



Roadmap

- Motivation
 - **Matrix tools**
 - Tensor tools
 - Case studies
- SVD, PCA
 - HITS, PageRank
 - CUR
 - Co-clustering





Examples of Matrices

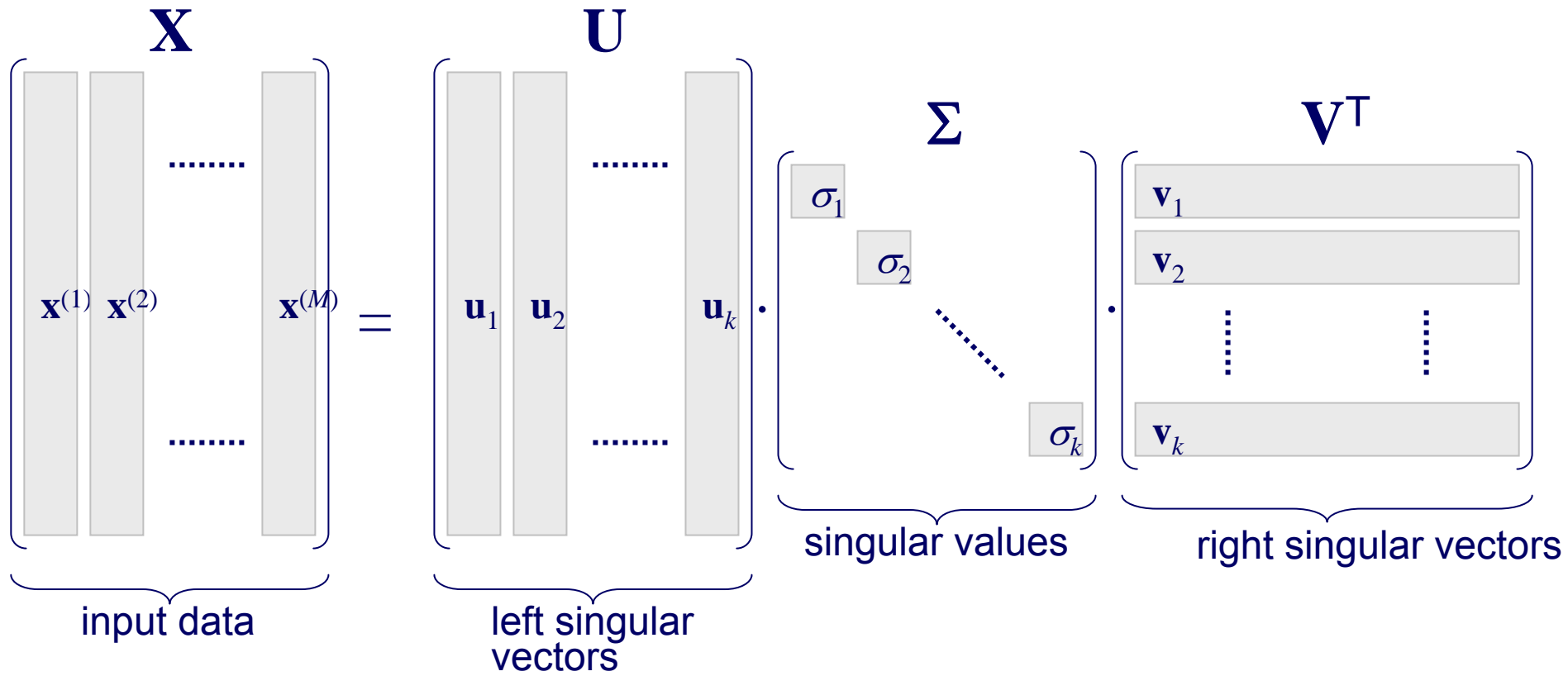
- Example/Intuition: Documents and terms
- Find patterns, groups, concepts

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



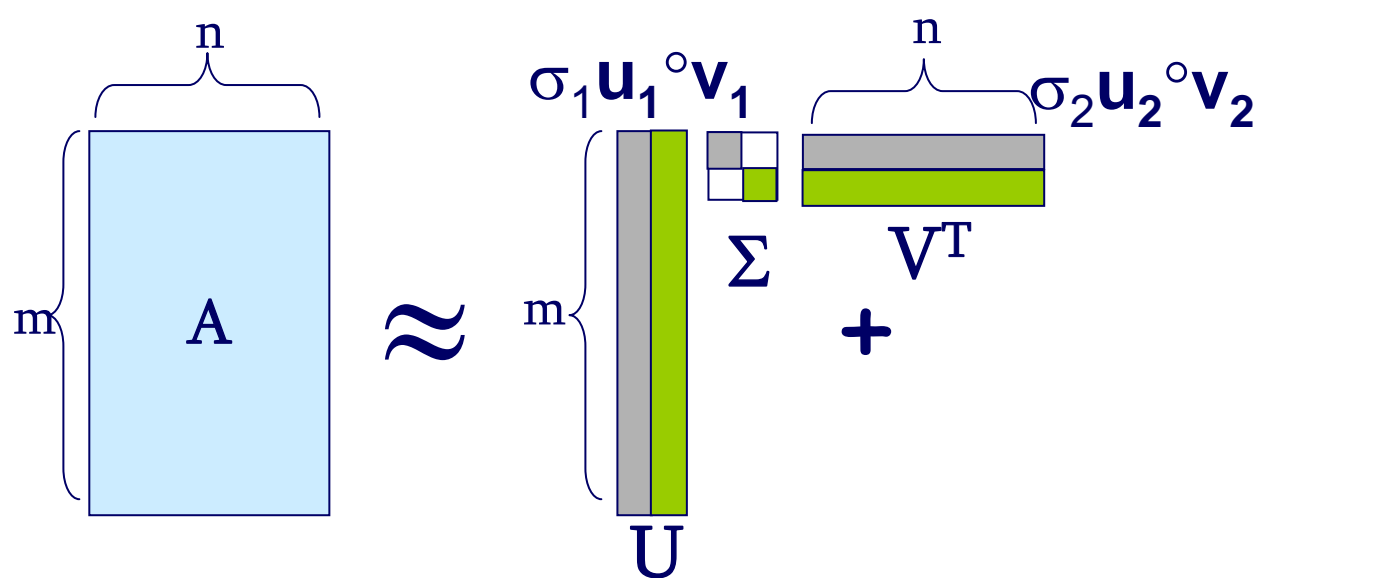
Singular Value Decomposition (SVD)

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





SVD as spectral decomposition

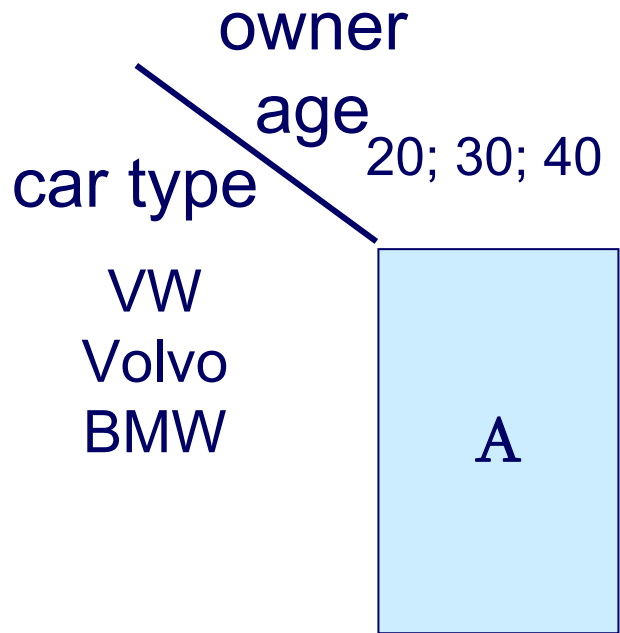
$$A \approx U \Sigma V^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i$$


The diagram illustrates the SVD decomposition of matrix A (size $m \times n$) into three matrices: U (size $m \times m$), Σ (size $m \times n$), and V^T (size $n \times n$). The matrix U is shown as a vertical rectangle with a green column and a grey column. The matrix Σ is shown as a small square with a green and grey block. The matrix V^T is shown as a horizontal rectangle with a green row and a grey row. The approximation is also shown as a sum of rank-1 matrices: $\sigma_1 \mathbf{u}_1 \circ \mathbf{v}_1 + \sigma_2 \mathbf{u}_2 \circ \mathbf{v}_2$.

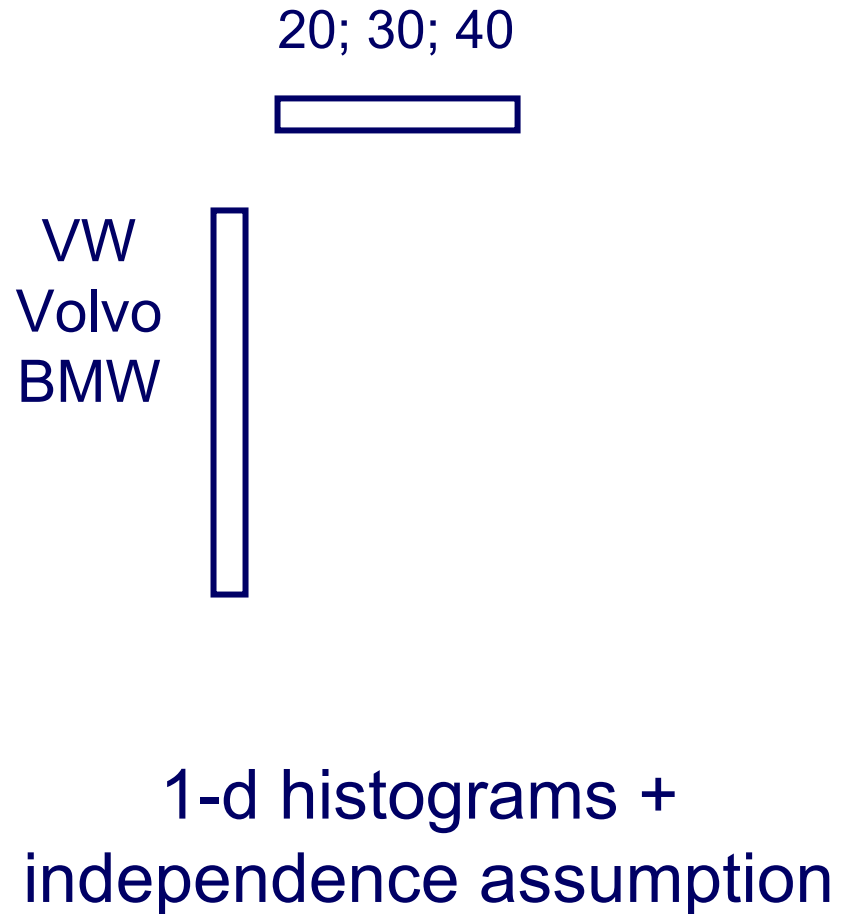
- Best rank- k approximation in L2 and Frobenius
- SVD only works for static matrices (a single 2nd order tensor)



Vector outer product – intuition:



2-d histogram





SVD - Example

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

retrieval

inf. ↓

data brain lung

$$\begin{array}{c} \uparrow \\ \text{CS} \\ \downarrow \\ \uparrow \\ \text{MD} \\ \downarrow \end{array}
 \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}
 =
 \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}
 \times
 \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}
 \times
 \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$



SVD - Example

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

retrieval
 inf. ↓ brain lung
 data

CS-concept
 MD-concept

$$\begin{array}{c} \uparrow \\ \text{CS} \\ \downarrow \\ \uparrow \\ \text{MD} \\ \downarrow \end{array}
 \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}
 =
 \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}
 \times
 \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}
 \times
 \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$



SVD - Example

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

doc-to-concept
similarity matrix

retrieval CS-concept
inf. ↓ brain lung MD-concept

data

↑ CS
↓
↑ MD
↓

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

The value 0.18 in the first row, first column of the middle matrix is circled in red. An arrow points from the label 'CS-concept' to this value.



SVD - Example

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

retrieval
inf. ↓ brain lung

data

‘strength’ of CS-concept

↑

CS

↓

↑

MD

↓

1	1	1	0	0	=	0.18	0	x	9.64	0	x	0.58	0.58	0.58	0	0		
2	2	2	0	0		0.36	0		0	5.29		0	0	0	0	0.71	0.71	
1	1	1	0	0		0.18	0		0	0		0	0	0	0	0	0	
5	5	5	0	0		0.90	0		0	0.53		0	0	0	0	0	0	0
0	0	0	2	2		0	0.80		0	0.27		0	0	0	0	0	0	0
0	0	0	3	3		0	0.27		0	0.27		0	0	0	0	0	0	0
0	0	0	1	1		0	0.27		0	0.27		0	0	0	0	0	0	0



SVD - Example

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

term-to-concept
similarity matrix

retrieval
inf. ↓
data brain lung

CS

MD

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

CS-concept



SVD - Example

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

term-to-concept
similarity matrix

↑ CS
↓
↑ MD
↓

data	inf.	retrieval	brain	lung					
	↓								
1	1	1	0	0	=	0.18	0		
2	2	2	0	0		0.36	0		
1	1	1	0	0		0.18	0		
5	5	5	0	0		0.90	0		
0	0	0	2	2		0	0.53		
0	0	0	3	3		0	0.80		
0	0	0	1	1		0	0.27		

CS-concept

\mathbf{X}	9.64	0		
	0	5.29		

\mathbf{X}

0.58	0.58	0.58	0	0
0	0	0	0.71	0.71

→



SVD - Interpretation

‘documents’, ‘terms’ and ‘concepts’:

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

A: term-to-term ($[m \times m]$) similarity matrix

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

A: document-to-document ($[n \times n]$) similarity matrix



SVD properties

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

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

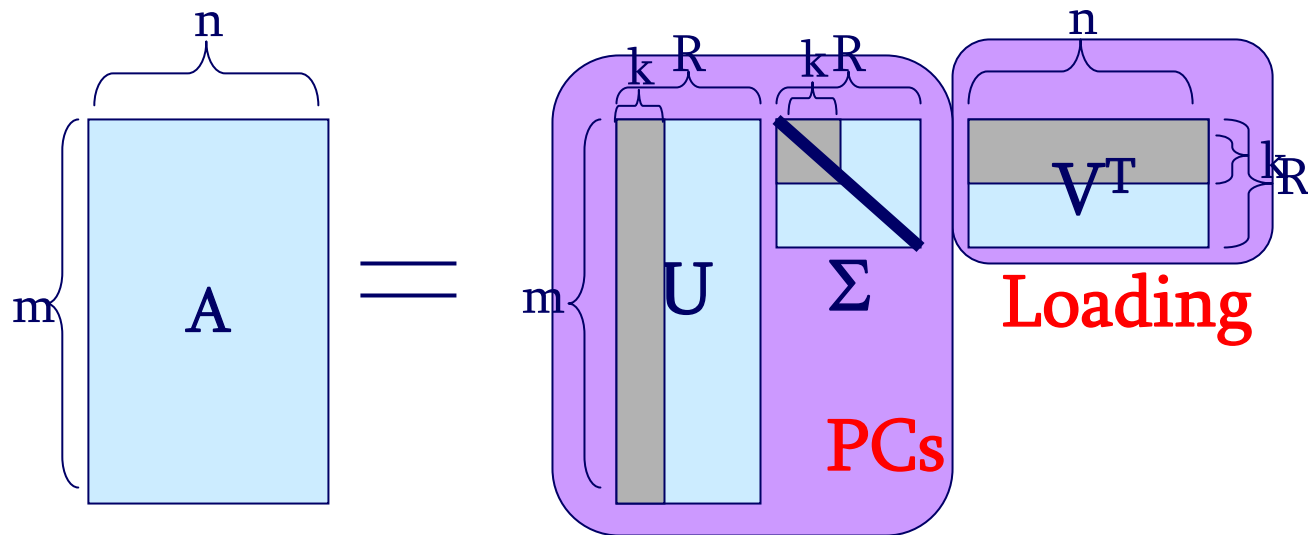
Further reading:

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



Principal Component Analysis (PCA)

- SVD $A = U\Sigma V^T$

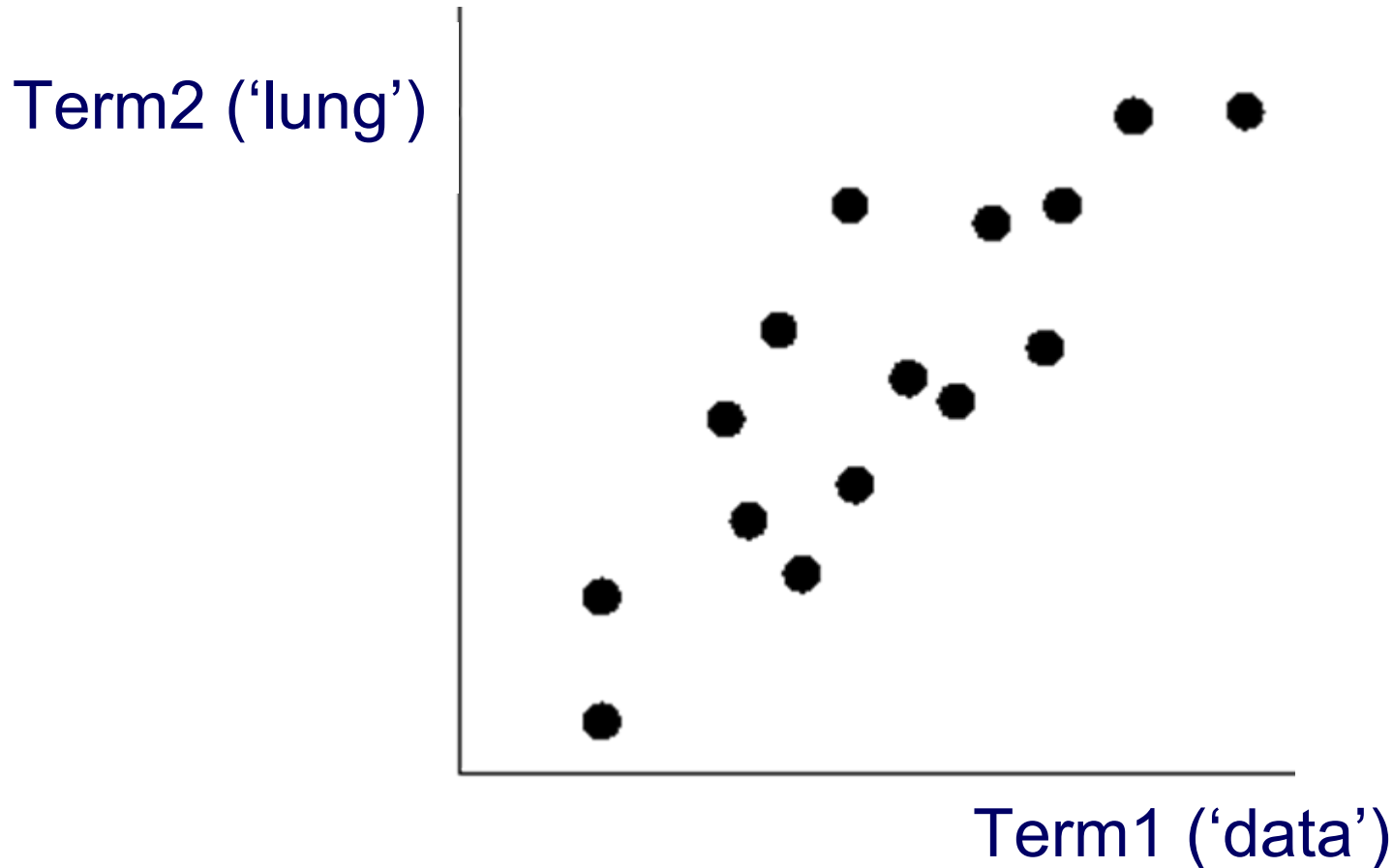


- PCA is an important application of SVD
- Note that U and V are dense and may have negative entries



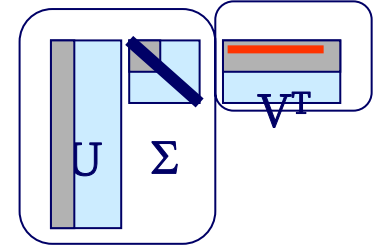
PCA interpretation

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



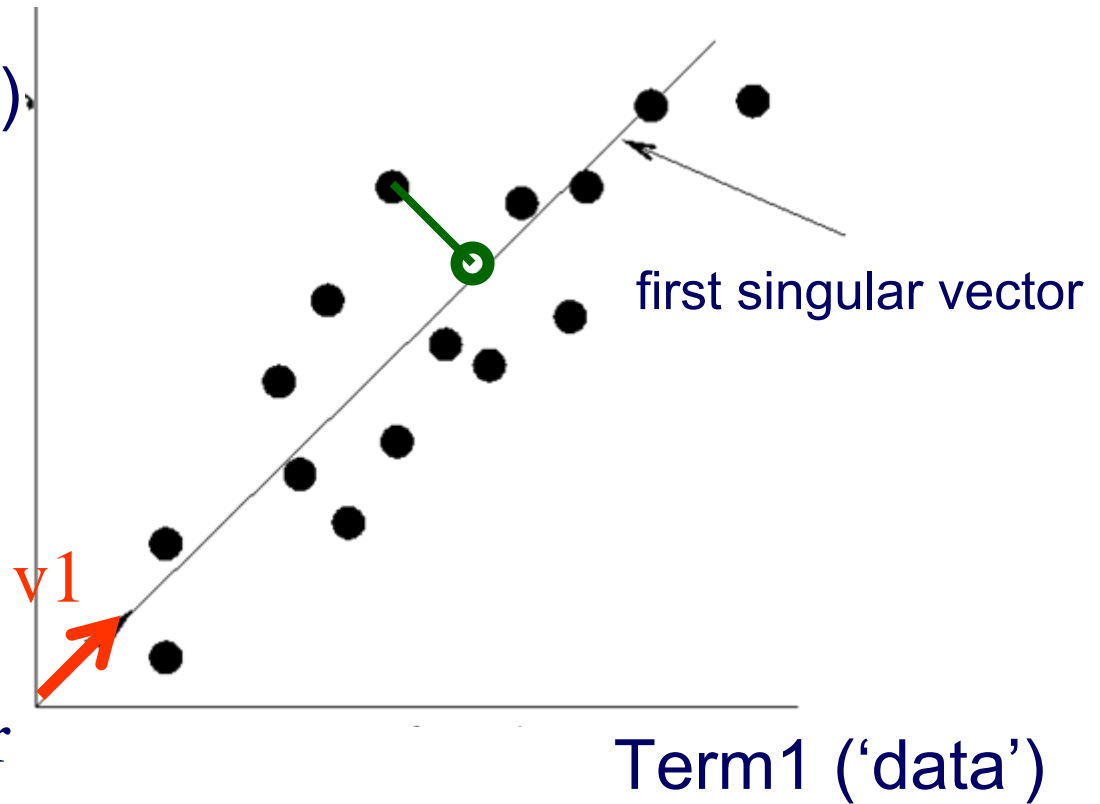


PCA - interpretation



Term2 ('retrieval')

PCA projects points
Onto the "best" axis



- minimum RMS error



Roadmap

- Motivation
 - **Matrix tools**
 - Tensor tools
 - Case studies
- SVD, PCA
 - **HITS, PageRank**
 - CUR
 - Co-clustering





Kleinberg's algorithm HITS

- Problem defn: 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

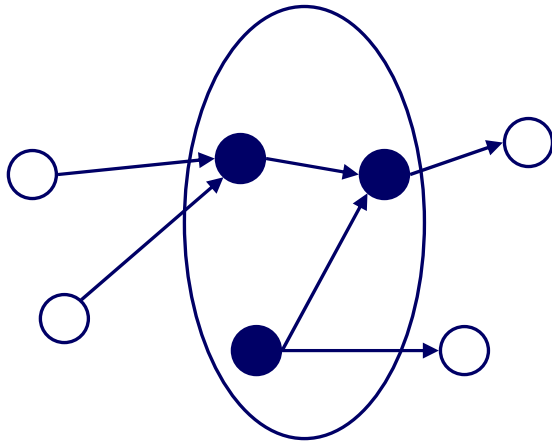
Further reading:

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



Kleinberg's algorithm HITS

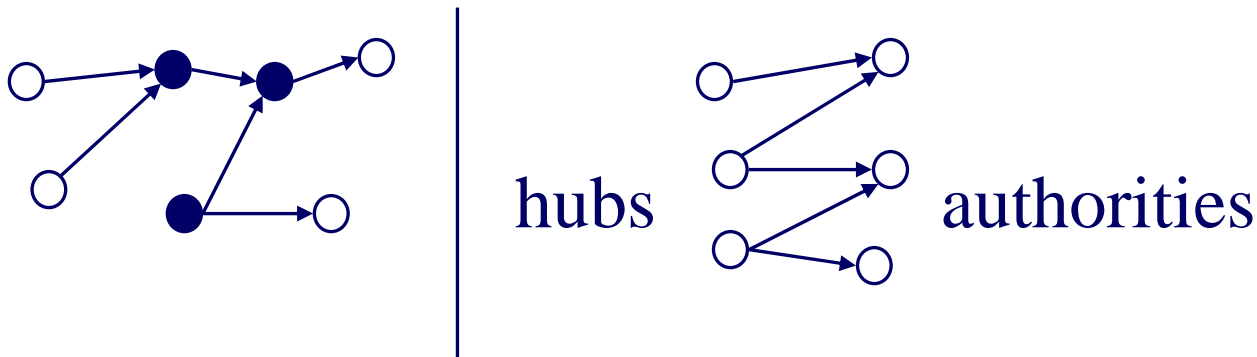
- Step 1: expand by one move forward and backward





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'





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



Kleinberg's algorithm: HITS

Let \mathbf{A} be the adjacency matrix:

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

Let \mathbf{h} and \mathbf{a} be $[n \times 1]$ vectors with the
'hubness' and 'authoritativeness' scores.

Then:



Kleinberg's algorithm: HITS

Then:

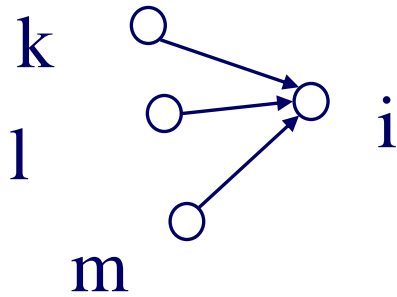
$$a_i = h_k + h_l + h_m$$

that is

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

or

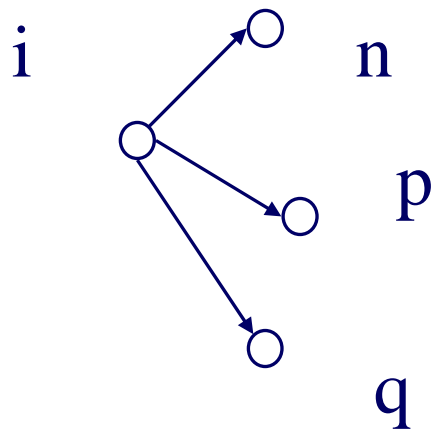
$$\mathbf{a} = \mathbf{A}^T \mathbf{h}$$





Kleinberg's algorithm: HITS

symmetrically, for the 'hubness':



$$h_i = a_n + a_p + a_q$$

that is

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

or

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



Kleinberg's algorithm: HITS

In conclusion, we want vectors \mathbf{h} and \mathbf{a} such that:

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

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

That is:

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



Kleinberg's algorithm: HITS

\mathbf{a} is a right singular vector of the adjacency matrix \mathbf{A} (by defn!), a.k.a the eigenvector of $\mathbf{A}^T \mathbf{A}$

Starting from random \mathbf{a}' and iterating, we'll eventually converge

Q: to which of all the eigenvectors? why?

A: to the one of the strongest eigenvalue,

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



Kleinberg's algorithm - discussion

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



Roadmap

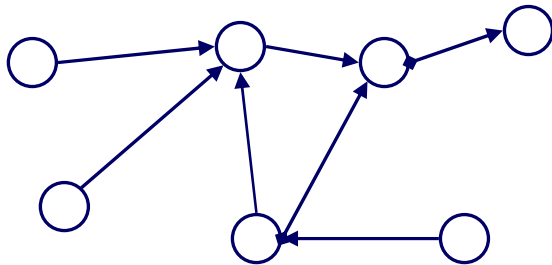
- Motivation
 - **Matrix tools**
 - Tensor tools
 - Case studies
- SVD, PCA
 - HITS, **PageRank**
 - CUR
 - Co-clustering





Motivating problem: PageRank

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



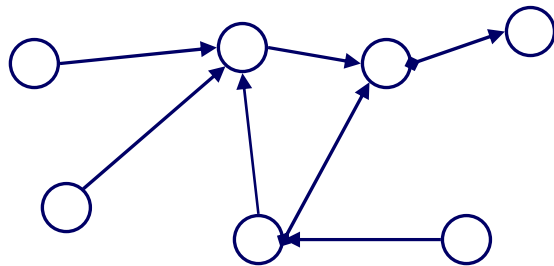
A node is important, if it is connected with important nodes (recursive, but OK!)



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))

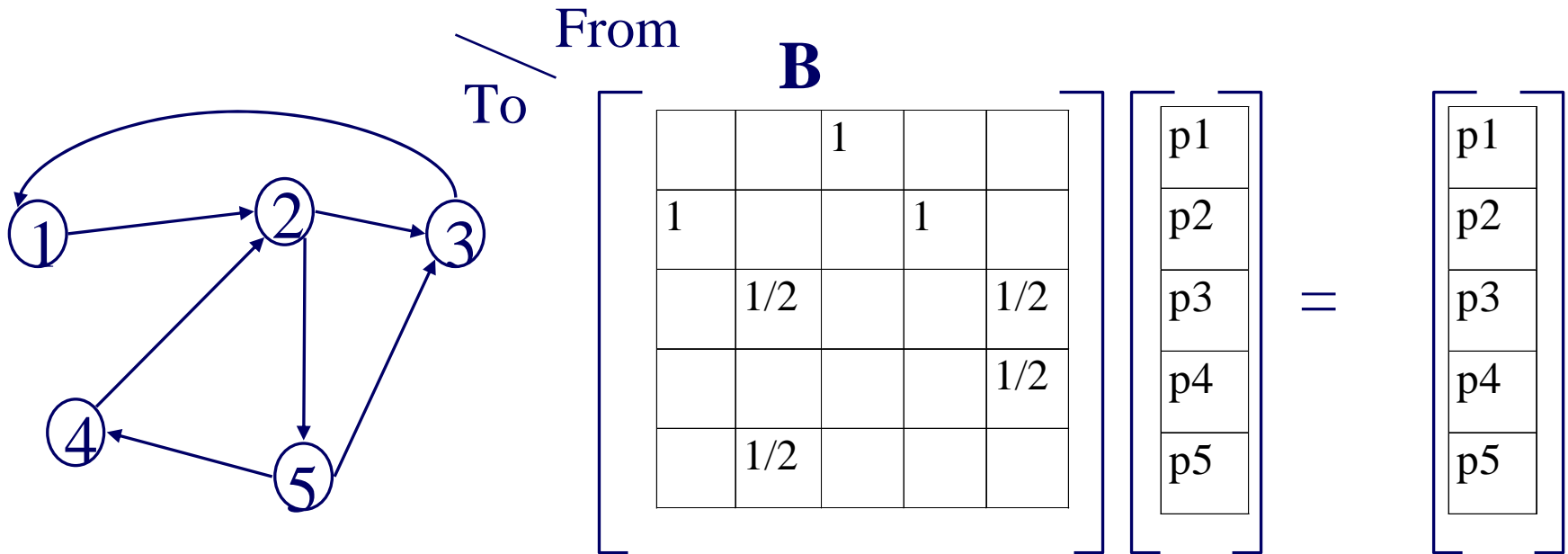


A node has high **ssp**, if it is connected with **high ssp** nodes (recursive, but OK!)



(Simplified) PageRank algorithm

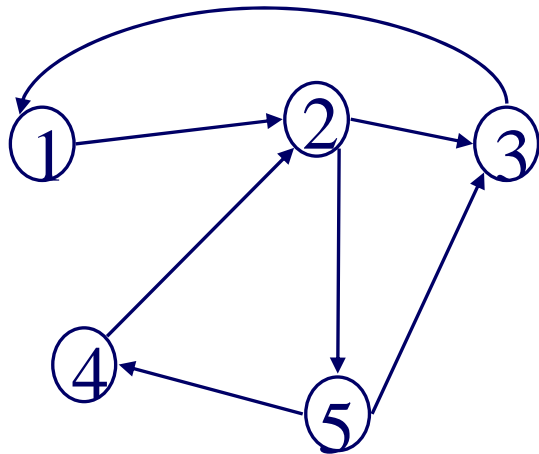
- Let \mathbf{A} be the transition matrix (= adjacency matrix); let \mathbf{B} be the transpose, column-normalized - then





(Simplified) PageRank algorithm

- $B p = p$



$$B p = p$$

		1		
1			1	
	1/2			1/2
				1/2
	1/2			

p1
p2
p3
p4
p5

$$=$$

p1
p2
p3
p4
p5



(Simplified) PageRank algorithm

- $\mathbf{B} \mathbf{p} = \mathbf{1} * \mathbf{p}$
- thus, \mathbf{p} is the **eigenvector** that corresponds to the highest eigenvalue (=1, since the matrix is column-normalized)
- Why does such a \mathbf{p} exist?
 - \mathbf{p} exists if \mathbf{B} is $n \times n$, nonnegative, irreducible [Perron–Frobenius theorem]



(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

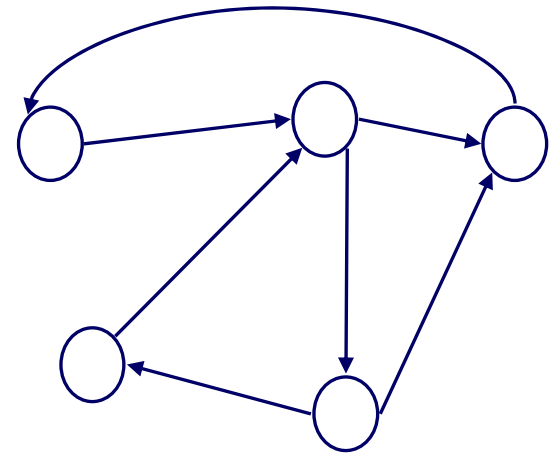
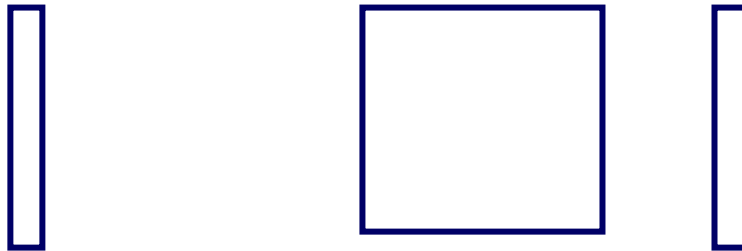


Full Algorithm

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

$$\mathbf{p} = c \mathbf{B} \mathbf{p} + (1-c)/n \mathbf{1} \Rightarrow$$

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





Roadmap

- Motivation
 - **Matrix tools**
 - Tensor tools
 - Case studies
- SVD, PCA
 - HITS, PageRank
 - **CUR**
 - Co-clustering





Motivation of CUR or CMD

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

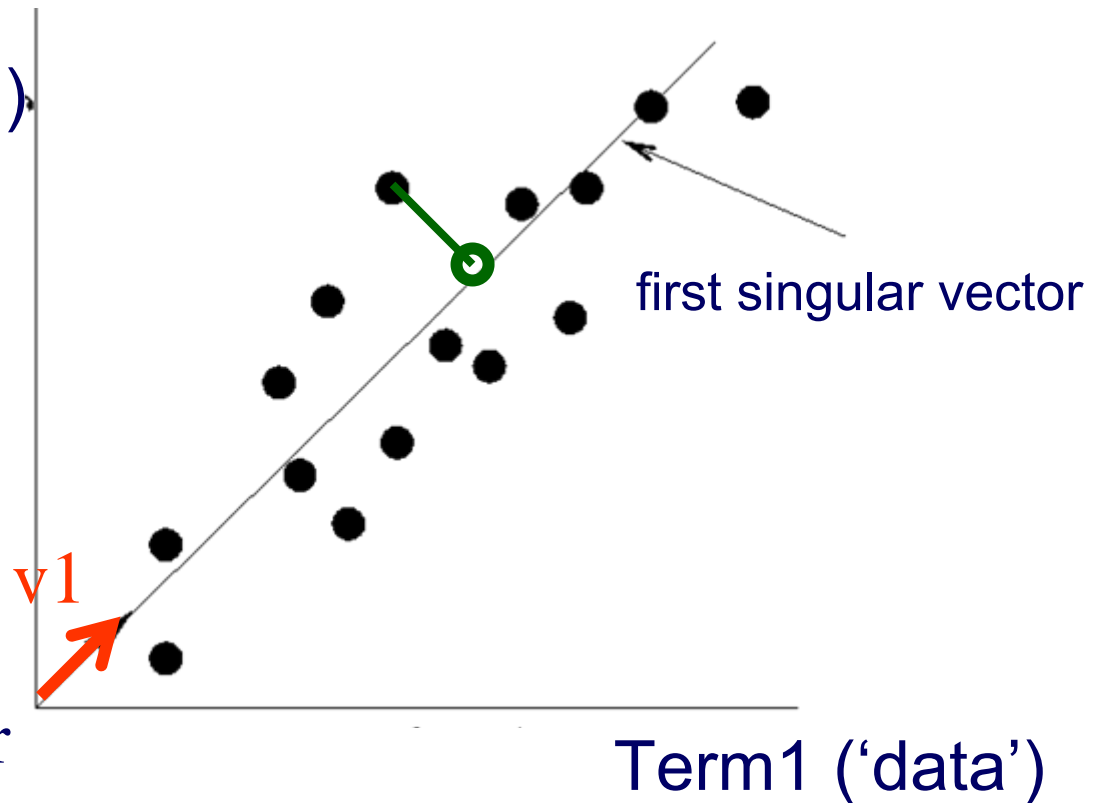


PCA - interpretation

Term2 ('retrieval')

PCA projects points
Onto the "best" axis

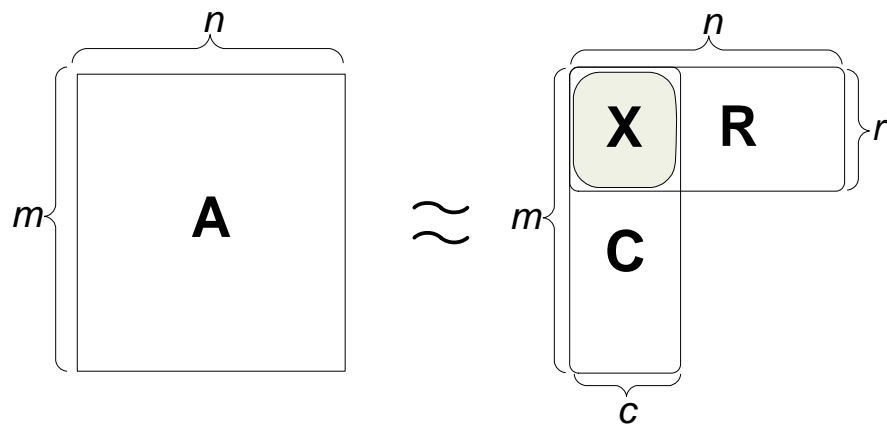
- minimum RMS error



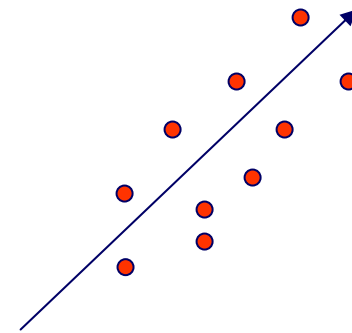


CUR

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



U is the pseudo-inverse of X

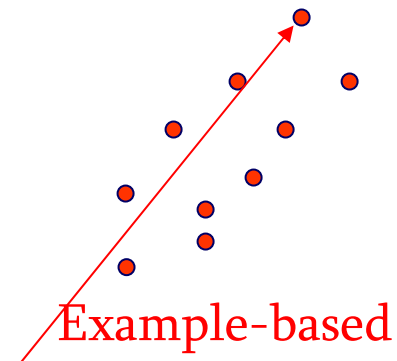
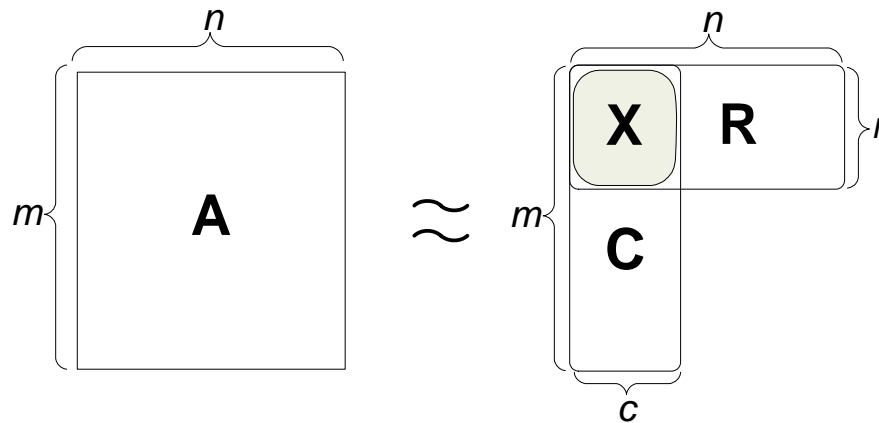


Orthogonal projection



CUR

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



U is the pseudo-inverse of X :

$$U = X^\dagger = (U^T U)^{-1} U^T$$



CUR (cont.)

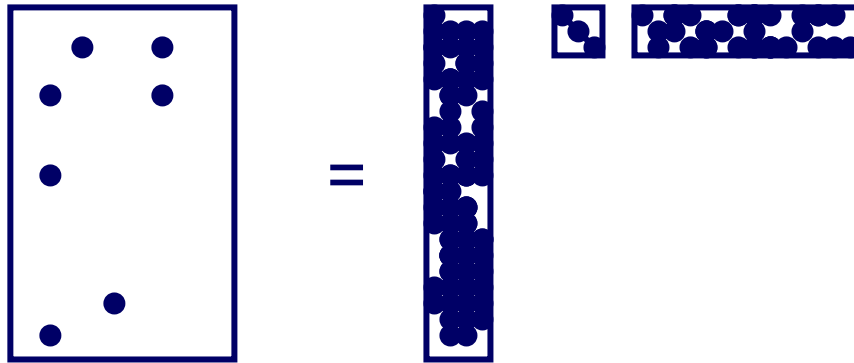
- Key question:
 - How to select/sample the columns and rows?
- Uniform sampling
- Biased sampling
 - CUR w/ absolute error bound
 - CUR w/ relative error bound

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.

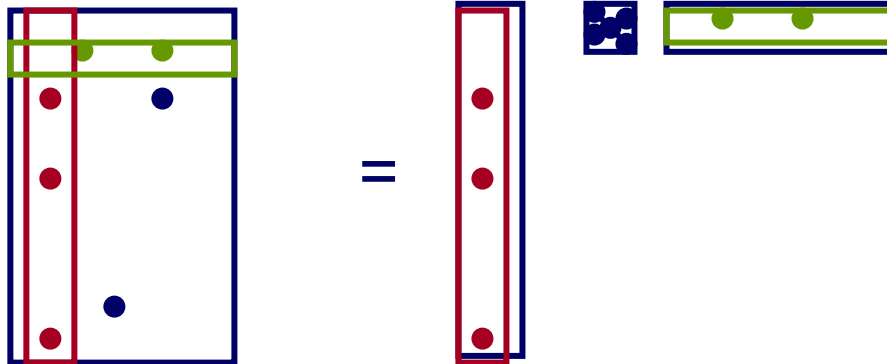


The sparsity property – pictorially:



SVD/PCA:
Destroys sparsity

$$U \Sigma V^T$$




CUR: maintains sparsity

$$C U R$$



The sparsity property




sparse and small

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

Big but sparse Big and dense

Detailed description: This diagram illustrates the SVD decomposition. A sad face emoji is in the top left. The equation $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ is centered. An arrow points from the text 'Big but sparse' to matrix \mathbf{A} . Another arrow points from the text 'Big and dense' to matrix \mathbf{U} . A third arrow points from the text 'Big and dense' to matrix \mathbf{V}^T . A fourth arrow points from the text 'sparse and small' to matrix $\mathbf{\Sigma}$.



dense but small

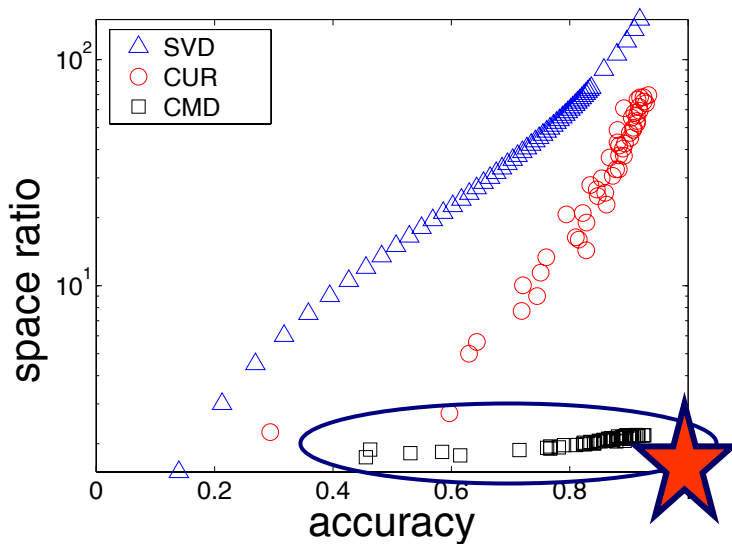
$$\text{CUR: } \mathbf{A} = \mathbf{C} \mathbf{U} \mathbf{R}$$

Big but sparse Big but sparse

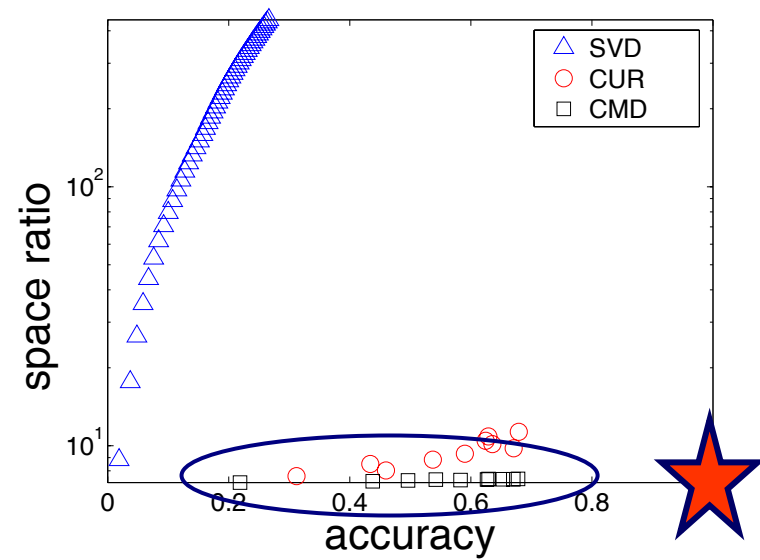
Detailed description: This diagram illustrates the CUR decomposition. A happy face emoji is in the top left. The equation $\mathbf{A} = \mathbf{C} \mathbf{U} \mathbf{R}$ is centered. An arrow points from the text 'Big but sparse' to matrix \mathbf{A} . Another arrow points from the text 'Big but sparse' to matrix \mathbf{C} . A third arrow points from the text 'Big but sparse' to matrix \mathbf{R} . A fourth arrow points from the text 'dense but small' to matrix \mathbf{U} .



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

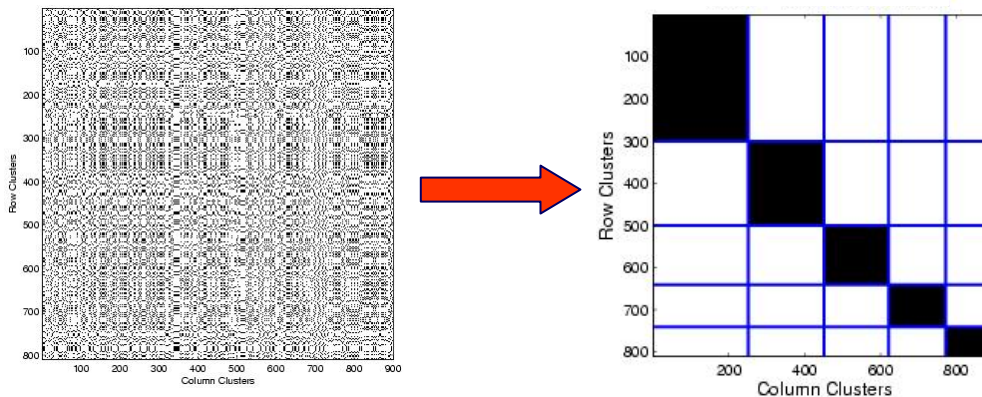
- Motivation
 - **Matrix tools**
 - Tensor tools
 - Case studies
- SVD, PCA
 - HITS, PageRank
 - CUR
 - **Co-clustering**





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





Co-clustering

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

Reference:

1. Dhillon et al. Information-Theoretic Co-clustering, KDD'03



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

$$\begin{bmatrix} .05 & .05 & .05 & 0 & 0 & 0 \\ .05 & .05 & .05 & 0 & 0 & 0 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ 0 & 0 & 0 & .05 & .05 & .05 \\ .04 & .04 & 0 & .04 & .04 & .04 \\ .04 & .04 & .04 & 0 & .04 & .04 \end{bmatrix}$$

med. terms

cs terms

common terms

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

term x term-group



Co-clustering

Observations

- uses KL divergence, instead of L2
- the middle matrix is **not** diagonal
 - we'll see that again in the Tucker tensor decomposition



Problem with Information Theoretic Co-clustering

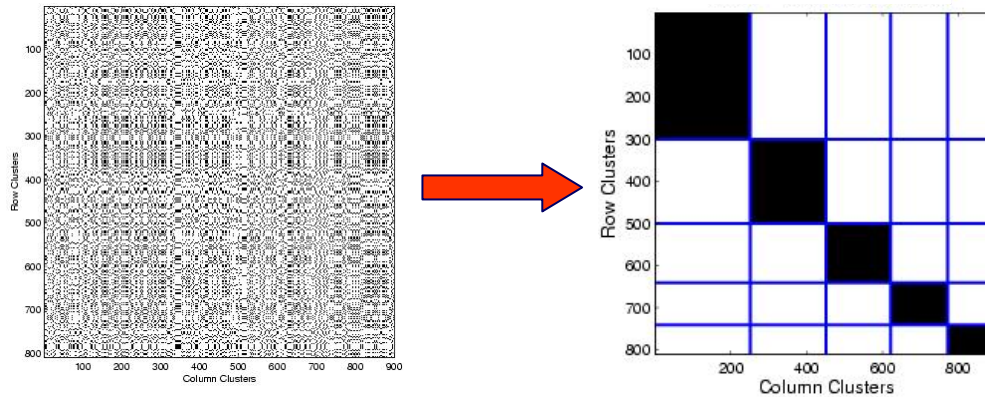
- Number of row and column groups must be specified

Desiderata:

- ✓ **Simultaneously discover** row and column groups
- ✗ **Fully Automatic:** No “magic numbers”
- ✓ **Scalable** to large graphs



Cross-association



Desiderata:

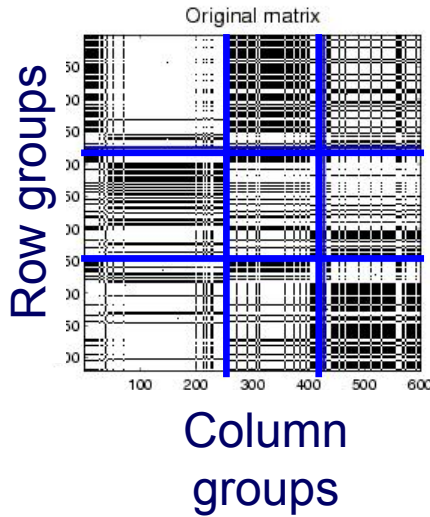
- ✓ **Simultaneously discover** row and column groups
- ✓ **Fully Automatic:** No “magic numbers”
- ✓ **Scalable** to large matrices

Reference:

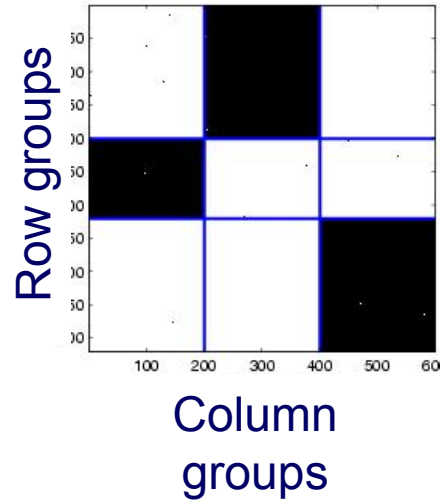
1. Chakrabarti et al. Fully Automatic Cross-Associations, KDD'04



What makes a cross-association “good”?



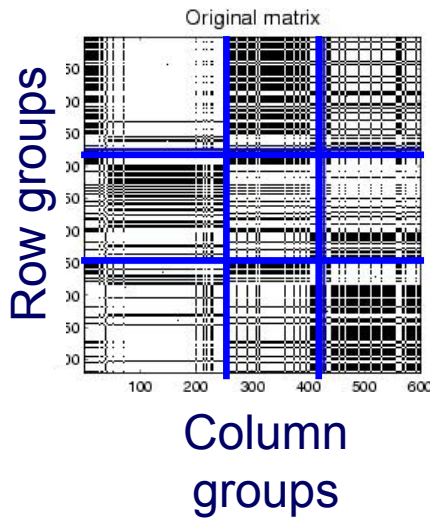
versus



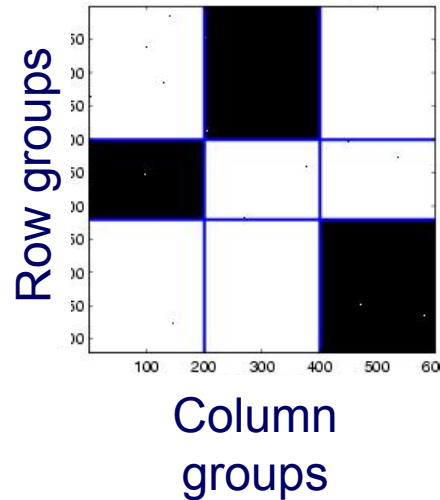
Why is this better?



What makes a cross-association “good”?



versus

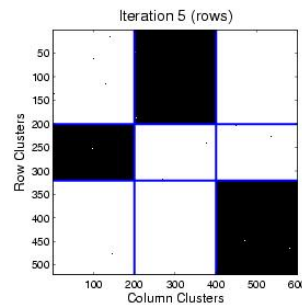
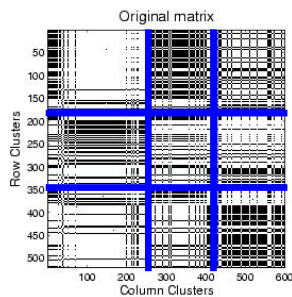


Why is this better?

simpler; easier to describe
easier to compress!



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



Main Idea



$$\text{Total Encoding Cost} = \underbrace{\sum_i \text{size}_i * H(x_i)}_{\text{Code Cost}} + \underbrace{\text{Cost of describing cross-associations}}_{\text{Description Cost}}$$

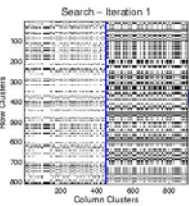
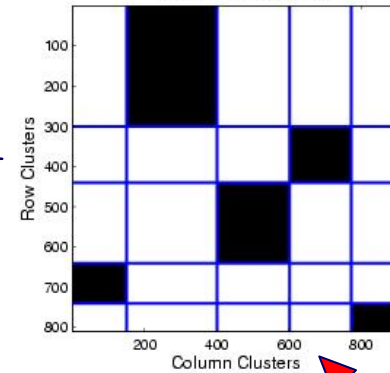
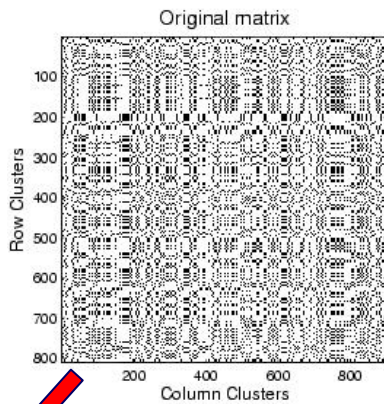
Minimize the total cost (# bits)
for lossless compression



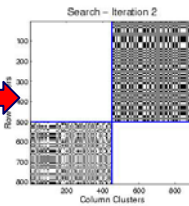
Algorithm

$l = 5$ col groups

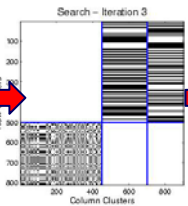
$k = 5$ row groups



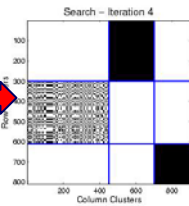
$k=1,$
 $l=2$



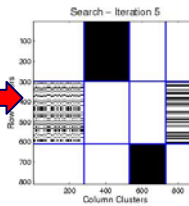
$k=2,$
 $l=2$



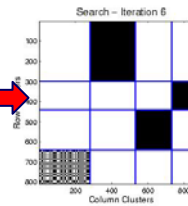
$k=2,$
 $l=3$



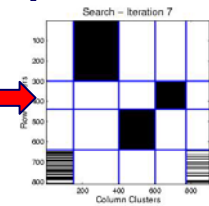
$k=3,$
 $l=3$



$k=3,$
 $l=4$



$k=4,$
 $l=4$



$k=4,$
 $l=5$



Algorithm

Code for cross-associations (matlab):

www.cs.cmu.edu/~deepay/mywww/software/CrossAssociations-01-27-2005.tgz

Variations and extensions:

- ‘Autopart’ [Chakrabarti, PKDD’04]
- www.cs.cmu.edu/~deepay



Matrix tools - summary

- SVD:
 - optimal for L2 – VERY popular (HITS, PageRank, Karhunen-Loeve, Latent Semantic Indexing, PCA, etc etc)
- C-U-R (CMD etc)
 - near-optimal; sparsity; interpretability