**CMU SCS**

Sandia National Laboratories

# Mining Large Time-evolving Data Using Matrix and Tensor Tools

| | |
|---|---|
| *Christos Faloutsos* | Carnegie Mellon Univ. |
| *Tamara G. Kolda* | Sandia National Labs |
| *Jimeng Sun* | Carnegie Mellon Univ. |

---

**CMU SCS**

Sandia National Laboratories

## About the tutorial

- Introduce **matrix and tensor tools** through **real mining applications**

- **Goal:** find **patterns, rules, clusters, outliers, …**
  - in matrices and
  - in tensors

**CMU SCS**

Sandia National Laboratories

# What is this tutorial about?

- Matrix tools
  - Singular Value Decomposition (SVD)
  - Principal Component Analysis (PCA)
  - Webpage ranking algorithms: HITS, PageRank
  - CUR decomposition
  - Co-clustering
  - Nonnegative Matrix Factorization (NMF)
- Tensor tools
  - Tucker decomposition
  - Parallel factor analysis (PARAFAC)
  - DEDICOM
  - Missing values
  - Nonnegativity
  - Incrementalization
- Applications, Software demo

Faloutsos, Kolda, Sun                                1-3

---

**CMU SCS**

Sandia National Laboratories

# What is this tutorial NOT about?

- **Classification methods**
- **Kernel methods**
- **Discriminative models**
  - Linear Discriminant Analysis (LDA)
  - Canonical Correlation Analysis (CCA)
- **Probabilistic latent variable models**
  - Probabilistic PCA
  - Probabilistic latent semantic indexing
  - Latent Dirichlet allocation

Faloutsos, Kolda, Sun                                1-4

**CMU SCS**

Sandia
National
Laboratories

# Motivation 1: Why "matrix"?

- Why matrices are important?

---

**CMU SCS**

Sandia
National
Laboratories

# Examples of Matrices:
## Graph - social network

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

**CMU SCS**

# Examples of Matrices:
## cloud of n-d points

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

Faloutsos, Kolda, Sun

1-7

**CMU SCS**

# 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   | ...  | ...   | ...   | ...  | ... |
| ...    | ...  | ...   | ...   | ...  | ... |

Faloutsos, Kolda, Sun

1-8

**CMU SCS**

# 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 | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... |

Faloutsos, Kolda, Sun

1-9

**CMU SCS**

# Examples of Matrices:
## Authors and terms

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

Faloutsos, Kolda, Sun

1-10

**CMU SCS**

Sandia
National
Laboratories

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

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

Faloutsos, Kolda, Sun

1-11

---

**CMU SCS**

Sandia
National
Laboratories

# Motivation 2: Why tensor?

- Q: what is a tensor?

Faloutsos, Kolda, Sun

1-12

**CMU SCS**

Sandia National Laboratories

# Motivation 2: Why tensor?

- A: N-D generalization of matrix:

ICML'07

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

Faloutsos, Kolda, Sun                                      1-13

---

**CMU SCS**

Sandia National Laboratories

# Motivation 2: Why tensor?

- A: N-D generalization of matrix:

ICML'05
ICML'06
ICML'07

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

Faloutsos, Kolda, Sun                                      1-14

**CMU SCS**

## Tensors are useful for 3 or more modes

Terminology: 'mode' (or 'aspect'):

3rd Mode

| data | mining | classif. | tree | ... |
|------|--------|----------|------|-----|
| 13 | 11 | 22 | 55 | ... |
| 5 | 4 | 6 | 7 | ... |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |

2nd Mode

1st Mode

Faloutsos, Kolda, Sun                                1-15

**CMU SCS**

## Motivating Applications

- Why matrices are important?
- Why tensors are useful?
  - P1: environmental sensors
  - P2: data center monitoring ('autonomic')
  - P3: social networks
  - P4: network forensics
  - P5: web mining
  - P6: face recognition

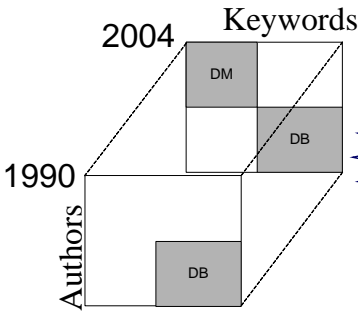Faloutsos, Kolda, Sun                                1-16

8

**CMU SCS**

Sandia National Laboratories

# P1: Environmental sensor monitoring

Data in three **modes**
(time, location, type)

Temperature

Light

Humidity

Voltage

Faloutsos, Kolda, Sun

5-17

---

**CMU SCS**

Sandia National Laboratories

# P2: Clusters/data center monitoring

Data in three **modes**
(time, machine, type)

Bytes Received — Unicast Packets Received — Bytes Sent — Unicast Packets Sent
Unprivileged CPU Utilization — Other CPU Utilization — Privileged CPU Utilization
CPU Idle Time

Hidden 1 — Hidden 2

- Monitor correlations of multiple measurements
- Automatically flag anomalous behavior
- Intemon: intelligent monitoring system
  - Prof. Greg Ganger and PDL
  - >100 machines in a data center
  - `warsteiner.db.cs.cmu.edu/demo/intemon.jsp`

Faloutsos, Kolda, Sun

1-18

**P3: Social network analysis**

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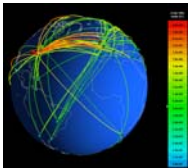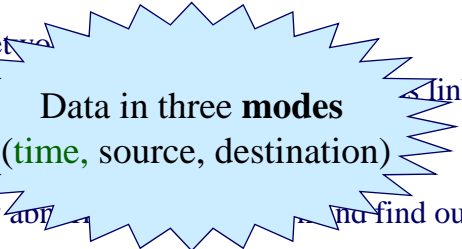- We plan to monitor the change of the community structure over time

Data in three **modes** (time, author, keyword)
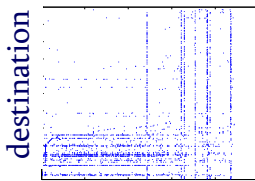


**P4: Network forensics**

- Directional net...
- A large ISP ... link capacity...
  - 450 GB/...
- Task: Identify ab... ...nd find out the cause
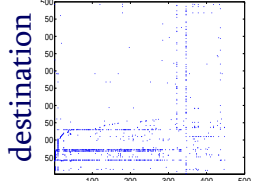
Data in three **modes** (time, source, destination)

abnormal traffic    normal traffic

**CMU SCS**

# P5: 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
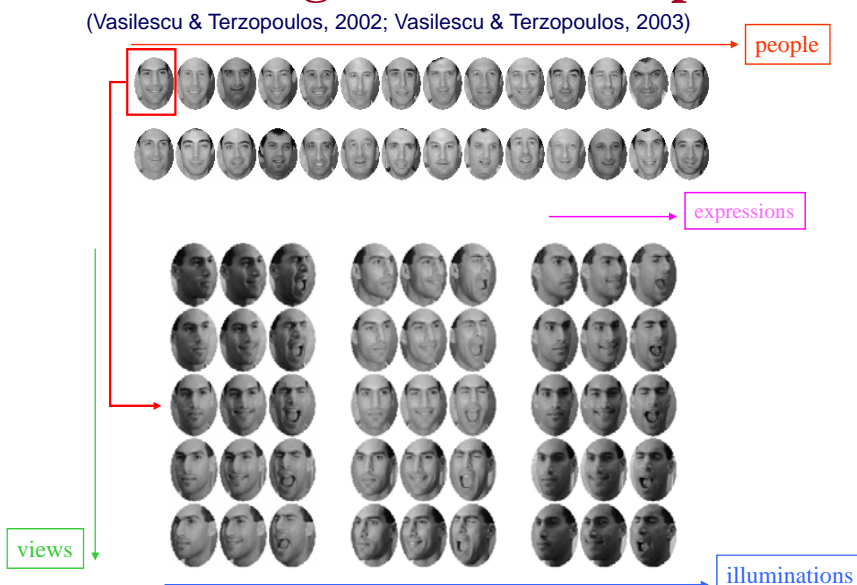
Data in three **modes**
(source, destination, text)

Faloutsos, Kolda, Sun

1-21

**CMU SCS**

# P6. Face recognition and compression

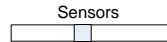(Vasilescu & Terzopoulos, 2002; Vasilescu & Terzopoulos, 2003)

people

expressions

views

illuminations

Faloutsos, Kolda, Sun

1-22

**CMU SCS**

Sandia National Laboratories
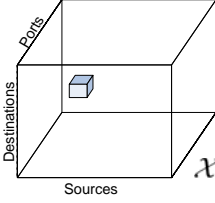
# Static Data model

- Tensor
  - Formally, $\mathcal{X} \in \mathbf{R}^{N_1 \times \ldots \times N_M}$
  - Generalization of matrices
  - Represented as multi-array, (~ data cube).

| Order | 1st | 2nd | 3rd |
|---|---|---|---|
| Correspondence | Vector | Matrix | 3D array |
| Example |  |  |  |

Faloutsos, Kolda, Sun

1-23

---

**CMU SCS**

Sandia National Laboratories

# Dynamic Data model

- Tensor Streams
  - A sequence of Mth order tensor

$$\mathcal{X}_1 \ldots \mathcal{X}_t \text{ where } \mathcal{X}_i \in \mathbf{R}^{N_1 \times \ldots \times N_M}$$

$t$ is increasing over time

| Order | 1st | 2nd | 3rd |
|---|---|---|---|
| Correspondence | Multiple streams | Time evolving graphs | 3D arrays |
| Example |  |  |  |

Faloutsos, Kolda, Sun

1-24

**CMU SCS**

Sandia
National
Laboratories

# Roadmap

- Motivation
- Matrix tools
- Tensor basics
- Tensor extensions
- Software demo
- Case studies



Faloutsos, Kolda, Sun                                            1-25

**CMU SCS**    Sandia National Laboratories

# Roadmap

- Motivation
- Matrix tools
- Tensor basics
- Tensor extensions
- Software demo
- Case studies

- SVD, PCA
- HITS, PageRank
- CUR
- Co-clustering
- Nonnegative Matrix factorization



Faloutsos, Kolda, Sun    2-1

---

**CMU SCS**    Sandia National Laboratories

# Singular Value Decomposition (SVD)

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\mathsf{T}$$



singular values    right singular vectors

input data    left singular vectors

Faloutsos, Kolda, Sun    2-4

**CMU SCS**

# SVD as spectral decomposition

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



– Best rank-k approximation in L2 and Frobenius
– SVD only works for static matrices (a single 2nd order tensor)

Faloutsos, Kolda, Sun                                    2-5

See also PARAFAC

---

**CMU SCS**

# SVD example

1st factor    2nd factor

$$\bullet \ \mathbf{A} = \mathbf{U}\ \Sigma\ \mathbf{V}^T = \sigma_1 \mathbf{u}_1 \circ \mathbf{v}_1 + \sigma_2 \mathbf{u}_2 \circ \mathbf{v}_2 + \ldots$$



6

6

**CMU SCS**

Sandia National Laboratories

# SVD properties

- **V** are the eigenvectors of the *covariance matrix* $\mathbf{X^T X}$, since

$$\mathbf{X^T X} = \left(\mathbf{U\Sigma V^T}\right)^T \left(\mathbf{U\Sigma V^T}\right) = \mathbf{V\Sigma^2 V^T}$$

- **U** are the eigenvectors of the *Gram (inner-product) matrix* $\mathbf{XX^T}$, since

$$\mathbf{XX^T} = \left(\mathbf{U\Sigma V^T}\right)\left(\mathbf{U\Sigma V^T}\right)^T = \mathbf{U\Sigma^2 U^T}$$

Further reading:
1. Ian T. Jolliffe, *Principal Component Analysis* (2nd ed), Springer, 2002.
2. Gilbert Strang, *Linear Algebra and Its Applications* (4th ed), Brooks Cole, 2005.

---

**CMU SCS**

Sandia National Laboratories

# SVD - Interpretation

'documents', 'terms' and 'concepts':

**Q:** if **A** is the document-to-term matrix, what is $\mathbf{A^T A}$?

A: term-to-term ([m x m]) similarity matrix

Q: $\mathbf{A A^T}$ ?

A: document-to-document ([n x n]) similarity matrix

Faloutsos, Kolda, Sun             2-8

**CMU SCS**

Sandia National Laboratories

# Principal Component Analysis (PCA)

- SVD $\quad \mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$



– PCA is an important application of SVD
– Note that U and V are dense and may have negative entries

Faloutsos, Kolda, Sun

2-9

---

**CMU SCS**

Sandia National Laboratories

# PCA interpretation

- best axis to project on: ('best' = min sum of squares of projection errors)



Term2 ('lung')

Term1 ('data')

2-10

**CMU SCS**

# PCA - interpretation

Term2 ('lung')

PCA projects points
Onto the "best" axis

v1

first singular vector

- minimum RMS error

Term1 ('data')

Faloutsos, Kolda, Sun                    2-11

---

**CMU SCS**

# Roadmap

- Motivation
- **Matrix tools**
- Tensor basics
- Tensor extensions
- Software demo
- Case studies

- SVD, PCA
- **HITS, PageRank**
- CUR
- Co-clustering
- Nonnegative Matrix factorization

Faloutsos, Kolda, Sun                    2-12

**CMU SCS**                                                                    Sandia National Laboratories

# Kleinberg's algorithm HITS

- Problem dfn: given the web and a query
- find the most 'authoritative' web pages for this query

Step 0: find all pages containing the query terms

Step 1: expand by one move forward and backward

Further reading:
1. J. Kleinberg. Authoritative sources in a hyperlinked environment. SODA 1998

---

**CMU SCS**                                                                    Sandia National Laboratories

# Kleinberg's algorithm HITS

- Step 1: expand by one move forward and backward



Faloutsos, Kolda, Sun                                          2-14

**CMU SCS**

# Kleinberg's algorithm HITS

- on the resulting graph, give high score (= 'authorities') to nodes that many important nodes point to
- give high importance score ('hubs') to nodes that point to good 'authorities'



hubs                    authorities

Faloutsos, Kolda, Sun                                          2-15

---

**CMU SCS**

# Kleinberg's algorithm HITS

observations

- recursive definition!
- each node (say, '$i$'-th node) has both an authoritativeness score $a_i$ and a hubness score $h_i$

Faloutsos, Kolda, Sun                                          2-16

**CMU SCS**

Sandia National Laboratories

# Kleinberg's algorithm: HITS

Let **A** be the adjacency matrix:

the ($i,j$) entry is 1 if the edge from $i$ to $j$ exists

Let **h** and **a** be [n x 1] vectors with the 'hubness' and 'authoritativiness' scores.

Then:

Faloutsos, Kolda, Sun                                    2-17

---

**CMU SCS**

Sandia National Laboratories

# Kleinberg's algorithm: HITS

Then:

$$a_i = h_k + h_l + h_m$$

k

l

m

i

that is

$a_i = $ Sum ($h_j$)    over all $j$ that
($j,i$) edge exists

or

**a** = **A**$^{\text{T}}$ **h**

Faloutsos, Kolda, Sun                                    2-18

**CMU SCS**

# Kleinberg's algorithm: HITS

symmetrically, for the 'hubness':

i

n

p

q

$$h_i = a_n + a_p + a_q$$

that is

$$h_i = \text{Sum } (q_j) \quad \text{over all } j \text{ that}$$
$$(i,j) \text{ edge exists}$$

or

$$\mathbf{h} = \mathbf{A} \, \mathbf{a}$$

Faloutsos, Kolda, Sun

2-19

---

**CMU SCS**

# Kleinberg's algorithm: HITS

In conclusion, we want vectors **h** and **a** such that:

$$\mathbf{h} = \mathbf{A} \, \mathbf{a}$$
$$\mathbf{a} = \mathbf{A}^{\mathrm{T}} \, \mathbf{h}$$

That is:

$$\mathbf{a} = \mathbf{A}^{\mathrm{T}}\mathbf{A} \, \mathbf{a}$$

Faloutsos, Kolda, Sun

2-20

**CMU SCS**                                                          Sandia National Laboratories

# Kleinberg's algorithm: HITS

**a** is a <u>right singular vector</u> of the adjacency matrix **A** (by dfn!), a.k.a the <u>eigenvector</u> of **A<sup>T</sup>A**

Starting from random **a'** and iterating, we'll eventually converge

Q: to which of all the eigenvectors? why?

A: to the one of the strongest eigenvalue,

$$(\mathbf{A}^T \mathbf{A})^k \, \mathbf{a} = \lambda_1^{\,k} \mathbf{a}$$

Faloutsos, Kolda, Sun                              2-21

---

**CMU SCS**                                                          Sandia National Laboratories

# Kleinberg's algorithm - discussion

- 'authority' score can be used to find 'similar pages' (how?)
- closely related to 'citation analysis', social networks / 'small world' phenomena

See also **TOPHITS**            Faloutsos, Kolda, Sun                2-22

**CMU SCS**

Sandia National Laboratories

# Roadmap

- Motivation
- **Matrix tools**
- Tensor basics
- Tensor extensions
- Software demo
- Case studies

- SVD, PCA
- HITS, **PageRank**
- CUR
- Co-clustering
- Nonnegative Matrix factorization

Faloutsos, Kolda, Sun

2-23

---

**CMU SCS**

Sandia National Laboratories

# Motivating problem: PageRank

Given a directed graph, find its most interesting/central node

A node is important, if it is connected with important nodes (recursive, but OK!)

Faloutsos, Kolda, Sun

2-24

# Motivating problem – PageRank solution

Given a directed graph, find its most interesting/central node

Proposed solution: Random walk; spot most 'popular' node (-> steady state prob. (ssp))

A node has high ssp, if it is connected with high ssp nodes (recursive, but OK!)

Faloutsos, Kolda, Sun

2-25

# (Simplified) PageRank algorithm

- Let **A** be the transition matrix (= adjacency matrix); let **A** become row-normalized - then



Faloutsos, Kolda, Sun

2-26

**CMU SCS**

# (Simplified) PageRank algorithm

- **A p = p**

$$A \qquad p \quad = \quad p$$



Faloutsos, Kolda, Sun

2-27

---

**CMU SCS**

# (Simplified) PageRank algorithm

- **A p = 1 * p**
- thus, **p** is the **eigenvector** that corresponds to the highest eigenvalue (=1, since the matrix is row-normalized)
- Why does it exist such a **p**?
  - **p** exists if A is nxn, nonnegative, irreducible [Perron–Frobenius theorem]

Faloutsos, Kolda, Sun

2-28

**CMU SCS**

Sandia National Laboratories

# (Simplified) PageRank algorithm

- In short: imagine a particle randomly moving along the edges
- compute its steady-state probabilities (ssp)

Full version of algo:  with occasional random jumps

Why? To make the matrix irreducible

Faloutsos, Kolda, Sun

2-29

---

**CMU SCS**

Sandia National Laboratories

# Full Algorithm

- With probability *1-c*, fly-out to a random node
- Then, we have

  $\mathbf{p} = c \, \mathbf{A} \, \mathbf{p} + (1\text{-}c)/n \, \mathbf{1} =>$

  $\mathbf{p} = (1\text{-}c)/n \, [\mathbf{I} - c \, \mathbf{A}]^{-1} \, \mathbf{1}$

Faloutsos, Kolda, Sun

2-30

**CMU SCS**

# Roadmap

- Motivation
- **Matrix tools**
- Tensor basics
- Tensor extensions
- Software demo
- Case studies

- SVD, PCA
- HITS, PageRank
- **CUR**
- Co-clustering
- Nonnegative Matrix factorization

Faloutsos, Kolda, Sun

2-31

---

**CMU SCS**

# Motivation of CUR or CMD

- SVD, PCA all transform data into some abstract space (specified by a set basis)
  - Interpretability problem
  - Loss of sparsity

Faloutsos, Kolda, Sun

2-32

**CMU SCS**

# Interpretability problem

- Each column of projection matrix $U_i$ is a linear combination of all dimensions along certain mode $U_i(:,1) = [0.5; -0.5; 0.5; 0.5]$
- All the data are projected onto the span of $U_i$
- It is hard to interpret the projections

Faloutsos, Kolda, Sun

2-33

**CMU SCS**

# The sparsity problem – pictorially:



SVD/PCA:
Destroys sparsity

$U \quad \Sigma \quad V^T$

CUR: maintains sparsity

C   U   R

2-34

---

**CMU SCS**

Sandia National Laboratories

# CUR

- Example-based projection: use actual rows and columns to specify the subspace
- Given a matrix $A \in R^{m \times n}$, find three matrices $C \in R^{m \times c}$, $U \in R^{c \times r}$, $R \in R^{r \times n}$, such that $\|A-CUR\|$ is small



U is the pseudo-inverse of X

Example-based
Orthogonal projection

Faloutsos, Kolda, Sun

2-35

---

**CMU SCS**

Sandia National Laboratories

# CUR (cont.)

- Key question:
    - How to select/sample the columns and rows?
- Uniform sampling [Williams & Seeger NIPS '00]
- Biased sampling
    - CUR w/ absolute error bound
    - CUR w/ relative error bound

Reference:
1. Tutorial: Randomized Algorithms for Matrices and Massive Datasets, SDM'06
2. Drineas et al. Subspace Sampling and Relative-error Matrix Approximation: Column-Row-Based Methods, ESA2006
3. Drineas et al., Fast Monte Carlo Algorithms for Matrices III: Computing a Compressed Approximate Matrix Decomposition, SIAM Journal on Computing, 2006.

**CMU SCS**

Sandia
National
Laboratories

# The sparsity property

SVD:  $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$

sparse and small

Big but sparse

Big and dense

CUR:  $\mathbf{A} = \mathbf{C} \, \mathbf{U} \, \mathbf{R}$

dense but small

Big but sparse

Big but sparse

2-37

---

**CMU SCS**

Sandia
National
Laboratories

# The sparsity property (cont.)



Network

DBLP

- CMD uses much smaller space to achieve the same accuracy
- CUR limitation: duplicate columns and rows
- SVD limitation: orthogonal projection densifies the data

Reference:
Sun et al. Less is More: Compact Matrix Decomposition for Large Sparse Graphs, SDM'07

# Roadmap

- Motivation
- **Matrix tools**
- Tensor basics
- Tensor extensions
- Software demo
- Case studies

- SVD, PCA
- HITS, PageRank
- CUR
- **Co-clustering etc**
- Nonnegative Matrix factorization

Faloutsos, Kolda, Sun                                                    2-39

---

# Co-clustering

- Given data matrix and the number of row and column groups $k$ and $l$
- Simultaneously
    - Cluster rows of $p(X, Y)$ into $k$ disjoint groups
    - Cluster columns of $p(X, Y)$ into $l$ disjoint groups

Faloutsos, Kolda, Sun                                                    2-40

# Co-clustering

- Let $X$ and $Y$ be discrete random variables
  - $X$ and $Y$ take values in $\{1, 2, ..., m\}$ and $\{1, 2, ..., n\}$
  - $p(X, Y)$ denotes the joint probability distribution—if not known, it is often estimated based on <u>co-occurrence</u> data
  - Application areas: <u>text mining</u>, market-basket analysis, analysis of browsing behavior, etc.
- Key Obstacles in Clustering Contingency Tables
  - High Dimensionality, Sparsity, Noise
  - Need for robust and scalable algorithms

<u>Reference:</u>
1. Dhillon et al. Information-Theoretic Co-clustering, KDD'03

---

$$p(x, y) = \quad m \begin{bmatrix} .05 & .05 & .05 & 0 & 0 & 0 \\ .05 & .05 & .05 & 0 & 0 & 0 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ .04 & .04 & 0 & .04 & .04 & .04 \\ .04 & .04 & .04 & 0 & .04 & .04 \end{bmatrix} \overset{n}{}$$

$$m \begin{bmatrix} .5 & 0 & 0 \\ .5 & 0 & 0 \\ 0 & .5 & 0 \\ 0 & .5 & 0 \\ 0 & 0 & .5 \\ 0 & 0 & .5 \end{bmatrix} k \begin{bmatrix} .3 & 0 \\ 0 & .3 \\ .2 & .2 \end{bmatrix} l \begin{bmatrix} .36 & .36 & .28 & 0 & 0 & 0 \\ 0 & 0 & 0 & .28 & .36 & .36 \end{bmatrix} = \begin{bmatrix} .054 & .054 & .042 & 0 & 0 & 0 \\ .054 & .054 & .042 & 0 & 0 & 0 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ 0 & 0 & 0 & .042 & .054 & .054 \\ .036 & .036 & 028 & .028 & .036 & .036 \\ .036 & .036 & .028 & .028 & .036 & .036 \end{bmatrix}$$

$$p(x \mid \hat{x}) \qquad p(\hat{x}, \hat{y}) \qquad p(y \mid \hat{y}) \qquad\qquad\qquad q(x, y)$$

#parameters that determine $q(x,y)$ are: $(m-k)+(kl-1)+(n-l)$

# Problem with Information Theoretic Co-clustering

- Number of row and column groups must be specified

Desiderata:

✓ Simultaneously discover row and column groups

✘ Fully Automatic: No "magic numbers"

✓ Scalable to large graphs

Faloutsos, Kolda, Sun

2-46

---

CMU SCS

# Cross-association



Desiderata:

✓ Simultaneously discover row and column groups

✓ Fully Automatic: No "magic numbers"

✓ Scalable to large matrices

Reference:
1. Chakrabarti et al. Fully Automatic Cross-Associations, KDD'04

**CMU SCS**

Sandia National Laboratories

# What makes a cross-association "good"?



Problem definition: given an encoding scheme
• decide on the # of col. and row groups *k* and *l*
• and reorder rows and columns,
• to achieve best compression

Faloutsos, Kolda, Sun

2-50

---

**CMU SCS**

details

Sandia National Laboratories

# Main Idea

| Good Compression | → | Better Clustering |
|---|---|---|

Total Encoding Cost = $\sum_i size_i * H(x_i)$ + Cost of describing cross-associations

Code Cost

Description Cost

Minimize the total cost (# bits)

for lossless compression

Faloutsos, Kolda, Sun

2-51

23

## Algorithm



I = 5 col groups

k = 5 row groups

| k=1, l=2 | k=2, l=2 | k=2, l=3 | k=3, l=3 | k=3, l=4 | k=4, l=4 | k=4, l=5 |

Faloutsos, Kolda, Sun

2-52

## Roadmap

- Motivation
- **Matrix tools**
- Tensor basics
- Tensor extensions
- Software demo
- Case studies

- SVD, PCA
- HITS, PageRank
- CUR
- Co-clustering, etc
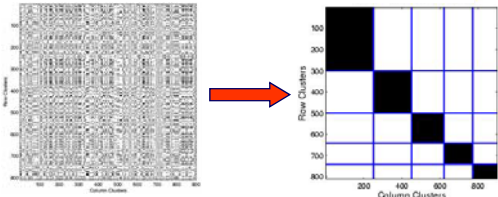- **Nonnegative Matrix factorization**



Faloutsos, Kolda, Sun

2-55

# Nonnegative Matrix Factorization

- Coming up soon with **nonnegative tensor factorization**

Faloutsos, Kolda, Sun                                              2-56

**CMU SCS**

# Roadmap

- Motivation
- Matrix tools
- Tensor basics
- Tensor extensions
- Software demo
- Case studies

- Tensor Basics
- Tucker
  - Tucker 1
  - Tucker 2
  - Tucker 3
- PARAFAC

3-1

---

**CMU SCS**

# Tensor Basics

**CMU SCS**

# A tensor is a multidimensional array

An $I \times J \times K$ tensor

$x_{1,1,1}$

$x_{ijk}$

3$^{rd}$ order tensor
mode 1 has dimension $I$
mode 2 has dimension $J$
mode 3 has dimension $K$

Note: Tutorial focus is on 3$^{rd}$ order, but everything can be extended to higher orders.

Column (Mode-1) Fibers

Row (Mode-2) Fibers     $x(1,:,5)$

Tube (Mode-3) Fibers

$x(:,5,1)$

Horizontal Slices          Lateral Slices          Frontal Slices

$X(:,:,1)$

**CMU SCS**

# Matricization : Converting a Tensor to a Matrix

Matricize (unfolding)     $(i,j,k)$ → $(i',j')$

Reverse Matricize     $(i',j')$ → $(i,j,k)$

$X_{(n)}$: The mode-**n** fibers are rearranged to be the columns of a matrix

$\mathcal{X}$          $X_{(3)}$

$$X_{(1)} = \begin{bmatrix} 1 & 3 & 5 & 7 \\ 2 & 4 & 6 & 8 \end{bmatrix}$$

$$\mathcal{X} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \begin{matrix} 5 & 7 \\ 6 & 8 \end{matrix}$$

$$X_{(2)} = \begin{bmatrix} 1 & 2 & 5 & 6 \\ 3 & 4 & 7 & 8 \end{bmatrix}$$

$$X_{(3)} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}$$

Vectorization

$$\mathbf{vec}(\mathcal{X}) = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \end{bmatrix}$$

3-4

2

CMU SCS

Sandia National Laboratories

# Tensor Mode-n Multiplication

$$\mathcal{X} \in \mathbb{R}^{I \times J \times K}, \ \mathbf{B} \in \mathbb{R}^{M \times J}, \ \mathbf{a} \in \mathbb{R}^{I}$$

- Tensor Times Matrix

$$\mathcal{Y} = \mathcal{X} \times_2 \mathbf{B} \in \mathbb{R}^{I \times M \times K}$$

$$y_{imk} = \sum_j x_{ijk} \, b_{mj}$$

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

Multiply each row (mode-2) fiber by **B**

- Tensor Times Vector

$$\mathbf{Y} = \mathcal{X} \, \bar{\times}_1 \, \mathbf{a} \in \mathbb{R}^{J \times K}$$

$$y_{jk} = \sum_i x_{ijk} \, a_i$$

Compute the dot product of $a$ and each column (mode-1) fiber

3-5

---

CMU SCS

Sandia National Laboratories

# Pictorial View of Mode-n Matrix Multiplication



Mode-1 multiplication (frontal slices)
$$\mathcal{Y} = \mathcal{X} \times_1 \mathbf{A}$$
$$\mathbf{Y}_{::k} = \mathbf{X}_{::k}\mathbf{A}^{\mathsf{T}}$$

Mode-2 multiplication (lateral slices)
$$\mathcal{Y} = \mathcal{X} \times_2 \mathbf{B}$$
$$\mathbf{Y}_{:j:} = \mathbf{X}_{:j:}\mathbf{B}^{\mathsf{T}}$$

Mode-3 multiplication (horizontal slices)
$$\mathcal{Y} = \mathcal{X} \times_3 \mathbf{C}$$
$$\mathbf{Y}_{i::} = \mathbf{X}_{i::}\mathbf{C}^{\mathsf{T}}$$

3-6

3

# Mode-n product Example

- Tensor times a matrix



3-7

# Mode-n product Example

- Tensor times a vector



3-8

# Outer, Kronecker, & Khatri-Rao Products

### 3-Way Outer Product

$$\mathfrak{X} = \mathbf{a} \circ \mathbf{b} \circ \mathbf{c}$$
$$x_{ijk} = a_i b_j c_k$$

Rank-1 Tensor

### Review: Matrix Kronecker Product

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1N}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2N}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{M1}\mathbf{B} & a_{M2}\mathbf{B} & \cdots & a_{MN}\mathbf{B} \end{bmatrix}$$

M x N    P x Q

MP x NQ

$$= \begin{bmatrix} \mathbf{a}_1 \otimes \mathbf{b}_1 & \mathbf{a}_1 \otimes \mathbf{b}_2 & \cdots & \mathbf{a}_N \otimes \mathbf{b}_Q \end{bmatrix}$$

### Matrix Khatri-Rao 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 x R   N x R                          MN x R

<u>Observe</u>: For two vectors **a** and **b**, **a** ○ **b** and **a** ⊗ **b** have the same elements, but one is shaped into a matrix and the other into a vector.

3-9

---

# Specially Structured Tensors

**CMU SCS**                                    Sandia National Laboratories

# Specially Structured Tensors

- Tucker Tensor

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

Our Notation

"core"

- Kruskal Tensor

$$\mathcal{X} = \sum_r \lambda_r\, \mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r$$
$$\equiv [\![\lambda\, ;\, \mathbf{U}, \mathbf{V}, \mathbf{W}]\!]$$

Our Notation



3-11

---

**CMU SCS**                                    Sandia National Laboratories

# Specially Structured Tensors

- Tucker Tensor

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

In matrix form:

$$\mathbf{X}_{(1)} = \mathbf{U}\mathbf{G}_{(1)}(\mathbf{W} \otimes \mathbf{V})^{\mathsf{T}}$$
$$\mathbf{X}_{(2)} = \mathbf{V}\mathbf{G}_{(2)}(\mathbf{W} \otimes \mathbf{U})^{\mathsf{T}}$$
$$\mathbf{X}_{(3)} = \mathbf{W}\mathbf{G}_{(3)}(\mathbf{V} \otimes \mathbf{U})^{\mathsf{T}}$$

$$\mathrm{vec}(\mathcal{X}) = (\mathbf{W} \otimes \mathbf{V} \otimes \mathbf{U})\mathrm{vec}(\mathcal{G})$$

- Kruskal Tensor

$$\mathcal{X} = \sum_r \lambda_r\, \mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r$$
$$\equiv [\![\lambda\, ;\, \mathbf{U}, \mathbf{V}, \mathbf{W}]\!]$$

In matrix form:

Let $\Lambda = \mathrm{diag}(\lambda)$

$$\mathbf{X}_{(1)} = \mathbf{U}\Lambda(\mathbf{W} \odot \mathbf{V})^{\mathsf{T}}$$
$$\mathbf{X}_{(2)} = \mathbf{V}\Lambda(\mathbf{W} \odot \mathbf{U})^{\mathsf{T}}$$
$$\mathbf{X}_{(3)} = \mathbf{W}\Lambda(\mathbf{V} \odot \mathbf{U})^{\mathsf{T}}$$

$$\mathrm{vec}(\mathcal{X}) = (\mathbf{W} \odot \mathbf{V} \odot \mathbf{U})\lambda$$

3-12

CMU SCS

# What is the HO Analogue of the Matrix SVD?

Matrix SVD:

$$\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{\mathsf{T}} = \quad = \quad \sigma_1 \boxed{\phantom{a}} + \sigma_2 \boxed{\phantom{a}} + \cdots + \sigma_R \boxed{\phantom{a}}$$

Tucker Tensor (finding bases for each subspace):

$$\mathbf{X} = \boldsymbol{\Sigma} \times_1 \mathbf{U} \times_2 \mathbf{V} = [\![ \boldsymbol{\Sigma} \, ; \mathbf{U}, \mathbf{V} ]\!]$$

Kruskal Tensor (sum of rank-1 components):

$$\mathbf{X} = \sum_{r=1}^{R} \sigma_r \, \mathbf{u}_r \circ \mathbf{v}_r = [\![ \sigma \, ; \mathbf{U}, \mathbf{V} ]\!]$$

3-13

CMU SCS

# Tensor Decompositions

**CMU SCS** Sandia National Laboratories

# Tucker Decomposition



$I$ x $J$ x $K$   $I$ x $R$   $J$ x $S$   $K$ x $T$   $R$ x $S$ x $T$

$$\mathcal{X} \approx [\![ \mathcal{G} \; ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!]$$

Given A, B, C, the optimal core is:

$$\mathcal{G} = [\mathcal{X} \; ; \mathbf{A}^\dagger, \mathbf{B}^\dagger, \mathbf{C}^\dagger]$$

- Proposed by Tucker (1966)
- AKA: Three-mode factor analysis, three-mode PCA, orthogonal array decomposition
- **A**, **B**, and **C** generally assumed to be orthonormal (generally assume they have full column rank)
- $\mathcal{G}$ is <u>not</u> diagonal
- Not unique

Recall the equations for converting a tensor to a matrix

$$\mathbf{X}_{(1)} = \mathbf{A}\mathbf{G}_{(1)}(\mathbf{C} \otimes \mathbf{B})^\mathsf{T}$$
$$\mathbf{X}_{(2)} = \mathbf{B}\mathbf{G}_{(2)}(\mathbf{C} \otimes \mathbf{A})^\mathsf{T}$$
$$\mathbf{X}_{(3)} = \mathbf{C}\mathbf{G}_{(3)}(\mathbf{B} \otimes \mathbf{A})^\mathsf{T}$$
$$\text{vec}(\mathcal{X}) = (\mathbf{C} \otimes \mathbf{B} \otimes \mathbf{A})\text{vec}(\mathcal{G})$$

3-15

---

**CMU SCS** Sandia National Laboratories

# Tucker Variations

See Kroonenberg & De Leeuw, Psychometrika,1980 for discussion.

- Tucker2

Identity Matrix



$I$ x $J$ x $K$   $I$ x $R$   $J$ x $S$   $R$ x $S$ x $K$

$$\mathcal{X} \approx [\![ \mathcal{G} \; ; \mathbf{A}, \mathbf{B}, \mathbf{I} ]\!]$$
$$\mathbf{X}_{(3)} \approx \mathbf{G}_{(3)}(\mathbf{B} \otimes \mathbf{A})^\mathsf{T}$$

- Tucker1



$I$ x $J$ x $K$   $I$ x $R$   $R$ x $J$ x $K$

$$\mathcal{X} \approx [\![ \mathcal{G} \; ; \mathbf{A}, \mathbf{I}, \mathbf{I} ]\!]$$
$$\mathbf{X}_{(1)} \approx \mathbf{A}\mathbf{G}_{(1)}$$

Finding principal components in only mode 1 can be solved via rank-R matrix SVD

3-16

# Solving for Tucker

$$\mathcal{X} \approx [\![ \mathcal{G} \, ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!]$$

Given A, B, C orthonormal, the optimal core is:

$$\mathcal{G} = [\![ \mathcal{X} \, ; \mathbf{A}^\mathsf{T}, \mathbf{B}^\mathsf{T}, \mathbf{C}^\mathsf{T} ]\!]$$

Tensor norm is the square root of the sum of all the elements squared

Eliminate the core to get:

$$\| \mathcal{X} - [\![ \mathcal{G} \, ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!] \|^2 = \| \mathcal{X} \|^2 - 2\langle \mathcal{X}, [\![ \mathcal{G} \, ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!] \rangle + \| \mathcal{G} \|^2$$

$$= \| \mathcal{X} \|^2 - \left\| [\![ \mathcal{X} \, ; \mathbf{A}^\mathsf{T}, \mathbf{B}^\mathsf{T}, \mathbf{C}^\mathsf{T} ]\!] \right\|^2$$

Minimize
s.t. **A,B,C** orthonormal

fixed          maximize this

If B & C are fixed, then we can solve for A as follows:

$$\left\| [\![ \mathcal{X} \, ; \mathbf{A}^\mathsf{T}, \mathbf{B}^\mathsf{T}, \mathbf{C}^\mathsf{T} ]\!] \right\| = \left\| \mathbf{A}^\mathsf{T} \mathbf{X}_{(1)} (\mathbf{C} \otimes \mathbf{B}) \right\|$$

Optimal **A** is R left leading singular vectors for $\boxed{\mathbf{X}_{(1)}(\mathbf{C} \otimes \mathbf{B})}$

3-17

$I \times J \times K$  $I \times R$  $J \times S$  $K \times T$  $R \times S \times T$

# Higher Order SVD (HO-SVD)

Not optimal, but often used to initialize Tucker-ALS algorithm.

$I \times J \times K$  $I \times R$  $J \times S$  $K \times T$  $R \times S \times T$

(Observe connection to Tucker1)

$\mathbf{A} = \text{leading } R \text{ left singular vectors of } \mathbf{X}_{(1)}$

$\mathbf{B} = \text{leading } S \text{ left singular vectors of } \mathbf{X}_{(2)}$

$\mathbf{C} = \text{leading } T \text{ left singular vectors of } \mathbf{X}_{(3)}$

$$\mathcal{G} = [\![ \mathcal{X} \, ; \mathbf{A}^\mathsf{T}, \mathbf{B}^\mathsf{T}, \mathbf{C}^\mathsf{T} ]\!]$$

De Lathauwer, De Moor, & Vandewalle, SIMAX, 1980          3-18

9

# Tucker-Alternating Least Squares (ALS)

*Successively solve for each component (**A**,**B**,**C**).*



$I$ x $J$ x $K$

$I$ x $R$

$J$ x $S$

$K$ x $T$

$R$ x $S$ x $T$

- Initialize
  - Choose R, S, T
  - Calculate **A**, **B**, **C** via HO-SVD
- Until converged do…
  - **A** = R leading left singular vectors of $\mathbf{X}_{(1)}(\mathbf{C}\otimes\mathbf{B})$
  - **B** = S leading left singular vectors of $\mathbf{X}_{(2)}(\mathbf{C}\otimes\mathbf{A})$
  - **C** = T leading left singular vectors of $\mathbf{X}_{(3)}(\mathbf{B}\otimes\mathbf{A})$
- Solve for core:

$$\mathcal{G} = [\![\,\mathcal{X}\,;\mathbf{A}^\mathsf{T},\mathbf{B}^\mathsf{T},\mathbf{C}^\mathsf{T}]\!]$$

Kroonenberg & De Leeuw, Psychometrika, 1980

3-19

# Tucker in Not Unique



$I$ x $J$ x $K$

$I$ x $R$

$J$ x $S$

$K$ x $T$

$R$ x $S$ x $T$

Tucker decomposition is <u>not</u> unique. Let Y be an RxR orthogonal matrix. Then…

$$\mathcal{X} \approx \mathcal{G}\times_1\mathbf{A}\times_2\mathbf{B}\times_3\mathbf{C} = \left(\mathcal{G}\times_1\mathbf{Y}^\mathsf{T}\right)\times_1(\mathbf{A}\mathbf{Y})\times_2\mathbf{B}\times_3\mathbf{C}$$

$$\mathbf{X}_{(1)} \approx \mathbf{A}\mathbf{G}_{(1)}(\mathbf{C}\otimes\mathbf{B})^\mathsf{T} = \mathbf{A}\mathbf{Y}\mathbf{Y}^\mathsf{T}\mathbf{G}_{(1)}(\mathbf{C}\otimes\mathbf{B})^\mathsf{T}$$

3-20

**CANDECOMP/PARAFAC Decomposition**

$I \times J \times K$   $I \times R$   $J \times R$   $K \times R$

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

$R \times R \times R$

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

3-21

---

**PARAFAC-Alternating Least Squares (ALS)**

*Successively solve for each component (**A**,**B**,**C**).*

$$\mathfrak{X} \approx [\![ \lambda ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!]$$

$$\mathbf{X}_{(1)} \approx \mathbf{A}\mathbf{\Lambda}(\mathbf{C} \odot \mathbf{B})^{\mathsf{T}}$$

$I \times J \times K$

**Find all the vectors in one mode at a time**

**KHATRI-RAO PRODUCT**
(column-wise Kronecker product)

$$\mathbf{C} \odot \mathbf{B} \equiv \begin{bmatrix} \mathbf{c}_1 \otimes \mathbf{b}_1 & \mathbf{c}_2 \otimes \mathbf{b}_2 & \cdots \mathbf{c}_R \otimes \mathbf{b}_R \end{bmatrix}$$

$$(\mathbf{C} \odot \mathbf{B})^{\dagger} \equiv (\mathbf{C}^{\mathsf{T}}\mathbf{C} * \mathbf{B}^{\mathsf{T}}\mathbf{B})^{\dagger}(\mathbf{C} \odot \mathbf{B})^{\mathsf{T}}$$

Hadamard Product

If **C**, **B**, and $\Lambda$ are fixed, the optimal A is given by:

$$\mathbf{A} = \mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})(\mathbf{C}^{\mathsf{T}}\mathbf{C} * \mathbf{B}^{\mathsf{T}}\mathbf{B})^{\dagger}\mathbf{\Lambda}^{-1}$$

*Repeat for **B**,**C**, etc.*   3-22

# PARAFAC is often unique

$I$ x $J$ x $K$

Assume PARAFAC decomposition is exact.

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

Sufficient condition for uniqueness (Kruskal, 1977):

$$2R + 2 \le k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}}$$

$k_{\mathbf{A}}$ = k-rank of **A** = max number k such that every set of k columns of **A** is linearly independent

3-23

# Tucker vs. PARAFAC Decompositions

- Tucker
  - Variable transformation in each mode
  - Core G may be dense
  - A, B, C generally orthonormal
  - Not unique

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

$I$ x $J$ x $K$    $I$ x $R$    $J$ x $S$    $K$ x $T$    $R$ x $S$ x $T$

$I$ x $J$ x $K$

3-24

12

**CMU SCS**

# Roadmap

- Motivation
- Matrix tools
- Tensor basics
- Tensor extensions
  - Other decompositions
  - Nonnegative PARAFAC
  - Handling missing values
- Software demo
- Case studies



4-1

---

**CMU SCS**

# Other Tensor Decompositions

CMU SCS

# Combining Tucker & PARAFAC

$$\mathcal{X} \approx [\![ \hat{\mathcal{X}} ; \mathbf{U}, \mathbf{V}, \mathbf{W} ]\!]$$

$M \times I$

$N \times J$

$P \times K$

$\mathbf{W}$

$\mathcal{X}$    $\approx$    $\mathbf{U}$    $\hat{\mathcal{X}}$    $\mathbf{V}$

$I \times J \times K$

$M \times N \times P$

$$\hat{\mathcal{X}} = [\![ \mathcal{X} ; \mathbf{U}^{\mathsf{T}}, \mathbf{V}^{\mathsf{T}}, \mathbf{W}^{\mathsf{T}} ]\!]$$

$$\hat{\mathcal{X}} \approx [\![ \hat{\lambda} ; \hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}} ]\!]$$

$$\mathcal{X} \approx [\![ \hat{\lambda} ; \mathbf{U}\hat{\mathbf{A}}, \mathbf{V}\hat{\mathbf{B}}, \mathbf{W}\hat{\mathbf{C}} ]\!] \equiv [\![ \lambda ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!]$$

Bro and Andersson, 1998

- Step 1: Choose orthonormal matrices U, V, W to compress tensor (Tucker tensor!)
  - Typically HO-SVD can be used
- Step 2: Run PARAFAC on smaller tensor

- Step 3: Reassemble result

4-3

---

CMU SCS

# 2-Way DEDICOM

$N \times N$          $N \times M$        $M \times N$

$\mathbf{X}$    $=$    $\mathbf{A}$    $\mathbf{R}$    $\mathbf{A^T}$

Dense, nonsymmetric M x M matrix

- 2-way DEDICOM introduced by Harshman (1978)
- X is a matrix of interactions between N entities
- Interactions can be nonsymmetric
- Assumes there are "M" roles
- Each entity has a weight for each role in A
- $R_{ij}$ = interaction weight for roles i & j

4-4

**CMU SCS**

# 3-Way DEDICOM



3-way DEDICOM

$$\mathbf{X}_{::k} = \mathbf{A}\, \mathbf{D}_{::k}\, \mathbf{R}\, \mathbf{D}_{::k}\, \mathbf{A}^{\mathsf{T}}$$

- 3-way DEDICOM due to Kiers (1993)
- Once again, X captures interactions among entities
- Third dimension can correspond to time
- Diagonal slices capture participation of each role at each time
- See Bader et al., SAND2006-7744 , for application to Enron email data

4-5

---

**CMU SCS**

# Nonnegativity

---

CMU SCS

Sandia
National
Laboratories

# Non-negative Matrix Factorization

$$\left\| \mathbf{X} - \mathbf{A}\mathbf{B}^{\mathsf{T}} \right\|$$  — Minimize subject to elements of **A** and **B** being positive.

Update formulas (do not increase objective function):

$$\mathbf{A} = \mathbf{A} * (\mathbf{X}\mathbf{B}) \oslash (\mathbf{A}\mathbf{B}^{\mathsf{T}}\mathbf{B})$$
$$\mathbf{B} = \mathbf{B} * (\mathbf{X}^{\mathsf{T}}\mathbf{A}) \oslash (\mathbf{B}\mathbf{A}^{\mathsf{T}}\mathbf{A})$$

Elementwise multiply          Elementwise divide
(Hadamard product)

Lee & Seung, Nature, 1999                                    4-7

---

CMU SCS

Sandia
National
Laboratories

# Non-negative 3-Way PARAFAC Factorization

$$\left\| \mathcal{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] \right\|$$  — Minimize subject to elements of **A**, **B** and **C** being positive.

Lee-Seung-like update formulas can be derived for 3D and higher:

$$\mathbf{A} = \mathbf{A} * (\mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})) \oslash (\mathbf{A}(\mathbf{C}^{\mathsf{T}}\mathbf{C} * \mathbf{B}^{\mathsf{T}}\mathbf{B}))$$
$$\mathbf{B} = \mathbf{B} * (\mathbf{X}_{(2)}(\mathbf{C} \odot \mathbf{A})) \oslash (\mathbf{B}(\mathbf{C}^{\mathsf{T}}\mathbf{C} * \mathbf{A}^{\mathsf{T}}\mathbf{A}))$$
$$\mathbf{C} = \mathbf{C} * (\mathbf{X}_{(3)}(\mathbf{B} \odot \mathbf{A})) \oslash (\mathbf{C}(\mathbf{B}^{\mathsf{T}}\mathbf{B} * \mathbf{A}^{\mathsf{T}}\mathbf{A}))$$

Elementwise multiply          Elementwise divide
(Hadamard product)

M. Mørup, L. K. Hansen, J. Parnas, S. M. Arnfred, *Decomposing the time-frequency representation of EEG using non-negative matrix and multi-way factorization*, 2006                    4-8

**CMU SCS**

# Handling Missing Data

---

**CMU SCS**

**Sandia National Laboratories**

# A Quick Overview on Handling Missing Data

- Consider sparse PARAFAC where $\mathcal{X}$ is missing data:

$$\mathcal{X} \approx [\![ \lambda \; ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!]$$

- Typically, missing values are just set to zero
- More sophisticated approaches for handling missing values:
  - Weighted least squares loss function
    - Ignore missing values
  - Data imputation
    - Estimate missing values
- See, e.g., Kiers, Psychometrika, 1997 and Srebro & Jaakkola, ICML 2003

4-10

# Weighted Least Squares

$$w_{ijk} = \begin{cases} 1 & x_{ijk} \text{ is known} \\ 0 & \text{otherwise} \end{cases}$$

Weight Tensor

- Weight the least squares problem so that the missing elements are ignored:

Weighted
Least Squares

$$\sum_i \sum_j \sum_k w_{ijk} \left( x_{ijk} - \sum_r \lambda_r a_{ir} b_{jr} c_{kr} \right)^2$$

- But this problem is often too hard to solve directly!

4-11

---

CMU SCS

# Missing Value Imputation

- Use the current estimate to fill in the missing values

$$\mathcal{E} = [\![ \lambda \, ; \mathbf{A}, \mathbf{B}, \mathbf{C} ]\!]$$

Current Estimate

- The tensor for the next iteration of the algorithm is:

Known Values    Estimates of Unknowns

$$\hat{\mathcal{X}} = \overbrace{\mathcal{W} * \mathcal{X}} + \overbrace{(1 - \mathcal{W}) * \mathcal{E}}$$
$$= \underbrace{\mathcal{X} - \mathcal{W} * \mathcal{E}} + \underbrace{\mathcal{E}}$$

Sparse!        Kruskal Tensor

- Challenge is finding a good initial estimate

4-12

**CMU SCS**

# Roadmap

- Motivation
- Matrix tools
- Tensor basics
- Tensor extensions
- Software demo
- Case studies



4-13

---

**CMU SCS**

# Computations with Tensors

Sandia National Laboratories

# Tensor Toolbox for MATLAB

**http://csmr.ca.sandia.gov/~tgkolda/TensorToolbox**

- Six object-oriented tensor classes
  - Working with tensors is easy
- Most comprehensive set of kernel operations in any language
  - E.g., arithmetic, logical, multiplication operations
- Sparse tensors are unique
  - Speed-ups of two orders of magnitude for smaller problems
  - Larger problems than ever before

| Workspace | | |
|---|---|---|
| Na... ▲ | Value | Class |
| a | \<4x3 double\> | double |
| b | \<4x3x2 tensor\> | tensor |
| c | \<5x4 double\> | double (sparse) |
| d | \<5x4x2 sptensor\> | sptensor |

- Free for research or evaluations purposes
- 297 unique registered users from all over the world (as of January 17, 2006)

4-15

Bader & Kolda, ACM TOMS 2006 & SAND2006-7592

---

Sandia National Laboratories

# Dense Tensors

- Largest tensor that can be stored on a laptop is 200 x 200 x 200

- Typically, tensor operations are reduced to matrix operations
  - Requires permuting and reshaping the tensor

- Example: Mode-n tensor-matrix multiply

$I$ x $J$ x $K$

Example: Mode-1 Matrix Multiply

$$\mathcal{Y} = \mathcal{X} \times_1 \mathbf{U}$$

$M$ x $J$ x $K$  $\quad I$ x $J$ x $K \quad$  $M$ x $I$

$$\mathbf{Y}_{(n)} = \mathbf{U}\mathbf{X}_{(n)}$$

$\qquad\qquad M$ x $I$

$M$ x $JK \qquad\qquad I$ x $JK$

4-16

# Sparse Tensors: Only Store Nonzeros

Example: Tensor-Vector Multiply (in all modes)

Store just the nonzeros of a tensor (assume coordinate format)

$$\alpha = \mathcal{X} \,\bar{\times}_1\, \mathbf{a} \,\bar{\times}_2\, \mathbf{b} \,\bar{\times}_3\, \mathbf{c}$$

$$= \sum_i \sum_j \sum_k x_{ijk}\, a_i\, b_j\, c_k$$

$$= \sum_p v_p\, a_{s(p,1)}\, b_{s(p,j)}\, c_{s(p,k)}$$

$p$th nonzero

1st subscript of $p$th nonzero

2nd subscript of $p$th nonzero

3rd subscript of $p$th nonzero

4-17

---

CMU SCS

# Tucker Tensors:
# Store Core & Factors

$$\mathcal{X} = \boxed{U}\ \boxed{\mathcal{G}}\ \boxed{V}\quad \boxed{W}$$

Tucker tensor stores the core (which can be dense, sparse, or structured) and the factors.

Example: Mode-3 Tensor-Vector Multiply

$$\mathbf{Y} = \mathcal{X} \,\bar{\times}_3\, \mathbf{z}$$

$$= (\mathcal{G} \times_1 \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W})\ \bar{\times}_3\, \mathbf{z}$$

$$= \mathcal{G} \times_1 \mathbf{U} \times_2 \mathbf{V} \,\bar{\times}_3\, \mathbf{W}^\mathsf{T}\mathbf{z}$$

$$= \underbrace{\mathcal{G} \,\bar{\times}_3\, \mathbf{W}^\mathsf{T}\mathbf{z}}_{\mathcal{H}} \times_1 \mathbf{U} \times_2 \mathbf{V} = [\![\mathcal{H}\,;\mathbf{U},\mathbf{V}]\!]$$

Result is a Tucker Tensor

4-18

**Kruskal Example: Store Factors**

Kruskal tensors store factor matrices and scaling vector.

Example: Norm

$$\|\mathcal{X}\|^2 = \|\,[\![\lambda\,; \mathbf{U}, \mathbf{V}, \mathbf{W}]\!]\,\|^2$$
$$= \|\,(\mathbf{W} \odot \mathbf{V} \odot \mathbf{U})\lambda\,\|^2$$
$$= \lambda^{\mathsf{T}}(\mathbf{W} \odot \mathbf{V} \odot \mathbf{U})^{\mathsf{T}}(\mathbf{W} \odot \mathbf{V} \odot \mathbf{U})\lambda$$
$$= \lambda^{\mathsf{T}}(\mathbf{W}^{\mathsf{T}}\mathbf{W} * \mathbf{V}^{\mathsf{T}}\mathbf{V} * \mathbf{U}^{\mathsf{T}}\mathbf{U})\lambda$$

4-19

**CMU SCS**

# Incrementalization

---

**CMU SCS**                                                    Sandia National Laboratories

## Incremental Tensor Decomposition

- Dynamic data model
  - Tensor Streams
- Dynamic Tensor Decomposition (DTA)
- Streaming Tensor Decomposition (STA)
- Window-based Tensor Decomposition (WTA)

**CMU SCS**

# Dynamic Tensor Stream

- Streams come with structure
  - (time, source, destination, port)
  - (time, author, keyword)
- How to summarize tensor streams effectively and incrementally?

Faloutsos, Kolda, Sun

5-3

**CMU SCS**

# Dynamic Data model

- Tensor Streams
  - A sequence of Mth order tensor

$$\mathcal{X}_1 \ldots \mathcal{X}_n \text{ where } \mathcal{X}_i \in \mathbf{R}^{N_1 \times \ldots \times N_M}$$

n is increasing over time

| Order | 1st | 2nd | 3rd |
|---|---|---|---|
| Correspondence | Multiple streams | Time evolving graphs | 3D arrays |
| Example | Sensors | keyword / author / time | Ports / Destinations / Sources $\mathcal{X}$ |

Faloutsos, Kolda, Sun

5-4

**CMU SCS**

# Incremental Tensor Decomposition

☺ Dynamic data model

- Tensor Streams
- Dynamic Tensor Decomposition (DTA)
- Streaming Tensor Decomposition (STA)
- Window-based Tensor Decomposition (WTA)

1. Jimeng Sun, Spiros Papadimitriou, Philip Yu. Window-based Tensor Analysis on High-dimensional and Multi-aspect Streams, *ICDM 2006*
2. Jimeng Sun, Dacheng Tao, Christos Faloutsos. Beyond Streams and Graphs: Dynamic Tensor Analysis, *KDD 2006*

Faloutsos, Kolda, Sun                                              5-5

---

**CMU SCS**

# Incremental Tensor Decomposition



Faloutsos, Kolda, Sun                                              5-6

**CMU SCS**

# 1st order DTA - problem

Given $x_1 \ldots x_n$ where each $x_i \in R^N$, find
$U \in R^{N \times R}$ such that the error e is
small: $e = \sum_{i=1}^{n} \| \mathbf{x}_i - \mathbf{x}_i UU^T \|_F^2$



Note that Y = XU

Faloutsos, Kolda, Sun                  5-7

---

**CMU SCS**

# 1st order Dynamic Tensor Analysis



<u>Input</u>: new data vector $x \in R^N$, old variance matrix $C \in R^{N \times N}$

<u>Output</u>: new projection matrix $U \in R^{N \times R}$

Algorithm:

1. update variance matrix $C_{new} = x^T x + C$
2. Diagonalize $U \Lambda U^T = C_{new}$
3. Determine the rank R and return U



Diagonalization has to be done for *every* new **x**!

Faloutsos, Kolda, Sun                  5-8

4

# M$^{th}$ order DTA



Reconstruct Variance Matrix

$\mathbf{U}_d^T$

$\mathbf{U}_d$  $\mathbf{S}_d$  $=$  $\mathbf{C}_d$

Construct Variance Matrix of Incremental Tensor

$\mathcal{X}$

Matricizing

Matricizing, Transpose

$\times$  $=$  $\mathbf{x}_{(d)}^T \mathbf{x}_{(d)}$

$\mathbf{X}_{(d)}^T$  $\mathbf{X}_{(d)}$

$\mathbf{C}_d$

Update Variance Matrix

Diagonalize Variance Matrix

$\mathbf{U}_d^T$

$\mathbf{U}_d$  $\mathbf{S}_d$

Faloutsos, Kolda, Sun                                    5-9

---

# M$^{th}$ order DTA – complexity

**Storage:**

$O(\prod N_i)$, i.e., size of an input tensor at a single timestamp

**Computation:**

$\sum N_i^3$ (or $\sum N_i^2$)    diagonalization of C

$+ \sum N_i \prod N_i$         matrix multiplication $X_{(d)}^T X_{(d)}$

For low order tensor(<3), diagonalization is the main cost

For high order tensor,  matrix multiplication is the main cost

Faloutsos, Kolda, Sun                                    5-10

Sandia National Laboratories

# Incremental Tensor Decomposition

☺ Dynamic data model
- Tensor Streams

☺ Dynamic Tensor Decomposition (DTA)

- Streaming Tensor Decomposition (STA)
- Window-based Tensor Decomposition (WTA)

1. Jimeng Sun, Spiros Papadimitriou, Philip Yu. Window-based Tensor Analysis on High-dimensional and Multi-aspect Streams, *ICDM 2006*
2. Jimeng Sun, Dacheng Tao, Christos Faloutsos. Beyond Streams and Graphs: Dynamic Tensor Analysis, *KDD 2006*

Faloutsos, Kolda, Sun

5-11

---

CMU SCS  Sandia National Laboratories

# 1st order Streaming Tensor Analysis (STA)

- Adjust U smoothly when new data arrive without diagonalization [VLDB05]
- For each new point x
  - Project onto current line
  - Estimate error
  - Rotate line in the direction of the error and in proportion to its magnitude

For each new point $x$ and for $i = 1, \ldots, k$ :

- $y_i := U_i^T x$    (proj. onto $U_i$)
- $d_i \leftarrow \lambda d_i + y_i^2$    (energy $\propto i$-th eigenval.)
- $e_i := x - y_i U_i$    (error)
- $U_i \leftarrow U_i + (1/d_i) y_i e_i$    (update estimate)
- $x \leftarrow x - y_i U_i$    (repeat with remainder)

Faloutsos, Kolda, Sun

error

Sensor 2

U

Sensor 1 5-12

6

# M$^{th}$ order STA

$\mathbf{X}^T_{(d)}$

Matricizing

$\mathcal{X}$

$x$

$U_1$ updated

$e_1$

$U_1$

$y_1$

- Run 1$^{st}$ order STA along each mode
- Complexity:
  - Storage: $O(\prod N_i)$
  - Computation: $\sum R_i \prod N_i$ which is smaller than DTA

Faloutsos, Kolda, Sun

5-13

---

# Incremental Tensor Decomposition

☺ Dynamic data model
- Tensor Streams

☺ Dynamic Tensor Decomposition (DTA)

☺ Streaming Tensor Decomposition (STA)

- **Window-based Tensor Decomposition (WTA)**

1. Jimeng Sun, Spiros Papadimitriou, Philip Yu. Window-based Tensor Analysis on High-dimensional and Multi-aspect Streams, *ICDM 2006*
2. Jimeng Sun, Dacheng Tao, Christos Faloutsos. Beyond Streams and Graphs: Dynamic Tensor Analysis, *KDD 2006*

Faloutsos, Kolda, Sun

5-14

**CMU SCS**

# Moving Window scheme (MW)

- Update the variance matrix $C_{(i)}$ **incrementally**
- Diagonalize **C(i) to find U(i)**

*A good and efficient initialization*

$$C_d^{old} - \square + \blacksquare = C_d^{new}$$

**Update variance matrix**

$U_{(d)}$    Diagonalize

Faloutsos, Kolda, Sun

5-17



**CMU SCS**

# Roadmap

- Motivation
- Matrix tools
- Tensor basics
- Tensor extensions
- Software demo
- **Case studies**

Faloutsos, Kolda, Sun

5-18

**CMU SCS**

# P1: sensor monitoring

2nd factor
Scaling factor 154

time          location          type

- 2nd factor captures an atypical trend:
  - Uniformly across all time
  - Concentrating on 3 locations
  - Mainly due to voltage
- Interpretation: two sensors have low battery, and the other one has high battery.

Faloutsos, Kolda, Sun

5-21



**CMU SCS**

# P3: Social network analysis

- Multiway latent semantic indexing (LSI)
  - Monitor the change of the community structure over time

2004    Keywords    Christos Faloutsos
DM                  $U_A$    Michael Stonebreaker
        DB
1990                        2004
Authors                     1990
        DB      $U_K$
            'Pattern'  'Query'

Faloutsos, Kolda, Sun

5-22

**CMU SCS**                                                    Sandia National Laboratories

# P3: Social network analysis (cont.)

| Authors | Keywords | Year |
|---|---|---|
| michael carey, michael stonebreaker, h. jagadish, hector garcia-molina | queri,parallel,optimization,concurr, ...ent | 1995 |
| surajit chaudhuri,mitch cherniack,michael stonebreaker,ugur etintemel | distribu,systems,view,storage,servic,process, cache | 2004 |
| jiawei han,jian pei,philip s. yu, jianyong wang,charu c. aggarwal | ...pattern,support,cluster, ...ner,queri | 2004 |

DB

DM

- Two groups are correctly identified: Databases and Data mining
- People and concepts are drifting over time

Faloutsos, Kolda, Sun                                           5-23

---

**CMU SCS**                                                    Sandia National Laboratories

# P4: Network anomaly detection



Abnormal traffic          Reconstruction error over time          Normal traffic

- Reconstruction error gives indication of anomalies.
- Prominent difference between normal and abnormal ones is mainly due to the unusual scanning activity (confirmed by the campus admin).

Faloutsos, Kolda, Sun                                           5-24

**CMU SCS**

# P5: Web graph mining

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

Faloutsos, Kolda, Sun                    5-25

---

**CMU SCS**

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

Sparse adjacency matrix and its SVD:

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

$$X \approx \sum_r \sigma_r \, h_r \circ a_r$$

authority scores for 1st topic

authority scores for 2nd topic

hub scores for 1st topic

hub scores for 2nd topic

Kleinberg, JACM, 1999

Faloutsos, Kolda, Sun                    5-26

**CMU SCS**

Sandia National Laboratories

# HITS Authorities on Sample Data

| 1st Principal Factor | |
|---|---|
| .97 | www.ibm.com |
| .24 | www.alphaw... |
| .08 | www-128.ibm... |
| .05 | www.develop... |
| .02 | www.research... |
| .01 | www.redbook... |
| .01 | news.com.c... |

| 2nd Principal Factor | |
|---|---|
| .99 | www.lehigh.edu |
| .11 | www2.lehigh.edu |
| .06 | www.lehigha... |
| .06 | www.lehighs... |
| .02 | www.bethleh... |
| .02 | www.adobe.... |
| .02 | lewisweb.cc.... |
| .02 | www.leo.lehi... |
| .02 | www.distanc... |
| .02 | fp1.cc.lehigh... |

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

| 4th Principal Factor | |
|---|---|
| .60 | www.pueblo.gsa.gov |
| .45 | www.whitehouse.gov |
| .35 | www.irs.gov |
| .31 | travel.state... |
| .22 | www.gsa.g... |
| .20 | www.ssa.g... |
| .16 | www.censu... |
| .14 | www.govbe... |
| .13 | www.kids.g... |
| .13 | www.usdoj... |

| 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 |

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

authority scores for 1st topic    authority scores for 2nd topic

$$\underset{\text{from}}{\text{to}} \quad = \quad | \quad + \quad | \quad + \quad ...$$

hub scores for 1st topic    hub scores for 2nd topic

Faloutsos, Kolda, Sun

5-27

---

**CMU SCS**

Sandia National Laboratories

# Three-Dimensional View of the Web

**Endangered Species**
Animals today are being threatened by a variety of environmental pressures. For example, the jaguar is losing prime habitat in the world. Zoos are trying to raise awareness of their plight.

**Jaguar FAQ**
Jaguars are an endangered species that live in the tropical rain forests of Central and South America. They live about 11 years in the wild and up to 22 years at a zoo.

**Rain Forest Zoo**
We have a new exhibit opening next month highlighting the endangered species of the Americas, including the jaguar.

**Online Atlas**
View maps of animal habitats from around the world, including those of endangered animals in North, South, and Central America.

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

Observe that this tensor is very sparse!

Website 1    jaguar    Website 2
endangered
species
zoo    endangered    jaguar    zoo
America
America
Website 3    Website 4

zoo
America
jaguar
species
endangered
1 2 3 4

Faloutsos, Kolda, Sun
Kolda, Bader, Kenny, ICDM05

5-28

**CMU SCS**

Sandia
National
Laboratories

# Topical HITS (TOPHITS)

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

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

authority scores
for 1st topic

authority scores
for 2nd topic

hub scores
for 1st topic

hub scores
for 2nd topic

Faloutsos, Kolda, Sun

5-29

---

**CMU SCS**

Sandia
National
Laboratories

# Topical HITS (TOPHITS)

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

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

term scores
for 1st topic

term scores
for 2nd topic

authority scores
for 1st topic

authority scores
for 2nd topic

hub scores
for 1st topic

hub scores
for 2nd topic

Faloutsos, Kolda, Sun

5-30

15

**CMU SCS**                                                                 Sandia National Laboratories

# TOPHITS Terms & Authorities on Sample Data

TOPHITS uses 3D analysis to find the dominant groupings of web pages and terms.

$$z_{ijk} = \begin{cases} \dfrac{1}{\log(w_k)+1} & \text{if } i \to j \text{ with term } k \\ 0 & \text{otherwise} \end{cases}$$

$w_k$ = # unique links using term k

**1st Principal Factor**

| | | | |
|---|---|---|---|
| .23 | JAVA | .86 | java.sun.com |
| .18 | SUN | | |
| .17 | PLATF | | |
| .16 | SOLAR | | |
| .16 | DEVEL | | |
| .15 | EDITIO | | |
| .15 | DOWN | | |
| .14 | INFO | | |
| .12 | SOFTW | | |
| .12 | NO-RE | | |

**2nd Principal Factor**

| | | | |
|---|---|---|---|
| .20 | NO-READABLE-TEXT | .99 | www.lehigh.edu |
| .16 | FACUL | | |
| .16 | SEARC | | |
| .16 | NEWS | | |
| .16 | LIBRA | | |
| .16 | COMP | | |
| .12 | LEHIG | | |

**3rd Principal Factor**

| | | | |
|---|---|---|---|
| .15 | NO-READABLE-TEXT | .97 | www.ibm.com |
| .15 | IBM | .18 | www.alphaworks.ibm.com |
| .12 | SERVI | | |
| .12 | WEBS | | |
| .12 | WEB | | |
| .11 | DEVEL | | |
| .11 | LINUX | | |
| .11 | RESOU | | |
| .11 | TECHN | | |
| .10 | DOWN | | |

**4th Principal Factor**

| | | | |
|---|---|---|---|
| .26 | INFORMATION | .87 | www.pueblo.gsa.gov |
| .24 | FEDERAL | .24 | www.irs.gov |
| .23 | CITIZE | | |
| .22 | OTHER | | |
| .19 | CENTE | | |
| .19 | LANGU | | |
| .25 | U.S | | |
| .15 | PUBLIC | | |
| .14 | CONSU | | |
| .13 | FREE | | |

**6th Principal Factor**

| | | | |
|---|---|---|---|
| .26 | PRESIDENT | .87 | www.whitehouse.gov |
| .25 | NO-READABLE-TEXT | .18 | www.irs.gov |
| .25 | BUSH | | |
| .25 | WELC | | |
| .17 | WHITE | | |
| .16 | U.S | | |
| .15 | HOUS | | |
| .13 | BUDG | | |
| .13 | PRES | | |
| .11 | OFFIC | | |

**12th Principal Factor**

| | | | |
|---|---|---|---|
| .75 | OPTIMIZATION | .35 | www.palisade.com |
| .58 | SOFTWARE | .35 | www.solver.com |
| .08 | DECIS | | |
| .07 | NEOS | | |
| .06 | TREE | | |
| .05 | GUIDE | | |
| .05 | SEARC | | |
| .05 | ENGIN | | |
| .05 | CONTI | | |
| .05 | ILOG | | |

**13th Principal Factor**

| | | | |
|---|---|---|---|
| .46 | ADOBE | .99 | www.adobe.com |
| .45 | READER | | |
| .45 | ACROI | | |

**16th Principal Factor**

| | | | |
|---|---|---|---|
| .50 | WEATHER | .81 | www.weather.gov |
| .24 | OFFICE | .41 | www.spc.noaa.gov |
| .23 | CENTER | .30 | lwf.ncdc.noaa.gov |
| .19 | NO-RE | | |
| .17 | ORGA | | |

**19th Principal Factor**

| | | | |
|---|---|---|---|
| .22 | TAX | .73 | www.irs.gov |
| .17 | TAXES | .43 | travel.state.gov |
| .15 | CHILD | .22 | www.ssa.gov |
| .15 | RETIREMENT | .08 | www.govbenefits.gov |
| .14 | BENEFITS | .06 | www.usdoj.gov |
| .14 | STATE | .03 | www.census.gov |
| .14 | INCOME | .03 | www.usmint.gov |
| .13 | SERVICE | .02 | www.nws.noaa.gov |
| .13 | REVENUE | .02 | www.gsa.gov |
| .12 | CREDIT | .01 | www.annualcreditreport.com |

(additional factor column)

| | |
|---|---|
| .15 | NWS |
| .15 | SEVER |
| .15 | FIRE |
| .15 | POLIC |
| .14 | CLIMA |

**Tensor PARAFAC**



term scores for 1st topic
term scores for 2nd topic
authority scores for 1st topic
authority scores for 2nd topic
hub scores for 1st topic
hub scores for 2nd topic

Faloutsos, Kolda, Sun

---

**CMU SCS**                                                                 Sandia National Laboratories

# Tensor faces

(Vasilescu & Terzopoulos, 2002; Vasilescu & Terzopoulos, 2003)



people

expressions

views

illuminations

Faloutsos, Kolda, Sun

5-32

**CMU SCS**

# Eigenfaces

- Facial images (identity change)

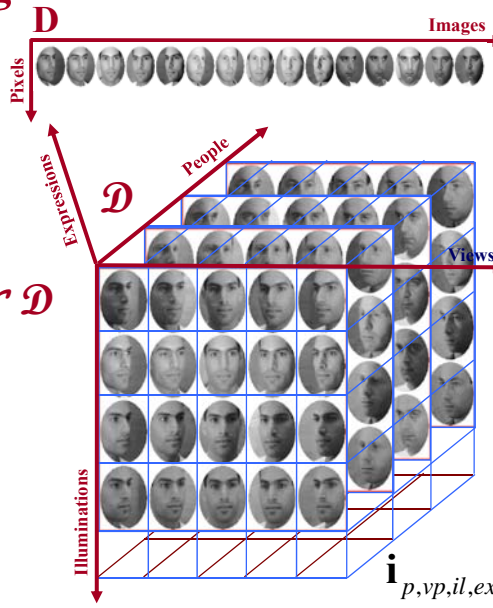- Eigenfaces bases vectors capture the variability in facial appearance (do not decouple pose, illumination, …)

Faloutsos, Kolda, Sun

5-33
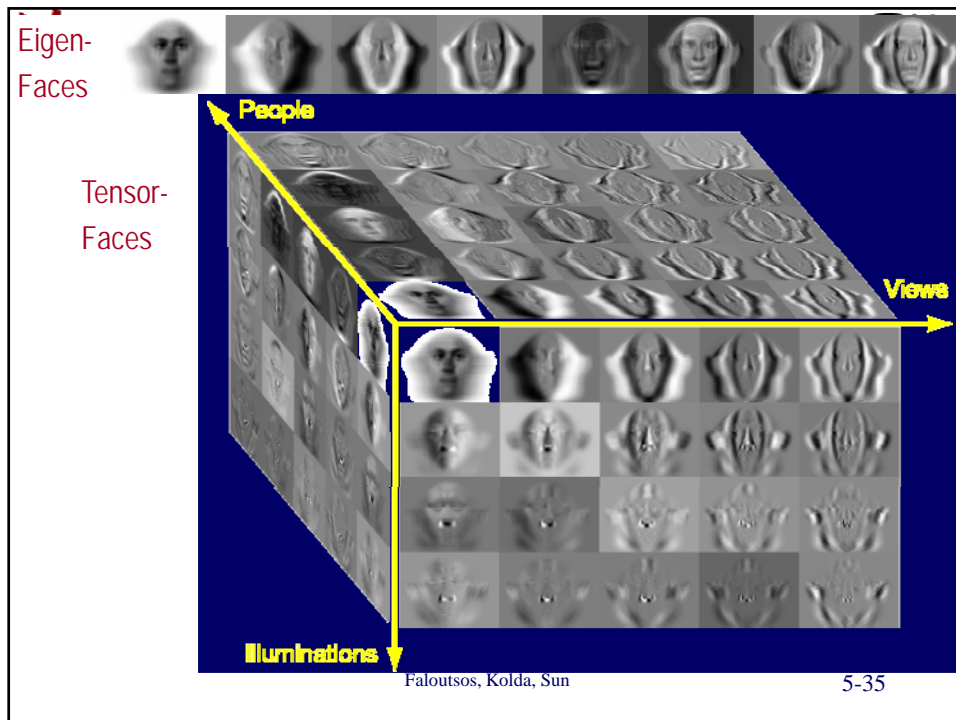
---

**CMU SCS**

# Data Organization

- Linear/PCA: **Data Matrix**

  - $\mathbb{R}^{\text{pixels x images}}$

  - a matrix of image vectors

- Multilinear: *Data Tensor $\mathcal{D}$*

  - $\mathbb{R}^{\text{people x views x illums x express x pixels}}$

  - N-dimensional matrix

  - 28 people, 45 images/person

  - 5 views, 3 illuminations,
    3 expressions per person

**D**

Pixels

Images

Expressions

$\mathcal{D}$

People

Views

Illuminations

$\mathbf{i}_{p,vp,il,ex}$

Faloutsos, Kolda, Sun

5-34

Eigen-
Faces

Tensor-
Faces

Faloutsos, Kolda, Sun                     5-35



CMU SCS

Sandia National Laboratories

# Strategic Data Compression = Perceptual Quality

- TensorFaces data reduction in illumination space primarily degrades illumination effects (cast shadows, highlights)

- PCA has *lower mean square error* but *higher perceptual error*

|  | TensorFaces | TensorFaces | PCA |
|---|---|---|---|
| Original |  | Mean Sq. Err. = 409.15 | Mean Sq. Err. = 85.75 |
|  | 6 illum + 11 people param. | 3 illum + 11 people param. | 33 parameters |
| 176 basis vectors | 66 basis vectors | 33 basis vectors | 33 basis vectors |

Faloutsos, Kolda, Sun                     5-36

**CMU SCS**                                             Sandia National Laboratories

## TensorFaces: An Application of the Tucker Decomposition

M.A.O. Vasilescu & D. Terzopoulos, CVPR'03

- Example: 7942 pixels x 16 illuminations x 11 subjects
- PCA (eigenfaces): SVD of 7942 x 176 matrix
- Tensorfaces: Tucker decomposition of 7942 x 16 x 11 tensor

7942 x 33      176 x 33

$$\mathbf{X} \approx \mathbf{E} \times_2 \mathbf{V}$$

eigenfaces      loadings

An image is represented by a linear combination of 33 eigenfaces.

7942 x 3 x 11      16 x 3      11 x 11

$$\mathcal{X} \approx \mathcal{T} \times_2 \mathbf{U}_{illum} \times_3 \mathbf{U}_{person}$$

tensorfaces      illumination      subjects

An image is represented by a multilinear combination of 33 tensorfaces using the outer product (or Kronecker product) of a length-3 illumination vector and a length-11 person vector.

|  | Original | PCA 11 Eigenfaces | TensorFaces 11 TensorFaces | PCA 22 Eigenfaces | TensorFaces 22 TensorFaces |
|---|---|---|---|---|---|
| RMSE: | | 14.62 | 33.47 | 9.26 | 20.22 |

5-37

---

**CMU SCS**                                             Sandia National Laboratories

## Summary

| Methods | Pros | Cons | Applications |
|---|---|---|---|
| SVD, PCA | Optimal in L2 and Frobenius | Dense representation, Negative entries | LSI, PageRank, HITS |
| CUR, CMD | Interpretability, sparse bases | Not optimal like SVD, dense core | DNA SNP data, network forensics |
| Co-clustering | Interpretability | Local minimum | Social networks, microarray data |
| Tucker | Flexible representation | Interpretability, non-uniqueness, dense core | TensorFaces |
| PARAFAC | Interpretability, efficient parse computation | Slow convergence | TOPHISTS |
| Incrementalization | Efficiency | Non-optimal | Tensor Streams |
| Nonnegativity | Interpretability, sparse results | Local minimum, non-uniqueness | Image segmentation |

**CMU SCS**

Sandia National Laboratories

# Conclusion

- Real data are often in high dimensions with multiple aspects (modes)
- Matrix and tensor provide elegant theory and algorithms for such data
- However, many problems are still open
    - skew distribution, anomaly detection, streaming algorithm, distributed/parallel algorithms, efficient out-of-core processing

Faloutsos, Kolda, Sun

5-39

**CMU SCS**

Sandia National Laboratories

# Thank you!

- **Christos Faloutsos**
  www.cs.cmu.edu/~christos

- Tamara Kolda
  csmr.ca.sandia.gov/~tgkolda

- Jimeng Sun
  www.cs.cmu.edu/~jimeng

Faloutsos, Kolda, Sun

5-40