# INTRO TO DATA SCIENCE
## LECTURE 13: DIMENSIONALITY REDUCTION

I. DIMENSIONALITY REDUCTION
II. PRINCIPAL COMPONENTS ANALYSIS
III. SINGULAR VALUE DECOMPOSITION
IV. OTHER METHODS

EXERCISE:
IV. DIMENSIONALITY REDUCTION IN SCIKIT-LEARN

# I. DIMENSIONALITY REDUCTION

*Q: What is dimensionality reduction?*

Q: What is dimensionality reduction?

A: A set of techniques for reducing the size (in terms of features, records, and/or bytes) of the dataset under examination.

Q: What is dimensionality reduction?

A: A set of techniques for reducing the size (in terms of features, records, and/or bytes) of the dataset under examination.

In general, the idea is to regard the dataset as a matrix and to decompose the matrix into simpler, meaningful pieces.

*Q: What is dimensionality reduction?*

*A: A set of techniques for reducing the size (in terms of features, records, and/or bytes) of the dataset under examination.*
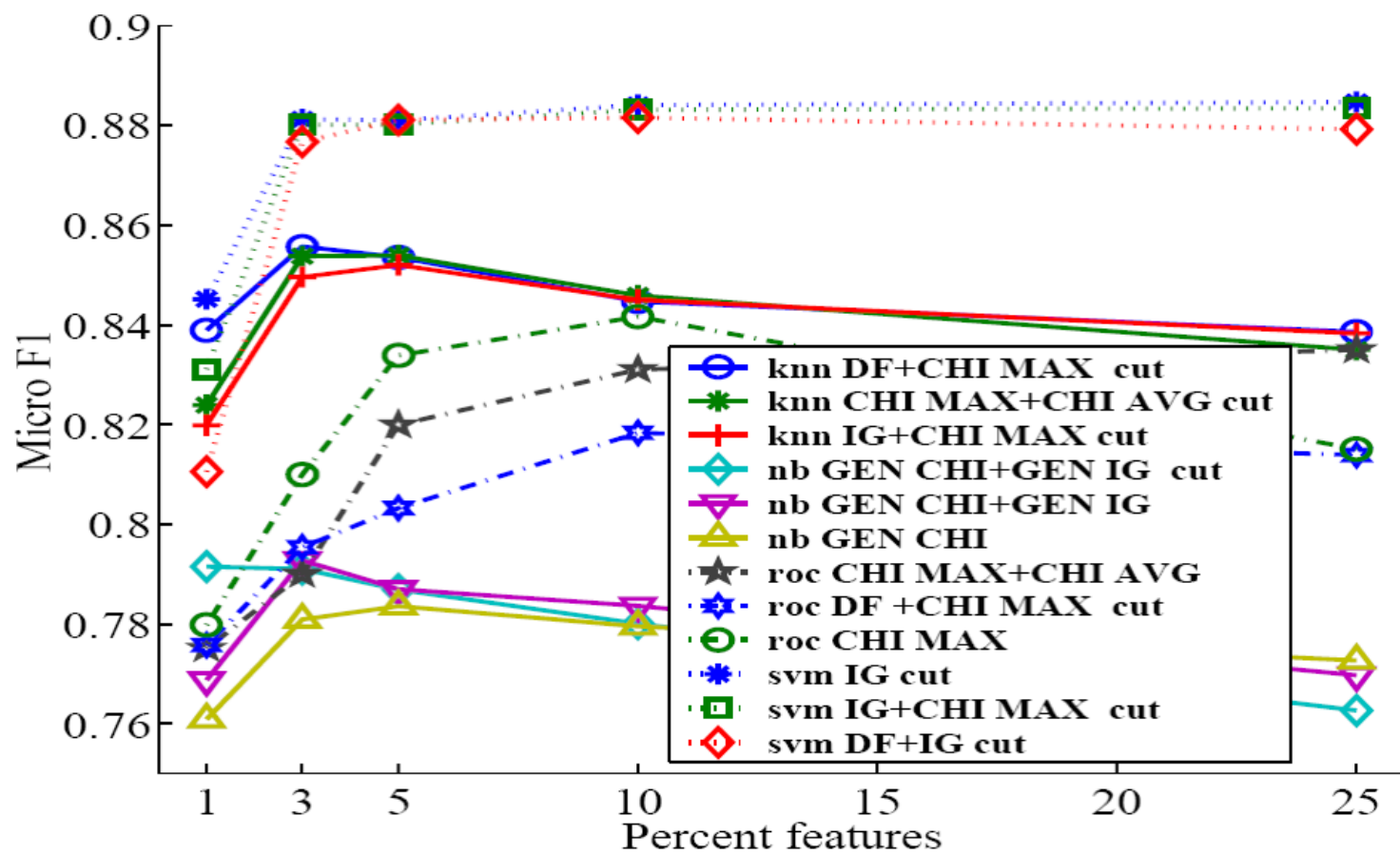
*In general, the idea is to regard the dataset as a matrix and to decompose the matrix into simpler, meaningful pieces.*

*Dimensionality reduction is frequently performed as a pre-processing step before another learning algorithm is applied.*

*Q: What are the goals of dimensionality reduction?*

*Q: What are the goals of dimensionality reduction?*

*– reduce computational expense*

*– reduce susceptibility to overfitting*

*– reduce noise in the dataset*

*– enhance our intuition*

*Q: How is dimensionality reduction performed?*

Q: How is dimensionality reduction performed?

A: There are two approaches: feature selection and feature extraction.

Q:  *How is dimensionality reduction performed?*
A:  *There are two approaches: feature selection and feature extraction.*

**feature selection** *– selecting a subset of features using an external criterion (filter) or the learning algo accuracy itself (wrapper)*

**feature extraction** *– mapping the features to a lower dimensional space*

# II. PRINCIPAL COMPONENT ANALYSIS

*Principal component analysis is a dimension reduction technique that can be used on a matrix of any dimensions.*
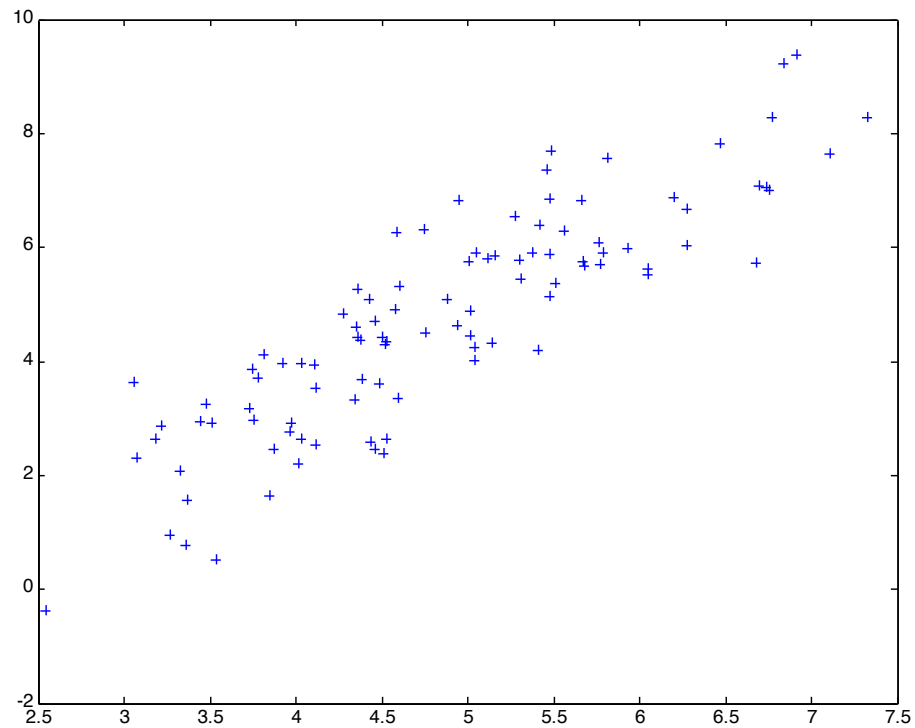
*Principal component analysis is a dimension reduction technique that can be used on a matrix of any dimensions.*

*This procedure produces a new basis, each of whose components retain as much variance from the original data as possible.*
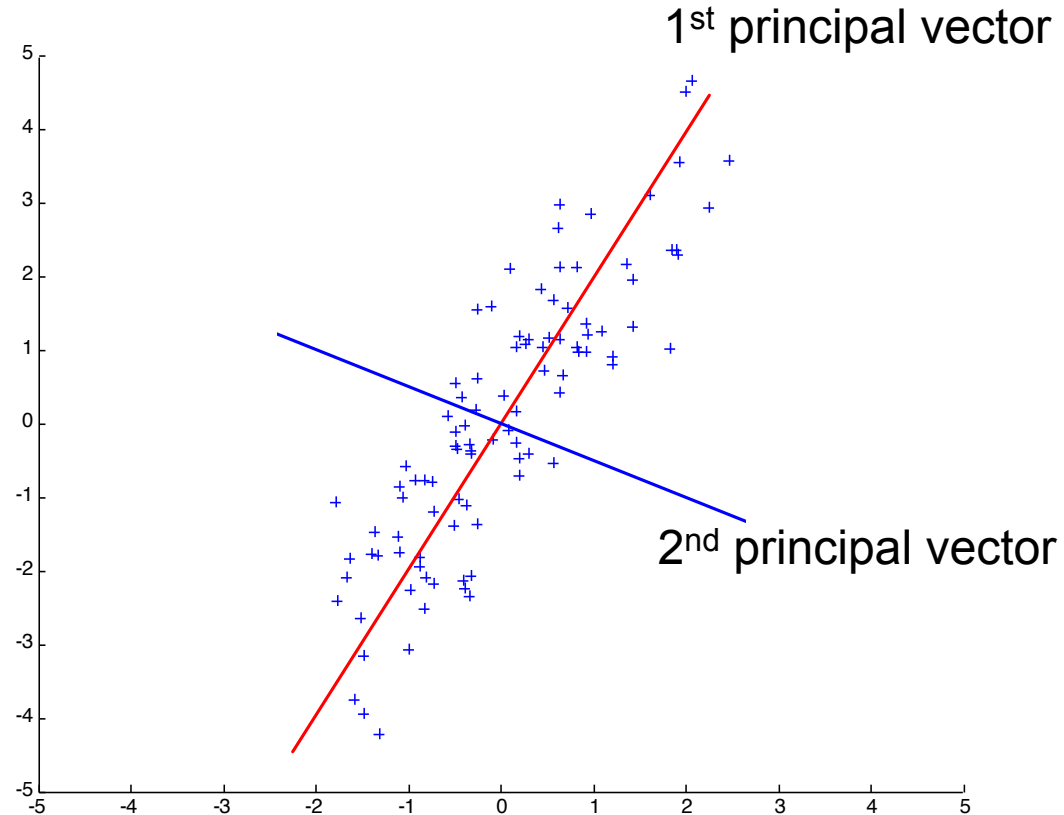
*Principal component analysis is a dimension reduction technique that can be used on a matrix of any dimensions.*
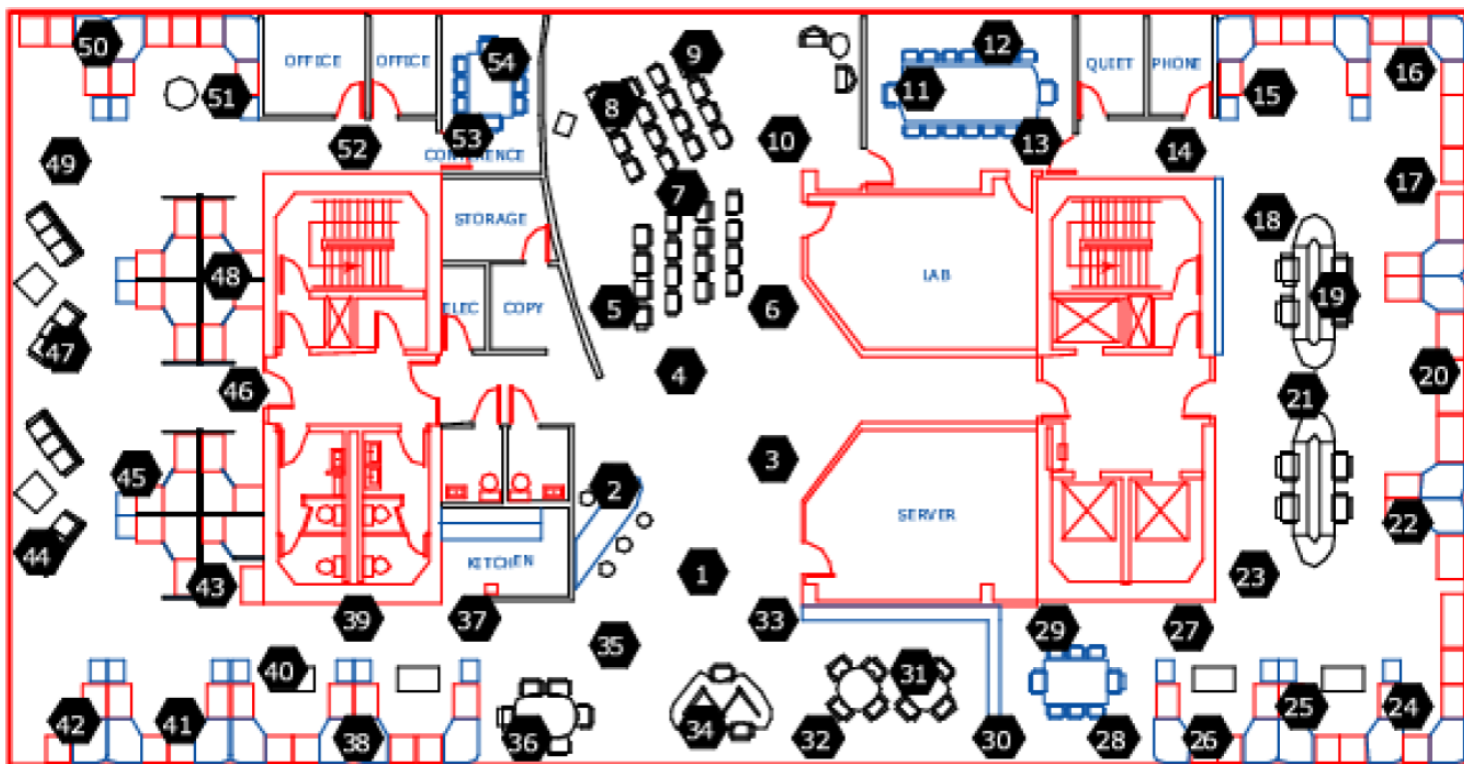
*This procedure produces a new basis, each of whose components retain as much variance from the original data as possible.*

*The PCA of a matrix A boils down to the* **eigenvalue decomposition** *of the* **covariance matrix** *of A.*

- Gives best axis to project
- Minimum RMS error
- Principal vectors are <span style="color:red">orthogonal</span>



1st principal vector

2nd principal vector

Sensors in Intel Berkeley Lab

Link quality vs. Distance between a pair of sensors
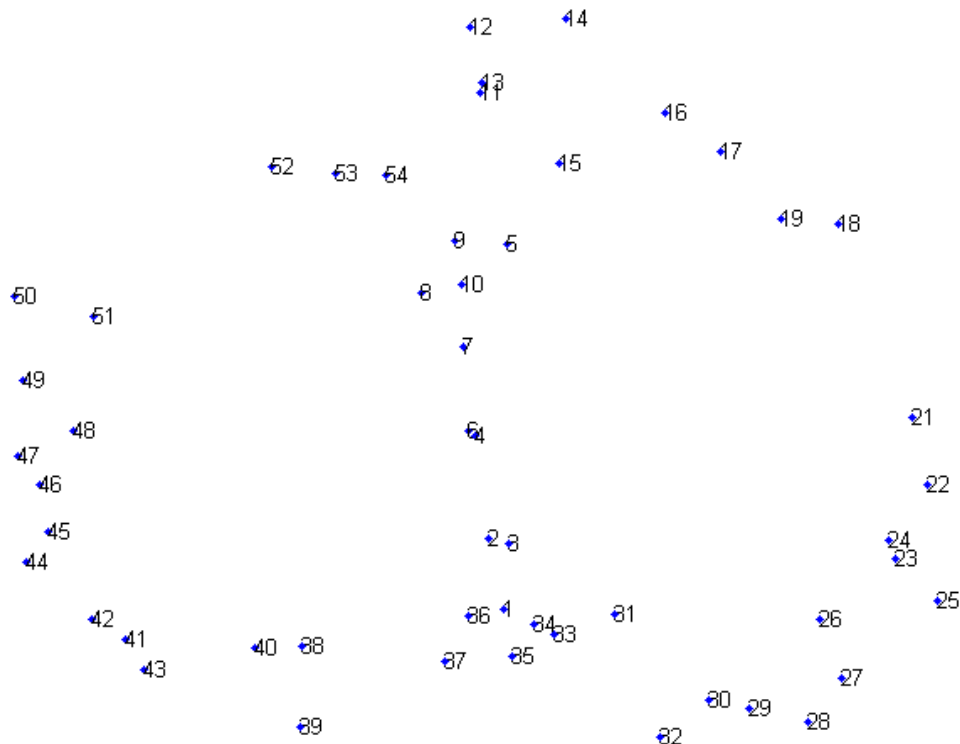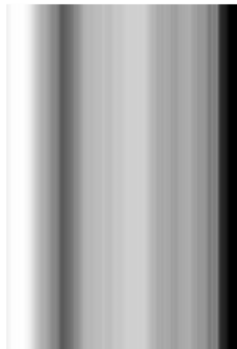
- Given a 54x54 matrix of pairwise link qualities
- Do PCA
- Project down to 2 principal dimensions
- PCA discovered the map of the lab

PCs # 0     PCs # 10     PCs # 20

PCs # 30     PCs # 40     PCs # 50

Source: http://glowingpython.blogspot.it/2011/07/pca-and-image-compression-with-numpy.html

- ## PCA algorithm:
  - 1. $X \leftarrow$ Create N x d data matrix, with one row vector $x_n$ per data point
  - 2. X  subtract mean $x$ from each row vector $x_n$ in X
  - 3. $\Sigma \leftarrow$ covariance matrix of X
  - Find eigenvectors and eigenvalues of $\Sigma$
  - PC's $\leftarrow$ the M eigenvectors with largest eigenvalues

- ## What if very large dimensional data?
  - e.g., Images ($d \geq 10^4$)
- ## Problem:
  - Covariance matrix $\Sigma$ is size ($d^2$)
  - $d=10_4 \rightarrow |\Sigma| = 10^8$

- ## Singular Value Decomposition (SVD)!
  - efficient algorithms available
  - some implementations find just top N eigenvectors

# III. SINGULAR VALUE DECOMPOSITION

# Singular Value Decomposition

- Problem:
  - #1: Find concepts in text
  - #2: Reduce dimensionality

| term / document | data | information | retrieval | brain | lung |
|---|---|---|---|---|---|
| CS-TR1 | 1 | 1 | 1 | 0 | 0 |
| CS-TR2 | 2 | 2 | 2 | 0 | 0 |
| CS-TR3 | 1 | 1 | 1 | 0 | 0 |
| CS-TR4 | 5 | 5 | 5 | 0 | 0 |
| MED-TR1 | 0 | 0 | 0 | 2 | 2 |
| MED-TR2 | 0 | 0 | 0 | 3 | 3 |
| MED-TR3 | 0 | 0 | 0 | 1 | 1 |

# SVD - Definition

$$A_{[n \times m]} = U_{[n \times r]} \Lambda_{[r \times r]} (V_{[m \times r]})^{\top}$$

- **A**: *n x m* matrix (e.g., n documents, m terms)
- **U**: *n x r* matrix (n documents, r concepts)
- $\Lambda$: *r x r* diagonal matrix (strength of each 'concept') (r: rank of the matrix)
- **V**: m x r matrix (m terms, r concepts)

# SVD - Properties

**THEOREM** [Press+92]: always possible to decompose matrix **A** into **A** = **U** $\Lambda$ **V**$^T$ , where

- **U,** $\Lambda$, **V**: unique (*)

- **U**, **V**: column orthonormal (ie., columns are unit vectors, orthogonal to each other)
    - **U$^T$U = I**; **V$^T$V = I (I:** identity matrix**)**

- $\Lambda$**:** singular value are positive, and sorted in decreasing order

# SVD - Interpretation

'documents', 'terms' and 'concepts':

- **U**: document-to-concept similarity matrix

- **V**: term-to-concept similarity matrix

- $\Lambda$: its diagonal elements: 'strength' of each concept

Projection:

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

# SVD - Example

- **A = U Λ V$^T$** - example:

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

Column labels: data, inf. retrieval, brain, lung

# SVD - Example

- **A** = **U** $\Lambda$ **V**$^T$ - example:

doc-to-concept similarity matrix

$$
\begin{array}{c}
\\
\text{data} \quad \substack{\text{inf.} \\ } \quad \substack{\text{retrieval} \\ \text{brain}} \quad \text{lung}
\end{array}
$$

CS-concept

MD-concept

$$
\begin{array}{c}
\text{CS} \\[3em]
\text{MD}
\end{array}
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 2 & 2 \\
0 & 0 & 0 & 3 & 3 \\
0 & 0 & 0 & 1 & 1
\end{bmatrix}
=
\begin{bmatrix}
0.18 & 0 \\
0.36 & 0 \\
0.18 & 0 \\
0.90 & 0 \\
0 & 0.53 \\
0 & 0.80 \\
0 & 0.27
\end{bmatrix}
\;x\;
\begin{bmatrix}
9.64 & 0 \\
0 & 5.29
\end{bmatrix}
\;x\;
\begin{bmatrix}
0.58 & 0.58 & 0.58 & 0 & 0 \\
0 & 0 & 0 & 0.71 & 0.71
\end{bmatrix}
$$

# SVD - Example

- **A = U Λ V$^T$** - example:

$$
\begin{array}{c}
\text{data} \\
\end{array}
\begin{array}{cc}
\text{inf.} & \text{retrieval} \\
& \text{brain} \quad \text{lung}
\end{array}
$$

'strength' of CS-concept

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

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

# SVD - Example

- **A** = **U** $\Lambda$ **V**$^T$ - example:

term-to-concept
similarity matrix

$$
\begin{array}{c}
\text{retrieval} \\
\text{inf.} \quad \text{lung} \\
\text{data} \quad \text{brain}
\end{array}
$$

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

CS-concept

# SVD – Dimensionality reduction

- Q: how exactly is dim. reduction done?
- A: set the smallest singular values to zero:

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

# SVD - Dimensionality reduction

$$
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 2 & 2 \\
0 & 0 & 0 & 3 & 3 \\
0 & 0 & 0 & 1 & 1
\end{bmatrix}
\sim
\begin{bmatrix}
0.18 \\
0.36 \\
0.18 \\
0.90 \\
0 \\
0 \\
0
\end{bmatrix}
\times
\begin{bmatrix}
9.64
\end{bmatrix}
\times
\begin{bmatrix}
0.58 & 0.58 & 0.58 & 0 & 0
\end{bmatrix}
$$

# SVD - Dimensionality reduction

$$
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 2 & 2 \\
0 & 0 & 0 & 3 & 3 \\
0 & 0 & 0 & 1 & 1
\end{bmatrix}
\sim
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0
\end{bmatrix}
$$

# LSI (latent semantic indexing)

Q1: How to do queries with LSI?

A: map query vectors into 'concept space' – how?

$$
\begin{array}{c}
\text{data} \quad \overset{\text{inf.}}{\underset{}{\text{retrieval}}} \quad \text{brain} \quad \text{lung} \\
\begin{array}{c} \text{CS} \\ \\ \\ \text{MD} \end{array}
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 2 & 2 \\
0 & 0 & 0 & 3 & 3 \\
0 & 0 & 0 & 1 & 1
\end{bmatrix}
=
\begin{bmatrix}
0.18 & 0 \\
0.36 & 0 \\
0.18 & 0 \\
0.90 & 0 \\
0 & 0.53 \\
0 & 0.80 \\
0 & 0.27
\end{bmatrix}
\times
\begin{bmatrix}
9.64 & 0 \\
0 & 5.29
\end{bmatrix}
\times
\begin{bmatrix}
0.58 & 0.58 & 0.58 & 0 & 0 \\
0 & 0 & 0 & 0.71 & 0.71
\end{bmatrix}
\end{array}
$$

# LSI (latent semantic indexing)

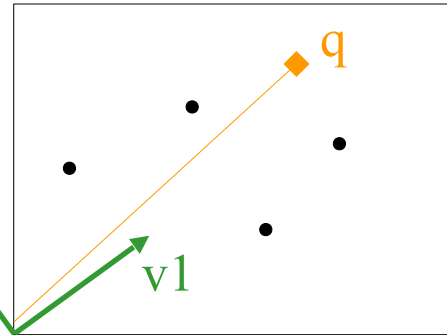Q: How to do queries with LSI?

A: map query vectors into 'concept space' – how?



$q=$

| | data | inf. | retrieval brain | lung |
|---|---|---|---|---|
| | 1 | 0 | 0 0 | 0 |

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

# LSI (latent semantic indexing)

compactly, we have:

$q_{concept}$ = q **V**

e.g.:

$$q = \begin{array}{c} \text{data} \quad \text{inf.} \quad \substack{\text{retrieval} \\ |} \quad \text{brain} \quad \text{lung} \end{array} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix} \overset{\text{CS-concept}}{=} \begin{bmatrix} 0.58 & 0 \end{bmatrix}$$

term-to-concept
similarities

# Multi-lingual IR
## (English query, on Spanish text?)

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

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

# Little example

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

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

Eg:

$$
d = \begin{array}{c} \overset{\text{data}}{} \overset{\text{inf.}}{\overset{\text{retrieval}}{|}} \overset{\text{brain}}{} \overset{\text{lung}}{} \\ \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \end{bmatrix} \end{array}
\begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix}
= \begin{bmatrix} 1.16 & 0 \end{bmatrix}
$$

CS-concept

term-to-concept similarities

# Little example

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

CS-concept

$$d= \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \end{bmatrix} \quad \cdots \cdots \cdots \quad \begin{bmatrix} 1.16 & 0 \end{bmatrix}$$

$$q= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \quad \cdots \cdots \cdots \quad \begin{bmatrix} 0.58 & 0 \end{bmatrix}$$

(columns: data, inf., retrieval, brain, lung)

# Multi-lingual IR

- Solution: ~ LSI

- Concatenate documents
- Do SVD on them
- Now when a new document comes project it into concept space
- Measure similarity in concept space

$$
\begin{array}{c}
\text{CS} \\
\\
\text{MD}
\end{array}
\left[
\begin{array}{ccccc|ccccc}
1 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 & 1 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 & 5 & 5 & 4 & 0 & 0 \\
0 & 0 & 0 & 2 & 2 & 0 & 0 & 0 & 2 & 2 \\
0 & 0 & 0 & 3 & 3 & 0 & 0 & 0 & 2 & 3 \\
0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 1
\end{array}
\right]
$$

data  inf.  retrieval  brain  lung     informacion  datos