

CMU SCS


15-826: Multimedia Databases and Data Mining

Lecture #22:
Independent Component Analysis (ICA)
Jia-Yu Pan and Christos Faloutsos

15-826 (c) C. Faloutsos and J-Y Pan (2017) #1

CMU SCS

Must-read Material



- *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 (2017) #2

CMU SCS

Outline

- Motivation
- Formulation
- PCA and ICA
- Example applications
- Conclusion


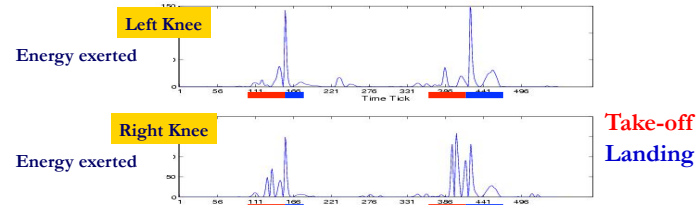
15-826 (c) C. Faloutsos and J-Y Pan (2017) #3

CMU SCS

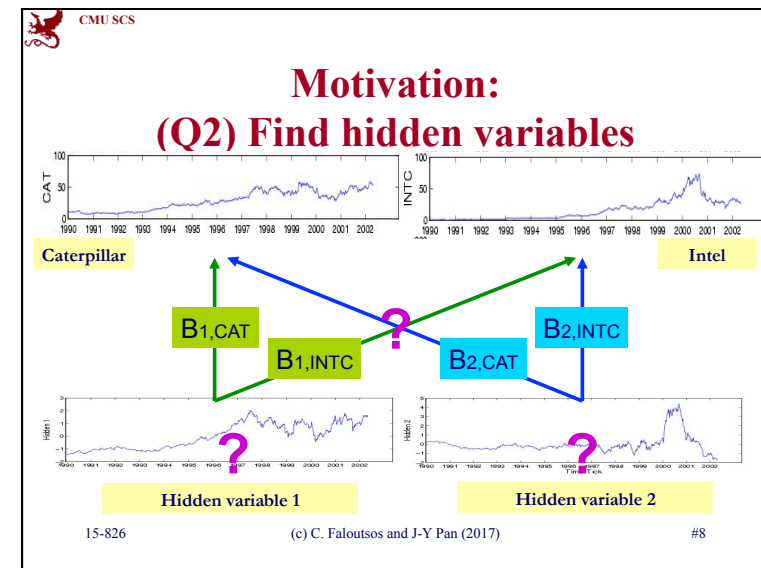
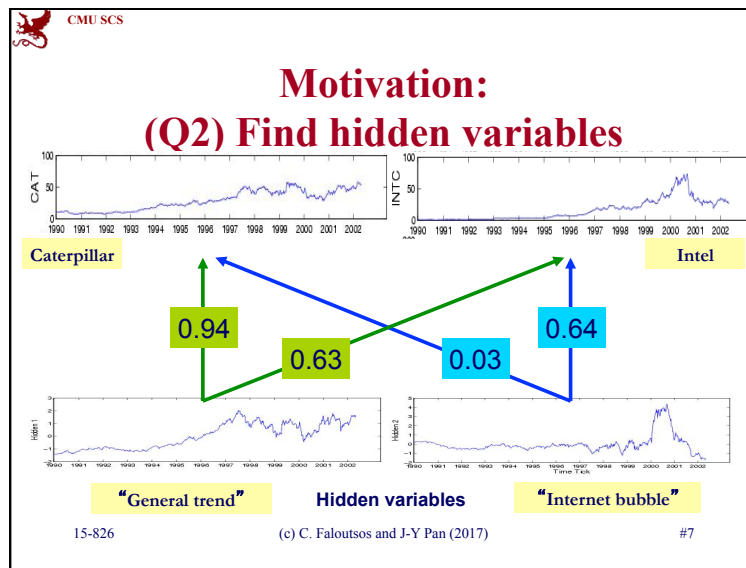
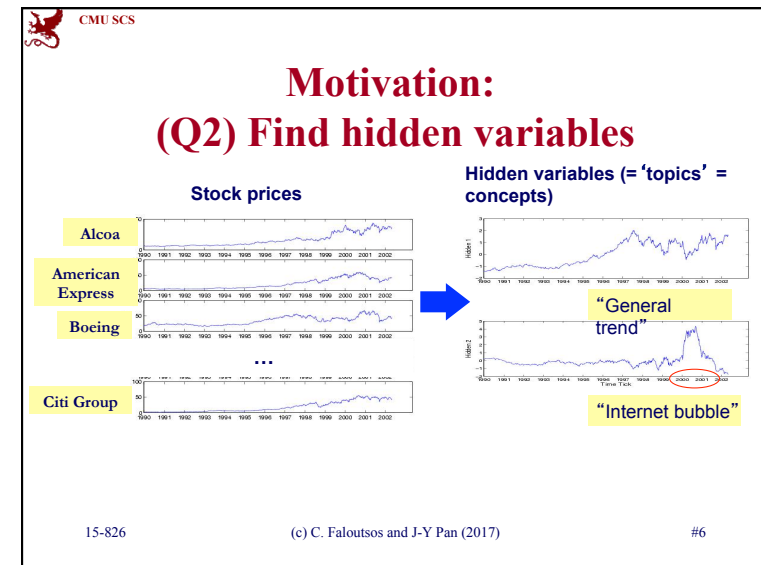
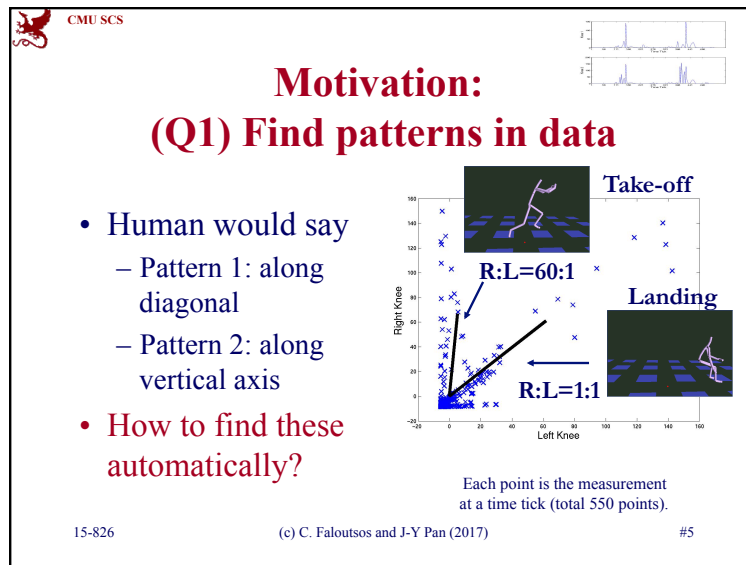
Motivation:

(Q1) Find patterns in data

- Motion capture data: broad jumps

15-826 (c) C. Faloutsos and J-Y Pan (2017) #4



CMU SCS

Motivation: Find hidden variables

- There are two sound sources in a cocktail party...

= "blind source separation"
(= we don't know the sources,
nor their mixing)

15-826 (c) C. Faloutsos and J-Y Pan (2017) #9

CMU SCS

Outline

- Motivation
- ➔ Formulation
- PCA and ICA
- Example applications
- Conclusion

15-826 (c) C. Faloutsos and J-Y Pan (2017) #10

CMU SCS

Formulation: Finding patterns

Given n data points,
each with m attributes.

Find patterns that describe
data properties the best.

15-826 (c) C. Faloutsos and J-Y Pan (2017) #11

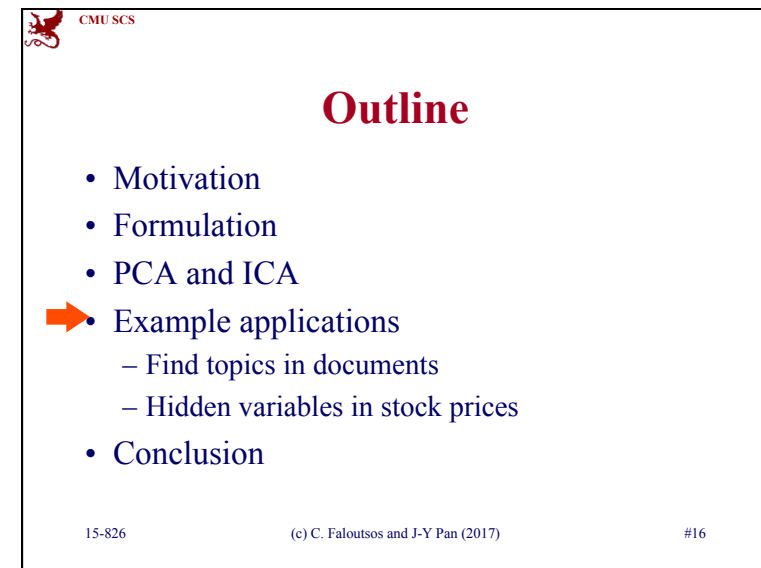
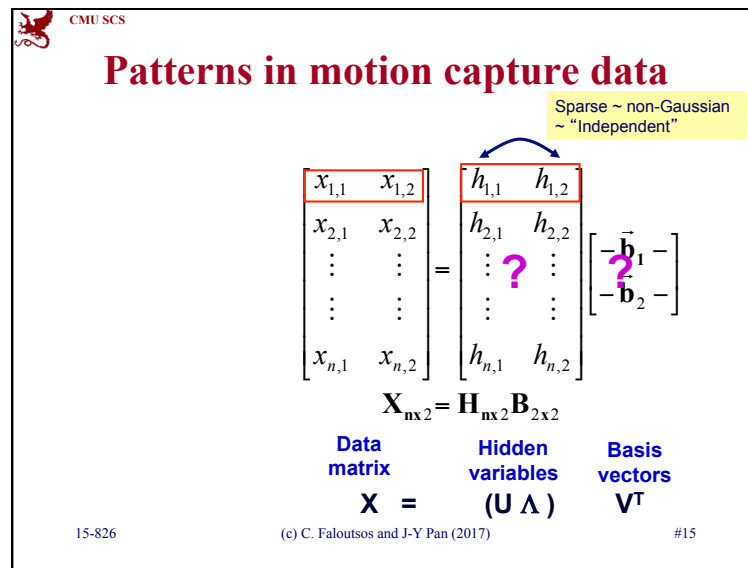
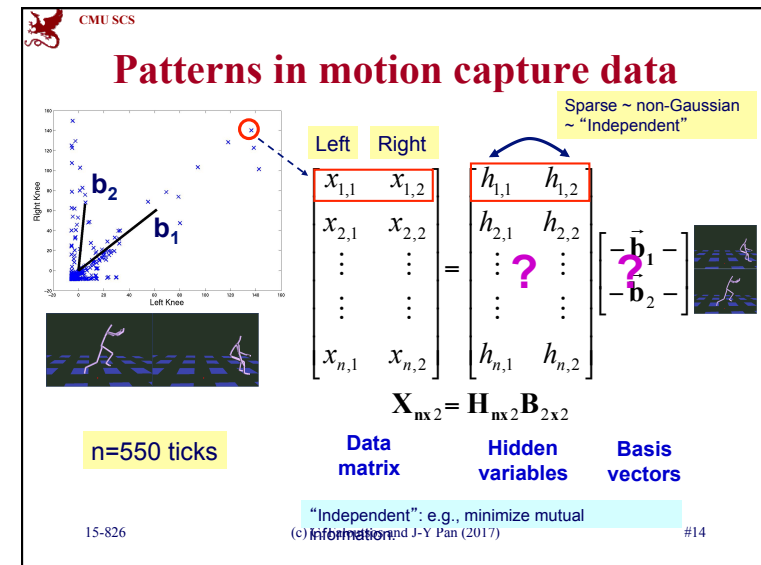
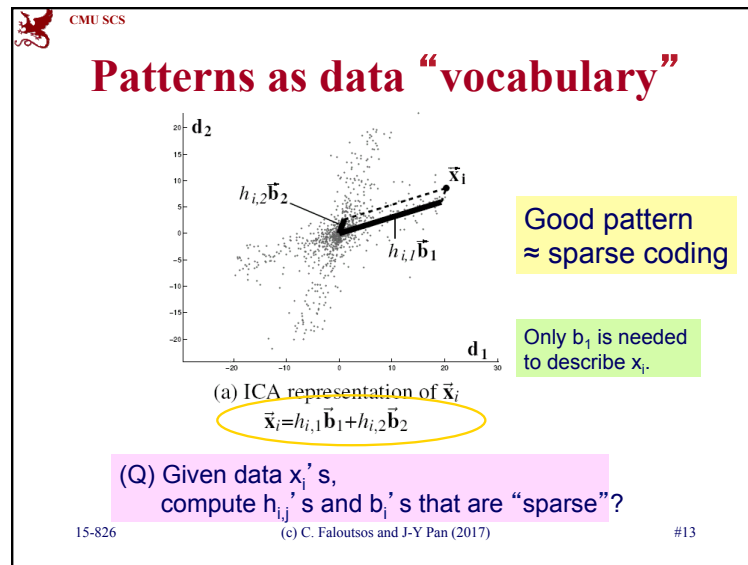
CMU SCS

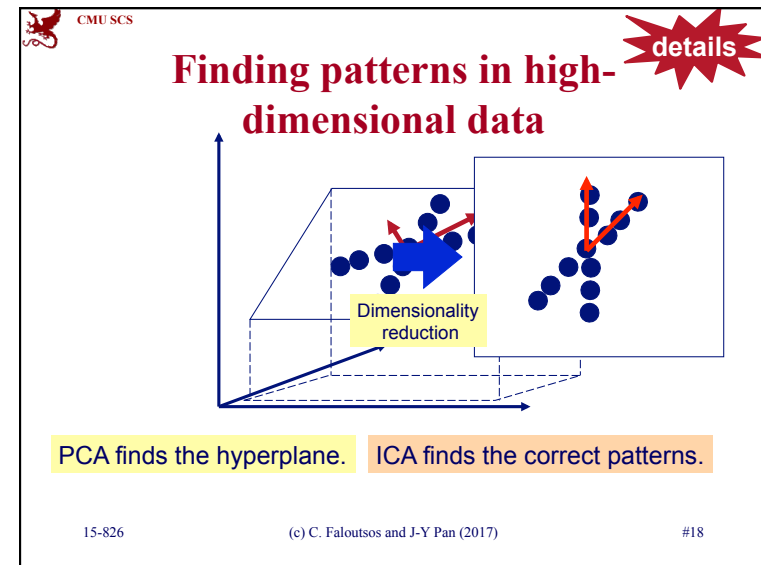
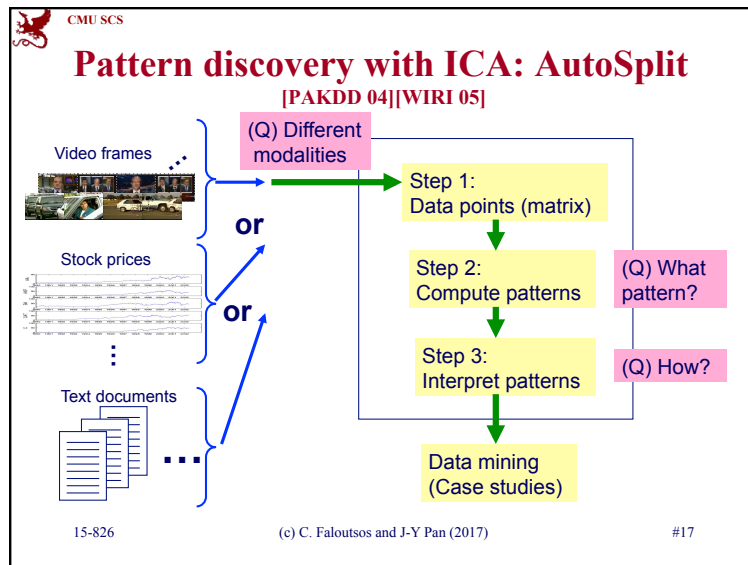
Linear representation

- Find vectors that describe the data set the best.
- Each point: linear combination of the vectors (patterns):

$$\vec{x}_i = h_{i,1} \vec{b}_1 + h_{i,2} \vec{b}_2$$

15-826 (c) C. Faloutsos and J-Y Pan (2017) #12

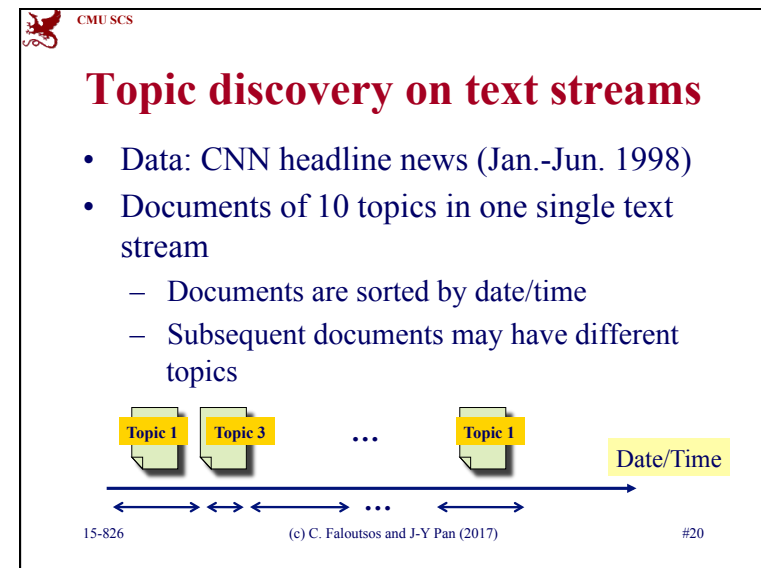




Outline

- Motivation
- Formulation
- PCA and ICA
- Example applications
 - ➔ Find topics in documents
 - Hidden variables in stock prices
 - Visual vocabulary for retinal images
- Conclusion

15-826 (c) C. Faloutsos and J-Y Pan (2017) #19



CMU SCS

Topic discovery on text streams

- Data: CNN headline news (Jan.-Jun. 1998)
- Documents of 10 topics in one single text stream
 - FIND: the document boundaries
 - AND: the terms of each topic

15-826 (c) C. Faloutsos and J-Y Pan (2017) #21

CMU SCS

Topic discovery on text streams

- Known: number of topics = 10
- Unknown: (1) topic of each document (2) topic description

15-826 (c) C. Faloutsos and J-Y Pan (2017) #22

CMU SCS

Topic discovery in documents

Step 1

New stories (n=1659) (30 words) → Windowing → $\begin{bmatrix} -\bar{x}_1 - \\ -\bar{x}_2 - \\ \vdots \\ -\bar{x}_n - \end{bmatrix} \mathbf{x}_{[n \times m]}$

$\mathbf{x}_i = [1, 5, \dots, 0]$

m=3887 (dictionary size)

Step 2

$\mathbf{X}_{[n \times m']} = \mathbf{H}_{[n \times m']}$ $\mathbf{B}_{[m' \times m']}$

(1) Find hyperplane ($m'=10$)
(2) Find patterns

Step 3

$\begin{bmatrix} -b'_1 - \\ -b'_2 - \\ \vdots \\ -b'_{10} - \end{bmatrix} \mathbf{b}'_i = [0, 0.7, \dots, 0.6]$

(Q) What does \mathbf{b}'_i mean?

15-826 (c) C. Faloutsos and J-Y Pan (2017) #23

CMU SCS

Step 3: Interpret the patterns

$\begin{bmatrix} -b'_1 - \\ -b'_2 - \\ \vdots \\ -b'_{10} - \end{bmatrix} \mathbf{b}'_i = [0, 0.7, \dots, 0.6]$

m=3887 (dictionary size)

Top words: "animal", "zoo", ...

A hidden topic!

Topics found

ID	Sorted word list				
A	Mckinne	Sergeant	sexual	Major	Armi
B	bomb	Rudolph	Clinic	Atlanta	Birmingham
C	Winfrei	Beef	Texa	Oprah	Cattl
D	Viagra	Drug	Impot	Pill	Doctor
E	Zamora	Graham	Kill	Former	Jone
H	Asia	Economi	Japan	Econom	Asian
I	Super	Bowl	Game	Team	Re
J	Peopl	Tornado	Florida	Re	bomb

General idea: related to the data attributes

15-826

CMU SCS

Step 3: Evaluate the patterns

ID	True Topic
1	Sgt. Gene McKinney is on trial for alleged sexual misconduct
2	A bomb explodes in a Birmingham, AL abortion clinic
3	The Cattle Industry in Texas sues Oprah Winfrey for defaming beef
4	New impotency drug Viagra is approved for use
5	Diane Zamora is convicted of helping to murder her lover's girlfriend

ID	Sorted word list				
A	mckinne	sergeant	sexual	major	armi
B	bomb	rudolph	clinic	atlanta	birmingham
C	winfrei	beef	texa	oprah	cattl
D	viagra	drug	Impot	pill	doctor
E	zamora	graham	kill	former	jone

AutoSplit finds correct topics.

15-826 #25

CMU SCS

Step 3: Evaluate the patterns

ID	AutoSplit				
A	mckinne	sergeant	sexual	major	armi
B	bomb	rudolph	clinic	atlanta	birmingham
C	winfrei	beef	texa	oprah	cattl
D	viagra	drug	Impot	pill	doctor
E	zamora	graham	kill	former	jone

ID	PCA				
A'	mckinne	bomb	women	sexual	sergeant
B'	bomb	mckinne	rudolph	clinic	atlanta
C'	winfrei	viagra	texa	beef	oprah
D'	viagra	winfrei	drug	texa	beef
E'	zamora	viagra	winfrei	graham	olymp

AutoSplit's topics are better than PCA.

15-826 (c) C. Faloutsos and J-Y Pan (2017) #26

CMU SCS

Step 3: Evaluate the patterns

AutoSplit				
A				
B				
C				
D				
E				

PCA				
A'				
B'				
C'				
D'				
E'				

PCA vectors mix the topics.

AutoSplit's topics are better than PCA.

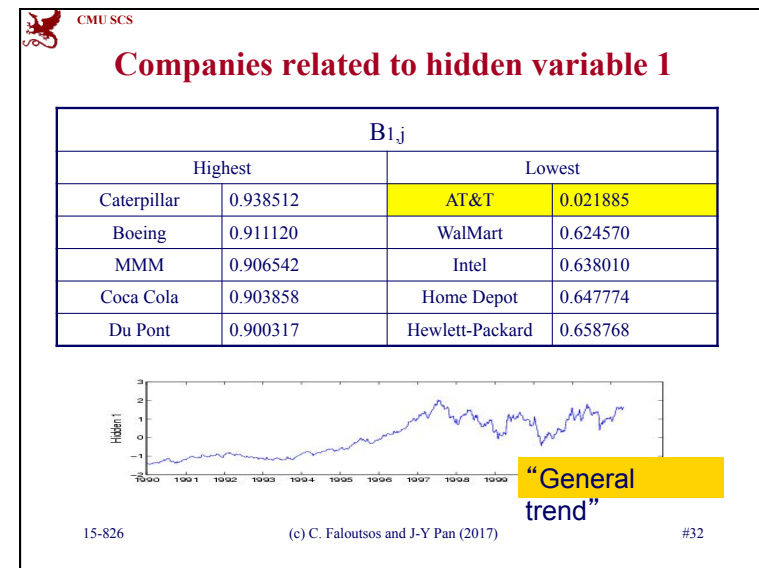
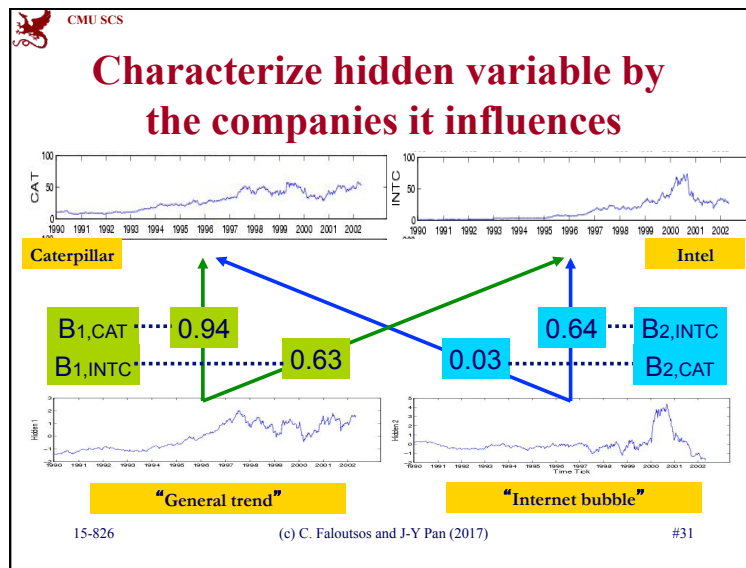
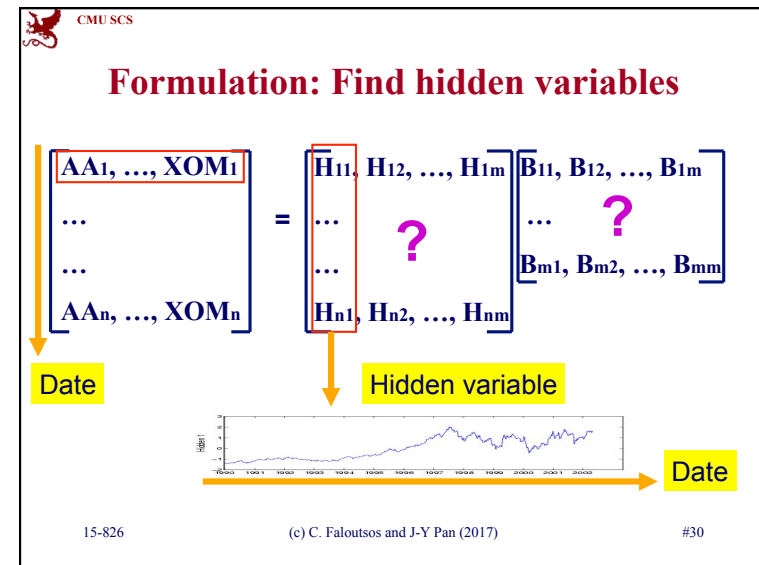
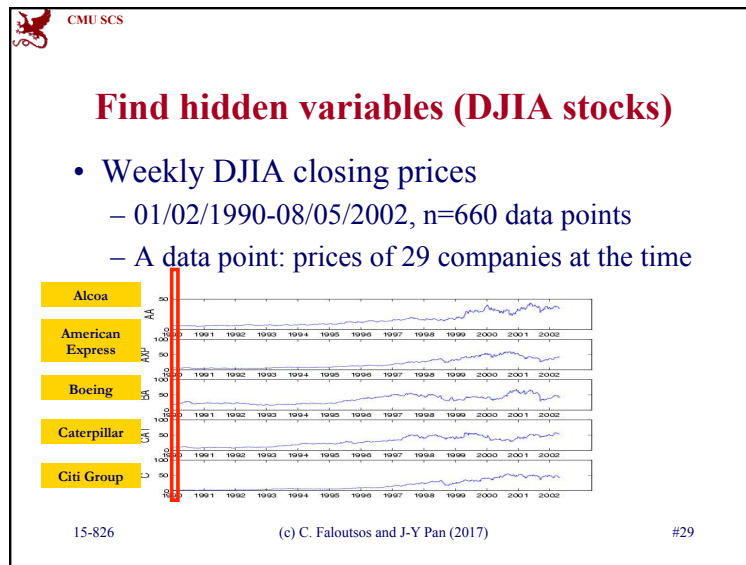
15-826 (c) C. Faloutsos and J-Y Pan (2017) #27

CMU SCS

Outline

- Motivation
- Formulation
- PCA and ICA
- Example applications
 - Find topics in documents
 - Hidden variables in stock prices
- Conclusion

15-826 (c) C. Faloutsos and J-Y Pan (2017) #28



CMU SCS

Companies related to hidden variable 1

B _{1,j}			
Highest		Lowest	
Caterpillar	0.938512	AT&T	0.021885
Boeing	0.911120	WalMart	0.624570
MMM	0.906542	Intel	0.638010
Coca Cola	0.903858	Home Depot	0.647774
Du Pont	0.900317	Hewlett-Packard	0.658768

All companies are affected by the “general trend” variable (with weights 0.6~0.9), except AT&T.

15-826 (c) C. Faloutsos and J-Y Pan (2017) #33

CMU SCS

General trend (and outlier)

15-826 (c) C. Faloutsos and J-Y Pan (2017) #34

CMU SCS

Companies related to hidden variable 2

B _{2,j}			
Highest		Lowest	
Intel	0.641102	Philip Morris	-0.194843
Hewlett-Packard	0.621159	International Paper	-0.089569
GE	0.509164	Caterpillar	0.031678
American Express	0.504871	Procter and Gamble	0.109576
Disney	0.490529	Du Pont	0.133337

Tech company

2000-2001 “Internet bubble”

15-826 (c) C. Faloutsos and J-Y Pan (2017) #35

CMU SCS

Companies related to hidden variable 2

B _{2,j}			
Highest		Lowest	
Intel	0.641102	Philip Morris	-0.194843
Hewlett-Packard	0.621159	International Paper	-0.089569
GE	0.509164	Caterpillar	0.031678
American Express	0.504871	Procter and Gamble	0.109576
Disney	0.490529	Du Pont	0.133337

Tech company

Companies affected by the “internet bubble” variable (with weights 0.5~0.6) are tech-related. Other companies are un-related (weights < 0.15).

15-826 (c) C. Faloutsos and J-Y Pan (2017) #36

CMU SCS

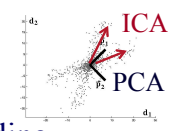
Outline

- Motivation
- Formulation
- PCA and ICA
- Example applications
 - Find topics in documents
 - Hidden variables in stock prices
 - Visual vocabulary for retinal images
- ➔ • Conclusion

15-826 (c) C. Faloutsos and J-Y Pan (2017) #37

CMU SCS

Conclusion




- ICA: more flexible than PCA in finding patterns.
- Many applications
 - Find topics and “vocabulary” for images
 - Find hidden variables in time series (e.g., stock prices)
 - Blind source separation
- Rule of thumb: plot after PCA;
 - if ‘chicken-feet’, try ICA

15-826 (c) C. Faloutsos and J-Y Pan (2017) #38

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 (2017) #39

CMU SCS

References

- Jia-Yu Pan, Andre Guilherme Ribeiro Balan, Eric P. Xing, Agma Juci Machado Traina, and Christos Faloutsos. Automatic Mining of Fruit Fly Embryo Images. In *Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2006.
- Arnab Bhattacharya, Vebjorn Ljosa, Jia-Yu Pan, Mark R. Verardo, Hyungjeong Yang, Christos Faloutsos, and Ambuj K. Singh. ViVo: Visual Vocabulary Construction for Mining Biomedical Images. In *Proceedings of the Fifth IEEE International Conference on Data Mining (ICDM)*, 2005.
- Masafumi Hamamoto, Hiroyuki Kitagawa, Jia-Yu Pan, and Christos Faloutsos. A Comparative Study of Feature Vector-Based Topic Detection Schemes for Text Streams. In *Proceedings of International Workshop on Challenges in Web Information Retrieval and Integration (WIRI)*, 2005, pp.125-130.
- Jia-Yu Pan, Hiroyuki Kitagawa, Christos Faloutsos, and Masafumi Hamamoto. AutoSplit: Fast and Scalable Discovery of Hidden Variables in Stream and Multimedia Databases. In *Proceedings of the The Eighth Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD)*, 2004.

15-826 (c) C. Faloutsos and J-Y Pan (2017) #40


 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 (2017) #41

 CMU SCS

Software

- Open source software: ‘fastICA’
<http://research.ics.tkk.fi/ica/fastica/>
- Or ‘autosplit’ :
www.cs.cmu.edu/~jypan/software/autosplit_cmu.tar.gz

15-826 (c) C. Faloutsos and J-Y Pan (2017) #42