

Carnegie Mellon

15-826: Multimedia Databases and Data Mining

Lecture #17: SVD – part II – applications

C. Faloutsos

1

Carnegie Mellon

Problems



- Q1: How to find ‘concepts’ in a document collection?
- Q2: how to answer queries in English, when documents are in Spanish?
- Q3: how to compress a customer x day matrix
- Q4: how to interpret the rules/concepts
- Q5: KL transform?



15-826


Copyright (c) 2019 C. Faloutsos

2



2

CarnegieMellon

Solutions



- Q1: How to find ‘concepts’ in a document collection?
- Q2: how to analyze a document in English, when document
- Q3: how to analyze a customer x day matrix
- Q4: how to interpret the rules/concepts
- Q5: KL transform?

15-826 Copyright (c) 2019 C. Faloutsos 3

3

CarnegieMellon

Must-read Material

- [MM Textbook Appendix D](#)

15-826 Copyright (c) 2019 C. Faloutsos 4

4

Carnegie Mellon

Outline

Goal: 'Find **similar / interesting** things'

- Intro to DB
- ➔ • Indexing - similarity search
- ↳ • Data Mining

15-826 Copyright (c) 2019 C. Faloutsos 5

5

Carnegie Mellon

Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
- spatial access methods
- fractals
- text
- ➔ • Singular Value Decomposition (SVD)
- multimedia
- ...

15-826 Copyright (c) 2019 C. Faloutsos 6

6

SVD - Detailed outline

- Motivation
- Definition - properties
- Interpretation
- Complexity
- ➔ • Case studies
- SVD properties
- Conclusions

7

SVD - Case studies

- ➔ • multi-lingual IR; LSI queries
- compression
- PCA - 'ratio rules'
- Karhunen-Lowe transform
- query feedbacks
- google/Kleinberg algorithms

8

Case study - LSI

Q1: How to do queries with LSI?

Q2: multi-lingual IR (english query, on spanish text?)

Case study - LSI

Q1: How to do queries with LSI?

Problem: Eg., find documents with 'data'

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

CarnegieMellon

Case study - LSI

Q1: How to do queries with LSI?
 A: map query vectors into 'concept space' – how?

$$q = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

data retrieval
 inf. ↓ brain lung

A: inner product (cosine similarity) with each 'concept' vector v_i

15-826
Copyright (c) 2019 C. Faloutsos
13

13

CarnegieMellon

Case study - LSI

Q1: How to do queries with LSI?
 A: map query vectors into 'concept space' – how?

$$q = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

data retrieval
 inf. ↓ brain lung

A: inner product (cosine similarity) with each 'concept' vector v_i

15-826
Copyright (c) 2019 C. Faloutsos
14

14

Case study - LSI

compactly, we have:

$$q \mathbf{V} = q_{\text{concept}}$$

Eg:

$$q = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

data retrieval
 inf. ↓ brain lung

$$\begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}$$

term-to-concept similarities

$$= \begin{bmatrix} 0.58 & 0 \end{bmatrix}$$

CS-concept
↓

15-826

Copyright (c) 2019 C. Faloutsos

15

15

Case study - LSI

Drill: how would the document ('information', 'retrieval') be handled by LSI?

15-826

Copyright (c) 2019 C. Faloutsos

16

16

CarnegieMellon

Case study - LSI

Drill: how would the document ('information' , 'retrieval') be handled by LSI? A: SAME:

$d_{\text{concept}} = d V$

Eg: $d = \begin{bmatrix} \text{data} & \text{inf.} & \text{retrieval} & \text{brain} & \text{lung} \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix}$

$V = \begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}$ (term-to-concept similarities)

$d_{\text{concept}} = \begin{bmatrix} 1.16 & 0 \end{bmatrix}$ (CS-concept)

15-826 Copyright (c) 2019 C. Faloutsos 17

17

CarnegieMellon

Case study - LSI

Observation: document ('information' , 'retrieval') will be retrieved by query ('data'), although it does not contain 'data' !!

$d = \begin{bmatrix} \text{data} & \text{inf.} & \text{retrieval} & \text{brain} & \text{lung} \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix}$


$q = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}$

$d_{\text{concept}} = \begin{bmatrix} 1.16 & 0 \\ 0.58 & 0 \end{bmatrix}$ (CS-concept)

15-826 Copyright (c) 2019 C. Faloutsos 18


18

CarnegieMellon



Solutions

- ✓ Q1: How to find ‘concepts’ in a document collection?
- Q2: how to answer queries in English, when documents are in Spanish?
- Q3: how to compress a customer x day matrix
- Q4: how to interpret the rules/concepts
- Q5: KL transform?



15-826 Copyright (c) 2019 C. Faloutsos 19

19

CarnegieMellon

Case study - LSI

- Q1: How to do queries with LSI?
- ➔ Q2: multi-lingual IR (english query, on spanish text?)

15-826 Copyright (c) 2019 C. Faloutsos 20

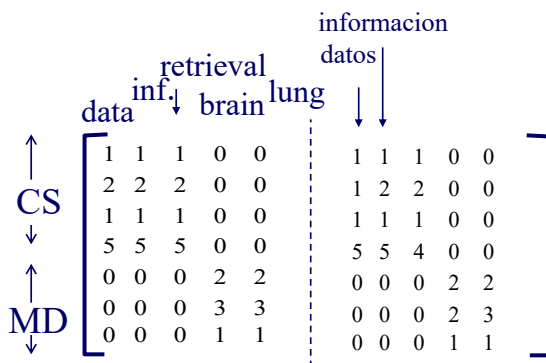
20

Case study - LSI

- Problem:
 - given many documents, translated to both languages (eg., English and Spanish)
 - answer queries across languages


Case study - LSI

- Solution: ~ LSI




CarnegieMellon

Solutions



- ✓ Q1: How to find ‘concepts’ in a document collection?
- ✓ Q2: how to answer queries in English, when documents are in Spanish?
- Q3: how to compress a customer x day matrix
- Q4: how to interpret the rules/concepts
- Q5: KL transform?



15-826 Copyright (c) 2019 C. Faloutsos 23

23

CarnegieMellon

Case study: compression

[Korn+97]

Problem:

- given a matrix
- compress it, but maintain ‘random access’
(surprisingly, its solution leads to data mining and visualization...)

Flip Korn, H. V. Jagadish, and Christos Faloutsos. *Efficiently supporting ad hoc queries in large datasets of time sequences*. SIGMOD '97, 289-300.

25

Carnegie Mellon

Problem - specs

- $\sim 10 \times 6$ rows; $\sim 10 \times 3$ columns; no updates;
- random access to any cell(s) ; small error: OK

customer	day	We 7/10/96	Th 7/11/96	Fr 7/12/96	Sa 7/13/96	Su 7/14/96
ABC Inc.		1	1	1	0	0
DEF Ltd.		2	2	2	0	0
GHI Inc.		1	1	1	0	0
KLM Co.		5	5	5	0	0
Smith		0	0	0	2	2
Johnson		0	0	0	3	3
Thompson		0	0	0	1	1

15-826
Copyright (c) 2019 C. Faloutsos
26

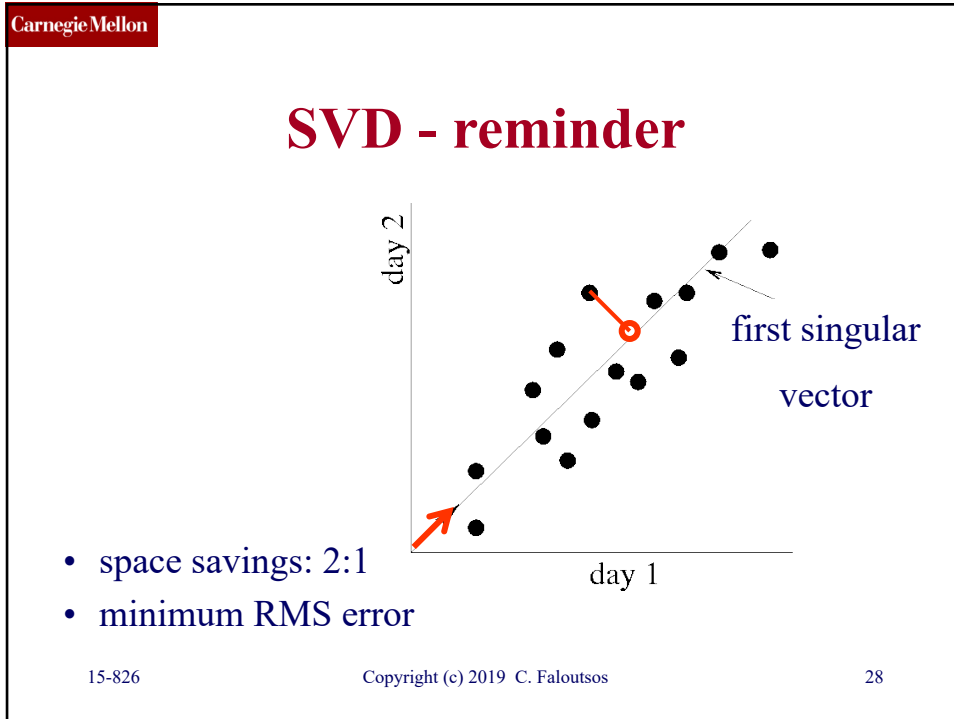
26

Carnegie Mellon

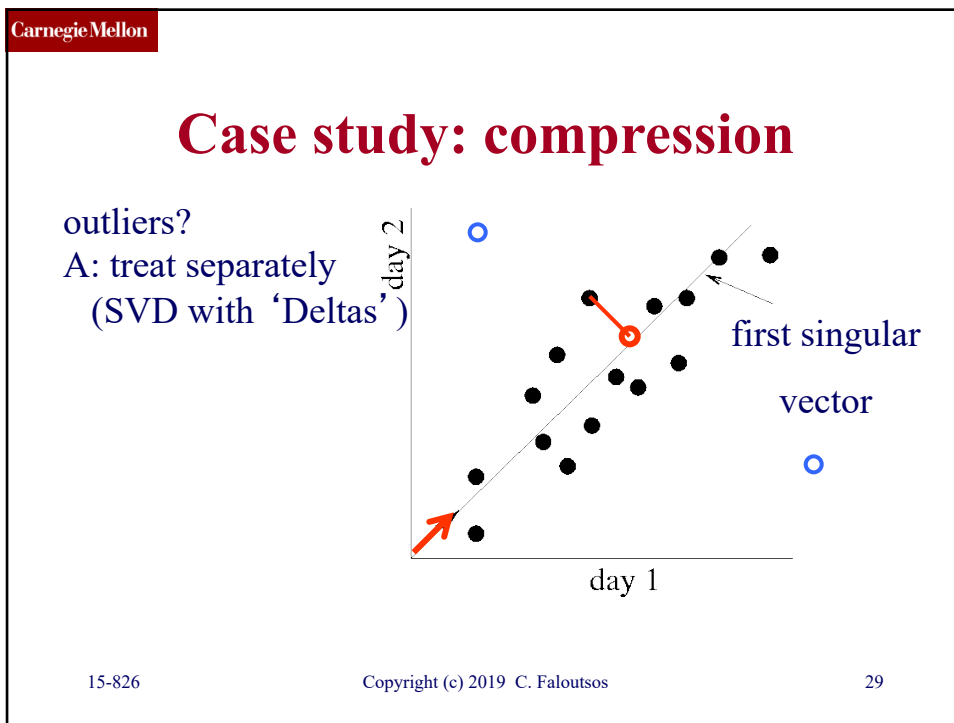
Idea

15-826
Copyright (c) 2019 C. Faloutsos
27

27



28

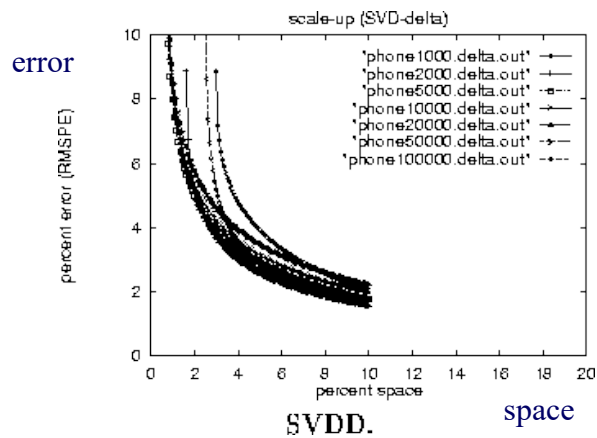


29

Compression - Performance

- 3 pass algo (-> scalability) (HOW?)
- random cell(s) reconstruction
- 10:1 compression with < 2% error

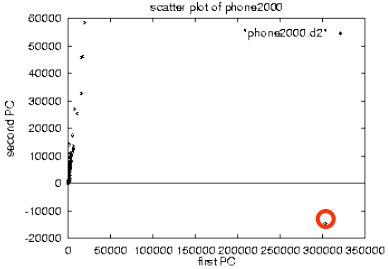
Performance - scaleup



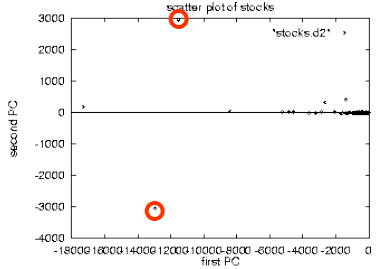
CarnegieMellon

Compression - Visualization

- no Gaussian clusters; Zipf-like distribution



(a) 'phone2000'




(b) 'stocks'

15-826 Copyright (c) 2019 C. Faloutsos 32


32

CarnegieMellon

Solutions



- ✓ Q1: How to find 'concepts' in a document collection?
- ✓ Q2: how to answer queries in English, when documents are in Spanish?
- ✓ Q3: how to compress a customer x day matrix
- Q4: how to interpret the rules/concepts
- Q5: KL transform?



15-826 Copyright (c) 2019 C. Faloutsos 34

34

CarnegieMellon

PCA - 'Ratio Rules'

[Korn+98]

Typically: 'Association Rules' (eg.,

{bread, milk} -> {butter}

But, can we discover more details? like:

\$-bread : \$-milk : \$-butter ~ \$2 : \$4 : \$3

Flip Korn, Alexandros Labrinidis, Yannis Kotidis, and Christos Faloutsos. *Ratio Rules: A New Paradigm for Fast, Quantifiable Data Mining*. (VLDB '98), 582-593.

35

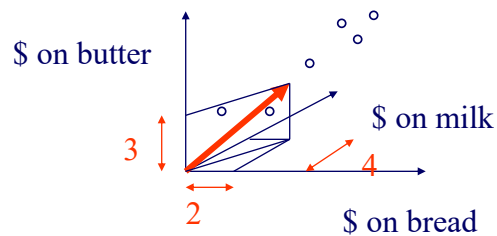
CarnegieMellon

PCA - 'Ratio Rules'

Idea: try to find 'concepts' :

- singular vectors dictate rules about ratios:

bread:milk:butter = 2:4:3



15-826

Copyright (c) 2019 C. Faloutsos

36

36

CarnegieMellon

PCA - 'Ratio Rules'

Identical to PCA = Principal Components Analysis

- ✓ – Q1: which set of rules is 'better' ?
- ✓ – Q2: how to reconstruct missing/corrupted values?
- ✓ – Q3: is there need for binary/bucketized values? **NO**
- ➔ – Q4: how to interpret the rules (= 'principal components')?

15-826

Copyright (c) 2019 C. Faloutsos

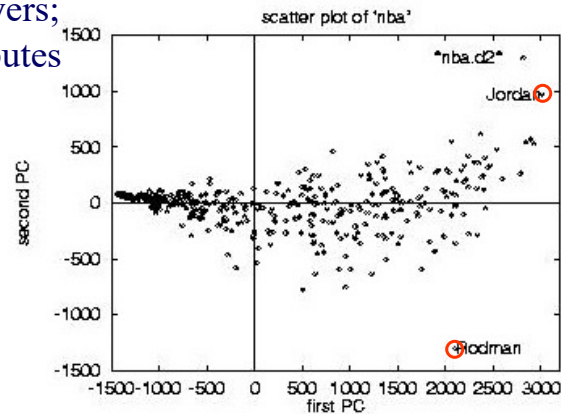
37

37

CarnegieMellon

PCA - Ratio Rules

NBA dataset
~500 players;
~30 attributes



15-826

Copyright (c) 2019 C. Faloutsos

38

38

PCA - Ratio Rules

- PCA: get singular vectors v_1, v_2, \dots
- ignore entries with small abs. value
- try to interpret the rest

PCA - Ratio Rules

NBA dataset - V matrix (term to 'concept' similarities)

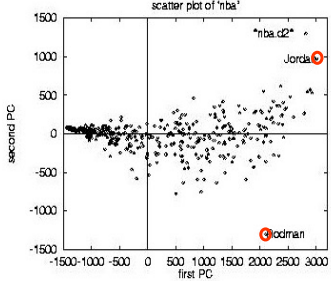
<i>field</i>	RR_1	RR_2	RR_3
minutes played	.808	-.4	
field goals			
goal attempts			
points	.406	.199	
total rebounds		-.489	.602
assists			-.486
steals			-.07

v_1

Carnegie Mellon

Ratio Rules - example

- RR1: minutes:points = 2:1
- corresponding concept?



→ v1

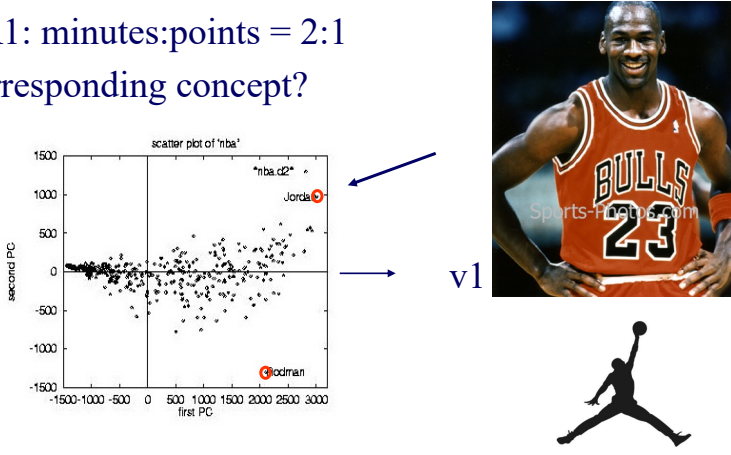
15-826 Copyright (c) 2019 C. Faloutsos 41

41

Carnegie Mellon

Ratio Rules - example

- RR1: minutes:points = 2:1
- corresponding concept?



→ v1

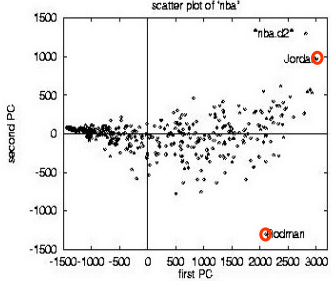

15-826 Copyright (c) 2019 C. Faloutsos 42

42

Carnegie Mellon

Ratio Rules - example

- RR1: minutes:points = 2:1
- CO
- CO

15-826 Copyright (c) 2019 C. Faloutsos

43

Carnegie Mellon

Ratio Rules - example

- RR1: minutes:points = 2:1
- corresponding concept?
- A: 'goodness' of player

15-826 Copyright (c) 2019 C. Faloutsos 44

44

CarnegieMellon

Ratio Rules - example

- RR2: points: rebounds negatively correlated(!)

<i>field</i>	RR_1	RR_2	RR_3
minutes played	.808	-.4	
field goals			
goal attempts			
points	.406	.199	
total rebounds		-.489	.602
assists			-.486
steals			-.07

15-826
Copyright (c) 2019 C. Faloutsos
45

45

CarnegieMellon

Ratio Rules - example

- RR2: points: rebounds negatively correlated(!) - concept?

v2

↑

15-826
Copyright (c) 2019 C. Faloutsos
46

46

Ratio Rules - example

- RR2: points: rebounds negatively correlated(!) - concept?
- A: position: offensive/defensive

Solutions

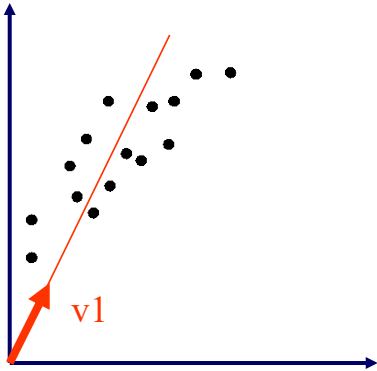


- ✓ Q1: How to find 'concepts' in a document collection?
- ✓ Q2: how to answer queries in English, when documents are in Spanish?
- ✓ Q3: how to compress a customer x day matrix
- ✓ Q4: how to interpret the rules/concepts
- Q5: KL transform?



Carnegie Mellon

K-L transform



[Duda & Hart]; [Fukunaga]

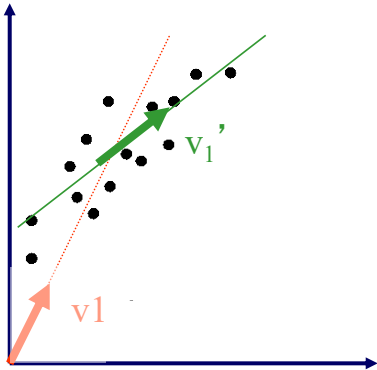
A subtle point:
SVD will give vectors that
go through the origin

15-826 Copyright (c) 2019 C. Faloutsos 50

50

Carnegie Mellon

K-L transform



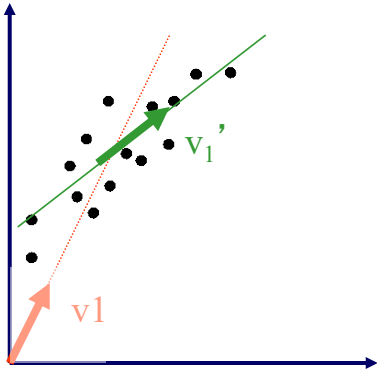
A subtle point:
SVD will give vectors that
go through the origin
Q: how to find v_1' ?

15-826 Copyright (c) 2019 C. Faloutsos 51

51

Carnegie Mellon

K-L transform



A subtle point:
SVD will give vectors that go through the origin
Q: how to find v_1' ?

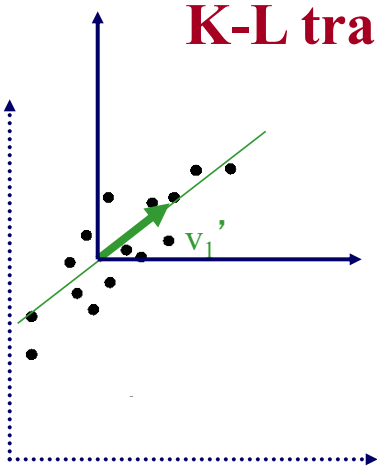
A: 'centered' PCA, ie.,
move the origin to center of gravity

15-826 Copyright (c) 2019 C. Faloutsos 52

52

Carnegie Mellon

K-L transform



A subtle point:
SVD will give vectors that go through the origin
Q: how to find v_1' ?

A: 'centered' PCA, ie.,
move the origin to center of gravity
and THEN do SVD

15-826 Copyright (c) 2019 C. Faloutsos 53

53

K-L transform


- How to ‘center’ a set of vectors (= data matrix)?
- What is the covariance matrix?
- A: see textbook
- (‘whitening transformation’)

Conclusions



- SVD: popular for dimensionality reduction / compression
- SVD is the ‘engine under the hood’ for PCA (principal component analysis)
- ... as well as the Karhunen-Lowe transform
- (and there is more to come ...)

CarnegieMellon

Solutions



- ✓ Q1: How to find ‘concepts’ in a document collection?
- ✓ Q2: how to analyze text in English, when documents are in other languages?
- ✓ Q3: how to analyze a customer x day matrix?
- ✓ Q4: how to interpret the rules/concepts?
- ✓ Q5: KL transform?

15-826 Copyright (c) 2019 C. Faloutsos 56

56

CarnegieMellon

References

- Duda, R. O. and P. E. Hart (1973). Pattern Classification and Scene Analysis. New York, Wiley.
- Fukunaga, K. (1990). Introduction to Statistical Pattern Recognition, Academic Press.
- Jolliffe, I. T. (1986). Principal Component Analysis, Springer Verlag.

15-826 Copyright (c) 2019 C. Faloutsos 57

57

References

- Korn, F., H. V. Jagadish, et al. (May 13-15, 1997). Efficiently Supporting Ad Hoc Queries in Large Datasets of Time Sequences. ACM SIGMOD, Tucson, AZ.
- Korn, F., A. Labrinidis, et al. (1998). Ratio Rules: A New Paradigm for Fast, Quantifiable Data Mining. VLDB, New York, NY.

References

- [Korn+, '00] Korn, F., A. Labrinidis, et al. (2000). "Quantifiable Data Mining Using Ratio Rules." VLDB Journal 8(3-4): 254-266.
- Press, W. H., S. A. Teukolsky, et al. (1992). Numerical Recipes in C, Cambridge University Press.