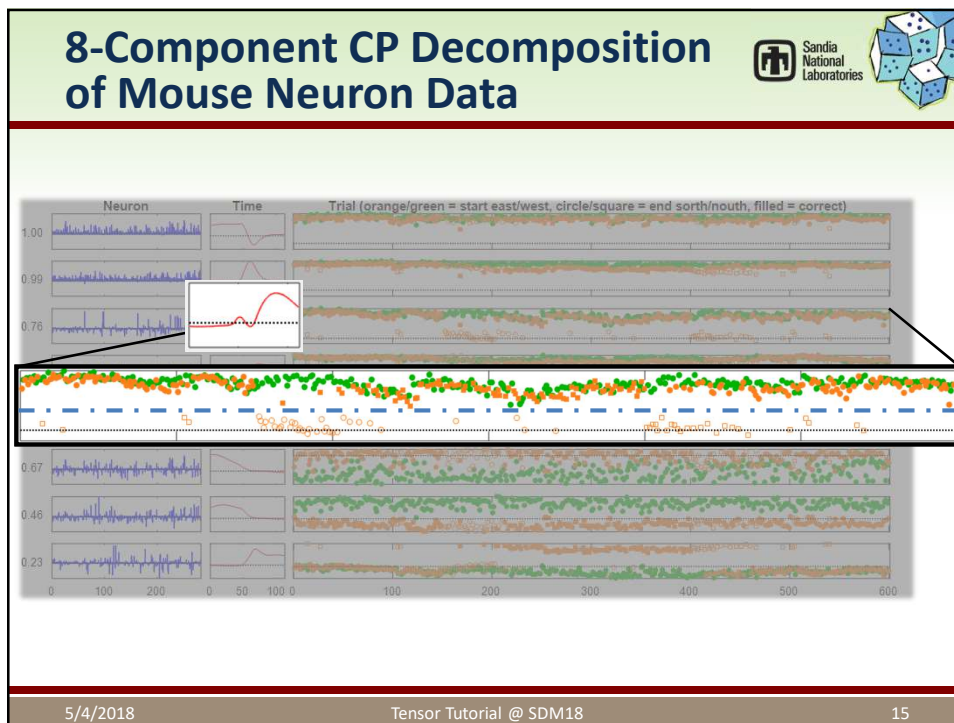
Past neuron-level work could only look at 2 factors at once: Time x Neuron or Trial x Neuron

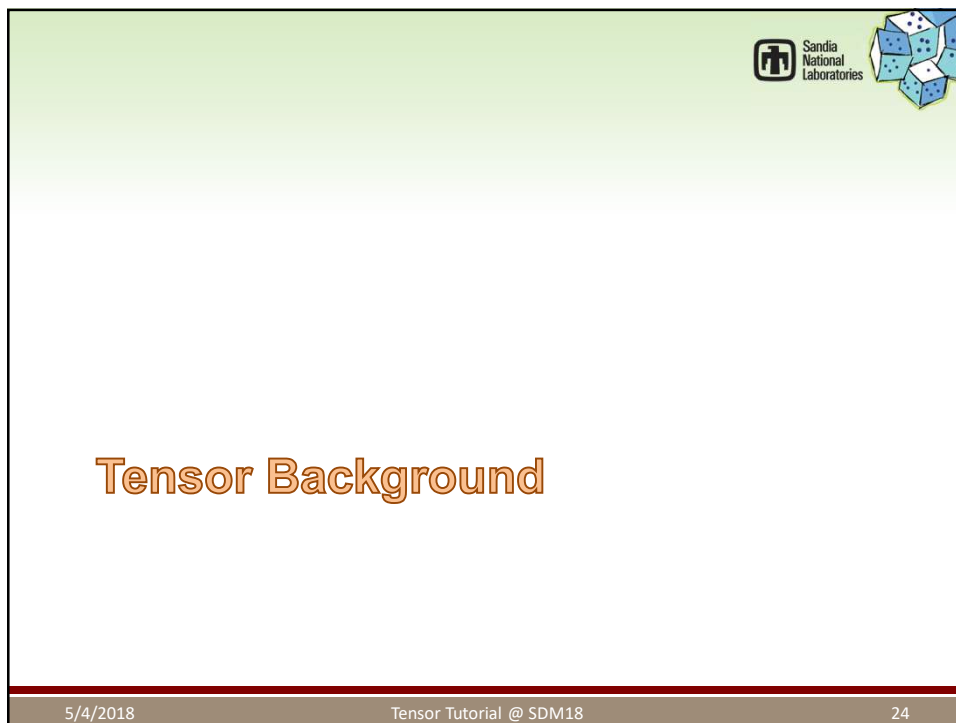
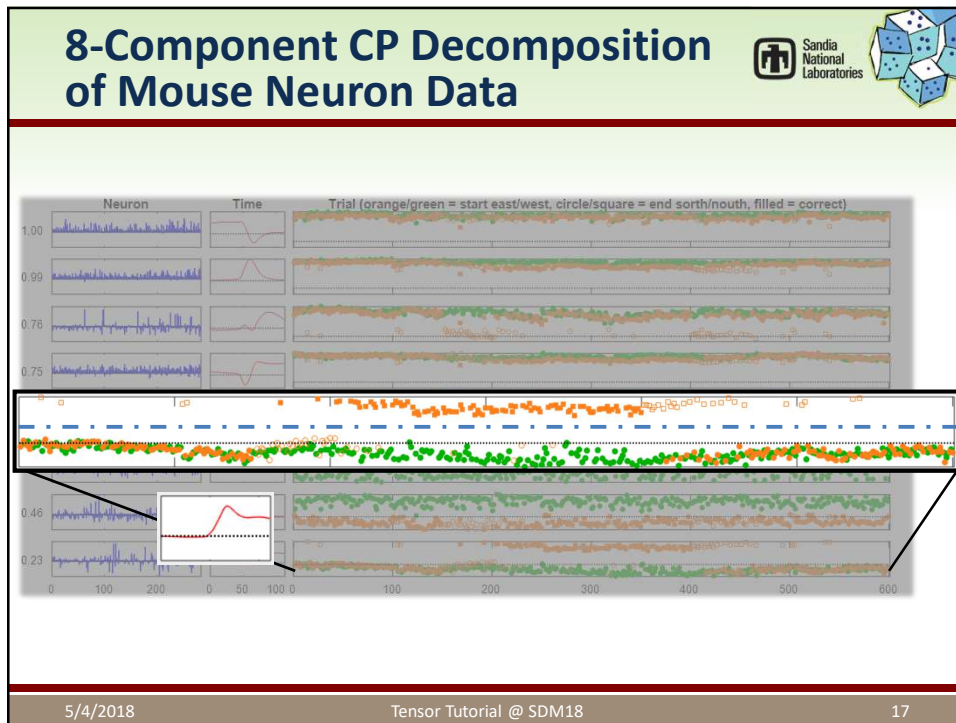
Prior tensor work in neuroscience for fMRI and EEG: Andersen and Rayens (2004), Mørup et al. (2004), Acar et al. (2007), De Vos et al. (2007), and more

Williams, Ganguli, Kolda, et al. 2017



5/4/2018
Tensor Tutorial @ SDM18
12



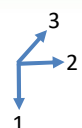




Working with Tensors

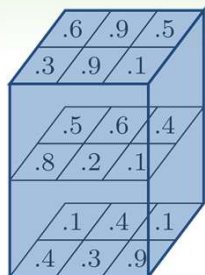



“Modes”



“Order” = Number of Modes or Ways

$d = 3$



\mathcal{X}

“Size”

$n_1 \times n_2 \times n_3$

$3 \times 3 \times 2$

“Frontal Slices”

$$\mathbf{X}(:, :, 1) = \begin{bmatrix} 0.3 & 0.9 & 0.1 \\ 0.8 & 0.2 & 0.1 \\ 0.4 & 0.3 & 0.9 \end{bmatrix}$$



$$\mathbf{X}(:, :, 2) = \begin{bmatrix} 0.6 & 0.9 & 0.5 \\ 0.5 & 0.6 & 0.4 \\ 0.1 & 0.4 & 0.1 \end{bmatrix}$$

This is how it's displayed in MATLAB.

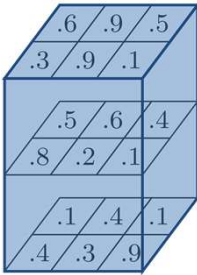
The term “dimension” is overloaded and may be confusing.

5/4/2018
Tensor Tutorial @ SDM18
25

Indexing a 3D Tensor

We use 1-based indexing throughout this talk to be consistent with MATLAB.



\mathcal{X}

Tensor Elements

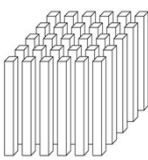
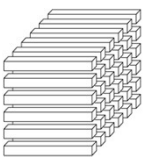
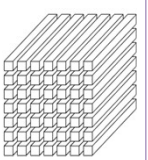
$x(1, 1, 1) = 0.3$

$x(2, 1, 1) = 0.8$

$x(1, 1, 2) = 0.6$

$x(3, 2, 1) = 0.3$

Tensor Fibers

Mode-1 Fibers Mode-2 Fibers Mode-3 Fibers



$$\mathbf{x}(:, 1, 1) = \begin{bmatrix} .3 \\ .8 \\ .4 \end{bmatrix}$$

$$\mathbf{x}(1, 1, :) = \begin{bmatrix} .3 \\ .6 \end{bmatrix}$$

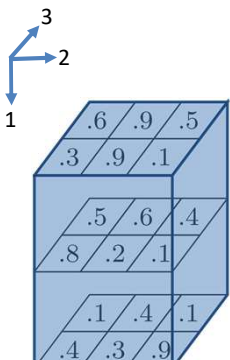
$$\mathbf{x}(1, :, 1) = \begin{bmatrix} .3 \\ .9 \\ .1 \end{bmatrix}$$

5/4/2018
Tensor Tutorial @ SDM18
26

Vectorizing a 3D Tensor (or How It's Stored in Memory!)

$(i_1, i_2, i_3) \rightarrow i' = (i_3 - 1)n_2n_1 + (i_2 - 1)n_1 + i_1$

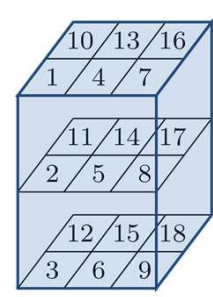


X
 $n_1 \times n_2 \times n_3$

\rightarrow **x** =

0.3
0.8
0.4
0.9
0.2
0.3
0.1
0.1
0.9
0.6
0.5
0.1
0.9
0.6
0.4
0.4
0.1

Ordering





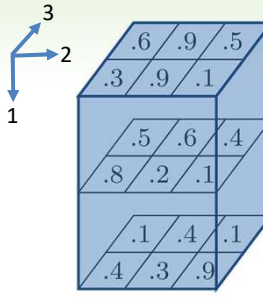
X

*This is MATLAB's default order.
It may not be everyone's.*

5/4/2018
Tensor Tutorial @ SDM18
27

Unfolding a 3D Tensor



X
 $n_1 \times n_2 \times n_3$

$(i_1, i_2, i_3) \rightarrow (i_1, i'_1), i'_1 = (i_3 - 1)n_2 + i_2$

$$\mathbf{X}_{(1)} = \begin{bmatrix} 0.3 & 0.9 & 0.1 & 0.6 & 0.9 & 0.5 \\ 0.8 & 0.2 & 0.1 & 0.5 & 0.6 & 0.4 \\ 0.4 & 0.3 & 0.9 & 0.1 & 0.4 & 0.1 \end{bmatrix} \quad n_1 \times n_2 n_3$$

$(i_1, i_2, i_3) \rightarrow (i_2, i'_2), i'_2 = (i_3 - 1)n_1 + i_1$



$$\mathbf{X}_{(2)} = \begin{bmatrix} 0.3 & 0.8 & 0.4 & 0.6 & 0.5 & 0.1 \\ 0.9 & 0.2 & 0.3 & 0.9 & 0.6 & 0.4 \\ 0.1 & 0.1 & 0.9 & 0.5 & 0.4 & 0.1 \end{bmatrix} \quad n_2 \times n_1 n_3$$

$(i_1, i_2, i_3) \rightarrow (i_3, i'_3), i'_3 = (i_2 - 1)n_1 + i_1$


$$\mathbf{X}_{(3)} = \begin{bmatrix} 0.3 & 0.8 & 0.4 & 0.9 & 0.2 & 0.3 & 0.1 & 0.1 & 0.9 \\ 0.6 & 0.5 & 0.1 & 0.9 & 0.6 & 0.4 & 0.5 & 0.4 & 0.1 \end{bmatrix} \quad n_3 \times n_1 n_2$$

Unfolding also known as "matricization".

5/4/2018
Tensor Tutorial @ SDM18
28

Working with d -way Tensors



\mathcal{X}

Size: $n_1 \times n_2 \times \dots \times n_d$

Single Element: $x(i_1, i_2, \dots, i_d)$ or x_i ← “multiindex”

Set of All Indices: $\mathcal{I} = \{i \equiv (i_1, i_2, \dots, i_d) \mid i_k \in \{1, \dots, n_k\} \text{ for } k = 1, \dots, d\}$



Geometric mean of the sizes: $n = \sqrt[d]{\prod_{k=1}^d n_k}$

Arithmetic mean of the sizes: $\bar{n} = \frac{1}{d} \sum_{k=1}^d n_k$

Total number of elements in the tensor:

$$|\mathcal{I}| = \prod_{k=1}^d n_k = n^d$$

5/4/2018
Tensor Tutorial @ SDM18
29

Vectorizing a d -way Tensor

Order: d

Size: $n_1 \times n_2 \times \dots \times n_d$

“subscripts” $(i_1, i_2, \dots, i_d) \in \mathcal{I} \rightarrow i' \in \{1, \dots, n^d\}$ “linear index”



$$i' = 1 + \sum_{k=1}^d (i_k - 1)n'_k$$

$$n'_k = \prod_{\ell=1}^{k-1} n_\ell$$

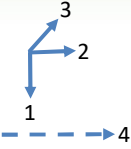
See MATLAB commands:
tt_sub2ind, tt_ind2sub

5/4/2018
Tensor Tutorial @ SDM18
30

Example: Ordering of Elements in a 4D Tensor

4th-order tensor of size $3 \times 4 \times 3 \times 2$



25	28	31	34
13	16	19	22
1	4	7	10

61	64	67	70
49	52	55	58
37	40	43	46

26	29	32	35
14	17	20	23
2	5	8	11



$\mathcal{X}(:, :, :, 1)$

62	65	68	71
50	53	56	59
38	41	44	47

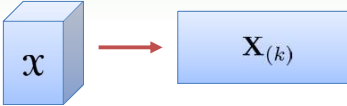
$\mathcal{X}(:, :, :, 2)$

5/4/2018
Tensor Tutorial @ SDM18
31

Unfolding a d -way Tensor

Matricization or Unfolding: Rearrange d -way array into 2-way array.



Size: $n_1 \times n_2 \times \dots \times n_d$ Size: $n_k \times (n^d/n_k)$

Mode- k unfolding: $(i_1, i_2, \dots, i_d) \rightarrow (i_k, i'_k)$



$$i'_k = 1 + \sum_{\ell=1}^{k-1} (i_\ell - 1)n'_\ell + \sum_{\ell=k+1}^d (i_\ell - 1)(n'_\ell/n_k)$$


$$n'_k = \prod_{\ell=1}^{k-1} n_\ell$$

$i'_k \in \{1, \dots, n^d/n_k\}$

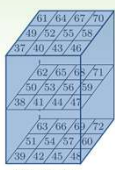
5/4/2018
Tensor Tutorial @ SDM18
32

Unfolding a 4D Tensor



$\mathbf{X}(:,:,1)$



$\mathbf{X}(:,:,2)$

Note the patterning!

block size = $n_k \times n'_k$

blocks = $n^d / n_k n'_k$

$n'_k = \prod_{\ell=1}^{k-1} n_\ell$

$\mathbf{X}_{(1)} =$

1	4	7	10	13	16	...	58	61	64	67	70
2	5	8	11	14	17	...	59	62	65	68	71
3	6	9	12	15	18	...	60	63	66	69	72

$\mathbf{X}_{(2)} =$

1	2	3	13	14	15	25	26	27	37	38	39	49	50	51	61	62	63
4	5	6	16	17	18	28	29	30	40	41	42	52	53	54	64	65	66
7	8	9	19	20	21	31	32	33	43	44	45	55	56	57	67	68	69
10	11	12	22	23	24	34	35	36	46	47	48	58	59	60	70	71	72

$\mathbf{X}_{(3)} =$

1	2	3	4	5	6	7	8	9	10	11	12	37	...	48
13	14	15	16	17	18	19	20	21	22	23	24	49	...	60
25	26	27	28	29	30	31	32	33	34	35	36	61	...	72



$\mathbf{X}_{(4)} =$

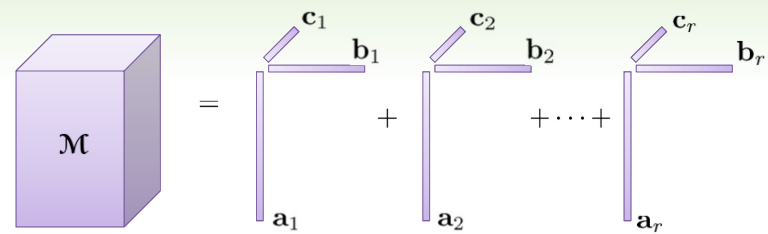
1	2	3	4	5	6	...	31	32	33	34	35	36
37	38	39	40	41	42	...	67	68	69	70	71	72

Patterning: Battaglini, Perros, Sub, & Vuduc 2015; Austin, Ballard, & Kolda 2016

5/4/2018
Tensor Tutorial @ SDM18
33

Ktensor: Sum of Outer Products




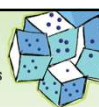
$\mathcal{M} = \mathbf{a}_1 \circ \mathbf{b}_1 \circ \mathbf{c}_1 + \mathbf{a}_2 \circ \mathbf{b}_2 \circ \mathbf{c}_2 + \dots + \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$

$$m(i_1, i_2, i_3) = \sum_{j=1}^r a(i_1, j) b(i_2, j) c(i_3, j)$$

Bader & Kolda 2007

5/4/2018
Tensor Tutorial @ SDM18
34

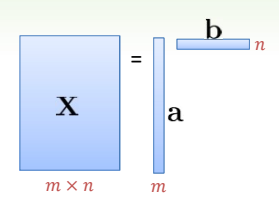
Vector Outer Product

Outer product of 2 vectors
 \Rightarrow rank-1 matrix

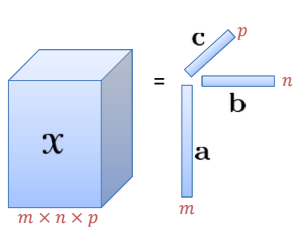
$$\mathbf{X} = \mathbf{a}\mathbf{b}^\top$$

$$\mathbf{X} = \mathbf{a} \circ \mathbf{b}$$

$$x(i_1, i_2) = a(i_1) b(i_2)$$



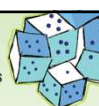
Outer product of 3 vectors
 \Rightarrow rank-1 tensor of order 3

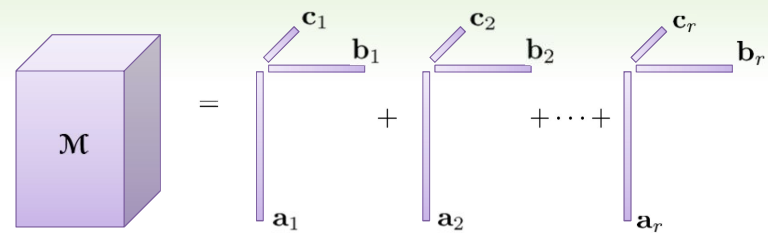
$$\mathcal{X} = \mathbf{a} \circ \mathbf{b} \circ \mathbf{c}$$

$$x(i_1, i_2, i_3) = a(i_1) b(i_2) c(i_3)$$


5/4/2018
Tensor Tutorial @ SDM18
35

Ktensor: Sum of Outer Products





$$\mathcal{M} = \mathbf{a}_1 \circ \mathbf{b}_1 \circ \mathbf{c}_1 + \mathbf{a}_2 \circ \mathbf{b}_2 \circ \mathbf{c}_2 + \dots + \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

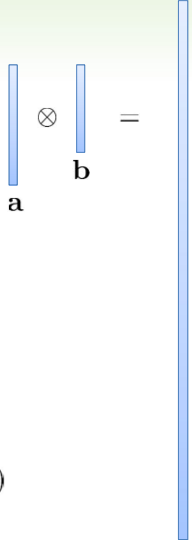
Bader & Kolda 2007

5/4/2018
Tensor Tutorial @ SDM18
36

Vector Kronecker Product



$$\mathbf{a} \otimes \mathbf{b} = \begin{bmatrix} a(1) \mathbf{b} \\ a(2) \mathbf{b} \\ \vdots \\ a(m) \mathbf{b} \end{bmatrix} = \begin{bmatrix} a(1) b(1) \\ a(1) b(2) \\ \vdots \\ a(1) b(n) \\ \vdots \\ a(m) b(1) \\ a(m) b(2) \\ \vdots \\ a(m) b(n) \end{bmatrix}$$



Associative: $(\mathbf{a} \otimes \mathbf{b}) \otimes \mathbf{c} = \mathbf{a} \otimes (\mathbf{b} \otimes \mathbf{c})$

5/4/2018
Tensor Tutorial @ SDM18
37

Connection between Kronecker Product and Outer Products

$$\mathbf{X} = \mathbf{a} \circ \mathbf{b} \iff \text{vec}(\mathbf{X}) = \mathbf{b} \otimes \mathbf{a}$$

$$\mathbf{X} = \mathbf{a} \circ \mathbf{b} \circ \mathbf{c} \iff \text{vec}(\mathbf{X}) = \mathbf{c} \otimes \mathbf{b} \otimes \mathbf{a}$$

$$\mathbf{X}_{(1)} = \mathbf{a}(\mathbf{c} \otimes \mathbf{b})'$$

$$\mathbf{X}_{(2)} = \mathbf{b}(\mathbf{c} \otimes \mathbf{a})'$$



$$\mathbf{X}_{(3)} = \mathbf{c}(\mathbf{b} \otimes \mathbf{a})'$$

$$\mathbf{X} = \mathbf{a}_1 \circ \mathbf{a}_2 \circ \dots \circ \mathbf{a}_d \iff \text{vec}(\mathbf{X}) = \mathbf{a}_d \otimes \dots \otimes \mathbf{a}_1$$

$$\mathbf{X}_{(k)} = \mathbf{a}_k(\mathbf{a}_d \otimes \dots \otimes \mathbf{a}_{k+1} \otimes \mathbf{a}_{k-1} \otimes \dots \otimes \mathbf{a}_1)'$$

5/4/2018
Tensor Tutorial @ SDM18
38

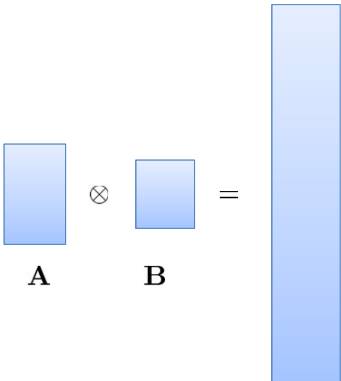
Khatri-Rao Product

Input two matrices with same number of columns
Columnwise-Kronecker Product



$$\mathbf{A} \odot \mathbf{B} = \begin{bmatrix} \mathbf{a}(:,1) \otimes \mathbf{b}(:,1) & \mathbf{a}(:,2) \otimes \mathbf{b}(:,2) & \cdots & \mathbf{a}(:,r) \otimes \mathbf{b}(:,r) \end{bmatrix}$$

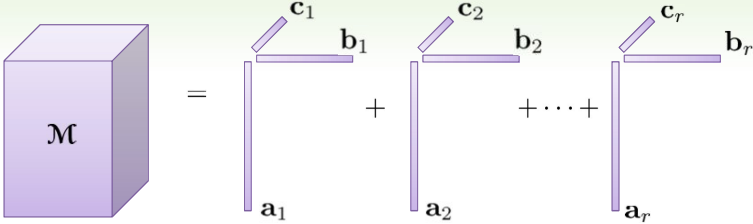
$m \times r$ $n \times r$ $mn \times r$



5/4/2018
Tensor Tutorial @ SDM18
39

Ktensors and Khatri-Rao Products



$$\mathcal{M} = \mathbf{a}_1 \circ \mathbf{b}_1 \circ \mathbf{c}_1 + \mathbf{a}_2 \circ \mathbf{b}_2 \circ \mathbf{c}_2 + \cdots + \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

$$\mathbf{M}_{(1)} = \mathbf{a}_1(\mathbf{c}_1 \otimes \mathbf{b}_1)' + \mathbf{a}_2(\mathbf{c}_2 \otimes \mathbf{b}_2)' + \cdots + \mathbf{a}_r(\mathbf{c}_r \otimes \mathbf{b}_r)'$$



$$\mathbf{M}_{(1)} = \underbrace{\begin{bmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_r \end{bmatrix}}_{\mathbf{A}} \underbrace{\begin{bmatrix} \mathbf{c}_1 \otimes \mathbf{b}_1 & \cdots & \mathbf{c}_r \otimes \mathbf{b}_r \end{bmatrix}'}_{(\mathbf{C} \odot \mathbf{B})'}$$

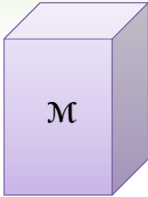
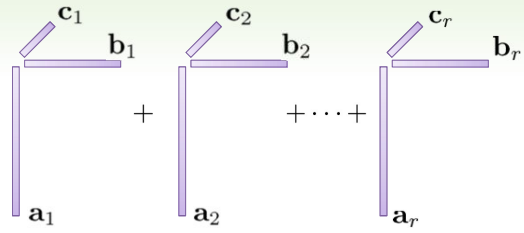
$$\mathbf{M}_{(1)} = \mathbf{A}(\mathbf{C} \odot \mathbf{B})'$$

Bader & Kolda 2007

5/4/2018
Tensor Tutorial @ SDM18
40

Ktensor Unfoldings


=


$$M_{(1)} = A(C \odot B)'$$

$$M_{(2)} = B(C \odot A)'$$

$$M_{(3)} = C(B \odot A)'$$

Factor Matrices

$$A = [a_1 \ a_2 \ \dots \ a_r]$$

$$B = [b_1 \ b_2 \ \dots \ b_r]$$



$$C = [c_1 \ c_2 \ \dots \ c_r]$$

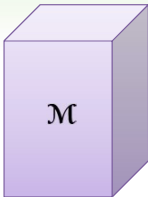
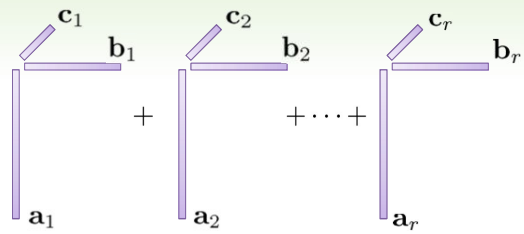
d-way case $M_{(k)} = A_k (A_d \odot \dots \odot A_{k+1} \odot A_{k-1} \odot \dots \odot A_1)'$

Bader & Kolda 2007

5/4/2018
Tensor Tutorial @ SDM18
41

Ktensor: Shorthand


=




$$\mathcal{M} = [A, B, C]$$

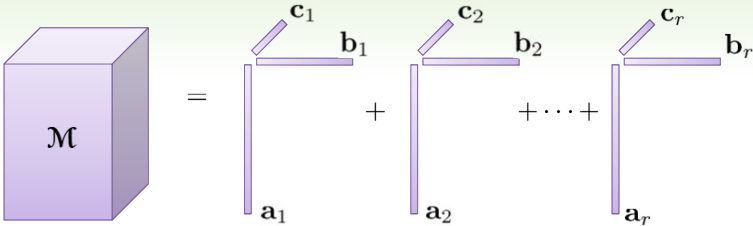
d-way case $\mathcal{M} = [A_1, A_2, \dots, A_d]$

Bader & Kolda 2007

5/4/2018
Tensor Tutorial @ SDM18
42

Scaling Ambiguity



$$\mathcal{M} = \mathbf{a}_1 \circ \mathbf{b}_1 \circ \mathbf{c}_1 + \mathbf{a}_2 \circ \mathbf{b}_2 \circ \mathbf{c}_2 + \dots + \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

$$\mathcal{M} = \mathbf{a}_1 \circ 2 \mathbf{b}_1 \circ \frac{1}{2} \mathbf{c}_1 + \mathbf{a}_2 \circ \mathbf{b}_2 \circ \mathbf{c}_2 + \dots + \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

$$\mathcal{M} = \lambda_1 \mathbf{a}_1 \circ \mathbf{b}_1 \circ \mathbf{c}_1 + \lambda_2 \mathbf{a}_2 \circ \mathbf{b}_2 \circ \mathbf{c}_2 + \dots + \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

Require all factors (i.e., columns of the factor matrices) to have unit norm, e.g., $\|\mathbf{a}_1\| = 1$



Fixed by ktensor/normalize function.

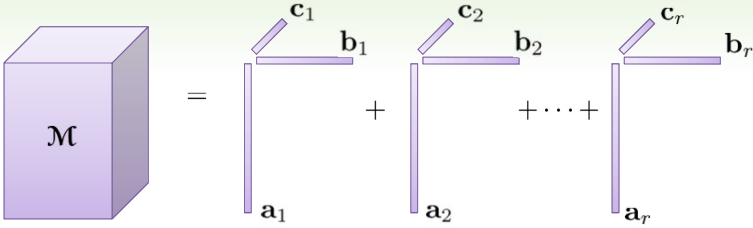
$$\mathcal{M} = \llbracket \boldsymbol{\lambda}, \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket$$

Bader & Kolda 2007

5/4/2018
Tensor Tutorial @ SDM18
43

Permutation Ambiguity in Ktensor








$$\mathcal{M} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket = \llbracket \mathbf{A}_1 \boldsymbol{\Pi}, \mathbf{A}_2 \boldsymbol{\Pi}, \dots, \mathbf{A}_d \boldsymbol{\Pi} \rrbracket$$

where $\boldsymbol{\Pi}$ is a $p \times p$ permutation matrix.

5/4/2018
Tensor Tutorial @ SDM18
44

Tensor Inner Product & Norm



$$\langle \mathcal{X}, \mathcal{Y} \rangle = \sum_{i_1=1}^{n_1} \sum_{i_2=1}^{n_2} \cdots \sum_{i_d=1}^{n_d} x(i_1, i_2, \dots, i_d) y(i_1, i_2, \dots, i_d)$$

Shorthand: $\langle \mathcal{X}, \mathcal{Y} \rangle = \sum_{i \in \mathcal{I}} x_i y_i$

$$\langle \mathcal{X}, \mathcal{Y} \rangle = \text{vec}(\mathcal{X})^T \text{vec}(\mathcal{Y})$$

$$\|\mathcal{X}\|^2 = \langle \mathcal{X}, \mathcal{X} \rangle = \sum_{i \in \mathcal{I}} x_i^2$$

5/4/2018
Tensor Tutorial @ SDM18
45

Matrix Hadamard Product



Input two same-sized matrices

$$\mathbf{A} * \mathbf{B} = \begin{bmatrix} a(1,1)b(1,1) & a(1,2)b(1,2) & \cdots & a(1,n)b(1,n) \\ a(2,1)b(2,1) & a(2,2)b(2,2) & \cdots & a(2,n)b(2,n) \\ \vdots & \vdots & \ddots & \vdots \\ a(m,1)b(m,1) & a(m,2)b(m,2) & \cdots & a(m,n)b(m,n) \end{bmatrix}$$

In MATLAB: A .* B

5/4/2018
Tensor Tutorial @ SDM18
46

Properties of Khatri-Rao Product

$$(\mathbf{A} \odot \mathbf{B}) \odot \mathbf{C} = \mathbf{A} \odot (\mathbf{B} \odot \mathbf{C})$$

$$(\mathbf{A} \odot \mathbf{B})^\top (\mathbf{A} \odot \mathbf{B}) = (\mathbf{A}^\top \mathbf{A} * \mathbf{B}^\top \mathbf{B})$$

$m \times r$ $n \times r$ $r \times r$ $r \times r$



$$(\mathbf{A} \odot \mathbf{B})^\dagger = (\mathbf{A}^\top \mathbf{A} * \mathbf{B}^\top \mathbf{B})^\dagger (\mathbf{A} \odot \mathbf{B})^\top$$

$mn \times r$ $r \times r$ $mn \times r$

↖ Moore-Penrose Pseudo-Inverse

5/4/2018
Tensor Tutorial @ SDM18
47

MTTKRP: Matricized-Tensor Times Khatri-Rao Product

Tensor of size $n_1 \times n_2 \times \dots \times n_d$: \mathbf{X}
 Matrices of size $n_k \times r$ for $k = 1, \dots, d$: $\mathbf{A}_1, \dots, \mathbf{A}_d$



$$\mathbf{B} = \mathbf{X}_{(k)} (\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)$$

$n_k \times r$ $n_k \times (n^d/n_k)$ $(n^d/n_k) \times r$

↖ Same size as \mathbf{A}_k



Bader & Kolda 2007

5/4/2018
Tensor Tutorial @ SDM18
48

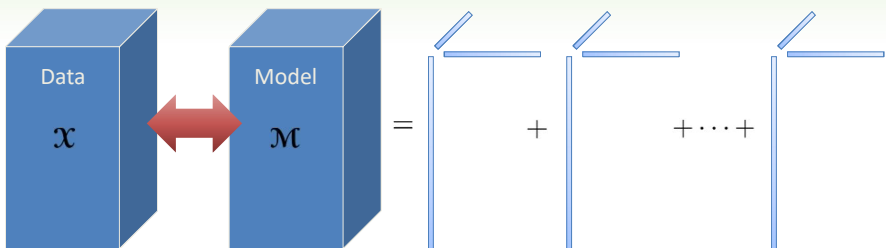



Computing the CP Decomposition

5/4/2018
Tensor Tutorial @ SDM18
50

CP Objective Function





$$\min \sum_i (x_i - m_i)^2 \quad \text{subject to} \quad m_i = \sum_{j=1}^r a_1(i_1, j) a_2(i_2, j) \cdots a_d(i_d, j)$$

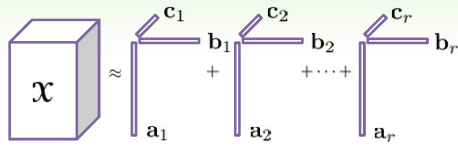
$$\min \|\mathcal{X} - \mathcal{M}\|^2 \quad \text{subject to} \quad \mathcal{M} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket$$

Hitchcock 1927, Harshman 1970, Carroll & Chang 1970

5/4/2018
Tensor Tutorial @ SDM18
51

CP Optimization Problem



Number of observations = n^d

Number of free parameters = $dr\bar{n}$

Typically, $dr\bar{n} \ll n^d$

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathcal{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket\|^2$$

General Form

$$\min_{\{\mathbf{A}_k\}} \|\mathcal{X} - \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket\|^2$$



Nonconvex

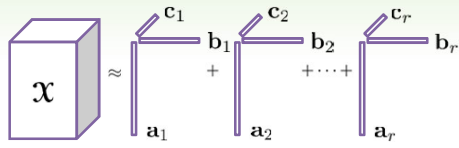
Permutation & Scaling Ambiguities

How to pick r ?

5/4/2018
Tensor Tutorial @ SDM18
52

CP-ALS: Fitting CP Model via Alternating Least Squares



*Convex (linear least squares)
subproblems can be solved exactly*

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \|\mathcal{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket\|^2$$

Repeat until convergence...



$$\min_{\mathbf{A}} \|\mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top\|^2$$

$$\min_{\mathbf{B}} \|\mathbf{X}_{(2)} - \mathbf{B}(\mathbf{C} \odot \mathbf{A})^\top\|^2$$

$$\min_{\mathbf{C}} \|\mathbf{X}_{(3)} - \mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top\|^2$$

5/4/2018
Tensor Tutorial @ SDM18
53

Special Structure of the Least Squares Problem

$$\min_{\mathbf{A}} \left\| \underbrace{\mathbf{X}_{(1)}}_{n_1 \times (n_2 n_3)} - \underbrace{\mathbf{A}}_{n_1 \times r} \underbrace{(\mathbf{C} \odot \mathbf{B})^\top}_{r \times (n_2 n_3)} \right\|^2$$

$$(\mathbf{C} \odot \mathbf{B}) \mathbf{A}^\top = \mathbf{X}_{(1)}^\top$$



$$\mathbf{A}^\top = (\mathbf{C} \odot \mathbf{B})^\dagger \mathbf{X}_{(1)}^\top$$

$$\mathbf{A}^\top = (\mathbf{C}^\top \mathbf{C} * \mathbf{B}^\top \mathbf{B})^\dagger (\mathbf{C} \odot \mathbf{B})^\top \mathbf{X}_{(1)}^\top$$

$$\mathbf{A} = \underbrace{\mathbf{X}_{(1)} (\mathbf{C} \odot \mathbf{B}) (\mathbf{C}^\top \mathbf{C} * \mathbf{B}^\top \mathbf{B})^\dagger}_{\text{MITKRP}}$$

5/4/2018
Tensor Tutorial @ SDM18
54

Special Structure of the Least Squares Problem: d -way



$$\min_{\mathbf{A}_k} \left\| \mathbf{X}_{(k)} - \mathbf{A}_k (\mathbf{A}_d \odot \cdots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \cdots \odot \mathbf{A}_1)^\top \right\|^2$$

$$\mathbf{A}_k \leftarrow \mathbf{X}_{(k)} (\mathbf{A}_d \odot \cdots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \cdots \odot \mathbf{A}_1) \mathbf{V}^\dagger$$

$$\mathbf{V} \leftarrow (\mathbf{A}_1^\top \mathbf{A}_1) * \cdots * (\mathbf{A}_{k-1}^\top \mathbf{A}_{k-1}) * (\mathbf{A}_{k+1}^\top \mathbf{A}_{k+1}) * \cdots * (\mathbf{A}_d^\top \mathbf{A}_d)$$

5/4/2018
Tensor Tutorial @ SDM18
55

CP-ALS Algorithm + Nuances



- 1: $\text{fit} \leftarrow 1 - (\|\mathcal{X} - \mathcal{M}\| / \|\mathcal{X}\|)$ ← *Roughly, the proportion of the data explained by the model.*
- 2: **for** $\ell = 1, 2, \dots$ **do**
- 3: $\text{oldfit} \leftarrow \text{fit}$
- 4: **for** $k = 1, 2, \dots, d$ **do**
- 5: $\mathbf{Z} \leftarrow \mathbf{X}_{(k)} (\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)$ {MTTKRP}
- 6: $\mathbf{V} \leftarrow (\mathbf{A}_1^T \mathbf{A}_1) * \dots * (\mathbf{A}_{k-1}^T \mathbf{A}_{k-1}) * (\mathbf{A}_{k+1}^T \mathbf{A}_{k+1}) * \dots * (\mathbf{A}_d^T \mathbf{A}_d)$
- 7: $\mathbf{A}_k \leftarrow \mathbf{Z} / \mathbf{V}$ {transposed backsolve}
- 8: $\boldsymbol{\lambda} \leftarrow$ column norms of \mathbf{A}_k *Important to normalize!*
- 9: $\mathbf{A}_k \leftarrow \mathbf{A}_k / \text{diag}(\boldsymbol{\lambda})$
- 10: **end for**
- 11: $\text{fit} \leftarrow 1 - (\|\mathcal{X} - \mathcal{M}\| / \|\mathcal{X}\|)$
- 12: **if** $|\text{oldfit} - \text{fit}| < \tau$ **then** *Stop when the fit stagnates.*
- 13: **break**
- 14: **end if**
- 15: **end for**

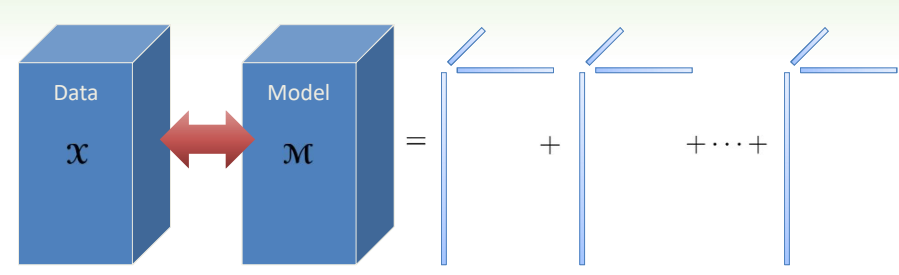
See MATLAB command:
cp_als

Harshman 1970, Carroll & Chang 1970

5/4/2018
Tensor Tutorial @ SDM18
56

CP Objective Function





$$\min \frac{1}{2} \|\mathcal{X} - \mathcal{M}\|^2 \text{ subject to } \mathcal{M} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket$$

Hitchcock 1927, Harshman 1970, Carroll & Chang 1970

5/4/2018
Tensor Tutorial @ SDM18
57

CP-OPT: Fitting CP Model via Direct Optimization

$$f = 1/2 \|\mathcal{X} - \mathcal{M}\|^2 = 1/2 \|\mathbf{X}_{(k)} - \mathbf{M}_{(k)}\|_F^2$$

$$\frac{\partial f}{\partial \mathbf{A}_k} = -(\mathbf{X}_{(k)} - \mathbf{M}_{(k)}) \frac{\partial \mathbf{M}_{(k)}}{\partial \mathbf{A}_k}$$


Recall: $\mathbf{M}_{(k)} = \mathbf{A}_k (\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)^\top$

$$\frac{\partial f}{\partial \mathbf{A}_k} = -(\mathbf{X}_{(k)} - \mathbf{M}_{(k)}) (\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)$$

$$= -\underbrace{\mathbf{X}_{(k)} (\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)}_{\text{MTTKRP}} + \mathbf{A}_k \mathbf{V}_k$$

where $\mathbf{V}_k = (\mathbf{A}_1^\top \mathbf{A}_1) * \dots * (\mathbf{A}_{k-1}^\top \mathbf{A}_{k-1}) * (\mathbf{A}_{k+1}^\top \mathbf{A}_{k+1}) * \dots * (\mathbf{A}_d^\top \mathbf{A}_d)$



Note
connect
to ALS!



Acar, Dunlavy, Kolda 2011

5/4/2018
Tensor Tutorial @ SDM18
58


CP-OPT

Calculating the function and gradient

- 1: for $k = 1, 2, \dots, d$ do
- 2: $\mathbf{Z}_k \leftarrow \mathbf{X}_{(k)} (\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)$ {MTTKRP}
- 3: $\mathbf{V}_k \leftarrow (\mathbf{A}_1^\top \mathbf{A}_1) * \dots * (\mathbf{A}_{k-1}^\top \mathbf{A}_{k-1}) * (\mathbf{A}_{k+1}^\top \mathbf{A}_{k+1}) * \dots * (\mathbf{A}_d^\top \mathbf{A}_d)$
- 4: $\Delta_k \leftarrow -\mathbf{Z}_k + \mathbf{A}_k \mathbf{V}_k$
- 5: end for
- 6: $f \leftarrow \|\mathcal{X}\|^2 + \mathbf{1}^\top [\mathbf{A}_d \mathbf{Z}_d] \mathbf{1} + \mathbf{1}^\top [\mathbf{A}_d^\top \mathbf{A}_d * \mathbf{V}_d] \mathbf{1}$

- Gradients not naturally in vector form.
- Need to convert before calling optimization method.





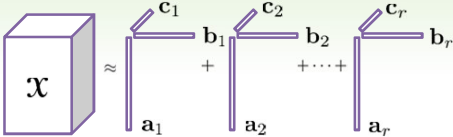
Scaling ambiguities can cause problems – may want to regularize to ameliorate

Acar, Dunlavy, Kolda 2011

5/4/2018
Tensor Tutorial @ SDM18
59

CP-OPT: Fitting CP via "All-at-once" Optimization

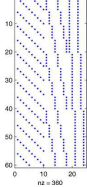





$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \|\mathcal{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket\|^2$$

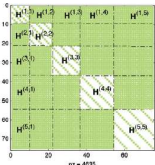
- CP-OPT (Acar et al.): 1st-order method, better accuracy than ALS when r is "too big"
- CP-NLS (Paatero, Tomasi & Bro): Damped Gauss-Newton, accurate but slow
- CP-Newton (Phan et al.): Newton method, superior to CP-OPT for high order

Structured
Jacobian

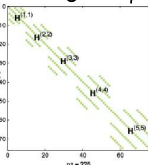


$n_2 = 250$

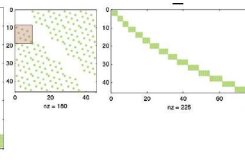
Structured Hessian can be written as
block diagonal plus low-rank correction



$n_2 = 4025$



$n_2 = 225$


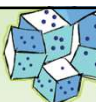


$n_2 = 225$

Paatero 1997; Tomasi & Bro 2005, 2006; Acar, Dunlavy, & Kolda 2011; Phan, Tichavský, & Cichocki 2013

5/4/2018
Tensor Tutorial @ SDM18
60

More on the CP Decomposition

More on the CP Decomposition

5/4/2018
Tensor Tutorial @ SDM18
61

Uniqueness \Rightarrow Interpretability

$k\text{-rank}(\mathbf{A})$ = maximum value k such that *any* k columns of \mathbf{A} are linearly independent

Model is **essentially unique** (i.e., up to permutation and scaling) under the condition that the sum of the factor matrix $k\text{-rank}$ values is $\geq 2r + d - 1$

$$k\text{-rank}(\mathbf{A}) + k\text{-rank}(\mathbf{B}) + k\text{-rank}(\mathbf{C}) \geq 2r + 2$$

Matrix factorization does not share this property!

Kruskal 1977, Sidiropoulos & Bro 2000

5/4/2018
Tensor Tutorial @ SDM18
62

Comparing two CP tensors



- Case 1: Single component
 - Normalize factors, absorbing weights into λ
 - Compute: score = $\rho(\mathbf{a}^T \hat{\mathbf{a}})(\mathbf{b}^T \hat{\mathbf{b}})(\mathbf{c}^T \hat{\mathbf{c}})$, $\rho = 1 - |\lambda - \hat{\lambda}| / \max\{\lambda, \hat{\lambda}\}$

Cosine of angle between vectors

 - score = 1 \Rightarrow perfect match
- Case 2: Multiple components
 - Compute: score = *average* of single component scores
 - But need to match the components to maximize the score!
 - Exact solution too expensive, so use greedy heuristic

5/4/2018
Tensor Tutorial @ SDM18
64

CP Rank

The (CP) **rank** of a tensor is the minimal value r such that the tensor can be expressed as the sum of exactly r rank-1 tensors and no fewer.

- The rank can be larger than the maximum dimension!

Tensor of size $2 \times 2 \times 2$:

$$\mathbf{X}(:, :, 1) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \mathbf{X}(:, :, 2) = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$



Rank is 3:

$$\mathbf{X} = \left[\underbrace{\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix}}_A, \underbrace{\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}}_B, \underbrace{\begin{bmatrix} 1 & 1 & 0 \\ -1 & 1 & 1 \end{bmatrix}}_C \right]$$

Kruskal 1977, Krusal 1983, Kruskal 1989, Bini et al. 1979; see also Kolda & Bader 2009

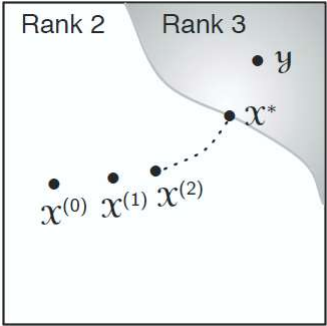
5/4/2018
Tensor Tutorial @ SDM18
65

CP Rank is NP-Hard



The (CP) **rank** of a tensor is the minimal value r such that the tensor can be expressed as the sum of exactly r rank-1 tensors and no fewer.

- Factorizations are **not nested** (Kolda 2001)
- The space of rank- r tensors is **not closed** except for $r = 1$ (Kruskal, Harshman, & Lundy 1989)
- The best rank- k approximation may **not exist** (de Silva & Lim 2006)
- Determining the rank of a tensor is **NP hard** (Håstad 1990, Hillar & Lim 2009)



Kolda & Bader 2009

5/4/2018
Tensor Tutorial @ SDM18
66



Tensor of Unknown Rank

- Specific $9 \times 9 \times 9$ tensor factorization problem
- Corresponds to being able to do fast matrix multiplication of two 3×3 matrices
- Rank is between 19 and 23 \rightarrow at most 621 variables

$x_{1,1,1} = 1$	$x_{4,2,1} = 1$	$x_{7,3,1} = 1$
$x_{1,4,2} = 1$	$x_{4,5,2} = 1$	$x_{7,6,2} = 1$
$x_{1,7,3} = 1$	$x_{4,8,3} = 1$	$x_{7,9,3} = 1$
$x_{2,1,4} = 1$	$x_{5,2,4} = 1$	$x_{8,3,4} = 1$
$x_{2,4,5} = 1$	$x_{5,5,5} = 1$	$x_{8,6,5} = 1$
$x_{2,7,6} = 1$	$x_{5,8,6} = 1$	$x_{8,9,6} = 1$
$x_{3,1,7} = 1$	$x_{6,2,7} = 1$	$x_{9,3,7} = 1$
$x_{3,4,8} = 1$	$x_{6,5,8} = 1$	$x_{9,6,8} = 1$
$x_{3,7,9} = 1$	$x_{6,8,9} = 1$	$x_{9,9,9} = 1$


Laderman 1976; Bini et al. 1979; Bläser 2003; Benson & Ballard, PPOPP'15

5/4/2018
Tensor Tutorial @ SDM18
67

Lab 1 (40 minutes)

1. Download labs and toolboxes from:
<https://tinyurl.com/tensor-tutorial-sdm18>
2. In MATLAB, change directory: `tensor_tutorial/labs`
3. Open and work through the following labs (in order):
 - a) `lab_setup.mlx`
 - b) `lab_intro_to_tensors.mlx`
 - c) `lab_cpals_cpopt.mlx`
 - d) `lab_cpals_cpopt_large_data.mlx`
4. You are free to continue working on the labs through the break.



Try it out.
Solutions Below.

When you see this in the labs, **STOP** and try out the exercises on your own.

5/4/2018
Tensor Tutorial @ SDM18
70







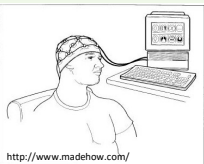

Missing Data

Tammy Kolda, Evrim Acar, Morten Mørup, and Danny Dunlavy

5/4/2018
Tensor Tutorial @ SDM18
71

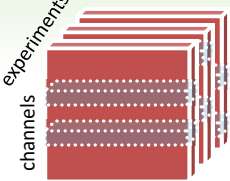
Tensor Factorizations with Missing Data?



<http://www.madehow.com/>

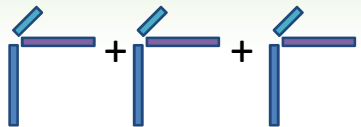
channels

experiments

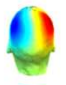


time-frequency

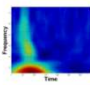
=



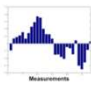
channel



time-freq



experiments



Biomedical signal processing


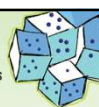
- EEG (electroencephalogram) signals can be recorded using electrodes placed on the scalp
- **Missing data problem** occurs when...
 - Electrodes get loose or disconnected, causing the signal to be unusable
 - Different experiments have overlapping but not identical channels

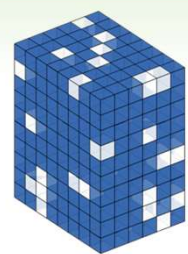
Can we still do this calculation if data are missing?

Acar, Dunlavy, Kolda & Mørup 2010 & 2011

5/4/2018
Tensor Tutorial @ SDM18
72

Approaches to Missing Data



$\Omega = \text{set of known entries (blue)}$

$$\min \sum_{i \in \Omega} (x_i - m_i)^2$$


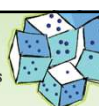
s.t. $\mathcal{M} = [\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]$

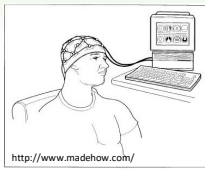
- Guess Missing Entries
- Expectation Maximization
 1. Generate initial model
 2. Estimate missing entries (using current model)
 3. Update model (data changed)
 4. Go back to step #2
- Ignore Missing Data
 - See equation at left
 - Closely related to ideas in matrix completion...except we haven't discussed about picking r

Acar, Dunlavy, Kolda & Mørup 2010 & 2011

5/4/2018
Tensor Tutorial @ SDM18
73

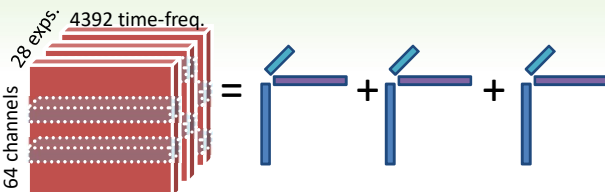
Brain dynamics can be captured even extensive missing channels



<http://www.madehow.com/>

28 expts. 4392 time-freq.





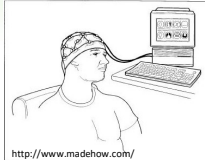
Number of Missing Channels	Replace Missing Entries with Mean
1	0.98
10	0.82
20	0.67
30	0.45
40	0.24

Acar, Dunlavy, Kolda & Mørup 2010 & 2011

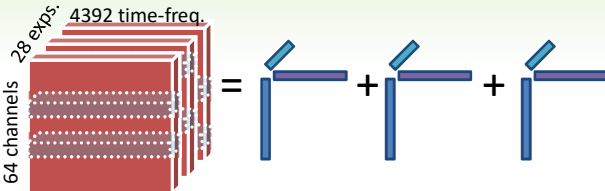
5/4/2018
Tensor Tutorial @ SDM18
74

Brain dynamics can be captured even extensive missing channels



64 channels 28 expts. 4392 time-freq.






Number of Missing Channels	Replace Missing Entries with Mean	Ignore Missing Entries
1	0.98	1.00
10	0.82	0.98
20	0.67	0.95
30	0.45	0.89
40	0.24	0.65

Acar, Dunlavy, Kolda & Mørup 2010 & 2011

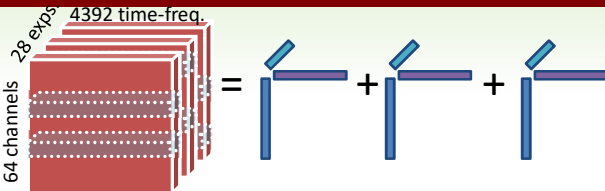
5/4/2018
Tensor Tutorial @ SDM18
75

Brain dynamics can be captured even extensive missing channels

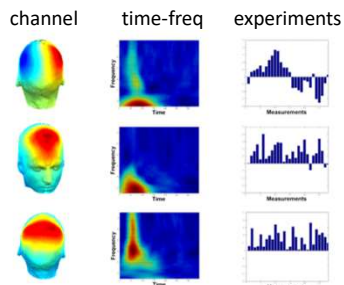





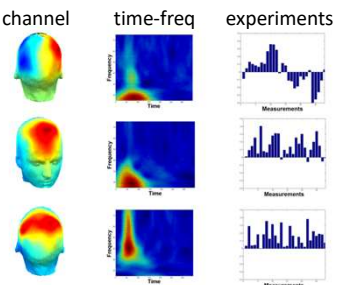
64 channels 28 expts. 4392 time-freq.



No Missing Data





30 Chan./Exp. Missing

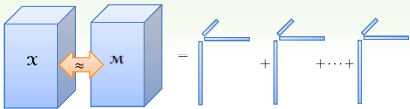


Acar, Dunlavy, Kolda & Mørup 2010 & 2011

5/4/2018
Tensor Tutorial @ SDM18
76


Review: Handling Missing Data



$$\mathcal{X} \approx \mathcal{M} = \sum_{j=1}^r \mathbf{a}_j \circ \mathbf{b}_j \circ \mathbf{c}_j = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$$

Original problem:

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \sum_i (x_i - m_i)^2 \text{ s.t. } \mathcal{M} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$$


$$F(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \sum_i (x_i - m_i)^2$$


$$F(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \sum_{i \in \Omega} (x_i - m_i)^2$$

Just compute sum at known values

Anderson-Bergman, Duersch, Hong, Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
82





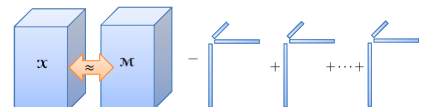
Generalized CP Decomposition

Tammy Kolda, Jed Duersch (Sandia), David Hong (Michigan), Cliff Anderson-Bergman (Sandia)

5/4/2018
Tensor Tutorial @ SDM18
83


Generalizing the Goodness-of-Fit Criteria



$$\mathcal{X} \approx \mathcal{M} = \sum_{j=1}^r \mathbf{a}_j \circ \mathbf{b}_j \circ \mathbf{c}_j = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$$

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \sum_i (x_i - m_i)^2 \text{ s.t. } \mathcal{M} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$$



$$F(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \sum_i (x_i - m_i)^2$$


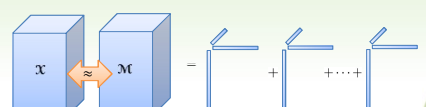
$$F(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \sum_i f(x_i, m_i)$$

Anderson-Bergman, Duersch, Hong, Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
84

"Standard" CP via Maximum Likelihood



Typically: Consider data to be low-rank plus "white noise"

$$x_i = m_i + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma)$$

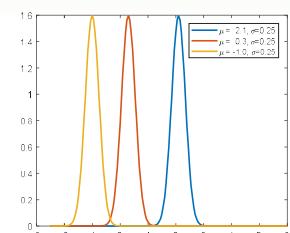
Equivalently, **Gaussian** with mean m_{ijk}

$$x_i \sim \mathcal{N}(m_i, \sigma)$$

Gaussian Probability Density Function (PDF)

$$\frac{e^{-(x-\mu)^2/2\sigma^2}}{\sqrt{2\pi\sigma^2}}$$

Probability Distribution Function:
Normal-distributed with constant σ



Minimize **negative log likelihood** with $\mu_i = m_i$ and σ constant for all entries:



$$-\log(\mathcal{L}(\mathcal{M})) = \sum_{ijk} \frac{(x_i - m_i)^2}{\cancel{\sigma^2}} + \cancel{1/\sigma^2} \cdot \cancel{\sqrt{2\pi\sigma^2}}$$

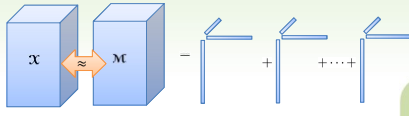
$$\min F(\mathcal{M}) = \sum_i (x_i - m_i)^2$$

Anderson-Bergman, Duersch, Hong, Kolda 2017

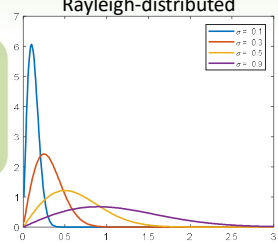
5/4/2018
Tensor Tutorial @ SDM18
85

“Rayleigh CP” with Linear Link



Probability Distribution Function:
Rayleigh-distributed



What if the data is nonnegative ($x_i \geq 0$)?
Assume data is Rayleigh-distributed.
 $x_i \sim \text{Rayleigh}(m_i)$

Rayleigh Probability Density Function (PDF)

$$\frac{x}{\sigma^2} e^{-x^2/(2\sigma^2)}$$

Minimize **negative log likelihood** with $\sigma_i = m_i$:

$$-\log(\mathcal{L}(\mathcal{M})) = \sum_i -\log \sigma_i + 2 \log m_i + \frac{x_i^2}{2m_i^2}$$

$$\min F(\mathcal{M}) = \sum_i 2 \log m_i + \frac{x_i^2}{2m_i^2}$$



Requires $m_i \geq 0$

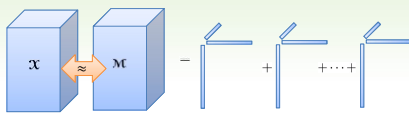
$\mathbb{E}(x_i) = m_i \sqrt{\frac{\pi}{2}}$

Anderson-Bergman, Duersch, Hong, Kolda 2017

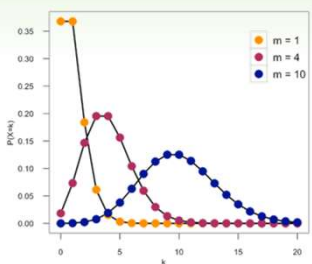
5/4/2018
Tensor Tutorial @ SDM18
86

Poisson CP: Identity Link



Probability Mass Function (PMF):



Consider data to be Poisson distributed with parameter m_i

$$x_i \sim \text{Poisson}(m_i)$$

$$\frac{e^{-\lambda} \lambda^x}{x!}$$

Minimizing **negative log likelihood** with $\lambda = m$:



$$\min F(\mathcal{X}, \mathcal{M}) = \sum_i m_i - x_i \log(m_i) + \log \frac{1}{x_i!}$$

Requires $m_i \geq 0$

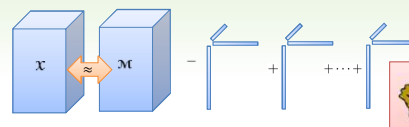
$\mathbb{E}(X_i) = m_i$

Chi & Kolda 2012, Anderson-Bergman, Duersch, Hong, and Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
87

“Boolean CP” with Odds Link



Random Coin Flip: Probability versus Odds

$p \in [0,1]$: probability of 1
 $r \geq 0$: odds ratio of 1

$$r = \frac{p}{1-p} \Leftrightarrow p = \frac{r}{1+r}$$

What if data is binary ($x_i \in \{0,1\}$)?

$m_i =$ odds ratio of $x_i = 1$.

$x_i \sim \text{Bernoulli}(m_i/(1 + m_i))$

Probability Mass Distribution (PMF) $p^x(1-p)^{1-x} \Leftrightarrow \left(\frac{r}{1+r}\right)^x \left(\frac{1}{1+r}\right)^{1-x}$



$$\mathbb{E}(x_i) = \frac{m_i}{1 + m_i}$$

Requires $m_i \geq 0$

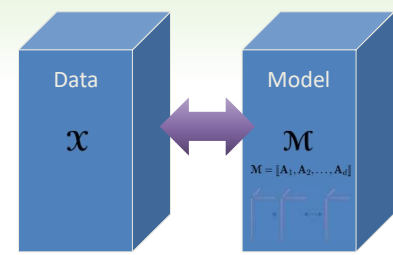
$$\min F(\mathcal{M}) = \sum_i \log(m_i + 1) - x_i \log m_i$$

Anderson-Bergman, Duersch, Hong, Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
88

Generalized CP



$$\min F(\mathcal{X}, \mathcal{M}) = \sum_i f(x_i, m_i)$$

s.t. $\mathcal{M} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket$

“Normal” CP uses:

$$f(x, m) = (x - m)^2$$

“Nonneg” CP uses :

$$f(x, m) = 2 \log m + x^2/(2m^2) \quad (x \geq 0, m \geq 0)$$

“Count” CP (Chi-Kolda 2012) uses:

$$f(x, m) = m - x \log m \quad (x \in \mathbb{N}, m \geq 0)$$

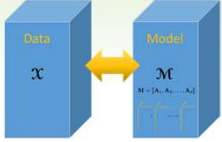
“Binary” CP uses:


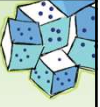
$$f(x, m) = \log(m + 1) - x \log m \quad (x \in \{0, 1\}, m \geq 0)$$

Anderson-Bergman, Duersch, Hong, Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
89

Generalized CP Gradient



$$\min F(\mathcal{X}, \mathcal{M}) = \sum_i f(x_i, m_i)$$

$$\text{s.t. } \mathcal{M} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket$$

Define a tensor \mathcal{G} such that $g(i_1, \dots, i_d) = g_i = \frac{\partial f}{\partial m}(x_i, m_i)$

Recall $m_i = \sum_{j=1}^r a_1(i_1, j) a_2(i_2, j) \dots a_d(i_d, j)$

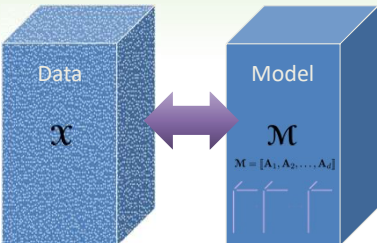
Then $\frac{\partial F}{\partial a_k(i_k, j)} = \underbrace{g(i_1, \dots, i_d)}_{\text{No dependency on model form (M)}} \underbrace{a_1(i_1, j) \dots a_{k-1}(i_{k-1}, j) a_{k+1}(i_{k+1}, j) \dots a_d(i_d, j)}_{\text{No dependency on function (f)}}$


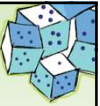
Matrix Version: $\frac{\partial F}{\partial \mathbf{A}_k} = \mathbf{G}^{(k)}(\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)$ MTTKRP!

Anderson-Bergman, Duersch, Hong, Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
90

Generalized CP Gradient with Missing Data



$$\min F(\mathcal{X}, \mathcal{M}) = \sum_{i \in \Omega} f(x_i, m_i)$$

$$\text{s.t. } \mathcal{M} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \rrbracket$$

Define a tensor \mathcal{G} such that



$$g_i = \begin{cases} \frac{\partial f}{\partial m}(x_i, m_i) & \text{if } i \in \Omega, \\ 0 & \text{if } i \notin \Omega. \end{cases}$$

Then (no change except G):

$$\frac{\partial F}{\partial \mathbf{A}_k} = \mathbf{G}^{(k)}(\mathbf{A}_d \odot \dots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \dots \odot \mathbf{A}_1)$$

Anderson-Bergman, Duersch, Hong, and Kolda 2017, Acar, Dunlavy, Kolda, Morup 2011

5/4/2018
Tensor Tutorial @ SDM18
91

Observations on Generalized CP



- Can use any optimization method to solve the optimization problem

$$\min F(\mathcal{X}, \mathcal{M}) = \sum_{i \in \Omega} f(x_i, m_i) \quad \text{s.t. } \mathcal{M} = [\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]$$
- Can easily add regularization

$$\min F(\mathcal{X}, \mathcal{M}) = \sum_{i \in \Omega} f(x_i, m_i) + \sum_k \rho_k(\mathbf{A}_k) \quad \text{s.t. } \mathcal{M} = [\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]$$
- If data tensor (X) is sparse, no guarantee gradient will be efficient
 - Standard CP is a special case that has exploitable structure
 - Scalability is potentially a major problem
- If known data is sparse (Ω), then gradient will be efficient
 - Doesn't require any sparsity in data tensor
 - Perhaps this can be exploited? Yes!

Anderson-Bergman, Duersch, Hong, Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
93

GCP Efficiency for Large Data

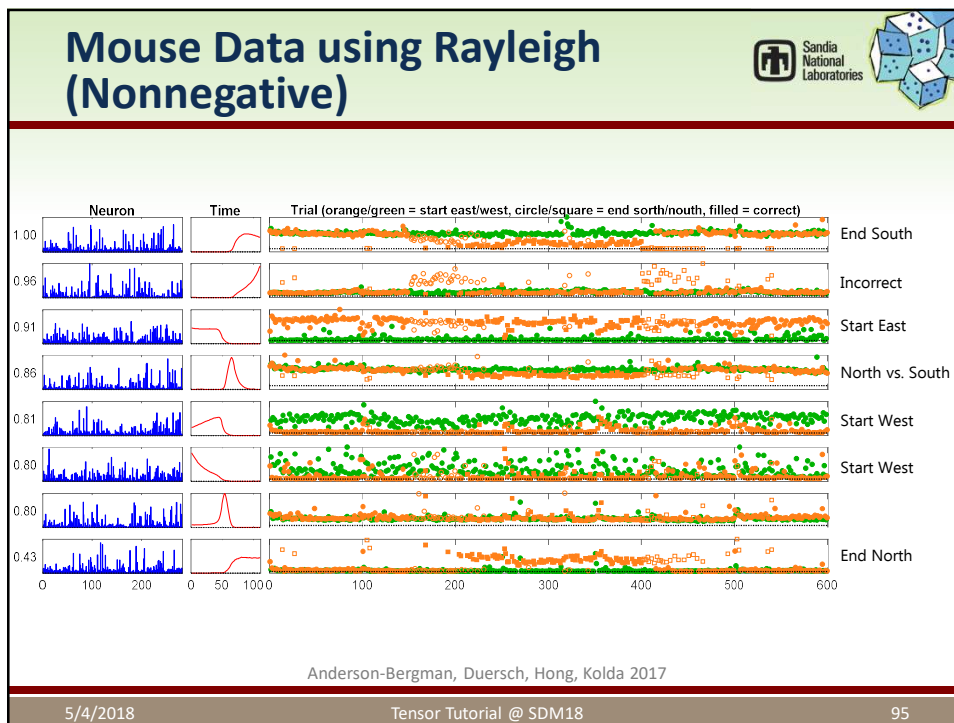
- If tensor is dense...

$$F(\mathcal{X}, \mathcal{M}) = \sum_{i \in \Omega} f(x_i, m_i)$$
 - Just leave out data at random!
 - Ω represents the data we keep, usually only 10% of the data or less
- If tensor is sparse...
 - Keep all the nonzeros but leave out all or most of the zeros
 - Add a penalty for the left out zero terms to encourage the model to stay small
$$F(\mathcal{X}, \mathcal{M}) = \sum_{i \in \Omega} f(x_i, m_i) + \underbrace{\sum_{i \notin \Omega} \alpha m_i + \beta m_i^2}_{\text{nearly free to compute}}$$



- Choices for α and β are up for debate

Anderson-Bergman, Duersch, Hong, Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
94




A Sparse Binary Dataset

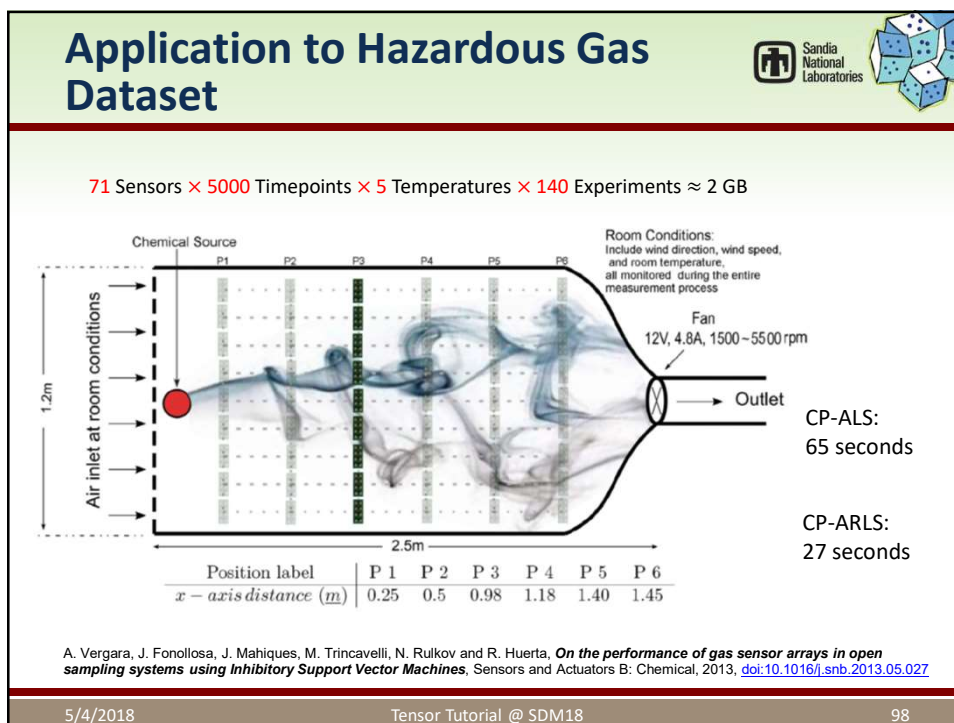
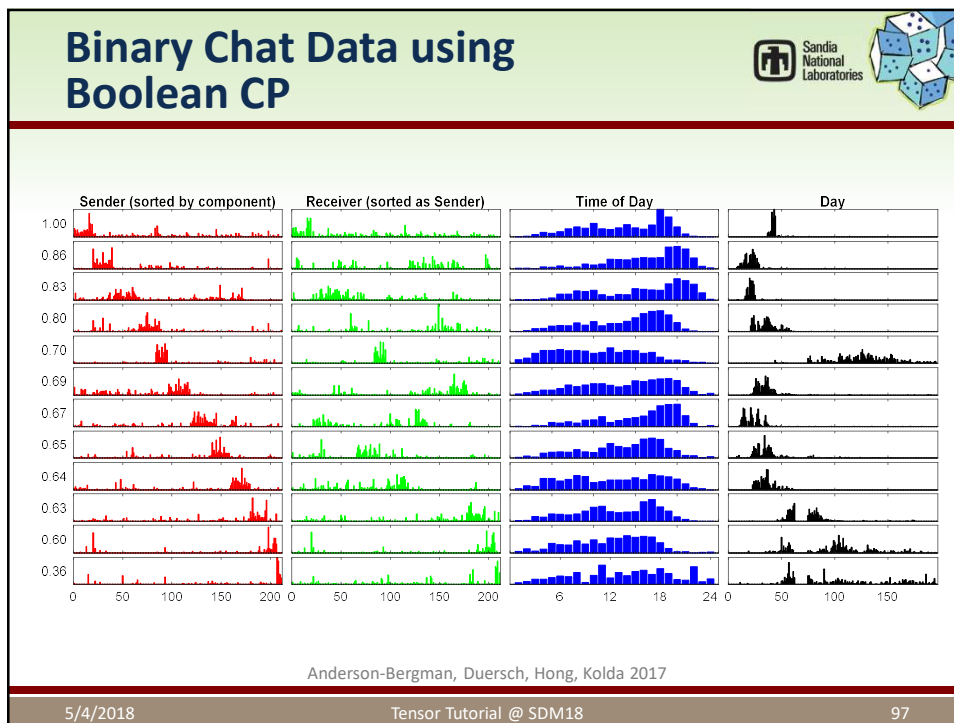
- UC Irvine Chat Network
 - 4-way binary tensor
 - Sender (211)
 - Receiver (211)
 - Hour of Day (24)
 - Day (196)
 - 14,849 nonzeros (very sparse)
- Goodness-of-fit (Boolean-odds):

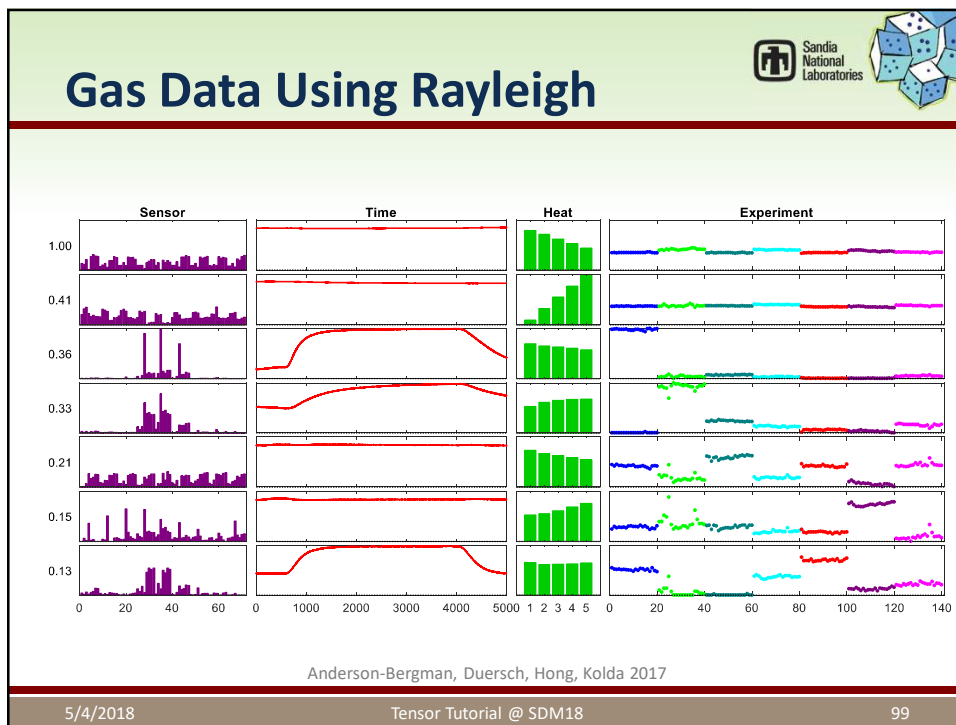
$$f(x, m) = \log(m + 1) - x \log m$$
- Use GCP to compute rank-12 decomposition





Opsahl, T., Panzarasa, P., 2009. Clustering in weighted networks. Social Networks 31 (2), 155-163, doi: 10.1016/j.socnet.2009.02.002

5/4/2018
Tensor Tutorial @ SDM18
96











Poisson Tensor Factorization


*Tammy Kolda, Eric Chi (NC State), Todd Plantenga (FireEye),
Sammy Hansen (Spotify)*

5/4/2018 Tensor Tutorial @ SDM18 107

Sparse Tensor Computations


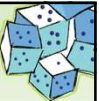
- Many real-world data analysis problems are naturally expressed as in terms of a *sparse* tensor
 - Computer traffic analysis
 - Term-document analysis
 - Email analysis
 - Link prediction
 - Web page analysis

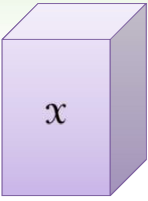


- But where do sparse tensors come from?
 - How are the entries generated?
 - How can we use that information in model fitting?

5/4/2018
Tensor Tutorial @ SDM18
108

Generating Sparse Test Data

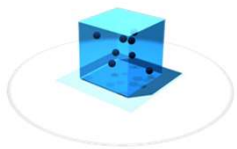


\mathcal{X}

$\sim \text{Poisson}$



$$\left(\begin{array}{c} \lambda_1 \text{ } \mathbf{c}_1 \\ \hline \mathbf{b}_1 \\ \hline \mathbf{a}_1 \end{array} + \begin{array}{c} \lambda_2 \text{ } \mathbf{c}_2 \\ \hline \mathbf{b}_2 \\ \hline \mathbf{a}_2 \end{array} + \dots + \begin{array}{c} \lambda_R \text{ } \mathbf{c}_R \\ \hline \mathbf{b}_R \\ \hline \mathbf{a}_R \end{array} \right)$$

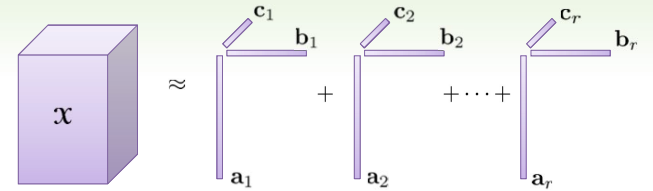
- Each “occurrence” generated as follows
- Choose component j proportional to λ_j
- Given component j :
 - Choose index i_1 proportional to \mathbf{a}_j
 - Choose index i_2 proportional to \mathbf{b}_j
 - Choose index i_3 proportional to \mathbf{c}_j
- Increment $x(i_1, i_2, i_3)$ by one
- Repeat



5/4/2018
Tensor Tutorial @ SDM18
109

Solving the Poisson Regression Problem





$$\min_{\mathcal{M}} \sum_i m_i - x_i \log m_i \text{ subject to } \mathcal{M} = [\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]$$

- Highly nonconvex problem!
 - Assume r is given
- Alternating Poisson regression
 - Assume $(d - 1)$ factor matrices are known and solve for the remaining one
 - Multiplicative updates like Lee & Seung (2000) for NMF, but improved
 - Typically assume data tensor \mathcal{X} is sparse and have special methods for this
 - Newton or Quasi-Newton method

Chi & Kolda 2012; Hansen, Plantenga, & Kolda 2015

5/4/2018
Tensor Tutorial @ SDM18
110

Alternating Poisson Regression (CP-APR)

Repeat until converged...

1. $\bar{\mathbf{A}} \leftarrow \arg \min_{\bar{\mathbf{A}} \geq 0} \sum_i m_i - x_i \log m_i \text{ s.t. } \mathcal{M} = \sum_j \bar{\mathbf{a}}_j \circ \mathbf{b}_j \circ \mathbf{c}_j$
2. $\lambda \leftarrow \mathbf{e}^T \bar{\mathbf{A}}; \mathbf{A} \leftarrow \bar{\mathbf{A}} \cdot \text{diag}(1/\lambda)$
3. $\bar{\mathbf{B}} \leftarrow \arg \min_{\bar{\mathbf{B}} \geq 0} \sum_i m_i - x_i \log m_i \text{ s.t. } \mathcal{M} = \sum_j \mathbf{a}_j \circ \bar{\mathbf{b}}_j \circ \mathbf{c}_j$
4. $\lambda \leftarrow \mathbf{e}^T \bar{\mathbf{B}}; \mathbf{B} \leftarrow \bar{\mathbf{B}} \cdot \text{diag}(1/\lambda)$
5. $\bar{\mathbf{C}} \leftarrow \arg \min_{\bar{\mathbf{C}} \geq 0} \sum_i m_i - x_i \log m_i \text{ s.t. } \mathcal{M} = \sum_j \mathbf{a}_j \circ \mathbf{b}_j \circ \bar{\mathbf{c}}_j$
6. $\lambda \leftarrow \mathbf{e}^T \bar{\mathbf{C}}; \mathbf{C} \leftarrow \bar{\mathbf{C}} \cdot \text{diag}(1/\lambda)$

} Fix **B,C**;
solve for **A**

} Fix **A,C**;
solve for **B**



} Fix **A,B**;
solve for **C**

Convergence Theory

Theorem: The CP-APR algorithm will **converge to a constrained stationary point** if the subproblems are strictly convex and solved exactly at each iteration.


Chi & Kolda 2012

5/4/2018
Tensor Tutorial @ SDM18
111

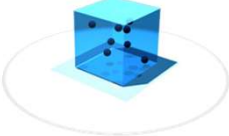



Solving the Subproblem

$$\min_{\mathbf{A} \geq 0} \sum_i m_i - x_i \log m_i \text{ s.t. } \mathcal{M} = \sum_j \bar{\mathbf{a}}_j \circ \mathbf{b}_j \circ \mathbf{c}_j$$



$$\min_{\mathbf{A} \geq 0} \sum_i \left(\sum_j \bar{a}_{i_1 j} b_{i_2 j} c_{i_3 j} \right) - x_i \log \left(\sum_j \bar{a}_{i_1 j} b_{i_2 j} c_{i_3 j} \right)$$





Lemma: The subproblems are **strictly convex** under mild conditions.

Chi & Kolda 2012

5/4/2018

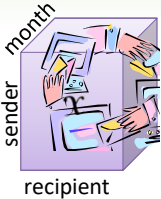
Tensor Tutorial @ SDM18

112

Motivating Example: Enron Email

- Emails from Enron FERC investigation
 - 8540 Messages
 - 28 Months (from Dec 1999 to Mar 2002)
 - 105 People (sent and received at least one email every month)
 - $x(i_1, i_2, i_3) = \# \text{ emails from sender } i_1 \text{ to recipient } i_2 \text{ in month } i_3$
 - $105 \times 105 \times 28 = 308,700$ possible entries
 - 8,500 nonzero counts
 - **3% dense**
- Questions: What can we learn about this data?
 - Each person labeled by Zhou et al. (2007); see also Owen and Perry (2010)
 - Seniority: 57% senior, 43% junior
 - Gender: 67% male, 33% female
 - Department: 24% legal, 31% trading, 45% other

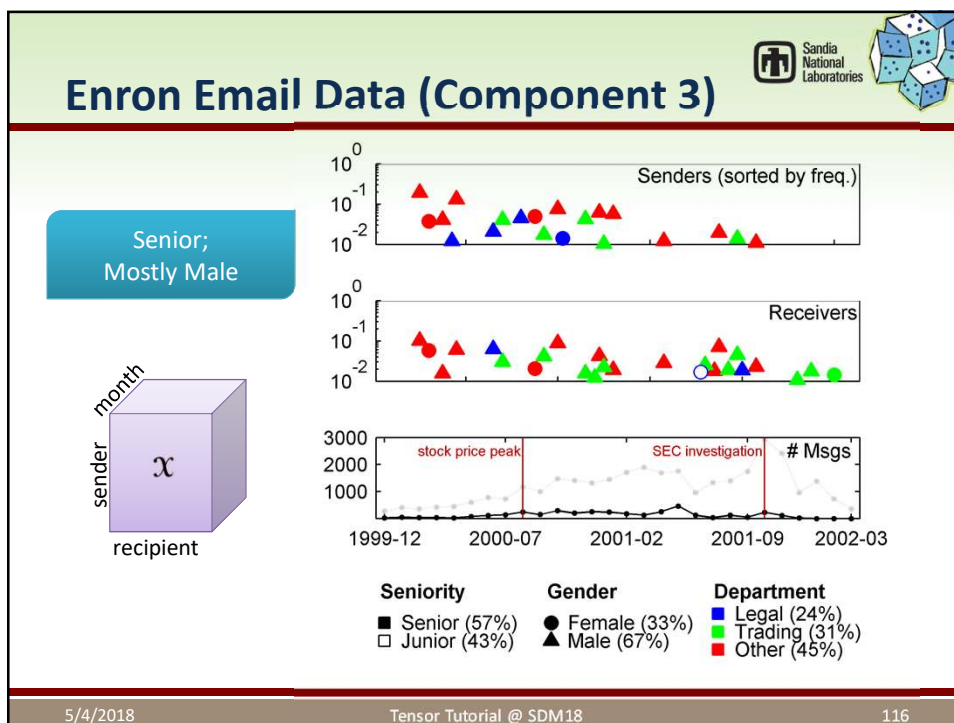
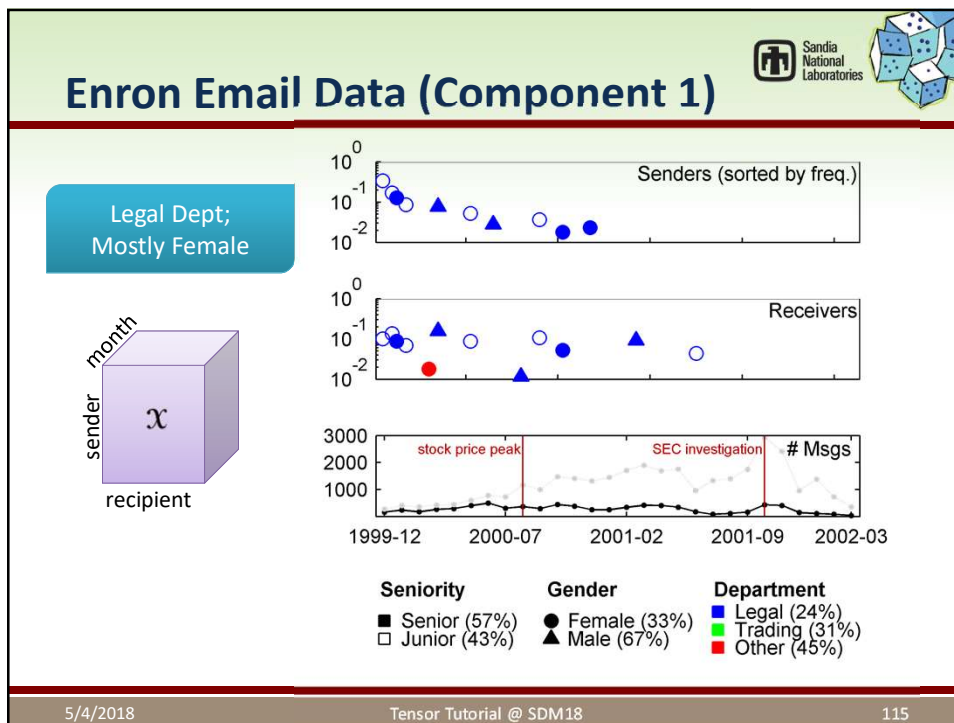


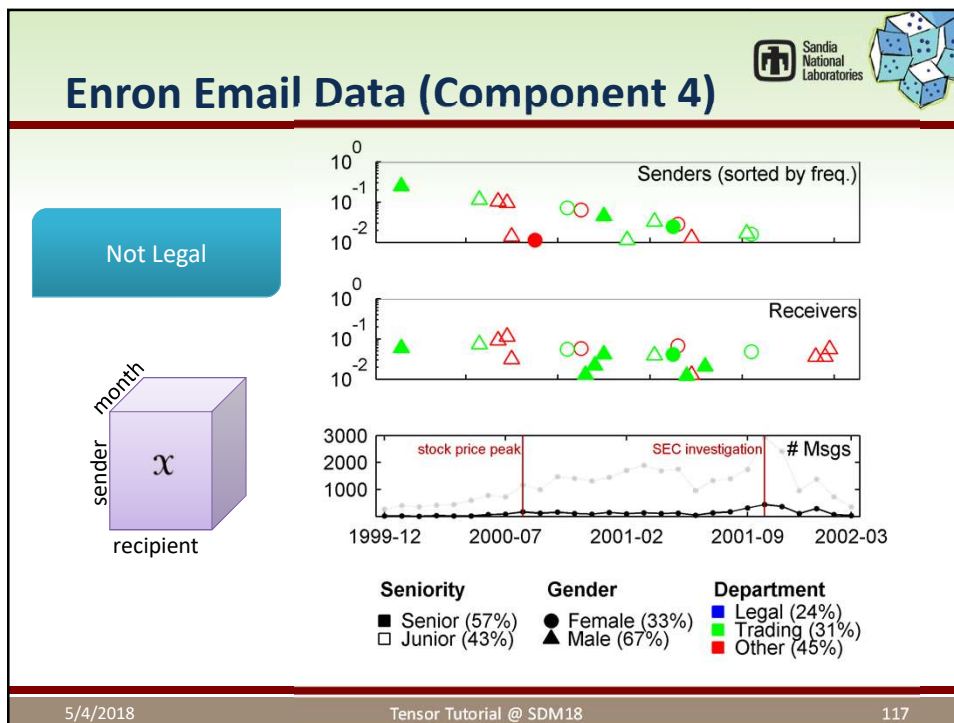
This information is not part of the tensor factorization


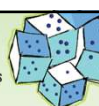
5/4/2018


Tensor Tutorial @ SDM18

114









Randomized Least Squares for CP Decomposition

Tammy Kolda, Casey Battaglini (GA Tech)
and Grey Ballard (Wake Forest)

5/4/2018
Tensor Tutorial @ SDM18
118

CP-ALS: Fitting CP Model via Alternating Least Squares

Repeat until fit stops changing...

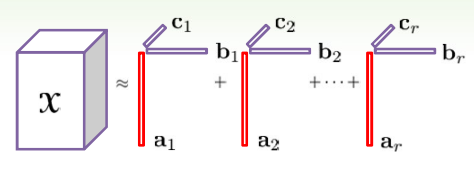
$$\min_{\mathbf{A}} \|\mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top\|^2$$

$$\min_{\mathbf{B}} \|\mathbf{X}_{(2)} - \mathbf{B}(\mathbf{C} \odot \mathbf{A})^\top\|^2$$

$$\min_{\mathbf{C}} \|\mathbf{X}_{(3)} - \mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top\|^2$$

\mathbf{X}

 \approx



$$\min_{\mathbf{A}} \|\mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top\|^2$$

“right hand sides”

 $\|\mathbf{X}_{(1)}\|$
 $n_k \times n^d/n_k$

\mathbf{A}
 $n_k \times r$

“matrix”

 $(\mathbf{C} \odot \mathbf{B})^\top$
 \vdots
 $(\mathbf{C}_r \odot \mathbf{B}_r)^\top$



Khatri-Rao Product

 $\|\mathbf{C} \odot \mathbf{B}\|_F^2$
 $r \times n^d/n_k$

Battaglino, Ballard, & Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
119

Recall: Special Structure of the Least Squares Problem

$$\min_{\mathbf{A}} \|\underbrace{\mathbf{X}_{(1)}}_{n_1 \times (n_2 n_3)} - \underbrace{\mathbf{A}}_{n_1 \times r} \underbrace{(\mathbf{C} \odot \mathbf{B})^\top}_{r \times (n_2 n_3)}\|^2$$

$$(\mathbf{C} \odot \mathbf{B})\mathbf{A}^\top = \mathbf{X}_{(1)}^\top$$

$$\mathbf{A}^\top = (\mathbf{C} \odot \mathbf{B})^\dagger \mathbf{X}_{(1)}^\top$$

$$\mathbf{A}^\top = (\mathbf{C}^\top \mathbf{C} * \mathbf{B}^\top \mathbf{B})^\dagger (\mathbf{C} \odot \mathbf{B})^\top \mathbf{X}_{(1)}^\top$$

$$\mathbf{A} = \mathbf{X}_{(1)} (\mathbf{C} \odot \mathbf{B}) (\mathbf{C}^\top \mathbf{C} * \mathbf{B}^\top \mathbf{B})^\dagger$$

The most expensive step is *not* the backsolve.


Rather, it's the formation of the Khatri-Rao product!


So, how will randomized methods help?


Battaglino, Ballard, & Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
120

CP Least Squares Problem

$\| \mathbf{X}_{(1)} \|^2_F$

 $n_k \times n^d / n_k$

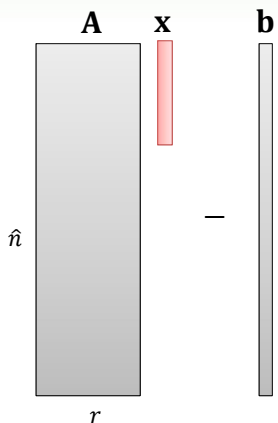
$- \mathbf{A}$

 $n_k \times r$

$[(\mathbf{C} \odot \mathbf{B})^\top]^2_F$

 $r \times n^d / n_k$

How to randomize this?

5/4/2018
Tensor Tutorial @ SDM18
121

Aside: Sketching for Standard Least Squares

$\min_x \| \mathbf{A} \mathbf{x} - \mathbf{b} \|^2$


\hat{n}

r

$\mathcal{O}(\hat{n}r^2)$



$\mathbf{x} \leftarrow \mathbf{A} \backslash \mathbf{b}$

Backslash causes MATLAB to automatically call the best solver (qr, etc.)

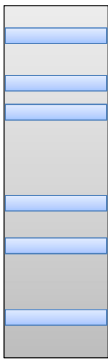
Sarlós 2006, Woodruff 2014

5/4/2018
Tensor Tutorial @ SDM18
122

Sampled Least Squares

Choose q rows, uniformly at random



\hat{n}

r

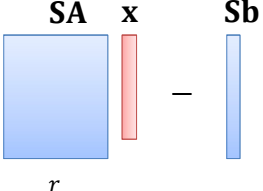
$\mathcal{O}(\hat{n}r^2)$

\mathbf{x}

\mathbf{b}

\mathbf{S}

approximate



q


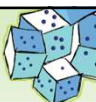
r


$\mathcal{O}(qr^2)$

Sampling only guaranteed to "work" if the A is incoherent.

5/4/2018
Tensor Tutorial @ SDM18
123

Enforcing Incoherence



\hat{n}

r

\mathbf{x}

\mathbf{b}



Mixing Matrix

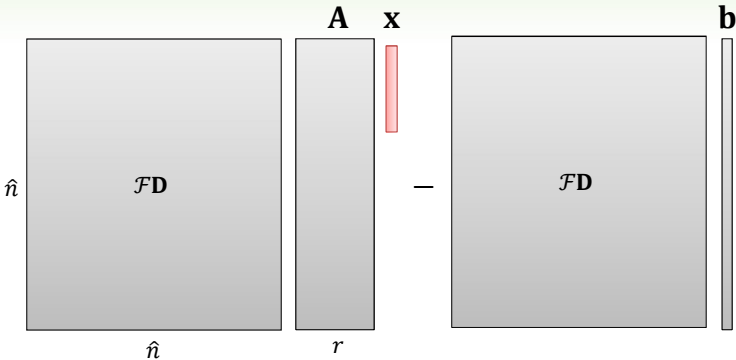
Mixing Matrix

- Many good choices of mixing matrix, such as a matrix with entries chosen from a uniform random distribution.
- But no reduction in cost!

5/4/2018
Tensor Tutorial @ SDM18
124

Fast Johnson-Lindenstrauss Transform (FJLT)





- Instead, use FFT (\mathcal{F}) followed by random diagonal with +/- 1 entries (\mathbf{D}).
- Costs only $r \log \hat{n}$ to apply
- Practical application in Blendenpik, yielding $\sim 4X$ speedup versus LAPACK

Sarlós 2006, Woodruff 2014, Ailon & Chazelle 2006, Avron, Maymounkov, & Toledo 2010

5/4/2018
Tensor Tutorial @ SDM18
125

Sampled/Mixed Least Squares

$$\min_x \|Ax - b\|$$

$$\min_x \|SAx - Sb\|$$

Sampling only,
No mixing.



$$\min_x \|S\mathcal{F}DAx - S\mathcal{F}Db\|$$

Sampling +
Mixing

Ailon & Chazelle 2006; Avron, Maymounkov, & Toledo 2010

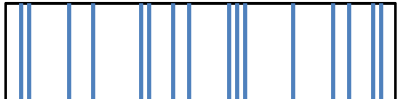
5/4/2018
Tensor Tutorial @ SDM18
126

CP-ARLS





$$\min_{\mathbf{A}} \|\mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^T\|^2$$


$$\mathbf{A} \leftarrow \mathbf{X}_{(1)} / (\mathbf{C} \odot \mathbf{B})^T$$



$\|\mathbf{X}_{(1)}\mathbf{S}^T$



- \mathbf{A}



$[(\mathbf{C} \odot \mathbf{B})^T\mathbf{S}^T] \|\|_F^2$

$$\mathbf{A} \leftarrow \mathbf{X}_{(1)}\mathbf{S}^T / (\mathbf{C} \odot \mathbf{B})^T\mathbf{S}^T$$



-

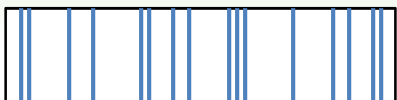
-

Battaglino, Ballard, & Kolda 2017


5/4/2018
Tensor Tutorial @ SDM18
127

CP-ARLS Trick #1: Avoid unfolding data tensor

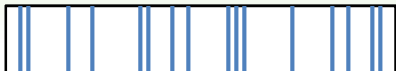





$\|\mathbf{X}_{(1)}\mathbf{S}^T$



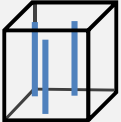
- \mathbf{A}




$[(\mathbf{C} \odot \mathbf{B})^T\mathbf{S}^T] \|\|_F^2$


Type equation here.


Trick 1:





...







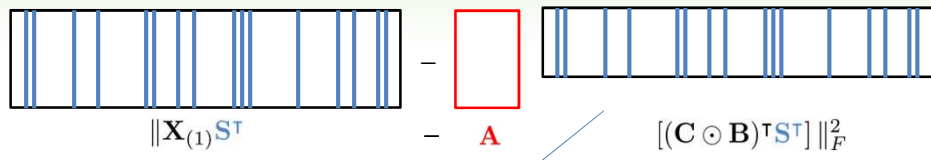
Data movement is often as expensive or more expensive than FLOPS.
Just move the minimum and no more.

Battaglino, Ballard, & Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
128

CP-ARLS Trick #2: Don't form Khatri-Rao Product



$\| \mathbf{X}_{(1)} \mathbf{S}^T - \mathbf{A} - [(\mathbf{C} \odot \mathbf{B})^T \mathbf{S}^T] \|_F^2$

Trick 2:

Each column in the sample is of the form:
 $(\mathbf{C}(\ell, :) .* \mathbf{B}(k, :))^T$



The Khatri-Rao product is actually the most expensive part of CP-ALS. Skip this and save lotsa time.

Mixing covered in the paper!

Battaglino, Ballard, & Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
129

Randomizing the Convergence Check

Repeat until fit stops changing...

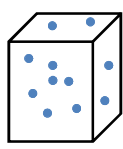
$\min_{\mathbf{A}} \| \mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^T \|^2$

$\min_{\mathbf{B}} \| \mathbf{X}_{(2)} - \mathbf{B}(\mathbf{C} \odot \mathbf{A})^T \|^2$

$\min_{\mathbf{C}} \| \mathbf{X}_{(3)} - \mathbf{C}(\mathbf{B} \odot \mathbf{A})^T \|^2$

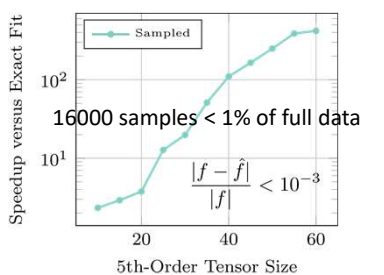
$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \sum_i (x_i - m_i)^2$

s.t. $\mathcal{M} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$



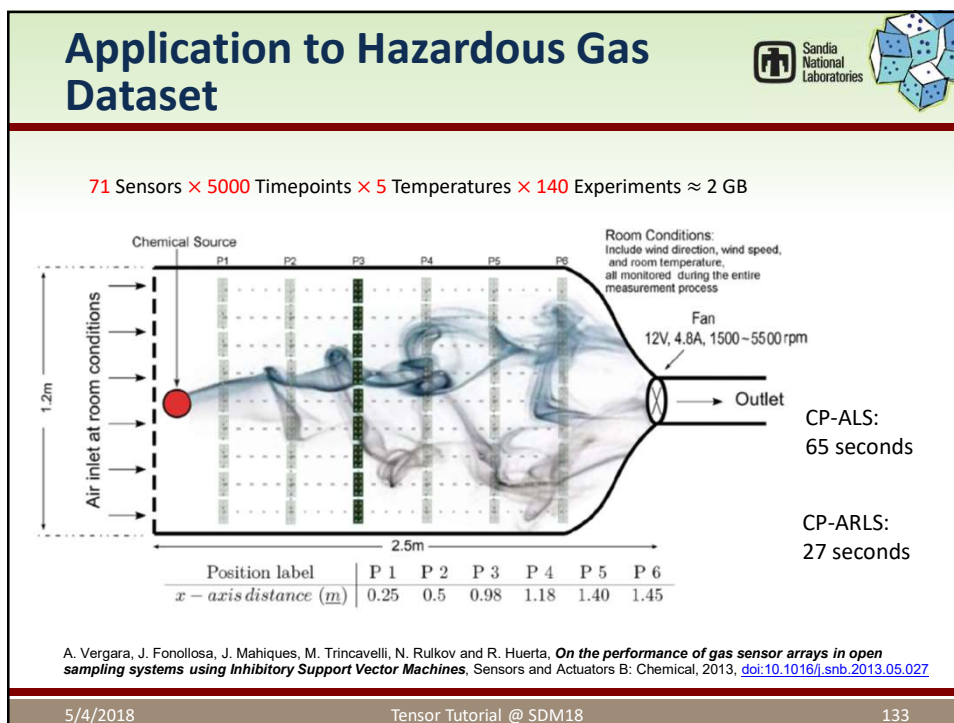
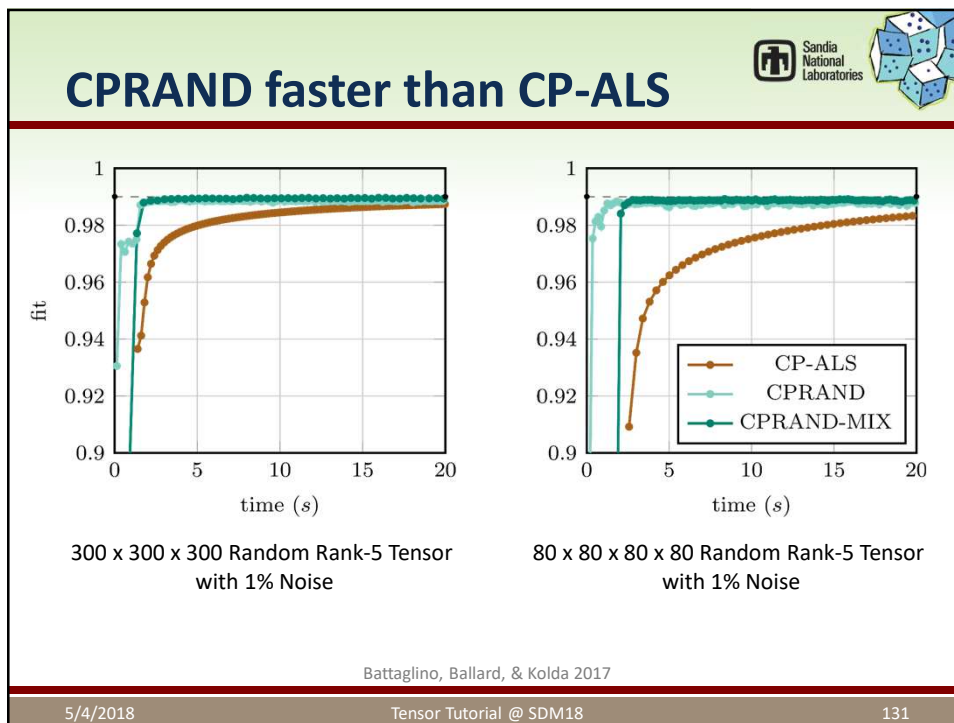
Estimate convergence of function values using small random subset of elements in function evaluation (use Chernoff-Hoeffding to bound accuracy)

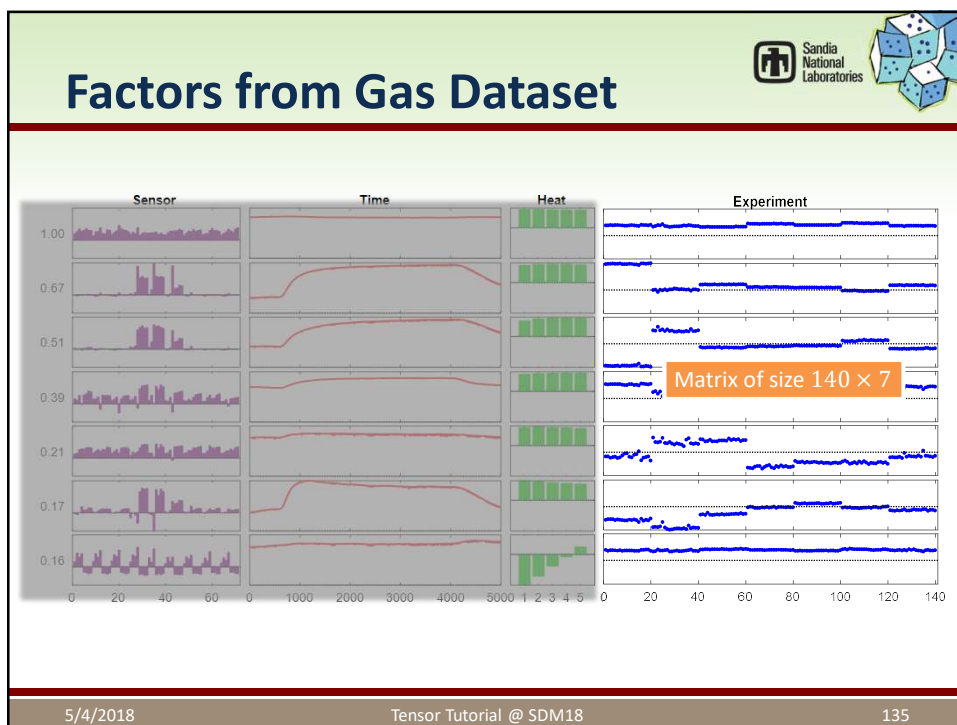
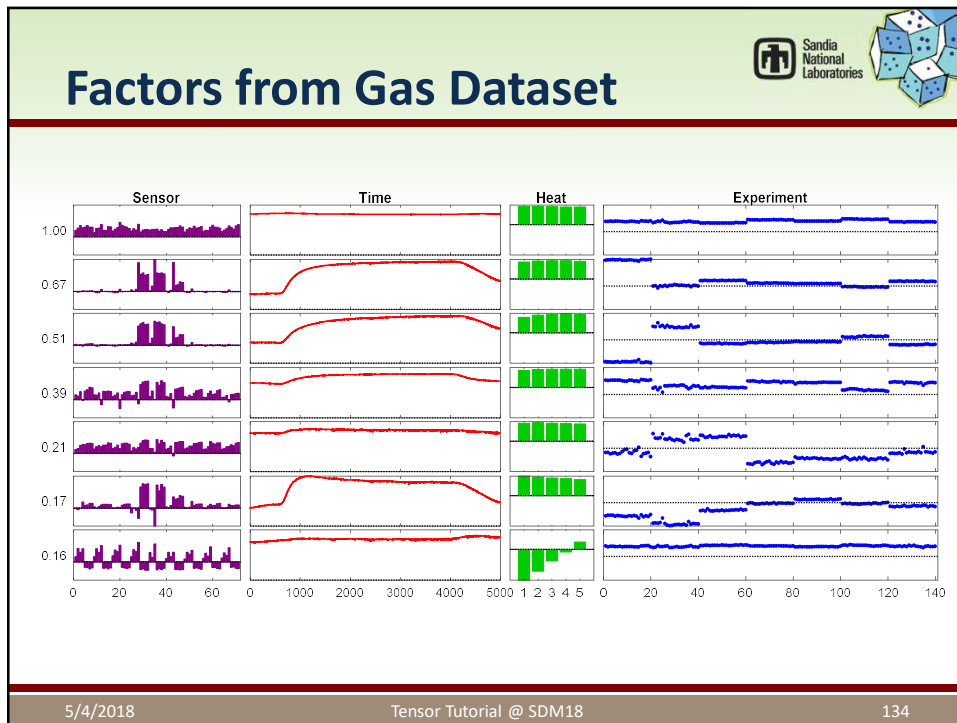
$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{n^d}{|\hat{\Omega}|} \sum_{i \in \hat{\Omega}} (x_i - m_i)^2$

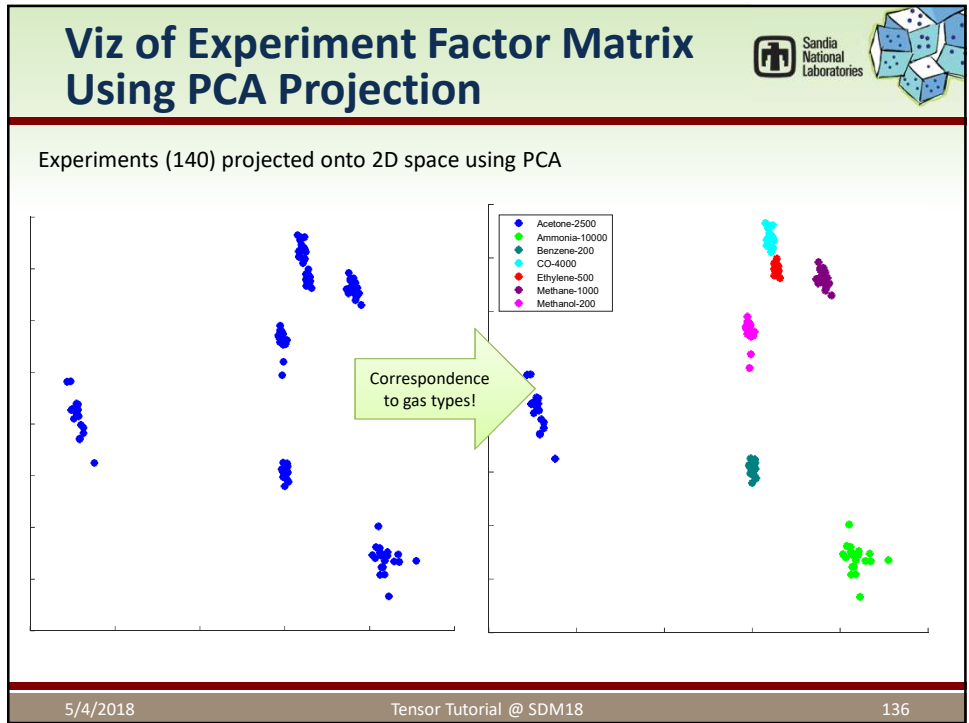


Battaglino, Ballard, & Kolda 2017

5/4/2018
Tensor Tutorial @ SDM18
130









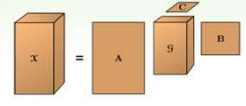
Summary

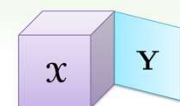
5/4/2018 Tensor Tutorial @ SDM18 138

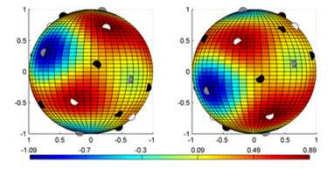



What We Didn't Cover

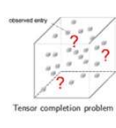
- Other tensor models
 - Tucker decomposition
 - Tensor train and hierarchical tensor decomposition
- Coupled factorizations
- Tensor eigenproblems
- Symmetric tensor decompositions
- Functional (i.e., continuous) tensor decomposition
- Tensor completion & tensor nuclear norm
- And much more!









$$f(x_1, \dots, x_d) = \sum_{j=1}^r \phi_j^{(1)}(x_1) \cdots \phi_j^{(d)}(x_d)$$



Tensor completion problem

5/4/2018
Tensor Tutorial @ SDM18
139

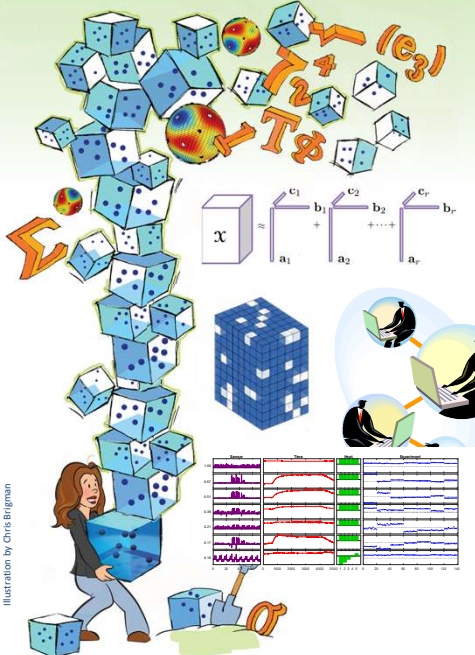


Illustration by Chris Bragman

Takeaways

- CP has many applications
- Nonconvex, difficult optimization problem
- Determining number of components is NP-hard but critical for real-world applications
- How to handle missing data? Ignore it!
- Alternative objective functions
 - Poisson tensor factorization
 - Generalized CP
- MTTKRP is key kernel
- Matrix sketching speeds up least squares subproblems
- Many open problems!
- Applications
 - Sensor data
 - Neuroscience data
 - Network data

5/4/2018
Tensor Tutorial @ SDM18
140



Acknowledgements

Current Collaborators

- Cliff Anderson-Bergman (Sandia)
- Gavin Baker (Sandia)
- Grey Ballard (Sandia/Wake Forest)
- Richard Barrett (Sandia)
- Casey Battaglini (Sandia/GA Tech)
- Jon Berry (Sandia)
- Karen Devine (Sandia)
- Jed Duersch (Sandia)
- Alex Gorodetsky (Sandia/Michigan)
- David Hong (Sandia/Michigan)
- Alicia Klinvex (Sandia)
- Hemanth Kolla (Sandia)
- Rich Lehoucq (Sandia)
- Jiajia Li (GA Tech)
- Eric Phipps (Sandia)
- Prashant Rai (Sandia)
- Keita Teranishi (Sandia)
- Rich Vuduc (GA Tech)
- Alex Williams (Sandia/Stanford)
- Kina Winoto (Sandia)
- Jeff Young (GA Tech)
- Plus many more previous collaborators!

Tensor Toolbox for MATLAB:
<https://www.tensortoolbox.org/>
 Bader, Kolda, Acar, Dunlavy, and others

Your Presenters Today

- Danny Dunlavy
- Tammy Kolda

Illustration by Chris Bringham

5/4/2018 Tensor Tutorial @ SDM18 141




References

- **Overall review:** T. G. Kolda and B. W. Bader, *Tensor Decompositions and Applications*, SIAM Review, 2009, [doi:10.1137/07070111X](https://doi.org/10.1137/07070111X)
- **Original MATLAB paper:** B. W. Bader and T. G. Kolda, *Algorithm 862: MATLAB Tensor Classes for Fast Algorithm Prototyping*, ACM Transactions on Mathematical Software, 2006, [doi:10.1145/1186785.1186794](https://doi.org/10.1145/1186785.1186794)
- **Special MATLAB classes:** B. W. Bader and T. G. Kolda, *Efficient MATLAB Computations with Sparse and Factored Tensors*, SIAM Journal on Scientific Computing, 2007, [doi:10.1137/060676489](https://doi.org/10.1137/060676489)
- **Application to Mouse Data:** A. Williams, T. G. Kolda, S. Ganguli, et al., *Unsupervised discovery of low-dimensional neural dynamics both within and across single trials through tensor analysis*, bioRxiv, 2017 (to appear in Neuron)
- **Application to Link Prediction:** D. M. Dunlavy, T. G. Kolda and E. Acar, *Temporal Link Prediction using Matrix and Tensor Factorizations*, ACM Transactions on Knowledge Discovery from Data, 2011, [doi:10.1145/1921632.1921636](https://doi.org/10.1145/1921632.1921636)
- **All-at-once Optimization Approach:** E. Acar, D. M. Dunlavy and T. G. Kolda, *A Scalable Optimization Approach for Fitting Canonical Tensor Decompositions*, Journal of Chemometrics, 2011, [doi:10.1002/cem.1335](https://doi.org/10.1002/cem.1335)
- **Missing data:** E. Acar, D. M. Dunlavy, T. G. Kolda, M. Mørup, *Scalable Tensor Factorizations for Incomplete Data*, Chemometrics and Intelligent Laboratory Systems, 2011, [doi:10.1016/j.chemolab.2010.08.004](https://doi.org/10.1016/j.chemolab.2010.08.004)
- **More missing data:** E. Acar, D. M. Dunlavy, T. G. Kolda, M. Mørup, *Scalable Tensor Factorizations with Missing Data*, SDM10: Proceedings of the 2010 SIAM International Conference on Data Mining, 2010, [doi:10.1137/1.9781611972801.61](https://doi.org/10.1137/1.9781611972801.61)
- **Randomized CP-ALS:** C. Battaglini, G. Ballard and T. G. Kolda, *A Practical Randomized CP Tensor Decomposition*, January 2017, [arXiv:1701.06600](https://arxiv.org/abs/1701.06600) (to appear in SIAM J. Matrix Analysis and Applications)
- **Poisson tensor factorization:** E. C. Chi and T. G. Kolda, *On Tensors, Sparsity, and Nonnegative Factorizations*, SIAM Journal on Matrix Analysis and Applications, 2012, [doi:10.1137/110859063](https://doi.org/10.1137/110859063)
- **More Poisson tensor factorization:** S. Hansen, T. Plantenga and T. G. Kolda, *Newton-Based Optimization for Kullback-Leibler Nonnegative Tensor Factorizations*, Optimization Methods and Software, 2015, [doi:10.1080/10556788.2015.1009977](https://doi.org/10.1080/10556788.2015.1009977)
- **Generalized CP:** Cliff Anderson-Bergman, J. Duersch, D. Hong, T. G. Kolda, *Generalized Canonical Polyadic Tensor Decomposition*, 2017 (coming soon)

Contact: Tammy Kolda, tgkolda@sandia.gov Danny Dunlavy, dmdunla@sandia.gov

5/4/2018 Tensor Tutorial @ SDM18 142

Lab 2 (20 minutes)



Try it out.
Solutions Below.

- `lab_missing_data.mlx`
- `lab_nonnegative_cp.mlx`
- `lab_boolean_cp.mlx`
- `lab_cross_validation.mlx`
- `lab_possion.mlx`
- `lab_randomization.mlx`

5/4/2018 Tensor Tutorial @ SDM18 143