# 15-826: Multimedia Databases and Data Mining

*Independent Component Analysis (ICA)*

Jia-Yu Pan and Christos Faloutsos

15-826                    (c) C. Faloutsos and J-Y Pan (2007)                    #1

---

# Outline

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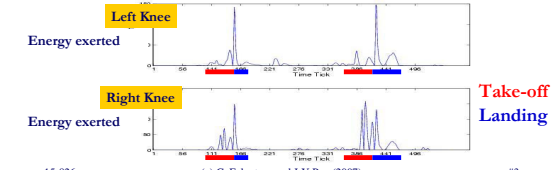- Formulation
- PCA and ICA
- Example applications
- Conclusion

15-826                    (c) C. Faloutsos and J-Y Pan (2007)                    #2

---

# Motivation:
# (Q1) Find patterns in data

- Motion capture data: broad jumps



**Left Knee**

Energy exerted

**Right Knee**

Energy exerted

**Take-off**
**Landing**

15-826                    (c) C. Faloutsos and J-Y Pan (2007)                    #3

# Motivation:
# (Q1) Find patterns in data

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

**Take-off**

R:L=60:1

**Landing**

R:L=1:1

Right Knee / Left Knee

Each point is the measurement
at a time tick (total 550 points).

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #4

---

# Motivation:
# (Q2) Find hidden variables

**Stock prices**                                **Hidden variables**

Alcoa

American Express

Boeing

...

Citi Group

"General trend"

"Internet bubble"

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #5

---

# Motivation:
# (Q2) Find hidden variables
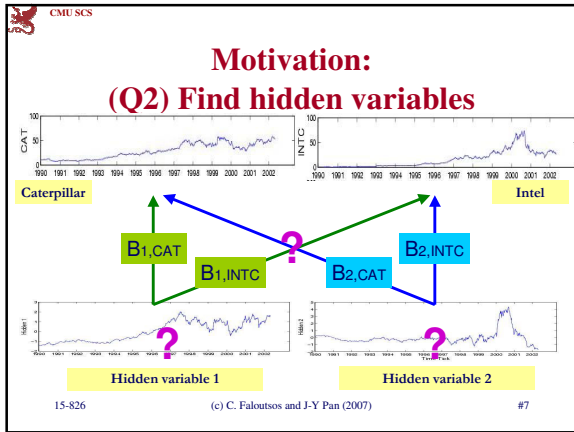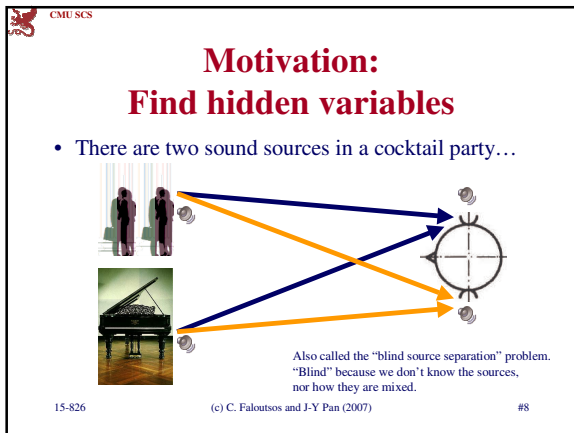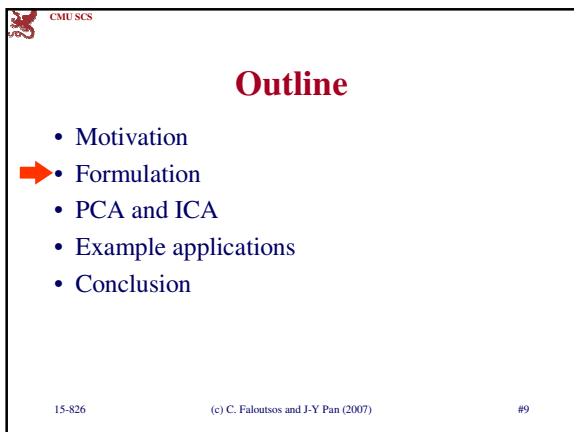
CAT

INTC

**Caterpillar**

**Intel**

0.94

0.63

0.03

0.64

"General trend"          **Hidden variables**          "Internet bubble"

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #6

# Motivation:
## (Q2) Find hidden variables

Caterpillar

Intel

$B_{1,CAT}$

$B_{1,INTC}$

$B_{2,CAT}$

$B_{2,INTC}$

**?**

**?**

**?**

Hidden variable 1

Hidden variable 2

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #7

---

# Motivation:
## Find hidden variables

- There are two sound sources in a cocktail party…

  Also called the "blind source separation" problem.
  "Blind" because we don't know the sources,
  nor how they are mixed.

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #8

---

# Outline

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

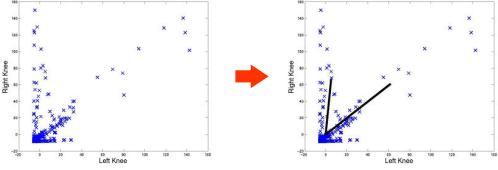15-826      (c) C. Faloutsos and J-Y Pan (2007)      #9
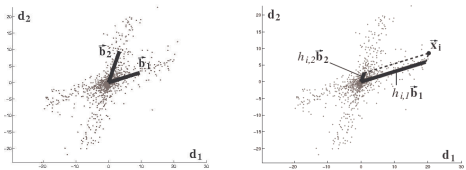
# Formulation: Finding patterns



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 (2007)        #10

---
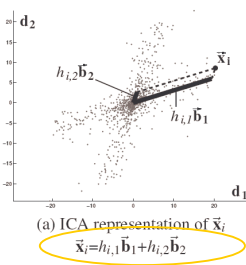
# Linear representation



- Find patterns that are vectors that describe the data set the best.
- Each point is described as a linear combination of the vectors (patterns):

$$\vec{\mathbf{x}}_\mathbf{i} = h_{i,1}\vec{\mathbf{b}}_\mathbf{1} + h_{i,2}\vec{\mathbf{b}}_\mathbf{2}$$

15-826        (c) C. Faloutsos and J-Y Pan (2007)        #11

---

# Patterns as data "vocabulary"



Good pattern
≈ sparse coding

Only $b_1$ is needed
to describe $x_i$.

(a) ICA representation of $\vec{x}_i$

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

(Q) Given data $x_i$'s,
compute $h_{i,j}$'s and $b_i$'s that are "sparse"?

15-826        (c) C. Faloutsos and J-Y Pan (2007)        #12

# Patterns in motion capture data

Sparse ~ non-Gaussian ~ "Independent"

Left | Right

$$\mathbf{X}_{nx2} = \mathbf{H}_{nx2}\mathbf{B}_{2x2}$$

$$\begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ \vdots & \vdots \\ \vdots & \vdots \\ x_{n,1} & x_{n,2} \end{bmatrix} = \begin{bmatrix} h_{1,1} & h_{1,2} \\ h_{2,1} & h_{2,2} \\ \vdots & \vdots \\ \vdots & \vdots \\ h_{n,1} & h_{n,2} \end{bmatrix}\,? \begin{bmatrix} -\vec{\mathbf{b}}_1 - \\ -\vec{\mathbf{b}}_2 - \end{bmatrix}?$$

n=550 ticks

**Data matrix**    **Hidden variables**    **Basis vectors**

"Independent": e.g., minimize mutual information.

15-826     (c) C. Faloutsos and J-Y Pan (2007)     #13

---

# Outline

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

15-826     (c) C. Faloutsos and J-Y Pan (2007)     #14

---

# Basis vectors and hidden variables
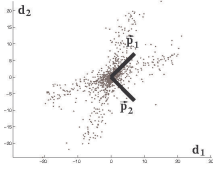
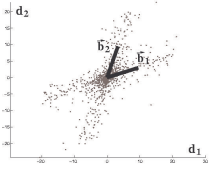- Goal: Knowing **X**, find **H** and **B**, where

  **X = H B**

- Problem: Under-constrained
  - Need additional assumptions/constraints

**X: data set**
**H: hidden variables**
**B: basis vectors**

15-826     (c) C. Faloutsos and J-Y Pan (2007)     #15

# PCA and ICA



PCA Vectors          ICA Vectors

- PCA vectors: major variations
  - Together = good "low-rank approximation"/dimensional reduction
  - Individually ≠ meaningful patterns
- Luckily, ICA detects the major meaningful patterns.

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #16

---

# PCA

- PCA (Principal Component Analysis)
  - Choose vectors which are orthonormal and
  - give smallest representation L2 error for dimensional reduction
- Matrices H and B can be solved by
  - SVD, neural networks, or many optimization methods

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #17

---

# PCA

- Extremely popular
  - Latent Semantic Indexing [Deerwester+90]
  - KL transform [Duda,Hart,Stork00]
  - EigenFace [Turk,Pentlend91]
  - Multiple time series correlation [Guha,Gunopulos,Koudas03]
- But, there is room for improvement.

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #18

# ICA

- ICA (Independent Component Analysis)
  - Make hidden variables $h_i$'s (columns of **H**) mutually independent.
- Many implementations
  - Many ways to define "independence"
  - Many ways to find the most independent **H**.
  - (**B** is found at the same time, since **X=HB**.)

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #19

---

# ICA

- Define "Independence": $p(h_i, h_j) = p(h_i) p(h_j)$
  - Zero mutual information
  - Non-Gaussianity, max. absolute Kurtosis
- To solve for H,B:
  - Neural networks, optimization methods (gradient ascent, fixed-point, …)

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #20

---

# An non-gaussian distribution: Laplacian pdf

$$P(x) = \frac{\lambda}{2} \exp(-\lambda |x|)$$

Sharper at 0,
and more heavy tail
than Gaussian pdf

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #21

# Maximizing non-gaussianity

- Assume $h_i$ ~ non-gaussian pdf (e.g. Laplacian pdf)
  - Fixed $h_i$ values, what is the most likely "**B**" ?
    - (data point **x** is given and fixed)
  - Find **B**, s.t. likelihood $P(\mathbf{x}|\mathbf{B})$ is maximized.

$$\bar{\mathbf{x}} = \bar{\mathbf{h}}\mathbf{B}$$

$$\Rightarrow P(\bar{\mathbf{x}}|\mathbf{B}) = \frac{P(\bar{\mathbf{h}})}{\det(\mathbf{B})}$$

$$P(\bar{\mathbf{h}}) = P_{Laplacian}(\bar{\mathbf{h}}) = P_{Laplacian}(\bar{\mathbf{x}}\mathbf{B}^{-1})$$

$$\Rightarrow P(\bar{\mathbf{x}}|\mathbf{B}) = \frac{P_{Laplacian}(\bar{\mathbf{x}}\mathbf{B}^{-1})}{\det(\mathbf{B})}$$

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #22

---

# Maximize likelihood

- Likelihood $P(\mathbf{x}|\mathbf{B})$ is a function of **B**, f(**B**)
- Gradient ascent
  - To find **B** which maximizes $P(\mathbf{x}|\mathbf{B})$

f(**B**)

**B***

**B**

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #23

---

# ICA by maximum likelihood
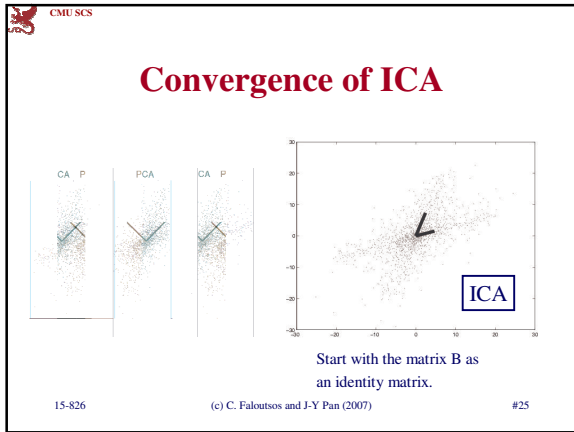
X: static

H

$$Z = -sign(H)$$
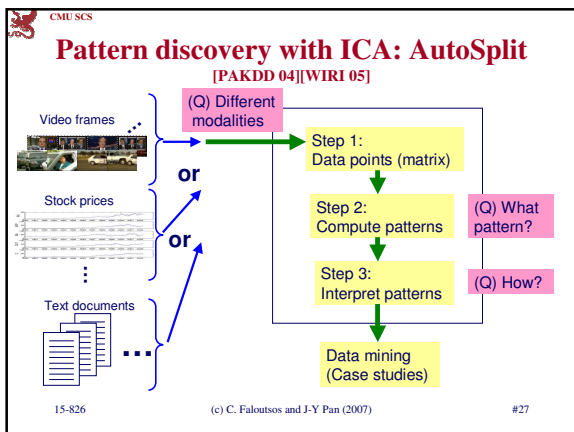$$\Delta B \propto -B^T Z^T H - nB^T$$
$$B = B + \varepsilon \Delta B$$

B

$$H = XB^{-1}$$

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #24

# Convergence of ICA

ICA

Start with the matrix B as
an identity matrix.

15-826 (c) C. Faloutsos and J-Y Pan (2007) #25

# 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 (2007) #26

# Pattern discovery with ICA: AutoSplit
## [PAKDD 04][WIRI 05]

Video frames

(Q) Different modalities

**or**

Stock prices

**or**

Text documents

Step 1:
Data points (matrix)

Step 2:
Compute patterns

(Q) What pattern?

Step 3:
Interpret patterns

(Q) How?

Data mining
(Case studies)

15-826 (c) C. Faloutsos and J-Y Pan (2007) #27

# Finding patterns in high-dimensional data

Dimensionality reduction

PCA finds the hyperplane.  ICA finds the correct patterns.

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #28

# 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 (2007)          #29

# Topic discovery on text streams

- Data: CNN headline news (Jan.-Jun. 1998)
- Documents of 10 topics in one single text stream
  - – Documents are sorted by date/time
  - – Subsequent documents may have different topics

Topic 1   Topic 3   …   Topic 1

Date/Time

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #30

# Topic discovery on text streams

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

| Topic 1 | Topic 3 | … | Topic 1 | Date/Time |
|---------|---------|---|---------|-----------|

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #31

---

# Topic discovery in documents

**Step 1**

New stories → Windowing (n=1659) (30 words) → $X_{[n \times m]}$

$$\begin{bmatrix} -\bar{x}_1- \\ -\bar{x}_2- \\ \vdots \\ -\bar{x}_n- \end{bmatrix}$$

aaron   zoo

$x_i = [1, 5, …, 0]$

m=3887 (dictionary size)

**Step 2**

$X_{[n \times m]} = H_{[n \times m]} \; (B_{[m \times m]})$

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

**Step 3**

$$\begin{bmatrix} -b_1- \\ -b_2- \\ \vdots \\ -b_{10}- \end{bmatrix}$$

aaron   animal   zoo

$b'_i = [0, 0.7, …, 0.6]$

**(Q) What does $b'_i$ mean?**

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #32

---

# Step 3: Interpret the patterns

$$\begin{bmatrix} -b_1- \\ -b_2- \\ \vdots \\ -b_{10}- \end{bmatrix}$$

aaron   animal   zoo

$b'_i = [0, 0.7, …, 0.6]$ → Top words : "animal", "zoo", …

A hidden topic!

m=3887 (dictionary size)

**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 |

**General idea: related to the data attributes**

| H | Asia | Economi | Japan | Econom | Asian |
|----|---------|---------|--------|--------|-------|
| I | Super | Bowl | Game | Team | Re |
| J | Peopl | Tornado | Florida | Re | bomb |

15-826

11

**Slide 1:**

# 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 #34

**Slide 2:**

# 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 (2007) #35

**Slide 3:**

# Step 3: Evaluate the patterns



PCA vectors mix the topics.

AutoSplit's topics are better than PCA.

15-826 (c) C. Faloutsos and J-Y Pan (2007) #36

# 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 (2007)                #37

---

## Find hidden variables (DJIA stocks)

- Weekly DJIA closing prices
  - 01/02/1990-08/05/2002, n=660 data points
  - A data point: prices of 29 companies at the time

| Alcoa |
| American Express |
| Boeing |
| Caterpillar |
| Citi Group |

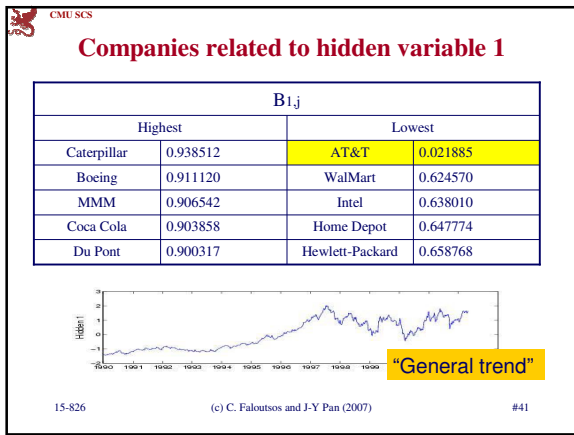15-826                (c) C. Faloutsos and J-Y Pan (2007)                #38
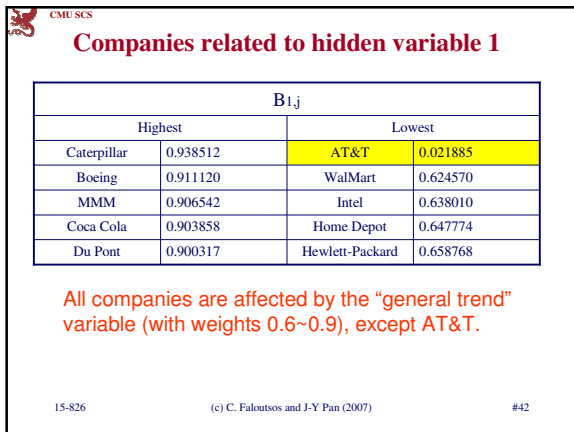
---

## Formulation: Find hidden variables

$$
\begin{bmatrix} AA_1, \ldots, XOM_1 \\ \ldots \\ \ldots \\ AA_n, \ldots, XOM_n \end{bmatrix} = \begin{bmatrix} H_{11}, H_{12}, \ldots, H_{1m} \\ \ldots \\ \ldots \\ H_{n1}, H_{n2}, \ldots, H_{nm} \end{bmatrix} \textbf{?} \begin{bmatrix} B_{11}, B_{12}, \ldots, B_{1m} \\ \ldots \\ B_{m1}, B_{m2}, \ldots, B_{mm} \end{bmatrix} \textbf{?}
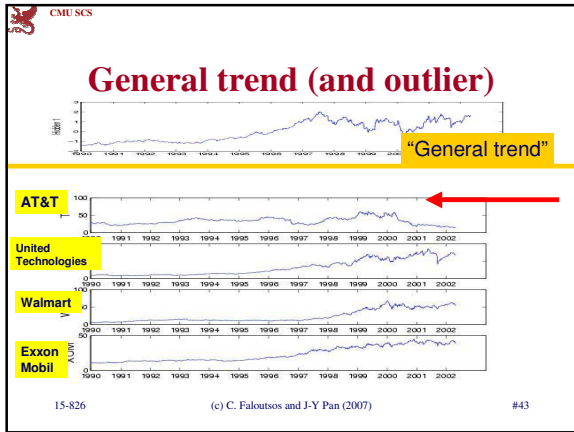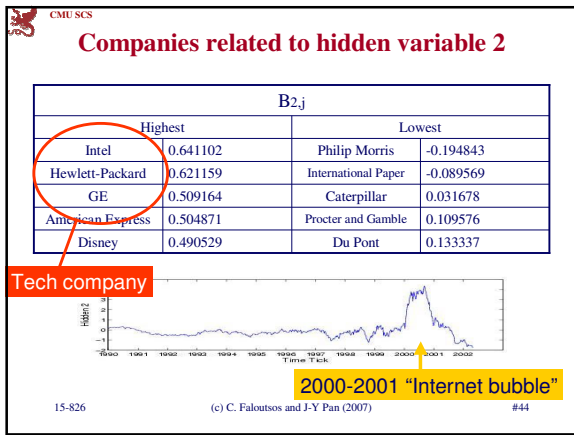$$

Date

Hidden variable

Date

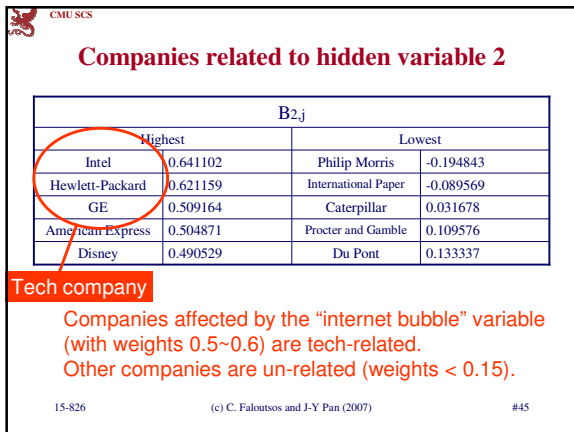15-826                (c) C. Faloutsos and J-Y Pan (2007)                #39

# Characterize hidden variable by the companies it influences



Caterpillar    Intel

$B_{1,CAT}$ — 0.94
$B_{1,INTC}$ — 0.63
0.64 — $B_{2,INTC}$
0.03 — $B_{2,CAT}$

"General trend"    "Internet bubble"

---

# 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 |



"General trend"

---

# 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.

## Slide 1

CMU SCS

# General trend (and outlier)

"General trend"

AT&T

United Technologies

Walmart

Exxon Mobil

(c) C. Faloutsos and J-Y Pan (2007)   #43

## Slide 2

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"

(c) C. Faloutsos and J-Y Pan (2007)   #44

## Slide 3

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).

(c) C. Faloutsos and J-Y Pan (2007)   #45

# 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 (2007)      #46

---

# Mining cat retinal images [ICDM 05]



**Retina**

Detachment Development

Distribution of 2 proteins

Normal | 1 day after detachment | 7 days after detachment | 28 days after detachment

Treatment

15-826    1h3dr    3d28dr    1d6dO$_2$    (c) C. Faloutso... ...an (2007)    #47

---

# "Vocabulary" for biomedical images?

- How to describe biomedical images?
- Analogy: the topics for text
  - Football reports: "touchdown", "punt", etc.
  - DB papers: "query", "optimization", etc.
- How to derive "visual vocabulary terms"?



Normal      7 days after detachment      "spongy"
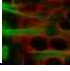
15-826      (c) C. Faloutsos and J-Y Pan (2007)      #48

Visual Vocabulary (ViVo) by AutoSplit

Visual vocabulary

Step 1: Tile image

8x12 tiles

Step 2: Extract tile features

Step 3: ViVo generation

*Feature 2*

V1

V2

*Feature 1*

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #49



Finding ViVos

ICA
PCA

$P_2$

$P_1$

$I_3$

$I_1$

$I_2$

PC 2

PC 1

Each point is a tile.
Projected to the 1st and 2nd PCA vectors.
(Feature vector: 512 color structure features.)

Red lines indicate ViVos.

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #50

Bio-mining with "ViVo"

✔ Visual Vocabulary for retinal images
 – using AutoSplit
➡ Evaluation
 – Qualitative: biological meanings of ViVos
 – Data mining: highlight "interesting" regions

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #51

17

## Biological interpretation of ViVos

| ID | ViVo | Description | Condition |
|---|---|---|---|
| V1 | | GFAP in inner retina (Müller cells) | Healthy |
| V10 | | Healthy outer segments of rod photoreceptors | Healthy |
| V8 | | Redistribution of rod opsin into cell bodies of rod photoreceptors | Detached |
| V11 | | Co-occurring processes: Müller cell hypertrophy and rod opsin redistribution | Detached |

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #52

---

## Biological interpretation of ViVos

| ID | ViVo | Description | Condition | ID | ViVo | Description | Condition |
|---|---|---|---|---|---|---|---|
| 2 | | GFAP in hypertrophy Müller cells | Morphological changes in inner retina | 6 | | Rod photoreceptor cell body | Background labeling |
| 3 | | GFAP in hypertrophy Müller cells | Morphological changes in inner retina | 7 | | GFAP in hypertrophy Müller cells | Morphological changes in inner retina |
| 4 | | GFAP in hypertrophy Müller cells | Morphological changes in inner retina | 9 | | Outer segment degeneration (rod opsin) | Detached |
| 5 | | Healthy outer segments of rod photoreceptors (rod opsin) | Healthy | 12 | | GFAP in hypertrophy Müller cells | Morphological changes in inner retina |

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #53

---

## Bio-mining with "ViVo"

✔ Visual Vocabulary for retinal images
  – using AutoSplit
• Evaluation
  ✔ Qualitative: biological meanings of ViVos
  ➡ Data mining: highlight "interesting" regions

15-826      (c) C. Faloutsos and J-Y Pan (2007)      #54

18

# Finding "distinguishing ViVos"

- Given: Images of two classes
  - Find the class-distinguishing ViVO ("DiVo")
  - Highlight distinguishing regions

Normal          Detached 3 days

15-826          (c) DiVo: "spongy" (2007)          #55

---

# Summary: a system viewpoint

Input          Output

Our system

V8: "spongy"

(1) Left: "n"; Right: "3d"

Accurate classification

DiVo analysis

ViVo interpretation

(2) Regions shown (V8):
"cells of rod photoreceptors"

(3) Description:
"Detachment occurs!"

"Rod opsin distributes from outer segment into cell bodies."

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #56

---

# 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 (2007)          #57

19

# 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

15-826       (c) C. Faloutsos and J-Y Pan (2007)       #58

---

# Vocabulary for embryo gene expressions



Vocabulary

with André Balan, Christos Faloutsos, Eric P. Xing

15-826       (c) C. Faloutsos and J-Y Pan (2007)       #59

---

# 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 (2007)       #60

**CMU SCS**

# Acknowledgement

- Prof. Tai Sing Lee
- Prof. Hiroyuki Kitagawa
- Prof. HyungJeong Yang
- Masafumi Hamamoto
- Prof. Nancy Pollard
- Prof. Jessica Hodgins

- Prof. Michael Lewicki
- Prof. Eric Xing
- CMU Informedia project
- UCSB DB Lab
- CMU bio-imaging center
- CMU graphics lab

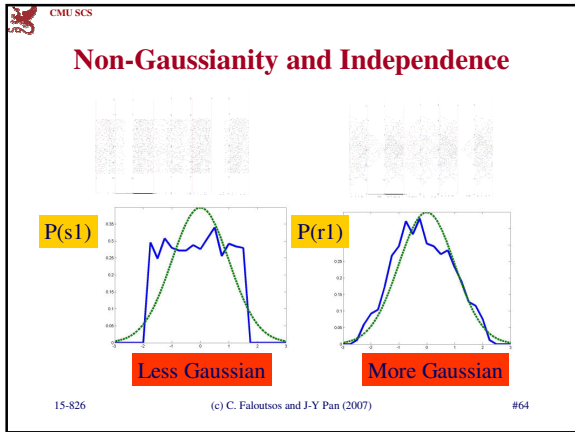15-826          (c) C. Faloutsos and J-Y Pan (2007)          #61
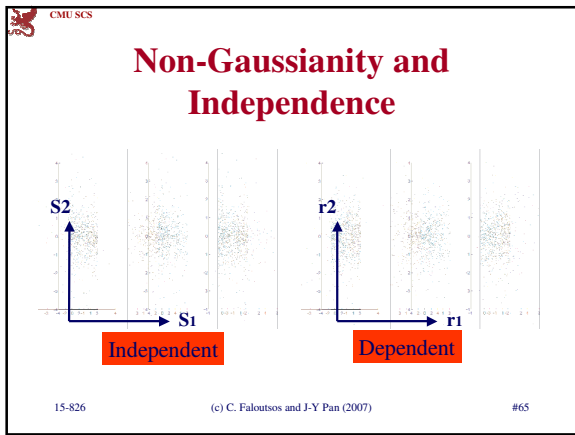
---

**CMU SCS**

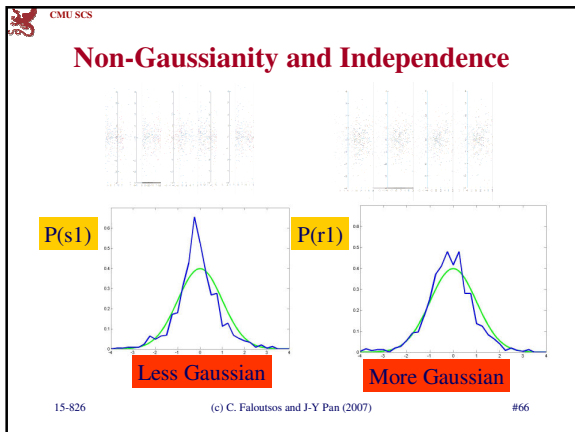15-826          (c) C. Faloutsos and J-Y Pan (2007)          #62

---

**CMU SCS**

# Independence



S2 / S1 — Independent

r2 / r1 — Dependent

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #63

## Non-Gaussianity and Independence



P(s1)  P(r1)

Less Gaussian   More Gaussian

## Non-Gaussianity and Independence



S2  r2

S1  r1

Independent   Dependent

## Non-Gaussianity and Independence



P(s1)  P(r1)

Less Gaussian   More Gaussian

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #67

---

# 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 (2007)          #68

---

# References

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

15-826          (c) C. Faloutsos and J-Y Pan (2007)          #69