

# CS665: Advanced Data Mining

**Lecture#18: Tensor**  
**U Kang**  
**KAIST**

## Most of slides by

- Dr. Tamara Kolda (Sandia N.L.)
- <http://csmr.ca.sandia.gov/~tgkolda>



- Dr. Jimeng Sun (Georgia Tech)
- <http://www.sunlab.org>



3h tutorial: [www.cs.cmu.edu/~christos/TALKS/SDM-tut-07/](http://www.cs.cmu.edu/~christos/TALKS/SDM-tut-07/)

# Outline

- ➡ ☐ **Motivation - Definitions**
- ☐ Tensor tools
- ☐ Case Studies
- ☐ Conclusion

# Motivation 0: Why “matrix”?

- Why matrices are important?

# Examples of Matrices:

## Graph - social network

	John	Peter	Mary	Nick	...
John	0	11	22	55	...
Peter	5	0	6	7	...
Mary	...	...	...	...	...
Nick	...	...	...	...	...
...	...	...	...	...	...

# Examples of Matrices:

## cloud of n-d points

	chol#	blood#	age	..	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary	...	...	...	...	...
Nick	...	...	...	...	...
...	...	...	...	...	...

# Examples of Matrices:

## Market basket

### ■ market basket as in Association Rules

	milk	bread	choc.	wine	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary	...	...	...	...	...
Nick	...	...	...	...	...
...	...	...	...	...	...

# Examples of Matrices:

## Documents and terms

	data	mining	classif.	tree	...
Paper#1	13	11	22	55	...
Paper#2	5	4	6	7	...
Paper#3	...	...	...	...	...
Paper#4	...	...	...	...	...
...	...	...	...	...	...



# Examples of Matrices:

## Authors and terms

	data	mining	classif.	tree	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary	...	...	...	...	...
Nick	...	...	...	...	...
...	...	...	...	...	...

# Examples of Matrices: sensor-ids and time-ticks

	temp1	temp2	humid.	pressure	...
t1	13	11	22	55	...
t2	5	4	6	7	...
t3	...	...	...	...	...
t4	...	...	...	...	...
...	...	...	...	...	...

# Motivation: Why tensors?

- Q: what is a tensor?

## Motivation 2: Why tensor?

- A: N-D generalization of matrix:

KDD'07

data            mining    classif.    tree            ...

John  
Peter  
Mary  
Nick  
...

13	11	22	55	...
5	4	6	7	...
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...

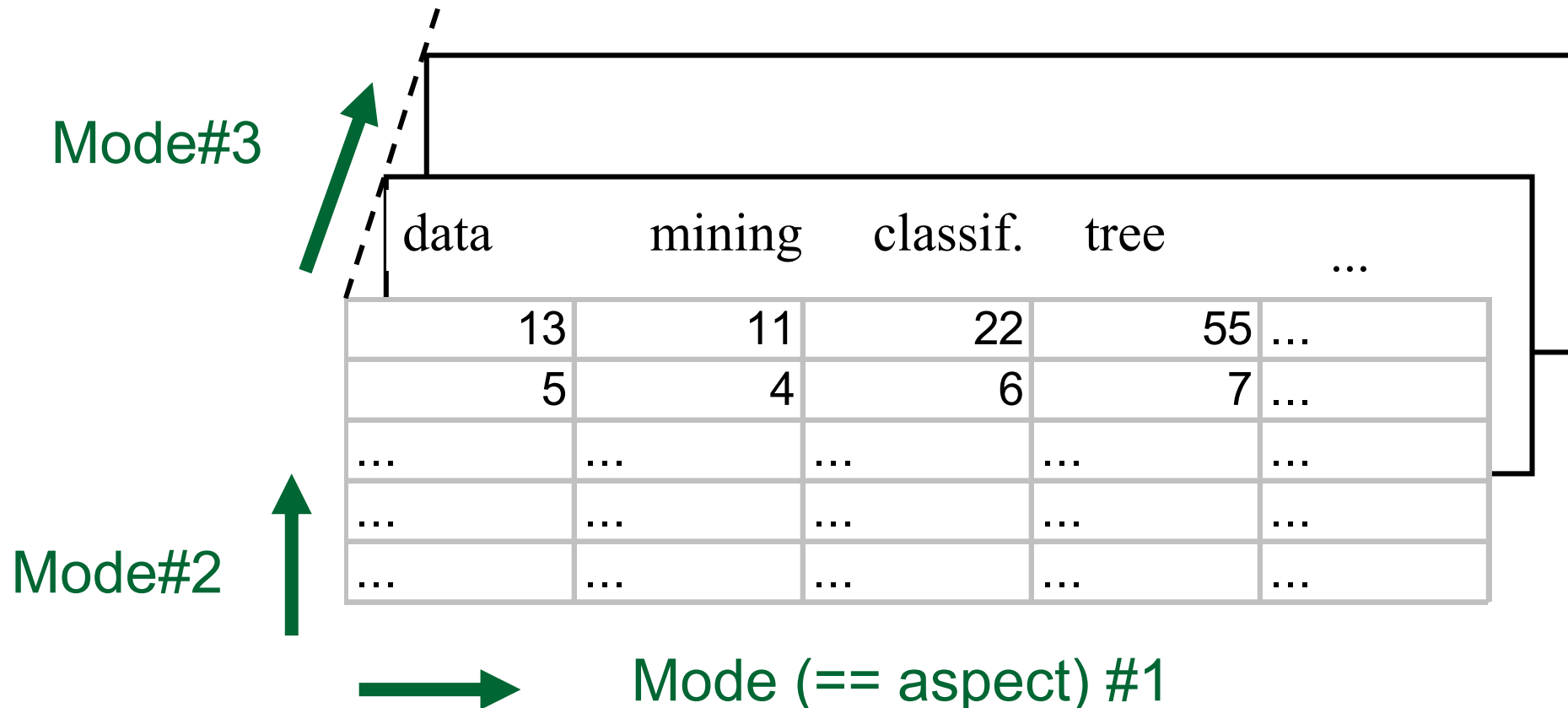
## Motivation 2: Why tensor?

- A: N-D generalization of matrix:

KDD'05					
KDD'06					
KDD'07	data	mining	classif.	tree	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary	...	...	...	...	...
Nick	...	...	...	...	...
...	...	...	...	...	...

# Tensors are useful for 3 or more modes

Terminology: ‘mode’ (or ‘aspect’):

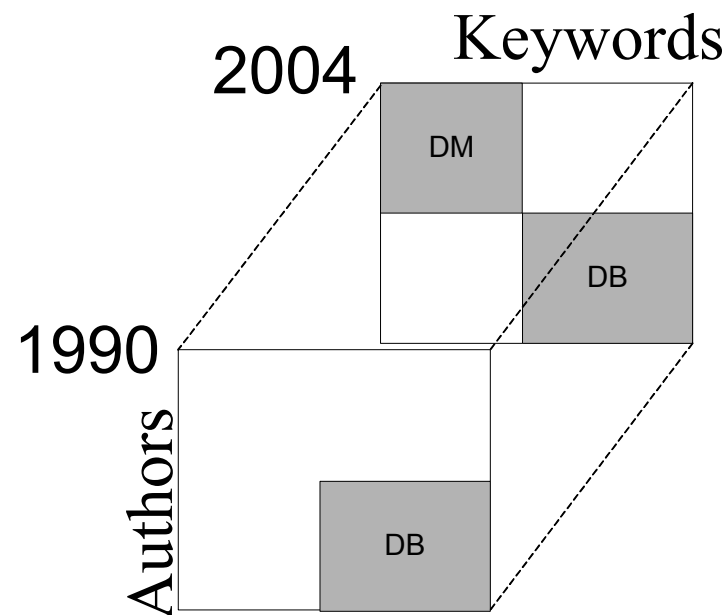


# Motivating Applications

- Why matrices are important?
- Why tensors are useful?
  - P1: social networks
  - P2: web mining

# P1: Social network analysis

- Previous research: people focus on static networks and find community structures
- We plan to monitor the change of the community structure over time





## P2: Web graph mining

- How to order the importance of web pages?
  - Kleinberg's algorithm HITS
  - PageRank
  - Tensor extension on HITS (**TOPHITS**)
    - context-sensitive hypergraph analysis

# Outline

☒ Motivation - Definitions

➡ ☐ **Tensor tools**

➡ Tensor basics

Parafac

Tucker

☐ Case Studies

☐ Conclusion

# Reminder: SVD

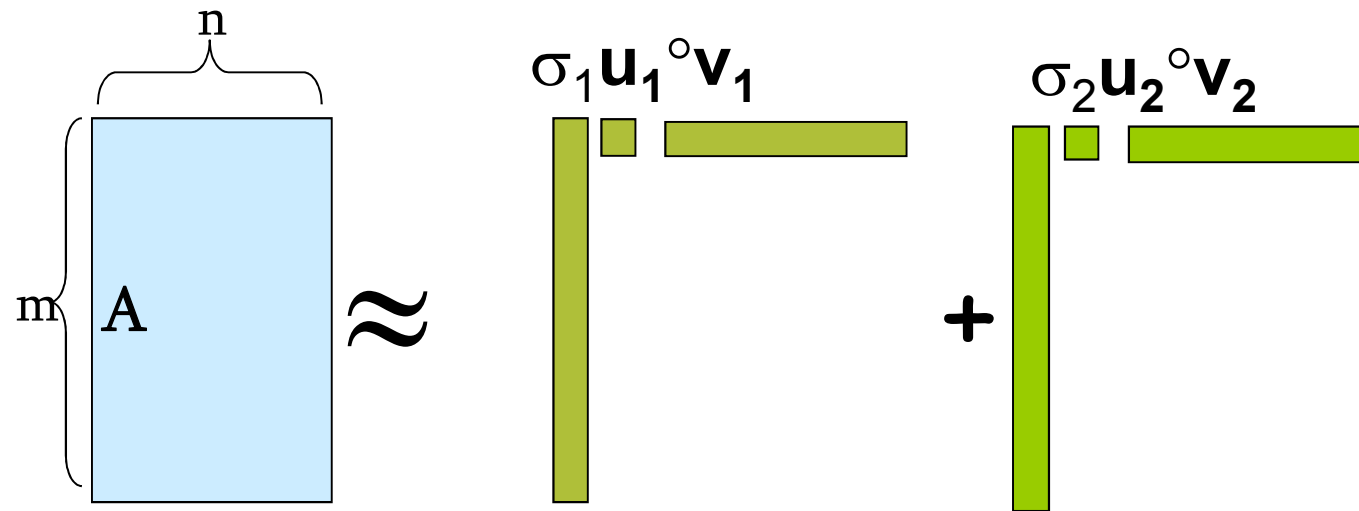
$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i$$

The diagram illustrates the SVD decomposition of matrix  $\mathbf{A}$ . Matrix  $\mathbf{A}$  is shown as a blue rectangle with dimensions  $m$  (rows) and  $n$  (columns). It is approximated by the product of three matrices:  $\mathbf{U}$  (green rectangle,  $m \times k$ ),  $\mathbf{\Sigma}$  (white square,  $k \times k$ ), and  $\mathbf{V}^T$  (green rectangle,  $k \times n$ ). The approximation is shown as  $\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . The resulting approximation is also shown as a sum of outer products:  $\sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i$ .

- Best rank- $k$  approximation in L2

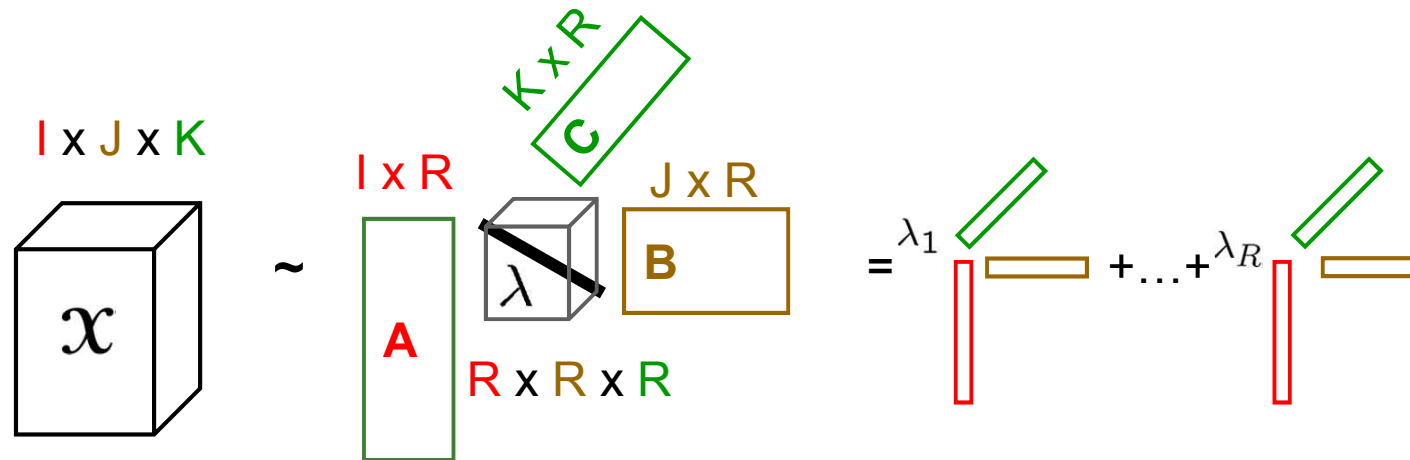
# Reminder: SVD

$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i$$



- Best rank- $k$  approximation in L2

## Goal: extension to $\geq 3$ modes



$$\mathcal{X} \approx [\lambda ; A, B, C] = \sum_r \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

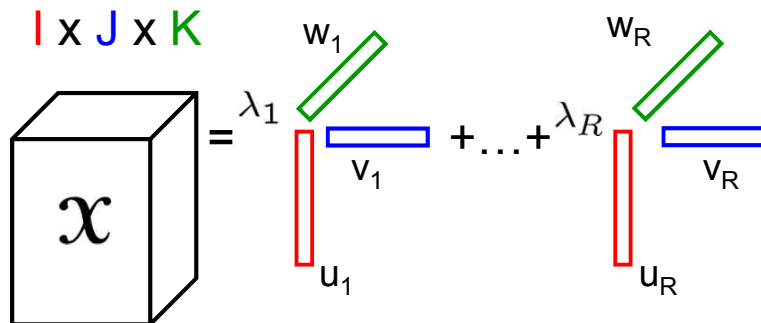
## Main points:

- 2 major types of tensor decompositions: PARAFAC and Tucker
- both can be solved with ``alternating least squares'' (ALS)
- Details follow

# Specially Structured Tensors

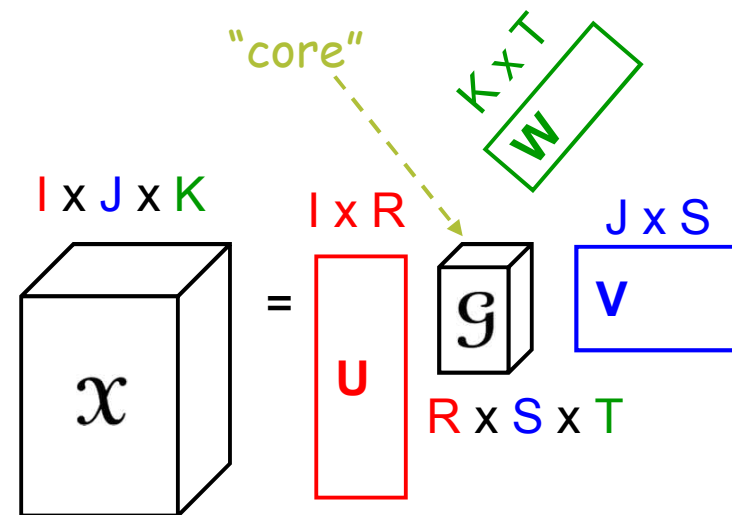
## ■ PARAFAC Tensor

$$\begin{aligned}\mathcal{X} &= \sum_r \lambda_r \mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r \\ &\equiv [\lambda; \mathbf{U}, \mathbf{V}, \mathbf{W}] \quad \left. \vphantom{[\lambda; \mathbf{U}, \mathbf{V}, \mathbf{W}]} \right\} \text{Our Notation}\end{aligned}$$



## ■ Tucker Tensor

$$\begin{aligned}\mathcal{X} &= \mathcal{G} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} \\ &= \sum_r \sum_s \sum_t g_{rst} \mathbf{u}_r \circ \mathbf{v}_s \circ \mathbf{w}_t \\ &\equiv [\mathcal{G}; \mathbf{U}, \mathbf{V}, \mathbf{W}] \quad \left. \vphantom{[\mathcal{G}; \mathbf{U}, \mathbf{V}, \mathbf{W}]} \right\} \text{Our Notation}\end{aligned}$$



# Outline

☒ Motivation - Definitions

➔ ☐ **Tensor tools**

Tensor basics

➔ Parafac

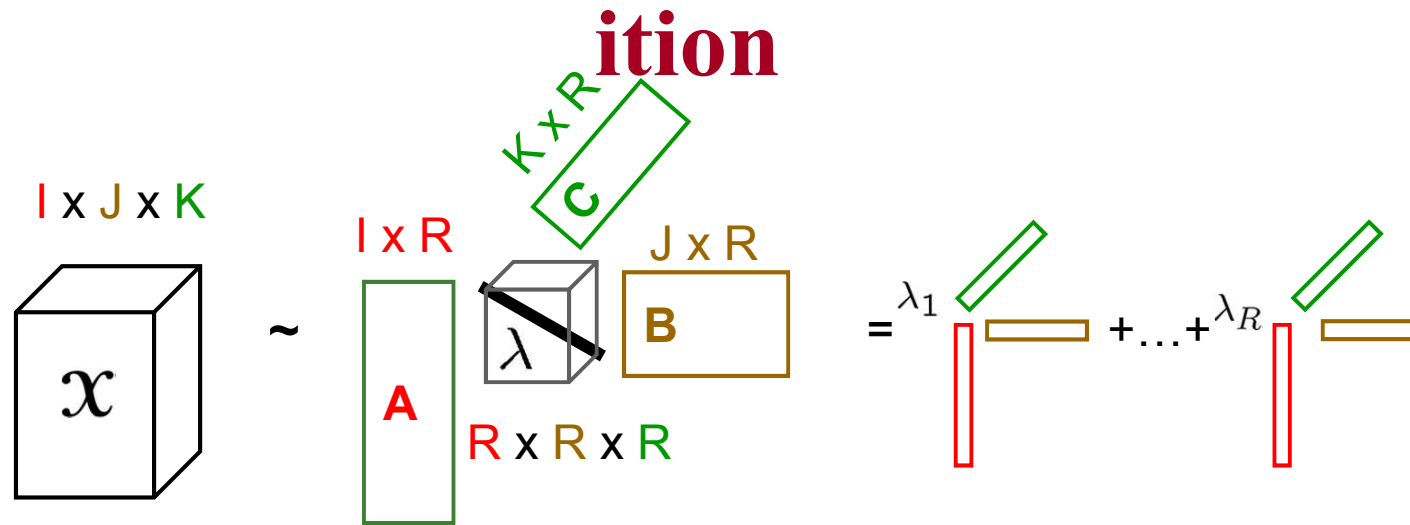
Tucker

☐ Case Studies

☐ Conclusion



# CANDECOMP/PARAFAC Decomposition



$$\mathcal{X} \approx [[\lambda; \mathbf{A}, \mathbf{B}, \mathbf{C}]] = \sum_r \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

- CANDECOMP = Canonical Decomposition (Carroll & Chang, 1970)
- PARAFAC = Parallel Factors (Harshman, 1970)
- Core is diagonal (specified by the vector  $\lambda$ )
- Columns of  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  are not orthonormal
- If  $R$  is minimum, then  $R$  is called the **rank** of the tensor (Kruskal 1977)
- Can have  $\text{rank}(\mathcal{X}) > \min\{I, J, K\}$

# Outline

☒ Motivation - Definitions

➡ ☐ **Tensor tools**

Tensor basics

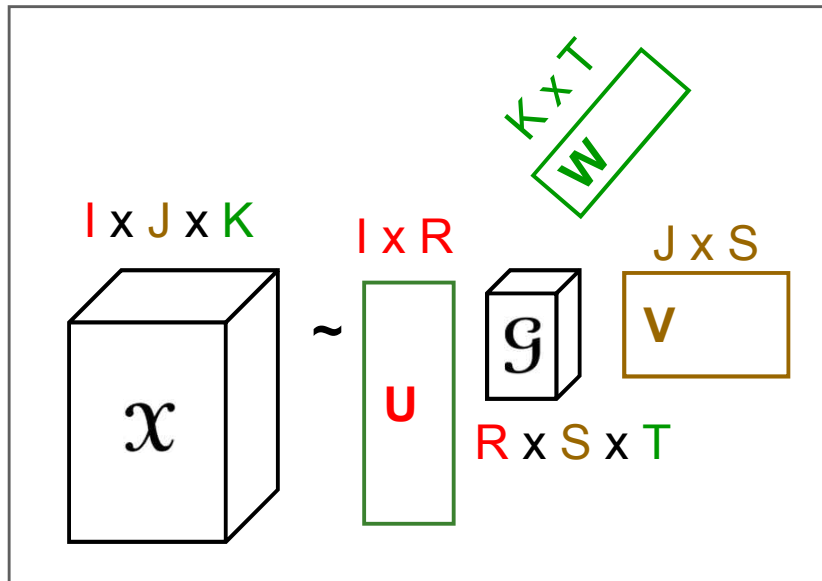
Parafac


➡ Tucker

☐ Case Studies

☐ Conclusion

# Tucker Decomposition - intuition



- author x keyword x conference
  - $\mathcal{U}$ : author x author-group
  - $\mathcal{V}$ : keyword x keyword-group
  - $\mathcal{W}$ : conf. x conf-group
  - $\mathcal{G}$ : how groups relate to each other
- Needs elaboration! 

# Intuition behind core tensor

- 2-d case: co-clustering
- [Dhillon et al. Information-Theoretic Co-clustering, KDD'03]

$$\begin{array}{c} \text{---} \\ n \\ \text{---} \end{array}$$

$$\begin{array}{c} m \\ \left[ \begin{array}{cccccc} .05 & .05 & .05 & 0 & 0 & 0 \\ .05 & .05 & .05 & 0 & 0 & 0 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ .04 & .04 & 0 & .04 & .04 & .04 \\ .04 & .04 & .04 & 0 & .04 & .04 \end{array} \right] \end{array}$$

eg, terms x documents

$$\begin{array}{c} m \\ \left[ \begin{array}{cc} k & l \\ \begin{array}{ccc} .5 & 0 & 0 \\ .5 & 0 & 0 \\ 0 & .5 & 0 \\ 0 & .5 & 0 \\ 0 & 0 & .5 \\ 0 & 0 & .5 \end{array} \end{array} \right] \begin{array}{c} k \\ \left[ \begin{array}{cc} l & n \\ \begin{array}{cc} .3 & 0 \\ 0 & .3 \\ .2 & .2 \end{array} \end{array} \right] \begin{array}{c} l \\ \left[ \begin{array}{cccccc} .36 & .36 & .28 & 0 & 0 & 0 \\ 0 & 0 & 0 & .28 & .36 & .36 \end{array} \right] \end{array} \end{array} = \begin{array}{c} \left[ \begin{array}{ccc|ccc} .054 & .054 & .042 & 0 & 0 & 0 \\ .054 & .054 & .042 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & .042 & .054 & .054 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ \hline .036 & .036 & .028 & .028 & .036 & .036 \\ .036 & .036 & .028 & .028 & .036 & .036 \end{array} \right] \end{array}$$

med. doc      cs doc

term group x  
doc. group

$$\begin{bmatrix} .5 & 0 & 0 \\ .5 & 0 & 0 \\ 0 & .5 & 0 \\ 0 & .5 & 0 \\ 0 & 0 & .5 \\ 0 & 0 & .5 \end{bmatrix}$$

$$\begin{bmatrix} .3 & 0 \\ 0 & .3 \\ .2 & .2 \end{bmatrix}$$

$$\begin{bmatrix} .36 & .36 & .28 & 0 & 0 & 0 \\ 0 & 0 & 0 & .28 & .36 & .36 \end{bmatrix} =$$

doc x  
doc group

| med. terms

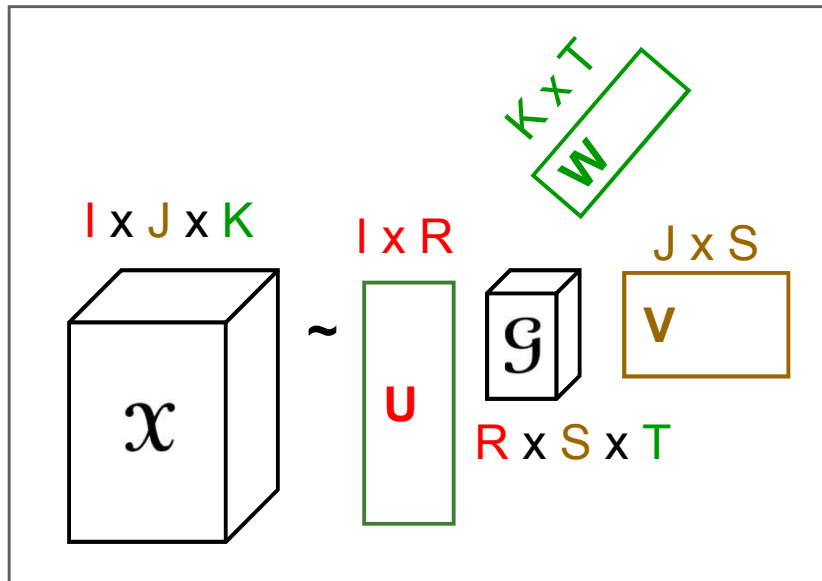
| cs terms

| common terms

$$\begin{bmatrix} .054 & .054 & .042 & 0 & 0 & 0 \\ .054 & .054 & .042 & 0 & 0 & 0 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ .036 & .036 & .028 & .028 & .036 & .036 \\ .036 & .036 & .028 & .028 & .036 & .036 \end{bmatrix}$$

term x  
term-group

# Tucker Decomposition



$$\begin{aligned}\mathcal{X} &= \mathcal{G} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} \\ &= \sum_r \sum_s \sum_t g_{rst} \mathbf{u}_r \circ \mathbf{v}_s \circ \mathbf{w}_t \\ &\equiv [\mathcal{G}; \mathbf{U}, \mathbf{V}, \mathbf{W}]\end{aligned}$$

- Proposed by Tucker (1966)
- AKA: Three-mode factor analysis, three-mode PCA, orthogonal array decomposition
- $\mathbf{U}$ ,  $\mathbf{V}$ , and  $\mathbf{W}$  generally assumed to be orthonormal (generally assume they have full column rank)
- $\mathcal{G}$  is not diagonal
- Not unique

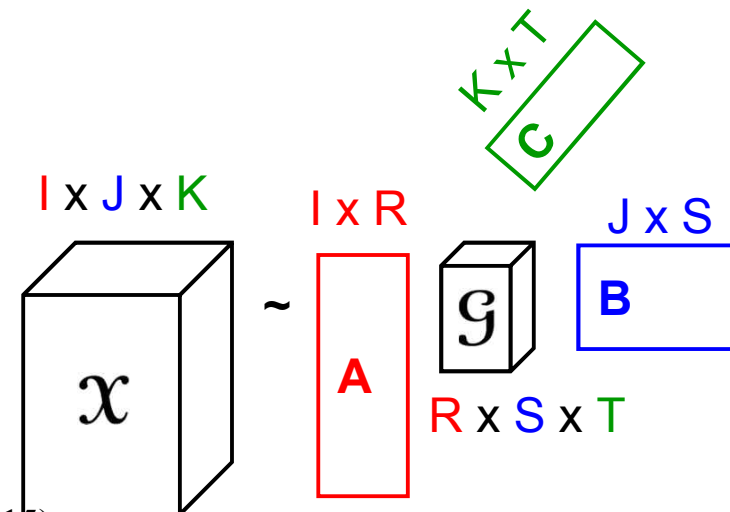
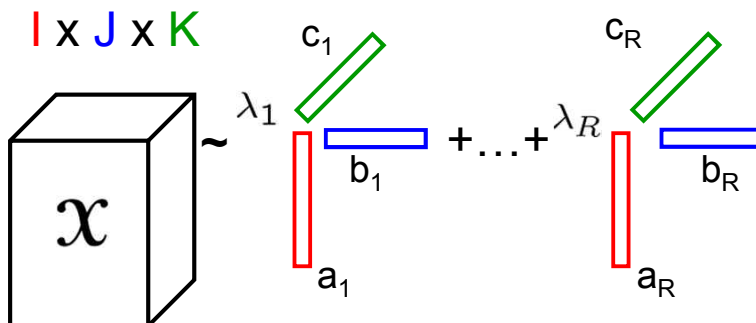
# PARAFAC vs. Tucker Decompositions

## ■ PARAFAC

- Sum of rank-1 components
- No core, i.e., superdiagonal core
- A, B, C may have linearly dependent columns
- Generally unique

## ■ Tucker

- Many interactions from groups
- Core G may be dense
- A, B, C generally orthonormal
- Not unique





# Tensor tools - summary

- Two main tools
  - PARAFAC
  - Tucker
- Both find row-, column-, tube-groups
  - but in PARAFAC the three groups are identical
- To solve: Alternating Least Squares, gradient descent, ...
  
- Toolbox: from Tamara Kolda:  
<http://csmr.ca.sandia.gov/~tgkolda/TensorToolbox/>

# Outline

☒ Motivation - Definitions

☒ Tensor tools

 ☐ **Case Studies**

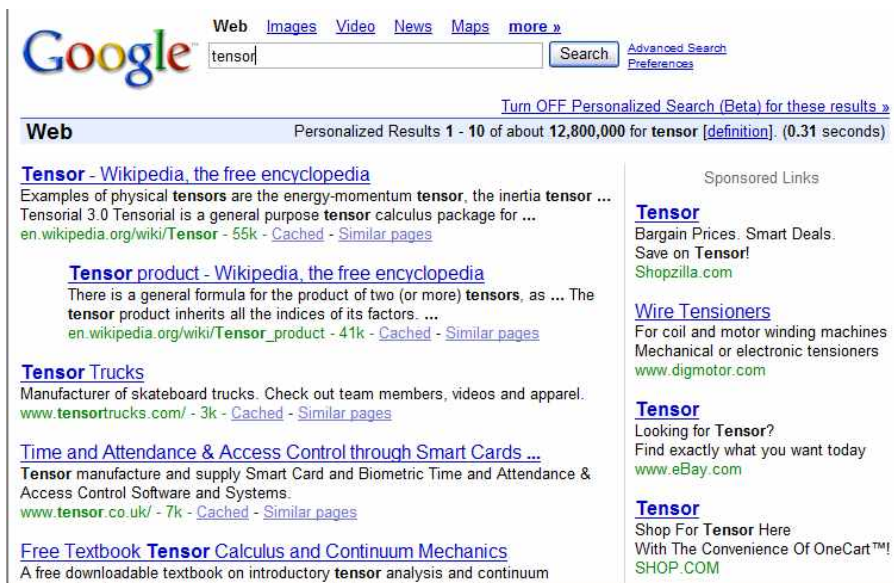
 HITS

GigaTensor

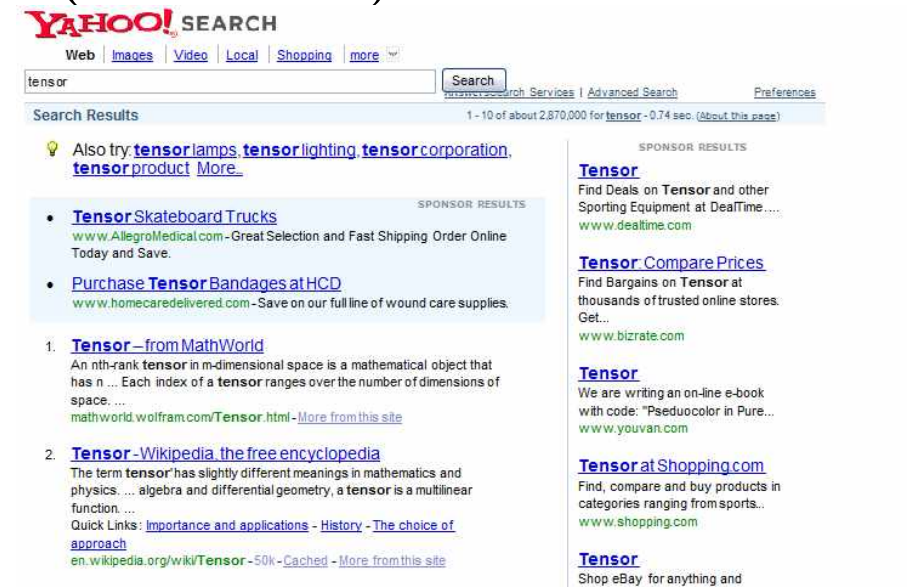
☐ Conclusion

# P1: Web graph mining

- How to order the importance of web pages?
  - Kleinberg's algorithm HITS
  - PageRank
  - Tensor extension on HITS (TOPHITS)

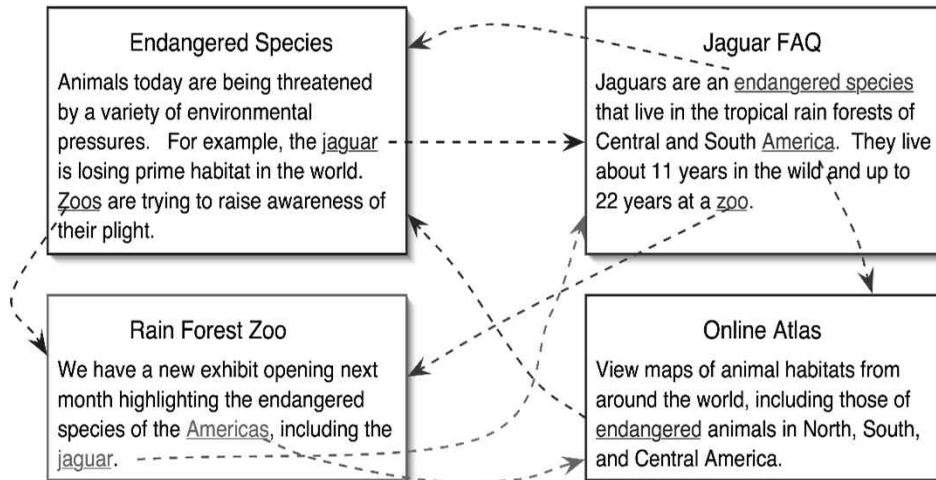


Google search results for "tensor". The search bar shows "tensor" and the search button is labeled "Search". Below the search bar, there are links for "Web", "Images", "Video", "News", "Maps", and "more »". The search results are displayed under the "Web" tab, showing personalized results 1 - 10 of about 12,800,000 for "tensor" [definition]. (0.31 seconds). The results include links to Wikipedia, Tensor product, Tensor Trucks, Time and Attendance & Access Control through Smart Cards, and Free Textbook Tensor Calculus and Continuum Mechanics.



Yahoo! search results for "tensor". The search bar shows "tensor" and the search button is labeled "Search". Below the search bar, there are links for "Web", "Images", "Video", "Local", "Shopping", and "more ». The search results are displayed under the "Search Results" tab, showing 1 - 10 of about 2,870,000 for "tensor" - 0.74 sec. (About this page). The results include sponsored links for TensorSkateboard Trucks, Purchase Tensor Bandages at HCD, and sponsored results for Tensor, Tensor Compare Prices, Tensor, and Tensor at Shopping.com.

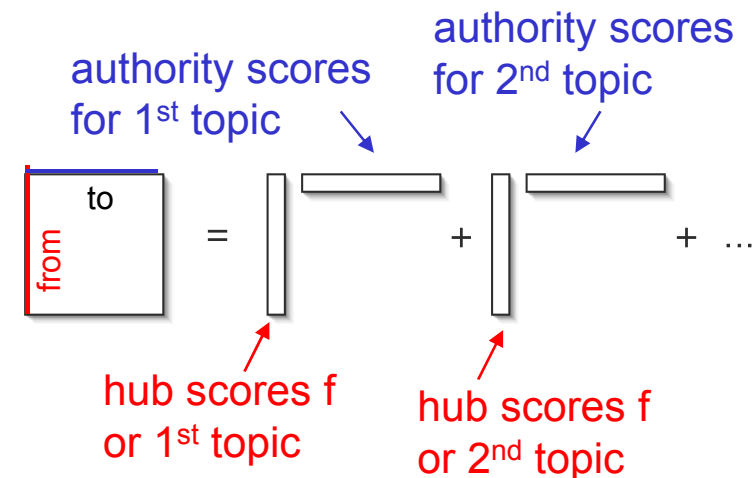
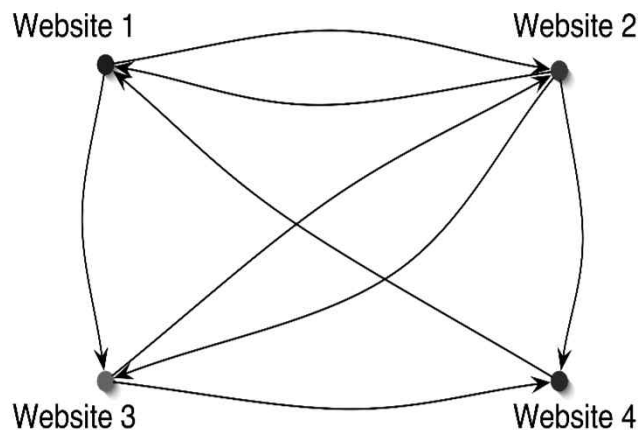
# Kleinberg's Hubs and Authorities (the HITS method)



Sparse adjacency matrix and its SVD:

$$x_{ij} = \begin{cases} 1 & \text{if page } i \text{ links to page } j \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbf{X} \approx \sum_r \sigma_r \mathbf{h}_r \circ \mathbf{a}_r$$



# HITS Authorities on Sample Data

1st Principal Factor	
.97	www.ibm.com
.24	www.alphaweb.com
.08	www-128.ibm.com
.05	www.developers.sun.com
.02	www.research.ibm.com
.01	www.redbook.ibm.com
.01	news.com.com

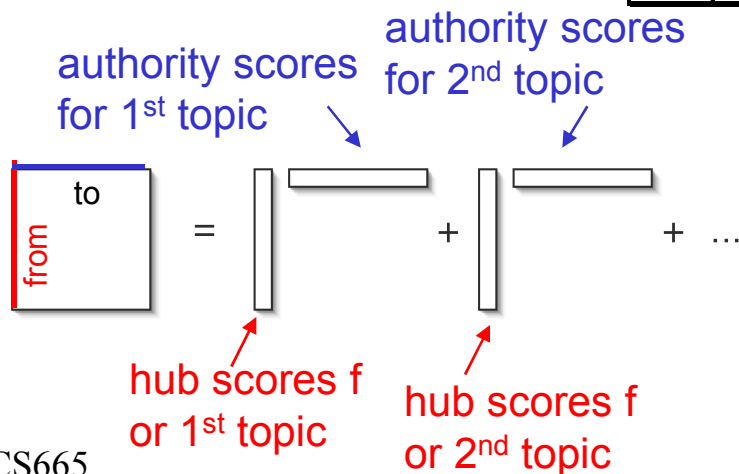
2nd Principal Factor	
.99	www.lehigh.edu
.11	www2.lehigh.edu
.06	www.lehigh.edu
.06	www.lehigh.edu
.02	www.bethlehem.edu
.02	www.adobe.com
.02	lewisweb.cc.lehigh.edu
.02	www.leo.lehigh.edu
.02	www.distance.edu
.02	fp1.cc.lehigh.edu

3rd Principal Factor	
.75	java.sun.com
.38	www.sun.com
.36	developers.sun.com
.24	see.sun.com
.16	www.samag.com
.13	docs.sun.com
.12	blogs.sun.com
.08	sunsolve.sun.com
.08	www.sun-catalog.com
.08	news.com.com

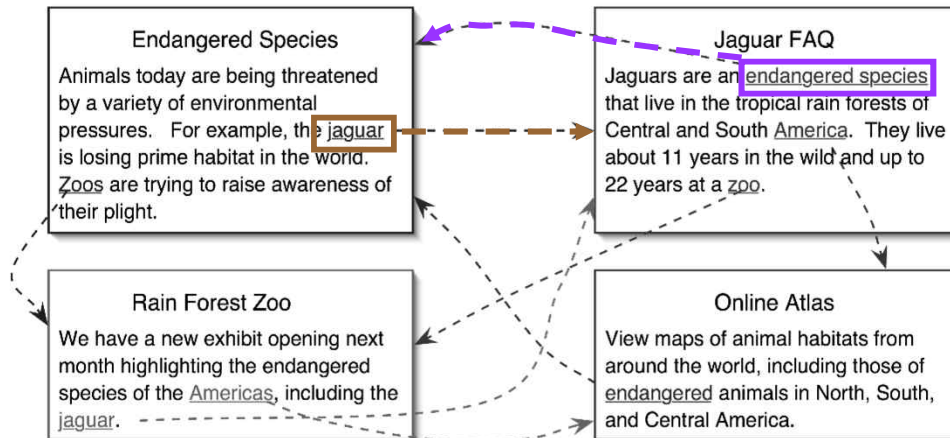
4th Principal Factor	
.60	www.pueblo.gsa.gov
.45	www.whitehouse.gov
.35	www.irs.gov
.31	travel.state.gov
.22	www.gsa.gov
.20	www.ssa.gov
.16	www.census.gov
.14	www.govbeat.com
.13	www.kids.gov
.13	www.usdoj.gov

We started our crawl from  
<http://www-neos.mcs.anl.gov/neos>,  
and crawled 4700 pages,  
resulting in 560  
cross-linked hosts.

6th Principal Factor	
.97	mathpost.asu.edu
.18	math.la.asu.edu
.17	www.asu.edu
.04	www.act.org
.03	www.eas.asu.edu
.02	archives.math.utk.edu
.02	www.geom.uiuc.edu
.02	www.fulton.asu.edu
.02	www.amstat.org
.02	www.maa.org



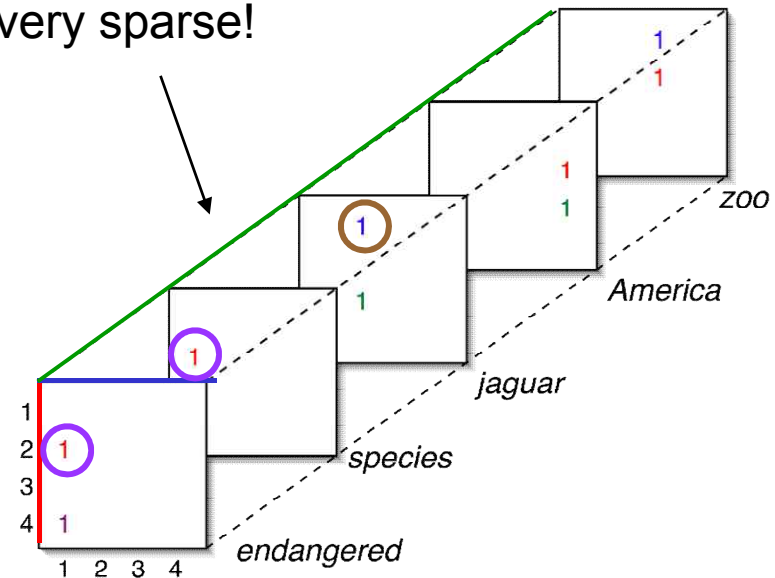
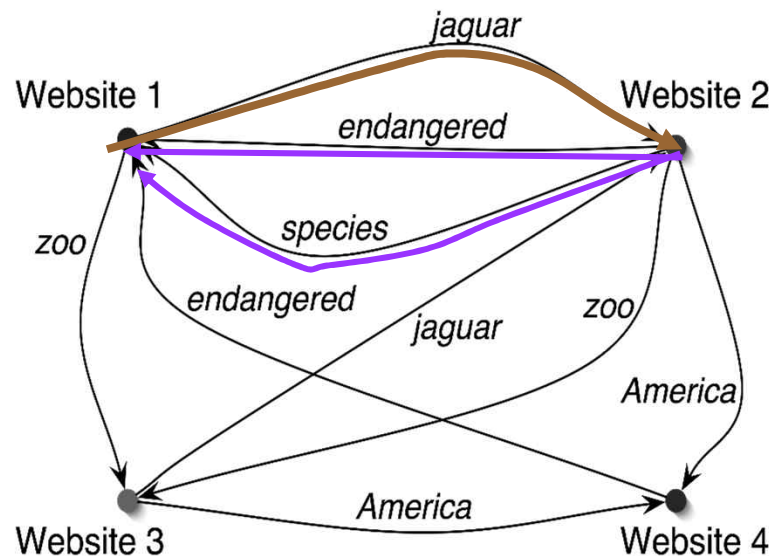
# Three-Dimensional View of the Web



[Kolda, Bader, Kenny, ICDM05]

$$x_{ijk} = \begin{cases} 1 & \text{if page } i \rightarrow \text{page } j \\ & \text{with term } k \\ 0 & \text{otherwise} \end{cases}$$

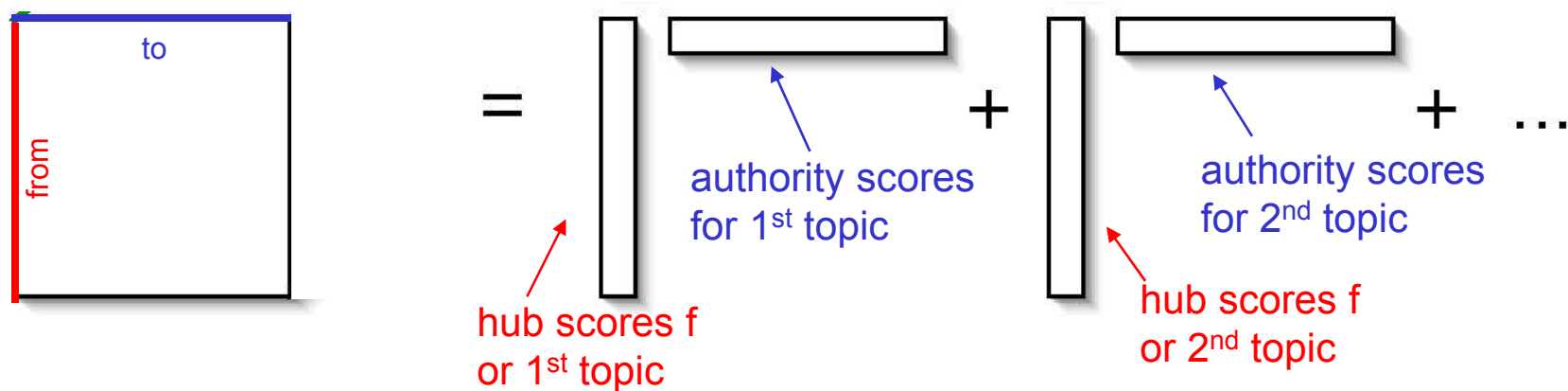
Observe that this tensor is very sparse!



# Topical HITS (TOPHITS)

**Main Idea:** Extend the idea behind the HITS model to incorporate term (i.e., topical) information.

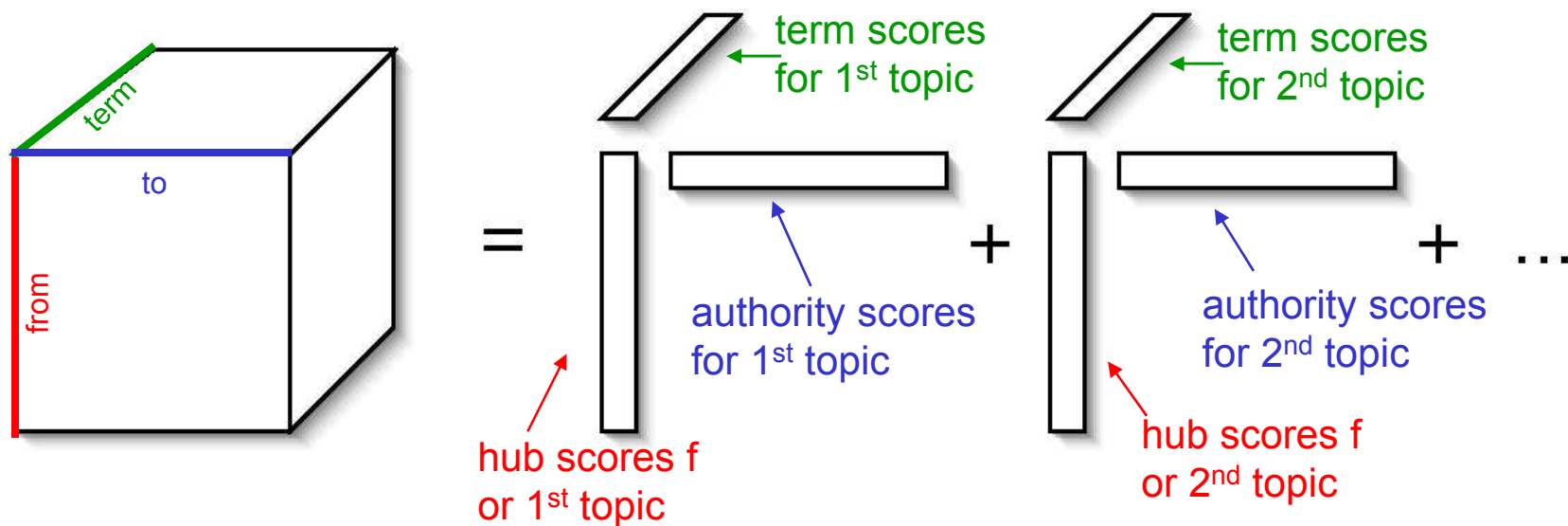
$$\mathbf{x} \approx \sum_{r=1}^R \lambda_r \mathbf{h}_r \circ \mathbf{a}_r$$



# Topical HITS (TOPHITS)

**Main Idea:** Extend the idea behind the HITS model to incorporate term (i.e., topical) information.

$$\mathbf{x} \approx \sum_{r=1}^R \lambda_r \mathbf{h}_r \circ \mathbf{a}_r \circ \mathbf{t}_r$$



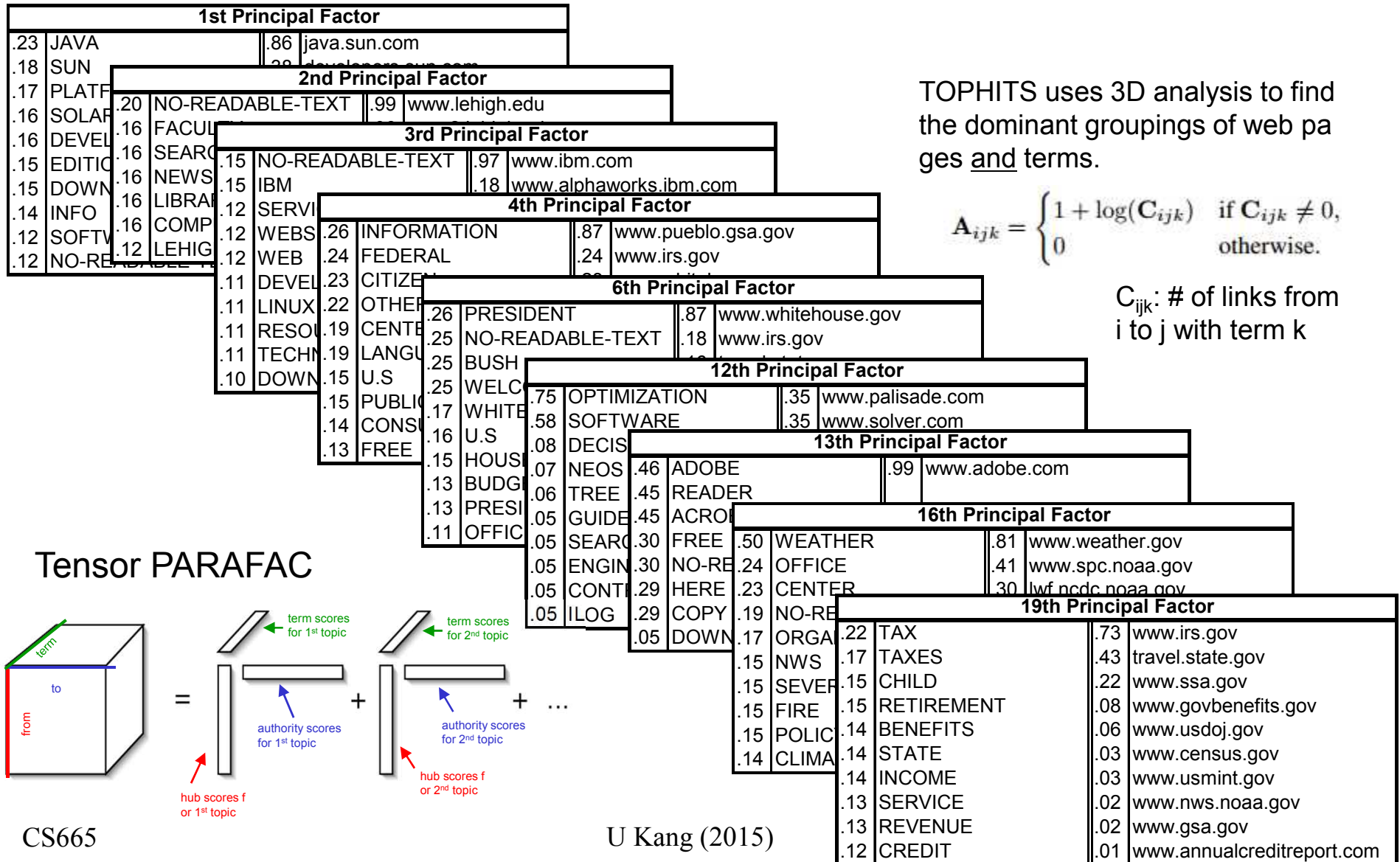


# TOPHITS Terms & Authorities

## on Sample Data

TERM

AUTHORITY



# Outline

☒ Motivation - Definitions

☒ Tensor tools

 ☐ **Case Studies**

HITS

 GigaTensor

☐ Conclusion

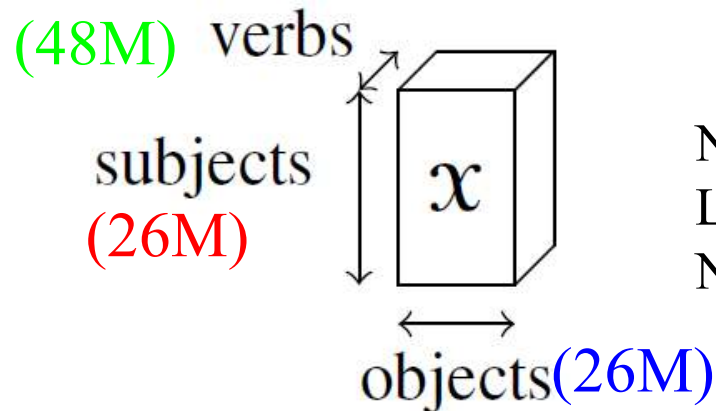
*U Kang, Evangelos Papalexakis, Abhay Harpale, and Christos Faloutsos. GigaTensor: Scaling Tensor Analysis Up By 100 Times - Algorithms and Discoveries, ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD) 2012, Beijing, China.*

## P2: N.E.L.L. analysis

- NELL: Never Ending Language Learner
  - Q1: dominant concepts / topics?
  - Q2: synonyms for a given new phrase?

“Eric Clapton plays  
guitar”

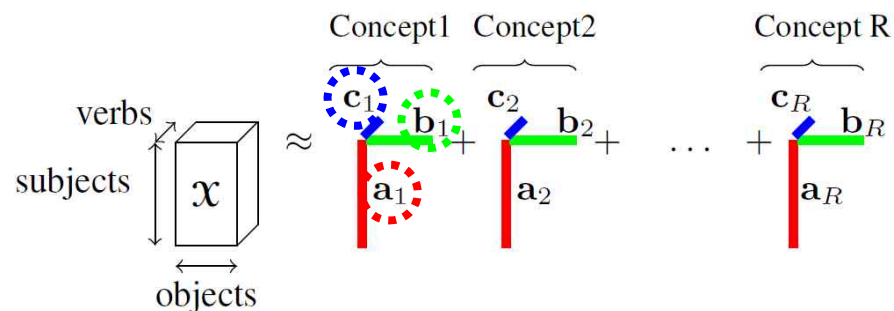
“Barrack Obama is the  
president of U.S.”



NELL (Never Ending  
Language Learner)  
Nonzeros = 144M

# A1: Concept Discovery

## ■ Concept Discovery in Knowledge Base



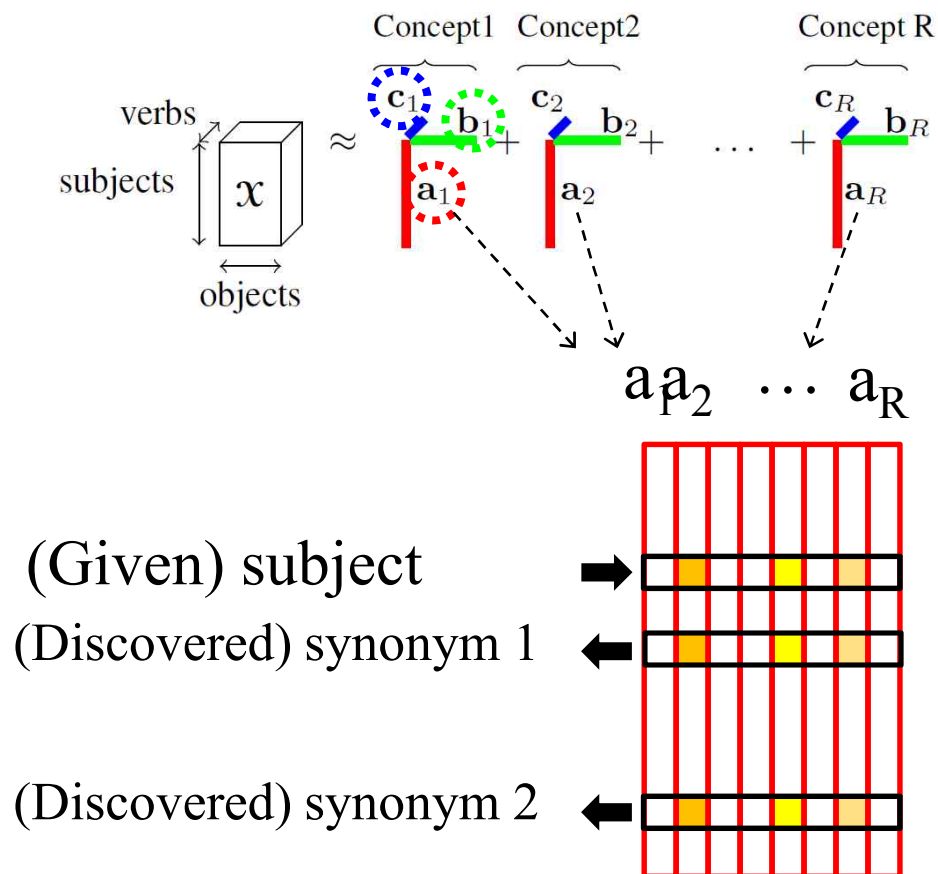
Noun Phrase 1	Noun Phrase 2	Context
<b>Concept 1: "Web Protocol"</b>		
internet	protocol	'np1' 'stream' 'np2'
file	software	'np1' 'marketing' 'np2'
data	suite	'np1' 'dating' 'np2'
<b>Concept 2: "Credit Cards"</b>		
credit	information	'np1' 'card' 'np2'
Credit	debt	'np1' 'report' 'np2'
library	number	'np1' 'cards' 'np2'
<b>Concept 3: "Health System"</b>		
health	provider	'np1' 'care' 'np2'
child	providers	'np' 'insurance' 'np2'
home	system	'np1' 'service' 'np2'
<b>Concept 4: "Family Life"</b>		
life	rest	'np2' 'of' 'my' 'np1'
family	part	'np2' 'of' 'his' 'np1'
body	years	'np2' 'of' 'her' 'np1'

# A1: Concept Discovery

Noun Phrase 1	Noun Phrase 2	Context
<b>Concept 1: "Web Protocol"</b>		
internet	protocol	'np1' 'stream' 'np2'
file	software	'np1' 'marketing' 'np2'
data	suite	'np1' 'dating' 'np2'
<b>Concept 2: "Credit Cards"</b>		
credit	information	'np1' 'card' 'np2'
Credit	debt	'np1' 'report' 'np2'
library	number	'np1' 'cards' 'np2'
<b>Concept 3: "Health System"</b>		
health	provider	'np1' 'care' 'np2'
child	providers	'np' 'insurance' 'np2'
home	system	'np1' 'service' 'np2'
<b>Concept 4: "Family Life"</b>		
life	rest	'np2' 'of' 'my' 'np1'
family	part	'np2' 'of' 'his' 'np1'
body	years	'np2' 'of' 'her' 'np1'

# A2: Synonym Discovery

## ■ Synonym Discovery in Knowledge Base



(Given) Noun Phrase	(Discovered) Potential Synonyms
pollutants	dioxin, sulfur dioxide, greenhouse gases, particulates, nitrogen oxide, air pollutants, cholesterol
disabilities	infections, dizziness, injuries, diseases, drowsiness, stiffness, injuries
vodafone	verizon, comcast
Christian history	European history, American history, Islamic history, history
disbelief	dismay, disgust, astonishment
cyberpunk	online-gaming
soul	body



# A2: Synonym Discovery

(Given) Noun Phrase	(Discovered) Potential Synonyms
pollutants	dioxin, sulfur dioxide, greenhouse gases, particulates, nitrogen oxide, air pollutants, cholesterol
disabilities	infections, dizziness, injuries, diseases, drowsiness, stiffness, injuries
vodafone	verizon, comcast
Christian history	European history, American history, Islamic history, history
disbelief	dismay, disgust, astonishment
cyberpunk	online-gaming
soul	body

# Outline

☒ Motivation - Definitions

☒ Tensor tools

☒ Case Studies

 ☐ **Conclusion**



# Conclusions

- Real data are often in high dimensions with multiple aspects (modes)
- Matrices and tensors provide elegant theory and algorithms

