

CMU SCS

15-826: Multimedia Databases and Data Mining

Independent Component Analysis (ICA)

C. Faloutsos


15-826
(c) C. Faloutsos and J-Y Pan (2005)
#1


CMU SCS


Citation

- AutoSplit: Fast and Scalable Discovery of Hidden Variables in Stream and Multimedia Databases, **Jia-Yu Pan**, Hiroyuki Kitagawa, Christos Faloutsos and Masafumi Hamamoto

PAKDD 2004, Sydney, Australia




15-826
(c) C. Faloutsos and J-Y Pan (2005)
#2


CMU SCS

Outline

- Motivation
- Related work: PCA, ICA
- Proposed methods
 - Batch-AutoSplit
 - AutoSplit
 - Clustering-AutoSplit
- Experimental results and discussion
- Conclusions

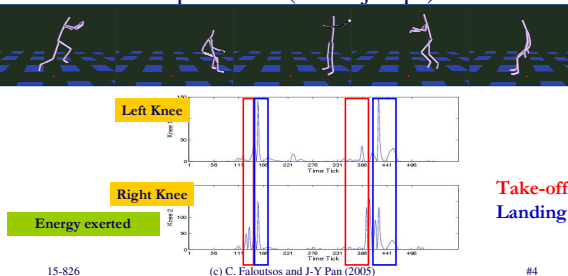
15-826
(c) C. Faloutsos and J-Y Pan (2005)
#3


CMU SCS


Motivation:

(Q1) Find patterns in data

- Motion capture data (broad jumps)



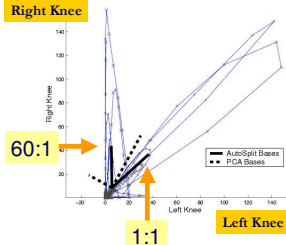
15-826
(c) C. Faloutsos and J-Y Pan (2005)
#4


CMU SCS


Motivation:

(Q1) Find patterns in data

- Human would say
 - Pattern 1: along diagonal
 - Pattern 2: along vertical axis
- How to find these automatically?

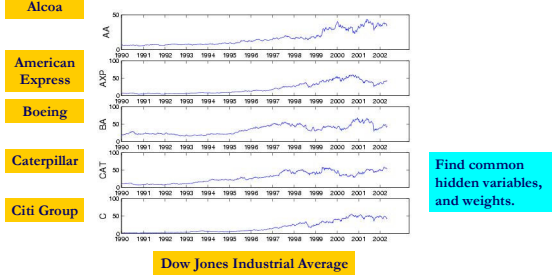


15-826
(c) C. Faloutsos and J-Y Pan (2005)
#5

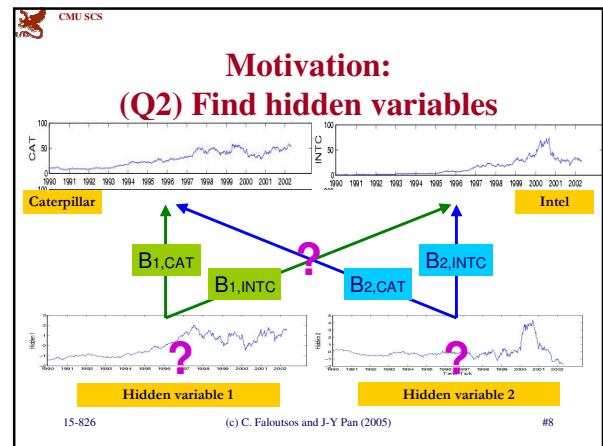
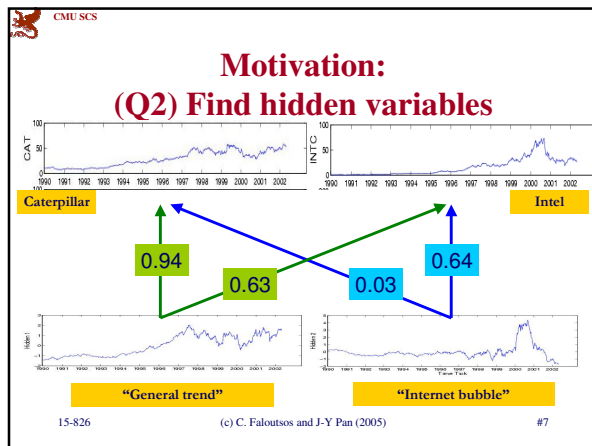

CMU SCS

Motivation:

(Q2) Find hidden variables



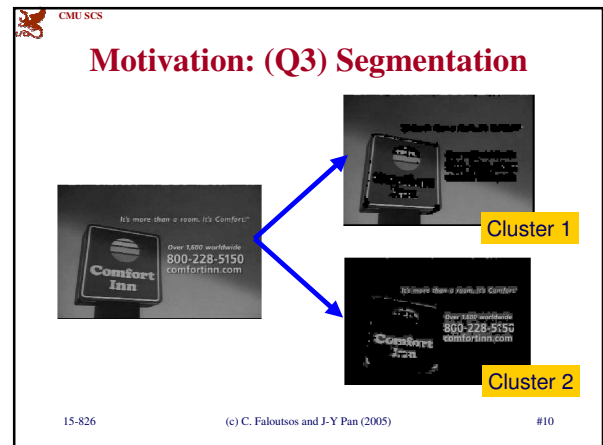
15-826
(c) C. Faloutsos and J-Y Pan (2005)
#6



Motivation:
Find hidden variables

- ICA: also known as ‘Blind Source Separation’
- ‘cocktail party problem’
 - in a party, we can hear two concurrent conversations,
 - but separate them (and tune-in on one of them only)
- http://www.cnl.salk.edu/~tewon/Blind/blind_audio.html
- (in stocks: one ‘discussion’ is the general economy trend; the other ‘discussion’ is the tech-stock boom)

15-826 (c) C. Faloutsos and J-Y Pan (2005) #9



Problem formulation

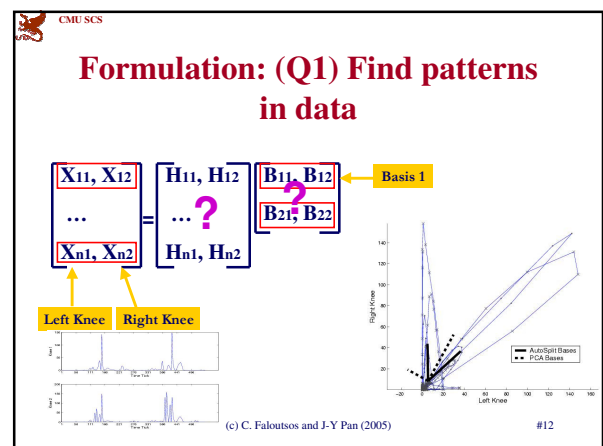
- Given n data items, each has m attributes
- Find the m hidden variables and the m bases

$$\begin{bmatrix} X_{11} & X_{12} & \dots & X_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{nm} \end{bmatrix} = \begin{bmatrix} H_{11} & H_{12} & \dots & H_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ H_{n1} & H_{n2} & \dots & H_{nm} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} & \dots & B_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ B_{m1} & B_{m2} & \dots & B_{mm} \end{bmatrix}$$

$X=HB$

Samples of the m -th hidden variable

15-826 (c) C. Faloutsos and J-Y Pan (2005) #11



CMU SCS

Formulation: (Q2) Find hidden variables

15-826 (c) C. Faloutsos and J-Y Pan (2005) #13

CMU SCS

Outline

- Motivation
- Related work: PCA, ICA
- Proposed methods
 - Batch-AutoSplit
 - AutoSplit
 - Clustering-AutoSplit
- Experimental results and discussion
- Conclusions

15-826 (c) C. Faloutsos and J-Y Pan (2005) #14

CMU SCS

Related Work: PCA

- Goal: Knowing \mathbf{X} , find \mathbf{H} and \mathbf{B} , where $\mathbf{X} = \mathbf{H}\mathbf{B}$
- Problem: Under-constrained

15-826 (c) C. Faloutsos and J-Y Pan (2005) #15

CMU SCS

Related Work: PCA

- PCA says
 - Choose bases/rows in \mathbf{B} which are **orthonormal**, and
 - Find such bases that give **smallest** representation **L2 error** (for dimensional reduction)
- Matrices \mathbf{H} and \mathbf{B} can be solved by
 - **SVD**, neural networks, or many optimization methods

15-826 (c) C. Faloutsos and J-Y Pan (2005) #16

CMU SCS

Related Work: PCA

- Extremely popular
 - Latent Semantic Indexing [Deerwester+90]
 - KL transform [Duda,Hart,Stork00]
 - EigenFace [Turk,Pentlind91]
 - Multiple time series correlation [Guha,Gunopulos,Koudas03]
- **But, there is room for improvement.**

15-826 (c) C. Faloutsos and J-Y Pan (2005) #17

CMU SCS

Related work: ICA

- ICA says
 - Let h_i 's (columns of \mathbf{H}) be mutually **independent**.
- **Many implementations**
 - “Independence = ?”:
 - No mutual information, Kurtosis, nonlinear decorrelation, ‘sparse coding’
 - To solve for \mathbf{H}, \mathbf{B} :
 - Neural networks, optimization methods (gradient ascent, fixed-point, ...)

15-826 (c) C. Faloutsos and J-Y Pan (2005) #18

CMU SCS

Outline

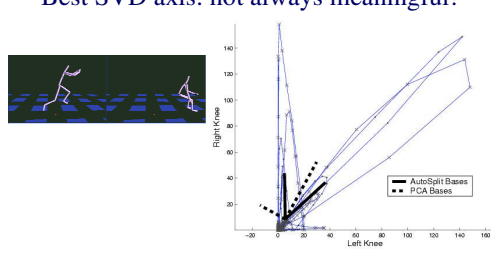
- Motivation
- Related work: PCA, ICA
- ➔ Proposed method: AutoSplit
 - Batch-AutoSplit
 - AutoSplit
 - Clustering-AutoSplit
- Experimental results and discussion
- Conclusions

15-826 (c) C. Faloutsos and J-Y Pan (2005) #19

CMU SCS

PCA sometimes misses essential features

- Best SVD axis: not always meaningful!



15-826 (c) C. Faloutsos and J-Y Pan (2005) #20

CMU SCS

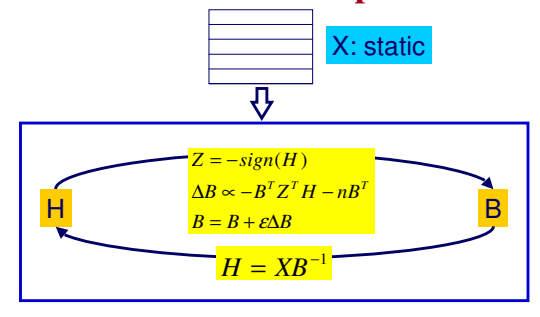
Batch-AutoSplit

- Assume hidden variables to be “independent”
- Intuition: Independent ~ less redundant
- ➔ Better data representation
- How to achieve this?

15-826 (c) C. Faloutsos and J-Y Pan (2005) #21

CMU SCS

Batch-AutoSplit



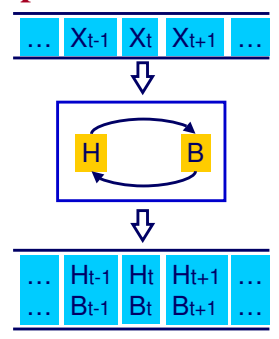
15-826 (c) C. Faloutsos and J-Y Pan (2005) #22

CMU SCS

AutoSplit

skip

- X: streaming
- X_t : t -th window
- Our solution:
 - H, B are estimated iteratively.
 - Adapt to newly arrived data

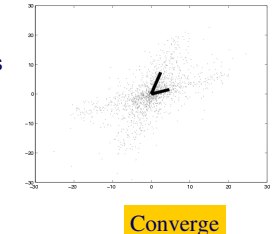


15-826 (c) C. Faloutsos and J-Y Pan (2005) #23

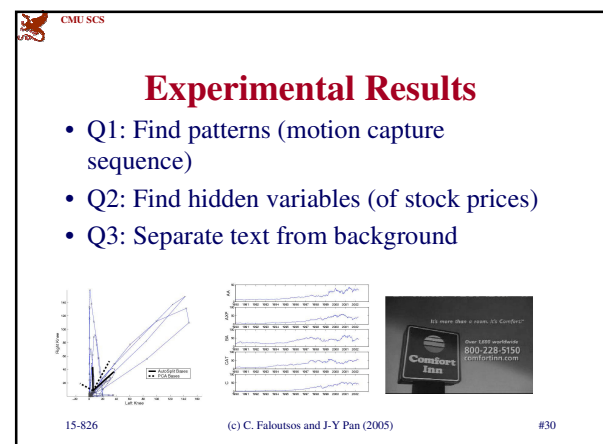
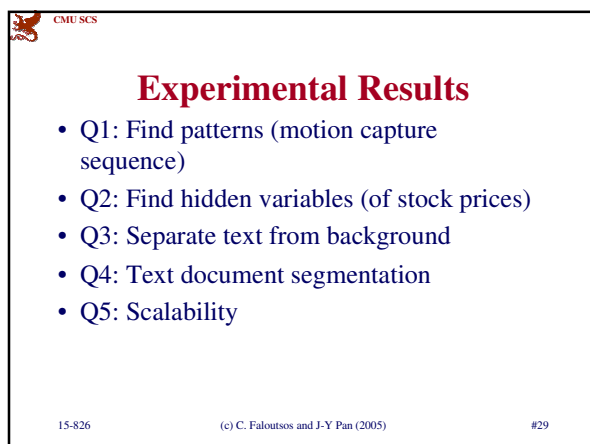
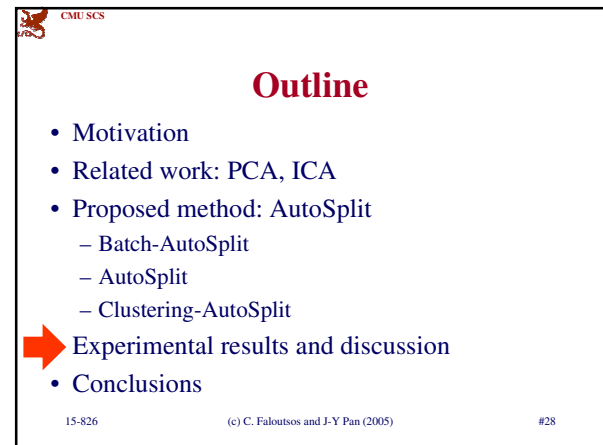
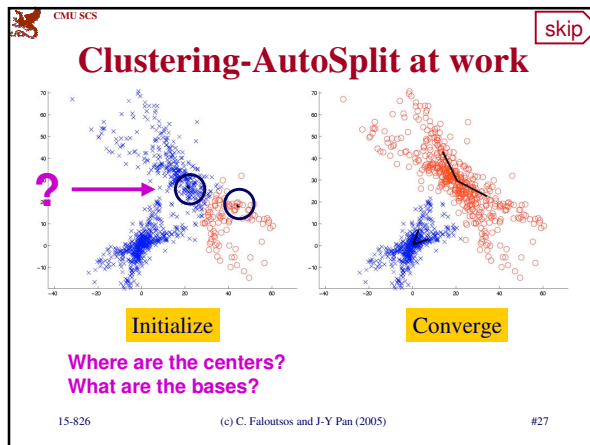
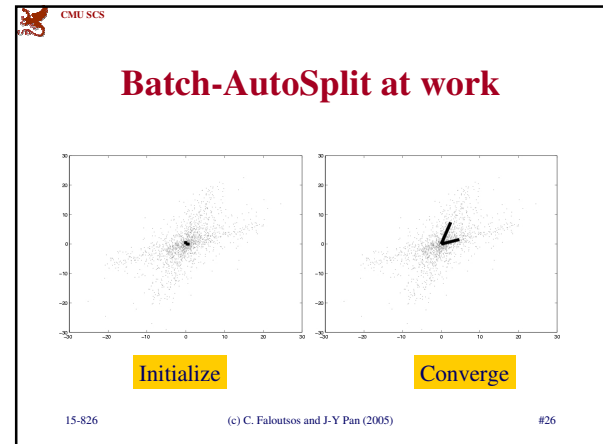
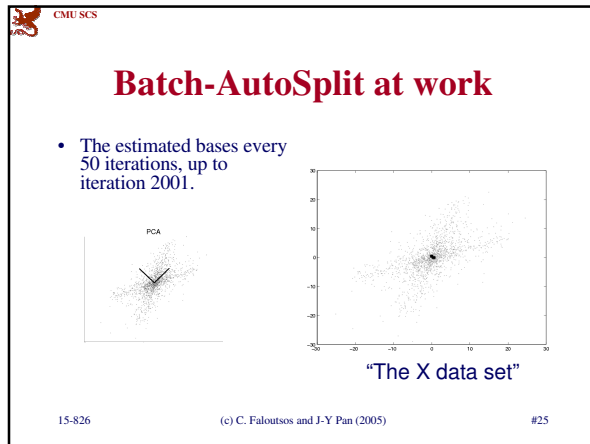
CMU SCS

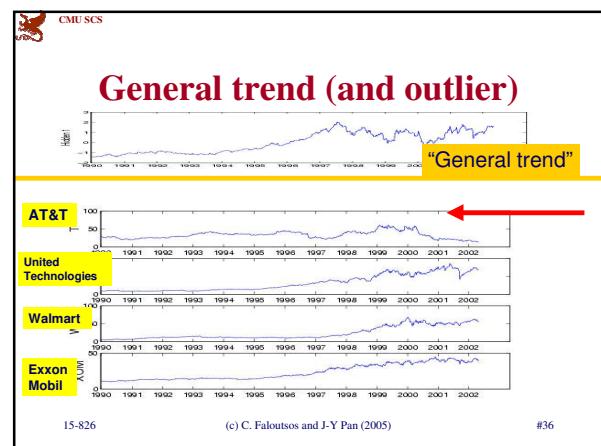
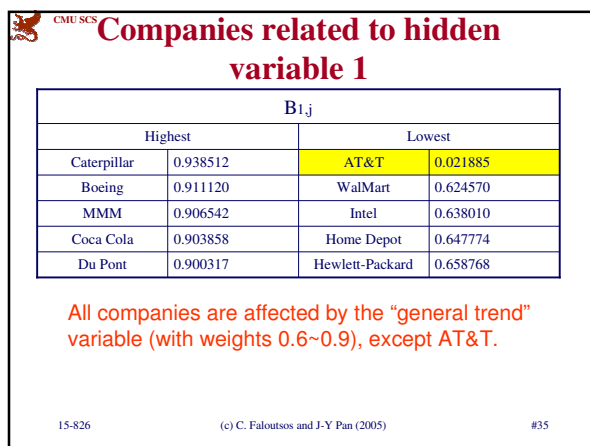
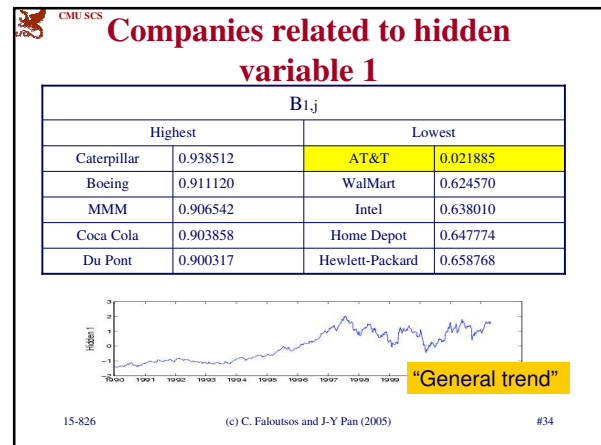
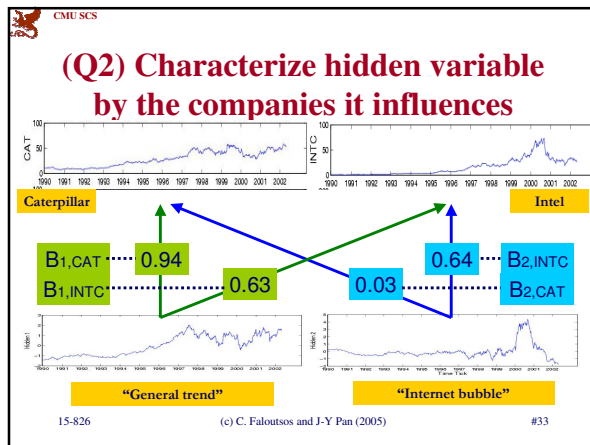
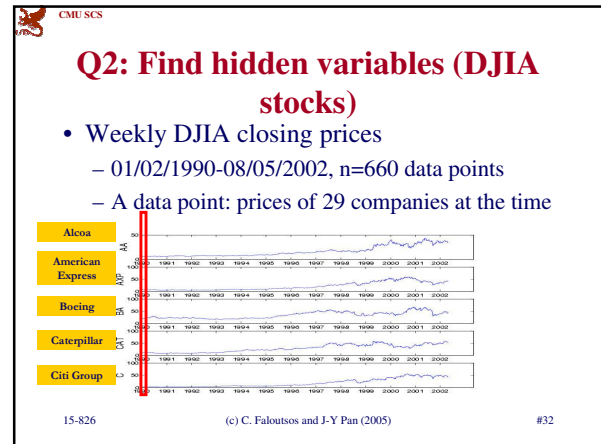
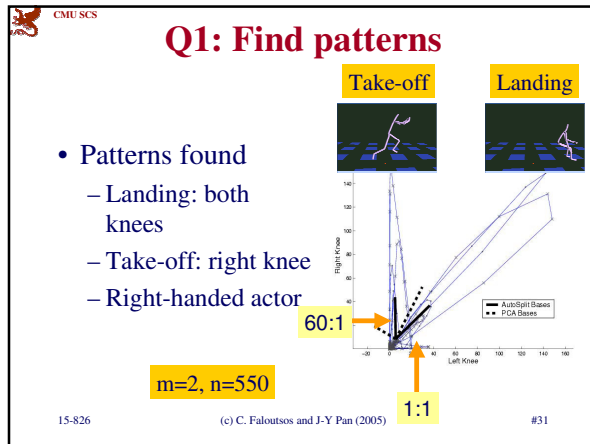
Batch-AutoSplit at work

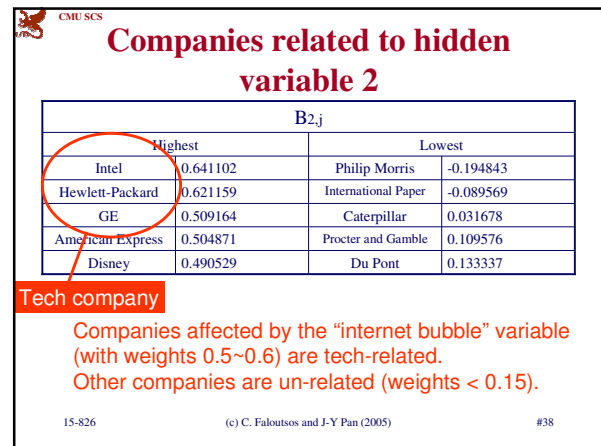
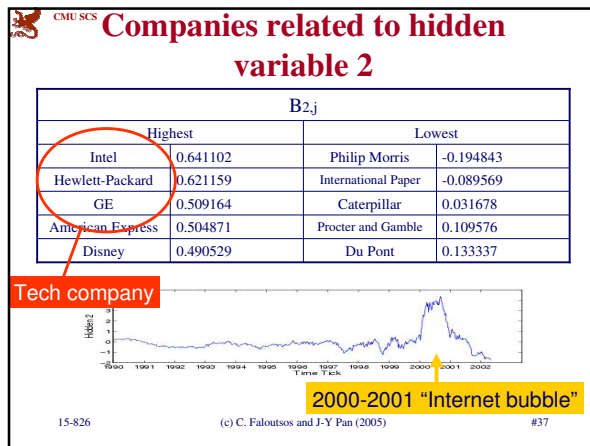
‘Sparse coding’: most points have ~zero for one or both coefficients



15-826 (c) C. Faloutsos and J-Y Pan (2005) #24







CMU SCS

Q3: Separate text from background

- Task: Separate text from background in video frames
- Data points: 6x6 patches from the video frame (50% overlap)
- Use Cluster-AutoSplit to cluster the patches into 2 groups.

15-826 (c) C. Faloutsos and J-Y Pan (2005) #39

CMU SCS

Q3: Separate text from background - Clustering-AutoSplit

Cluster 1 Cluster 2

15-826 (c) C. Faloutsos and J-Y Pan (2005) #40

CMU SCS

Q3: Separate text from background

- Cluster-AutoSplit gives cleaner separation
 - Does not get confused by the background edges

15-826 (c) C. Faloutsos and J-Y Pan (2005) #41

CMU SCS

Q4: Application on text streams

- Task: Segmentation and topic discovery
- Data: CNN headline news (Jan.-Jun. 1998)
- Documents of 10 topics in one single text stream (sorted by date)

15-826 (c) C. Faloutsos and J-Y Pan (2005) #47

CMU SCS

Problem set-up

oprah ... beef drug

← 30 words 30 words → ...

Each 30-word window:
a vector in V -dimensional space
 V : vocabulary size
then,
• we do PCA to 10 eigenvectors/topics,
• and THEN we do ICA:

15-826 (c) C. Faloutsos and J-Y Pan (2005) #48

CMU SCS

The 10 topics

ID	Topic	#Articles
TP_1	Sgt. Gene McKinney is on trial for alleged sexual misconduct	60
TP_2	A bomb explodes in a Birmingham, AL abortion clinic	18
TP_3	The Cattle Industry in Texas sues Oprah Winfrey for defaming beef	45
TP_4	New impotency drug Viagra is approved for use	52
TP_5	Diane Zamora is convicted of helping to murder her lover's girlfriend	22
TP_6	1998 Winter Olympic games	20
TP_7	The Pope's historic visit to Cuba in Winter 1998	39
TP_8	The economic crisis in Asia	69
TP_9	Superbowl XXXII	23
TP_{10}	Tornado in Florida	38

15-826 #49

CMU SCS

Process

- Windows of 30 words are taken
 - 50% overlap, 1659 windows
- Word vectors (3887-Dim)
- PCA to 10-Dim (n=10, T=1659)
 - Extract PCA bases for comparison
- ICA on the 10-Dim data vectors
 - Extract ICA bases
- Note: 10-Dim bases are mapped back to original space to identify member terms.

15-826 (c) C. Faloutsos and J-Y Pan (2005) #50

CMU SCS

Topics found (PCA)

pc_1	mckinne	bomb	women	sexual	sergeant
pc_2	bomb	mckinne	rudolph	clinic	atlanta
pc_3	winfrei	viagra	texa	beef	oprah
pc_4	viagra	winfrei	drug	texa	beef
pc_5	zamora	viagra	winfrei	graham	olymp
pc_6	zamora	graham	kill	viagra	jone
pc_7	pope	cuba	medal	olymp	castro
pc_8	asia	economi	japan	econom	asian
pc_9	bowl	super	re	peopl	medal
pc_{10}	peopl	tornado	super	bowl	florida

15-826 (c) C. Faloutsos and J-Y Pan (2005) #51

CMU SCS

Topics found (PCA)

pc_1	mckinne	bomb	women	sexual	sergeant
pc_2	bomb	mckinne	rudolph	clinic	atlanta
pc_3	winfrei	viagra	texa	beef	oprah
pc_4	viagra	winfrei	drug	texa	beef
pc_5	zamora	viagra	winfrei	graham	olymp
pc_6	zamora	graham	kill	viagra	jone
pc_7	pope	cuba	medal	olymp	castro
pc_8	asia	economi	japan	econom	asian
pc_9	bowl	super	re	peopl	medal
pc_{10}	peopl	tornado	super	bowl	florida

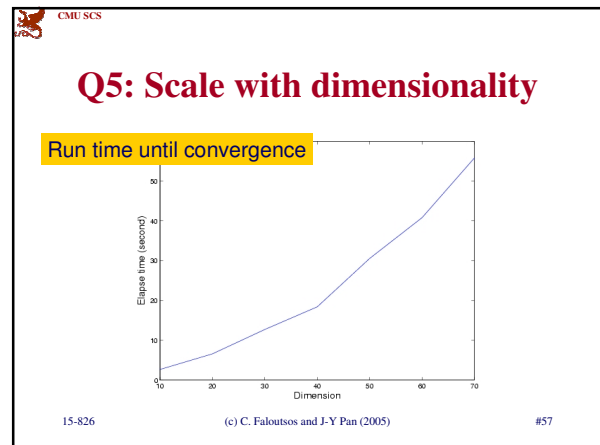
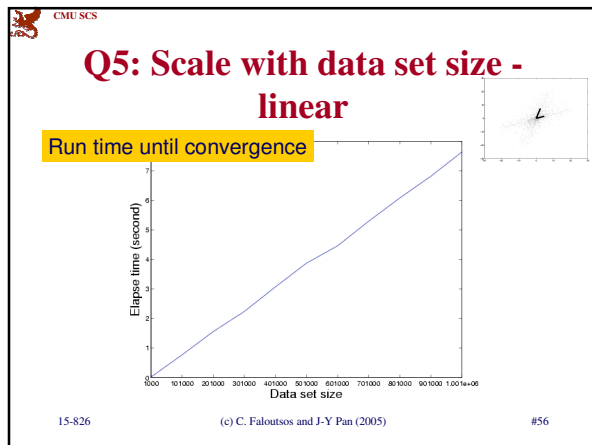
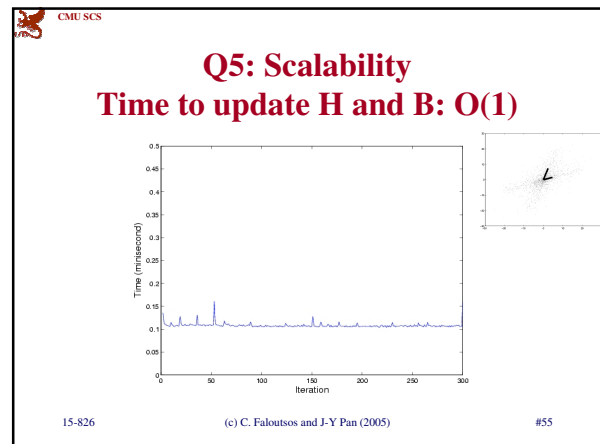
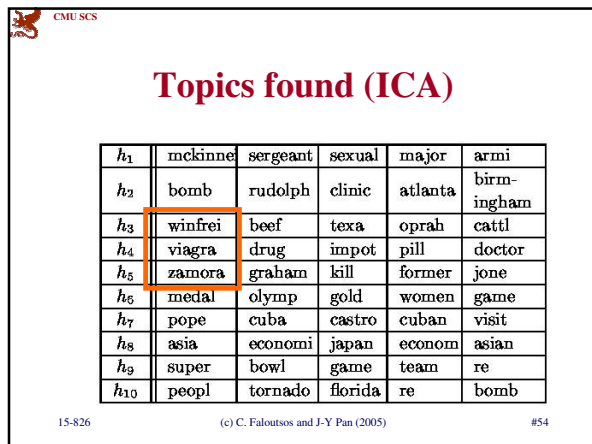
15-826 (c) C. Faloutsos and J-Y Pan (2005) #52

CMU SCS

Topics found (ICA)

h_1	mckinne	sergeant	sexual	major	armi
h_2	bomb	rudolph	clinic	atlanta	birm-ingham
h_3	winfrei	beef	texa	oprah	cattl
h_4	viagra	drug	impot	pill	doctor
h_5	zamora	graham	kill	former	jone
h_6	medal	olymp	gold	women	game
h_7	pope	cuba	castro	cuban	visit
h_8	asia	economi	japan	econom	asian
h_9	super	bowl	game	team	re
h_{10}	peopl	tornado	florida	re	bomb

15-826 (c) C. Faloutsos and J-Y Pan (2005) #53



CMU SCS

Conclusions

- AutoSplit (=ICA = sparse coding = Blind Source Separation)
 - Can find patterns (Q1)
 - Can find hidden variables (Q2)
 - **Better than PCA** (Q4)
 - Suitable for streams
 - Scalable (Q5)

15-826 (c) C. Faloutsos and J-Y Pan (2005) #58

CMU SCS

References

- Aapo Hyvärinen, Juha Karhunen, Erkki Oja: *Independent Component Analysis*, John Wiley & Sons, 2001

15-826 (c) C. Faloutsos and J-Y Pan (2005) #59