



Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

# Processamento e Recuperação de Informação

## Document Clustering and Dimensionality Reduction

Departamento de Engenharia Informática  
Instituto Superior Técnico

1º Semestre  
2018/2019



# Bibliography

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

- Bing Liu, Web Data Mining - Exploring Hyperlinks, Contents, and Usage Data. Chapter 4.
- Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Chapters 16 and 17.
- Ricardo Baeza-Yates and Berthier Ribeiro-Neto, Modern Information Retrieval, Chapter 2.
- Jure Leskovec, Anand Rajaraman, and Jeff Ullman, Mining of Massive Datasets, Chapters 7 and 11



# Outline

Processamento  
e Recuperação  
de Informação

## Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- 1 Introduction
  - Motivation
  - Basic Concepts
- 2 Partitioning Approaches
- 3 Dimensionality Reduction



# Motivation

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- **Problem:** Query terms can be ambiguous
  - E.g., query "star" retrieves documents about astronomy, animals, etc.



# Motivation

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- **Problem:** Query terms can be ambiguous
  - E.g., query "star" retrieves documents about astronomy, animals, etc.
  - Solution: **clustering** document responses to queries along lines of different topics



# Motivation

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- **Problem:** Query terms can be ambiguous
  - E.g., query "star" retrieves documents about astronomy, animals, etc.
  - Solution: **clustering** document responses to queries along lines of different topics
- **Problem:** Manual construction of topic hierarchies and taxonomies



# Motivation

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- **Problem:** Query terms can be ambiguous
  - E.g., query "star" retrieves documents about astronomy, animals, etc.
  - Solution: **clustering** document responses to queries along lines of different topics
- **Problem:** Manual construction of topic hierarchies and taxonomies
  - Solution: preliminary **clustering** of large samples of web documents



# Motivation

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- **Problem:** Query terms can be ambiguous
  - E.g., query "star" retrieves documents about astronomy, animals, etc.
  - Solution: **clustering** document responses to queries along lines of different topics
- **Problem:** Manual construction of topic hierarchies and taxonomies
  - Solution: preliminary **clustering** of large samples of web documents
- **Problem:** Speeding up similarity search





# Motivation

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- **Problem:** Query terms can be ambiguous
  - E.g., query "star" retrieves documents about astronomy, animals, etc.
  - Solution: **clustering** document responses to queries along lines of different topics
- **Problem:** Manual construction of topic hierarchies and taxonomies
  - Solution: preliminary **clustering** of large samples of web documents
- **Problem:** Speeding up similarity search
  - Solution: restrict the search for documents similar to a query to most representative **cluster(s)**



# An Example

Processamento  
e Recuperação  
de Informação

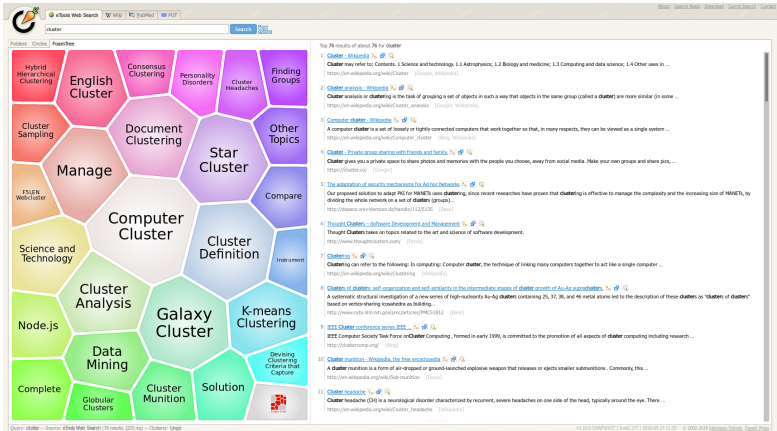
Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction



<http://search.carrot2.org>



# Another Example

Processamento  
e Recuperação  
de Informação

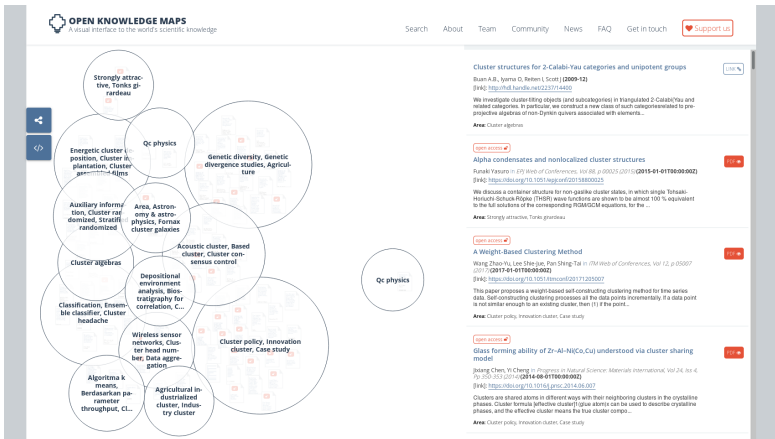
Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction



<https://openknowledgemaps.org/>



# Clustering

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

## Cluster Hypothesis:

- Given a 'suitable' clustering of a collection, if the user is interested in document  $d$  (or term  $t$ ), he is likely to be interested in other members of the cluster to which  $d$  ( $t$ ) belongs

## Clustering Task:

- Use measures of similarity to **cluster** a collection of documents/terms into groups, so that **similarity within a cluster is larger than across clusters**



# Clustering Concepts

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- Clustering paradigms:
  - **Bottom-up** agglomerative clustering
  - **Top-down** partitioning
- Dimensionality reduction:
  - Embedding of corpus in a low-dimensional space
  - Many different approaches, based on heuristics, linear algebra, probabilistic models, ...



# Clustering Approaches

Processamento  
e Recuperação  
de Informação

Introduction

Motivation

Basic Concepts

Partitioning  
Approaches

Dimensionality  
Reduction

- Partitioning Approaches
  - Bottom-up clustering
  - Top-down clustering
- Geometric Embedding Approaches
  - Self-organization map
  - Latent semantic indexing
- Generative models and probabilistic approaches
  - Single topic per document
  - Documents correspond to mixtures of multiple topics



# Outline

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction

## 1 Introduction

## 2 Partitioning Approaches

- Bottom-Up Clustering
- Top-Down Clustering

## 3 Dimensionality Reduction



# Partitioning Approaches

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction

- Partition document collection into a set  $G$  of  $k$  clusters  $[C_1, C_2, \dots, C_k]$
- Choices:
  - Minimize intra-cluster distance:  $\sum_i \sum_{d_1, d_2 \in C_i} \delta(d_1, d_2)$
  - Maximize intra-cluster semblance:  $\sum_i \sum_{d_1, d_2 \in C_i} \rho(d_1, d_2)$
- If cluster representations  $\vec{C}_i$  are available
  - Minimize  $\sum_i \sum_{d \in C_i} \delta(d, \vec{C}_i)$
  - Maximize  $\sum_i \sum_{d \in C_i} \rho(d, \vec{C}_i)$
- Soft clustering
  - $d$  assigned to  $C_i$  with confidence  $z_{d,i}$
  - Find  $z_{d,i}$  so as to minimize  $\sum_i \sum_{d \in C_i} z_{d,i} \delta(d, \vec{C}_i)$  or maximize  $\sum_i \sum_{d \in C_i} z_{d,i} \rho(d, \vec{C}_i)$
- Two ways to get partitions: **bottom-up clustering** and **top-down clustering**





# Hierarchical Agglomerative Clustering

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction

- Consider that  $G$  represents a set of clusters
- Initially  $G$  is a collection of **singleton groups**, each with one document  $d$
- Repeat
  - Find  $\Gamma, \Delta \in G$  with max similarity measure,  $s(\Gamma \cup \Delta)$
  - Merge group  $\Gamma$  with group  $\Delta$
- Use above info to plot the hierarchical merging process (**dendrogram**)
- To get desired number of clusters: cut across any level of the dendrogram



# Example Dendrogram

Processamento  
e Recuperação  
de Informação

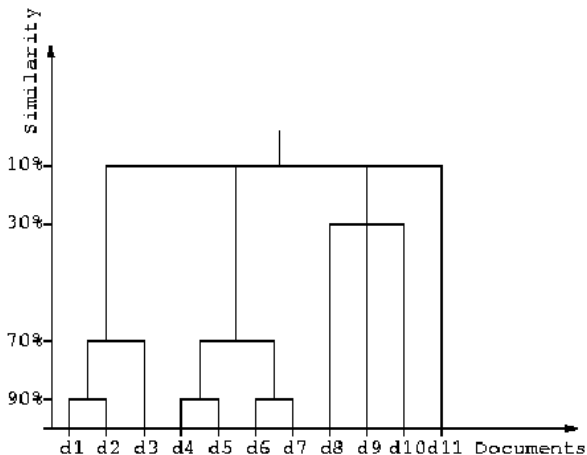
Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction



(b)

(example)



# Similarity Measures

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction

- Self-Similarity
  - Consider that  $\Phi$  represents a cluster
  - Average pairwise similarity between documents in  $\Phi$

$$s(\Phi) = \frac{1}{\binom{|\Phi|}{2}} \sum_{d_1, d_2 \in \Phi} s(d_1, d_2)$$

- $s(d_1, d_2)$  is the inter-document similarity measure (e.g., cosine of *TF-IDF* vectors)
- Other criteria:
  - Maximum/minimum pairwise similarity between documents in the clusters
- Complexity:  $O(n^2 \log n)$  with  $n^2$  space



# Top-Down Clustering

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction

- Use an internal representation for documents as well as clusters (**centroids**)
- Partition documents into  $k$  clusters
- 2 variants
  - **Hard**: 0/1 assignment of documents to clusters
  - **Soft**: documents belong to clusters with fractional scores
- Termination
  - When assignment of documents to clusters ceases to change much, or
  - when cluster centroids move negligibly over successive iterations



# The $k$ -Means Algorithm

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction

## Hard $k$ -Means

- Choose  $k$  arbitrary centroids
- Assign each document to nearest centroid
- Recompute centroids
- Contribution for updating cluster centroid

$$\Delta\mu_c = \sum_d \begin{cases} \eta(d - \mu_c) & \text{if } \mu_c \text{ is closest to } d \\ 0 & \text{otherwise} \end{cases}$$

$$\mu_c \leftarrow \mu_c + \Delta\mu_c$$



# The $k$ -Means Algorithm

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction

## Soft $k$ -Means

- Don't break close ties between document assignments to clusters and don't make documents contribute to a single cluster which wins narrowly
  - The contribution for updating cluster centroid  $\mu_c$  from document  $d$  is related to the current similarity between both

$$\Delta\mu_c = \eta \frac{1/|d - \mu_c|^2}{\sum_{\gamma} 1/|d - \mu_{\gamma}|^2} (d - \mu_c)$$



# An Example

Processamento  
e Recuperação  
de Informação

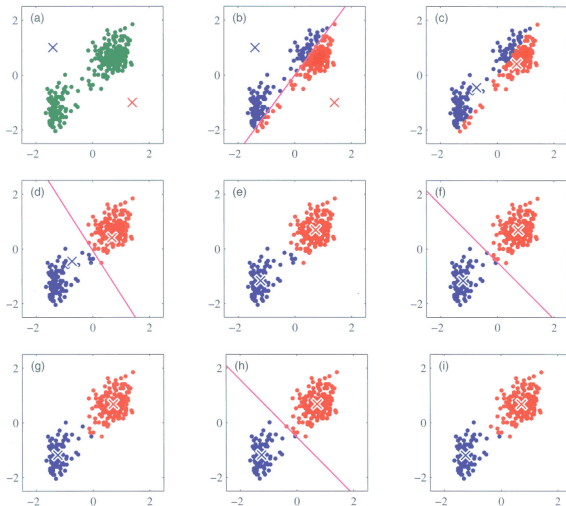
Introduction

Partitioning  
Approaches

Bottom-Up  
Clustering

Top-Down Clustering

Dimensionality  
Reduction



(online demo)



# Outline

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

## 1 Introduction

## 2 Partitioning Approaches

## 3 Dimensionality Reduction

- Self-Organizing Maps
- Latent Semantic Indexing





# Dimensionality Reduction

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

- **Goal:** Embedding of corpus in a low-dimensional space
- Self-Organizing Map (SOM)
  - Technique related to  $k$ -means
- Latent Semantic Indexing (LSI)
  - Linear transformations to reduce number of dimensions
- Other techniques:
  - Multidimensional scaling (MDS): Minimize the distortion of interpoint distances in the low-dimensional embedding as compared to the dissimilarity given in the input data
  - NN-based word embeddings: Use neural networks to learn a mapping from the high dimensional vocabulary space to a lower dimension concept-based space



# Self-Organizing Maps

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

- Like soft  $k$ -means
  - Determine association between clusters and documents
  - Associate a representative vector with each cluster and iteratively refine
- Unlike  $k$ -means
  - Embed the clusters in a low-dimensional space right from the beginning
  - A large number of clusters can be initialized even if eventually many are to remain devoid of documents
- Each cluster can be a slot in a square/hexagonal grid.
  - The grid structure defines the neighborhood  $N(c)$  for each cluster  $c$
- Also involves a proximity function  $h(\gamma, c)$  between clusters



# Update Rule

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

- Data item  $d$  activates node (closest cluster) as well as the neighborhood nodes
  - E.g., all nodes within  $n$  hops
- Update rule for node under the influence of  $d$  is:

$$\mu_{\gamma} \leftarrow \mu_{\gamma} + \eta h(\gamma, c_d)(d - \mu_{\gamma})$$

# A Self Organizing Map

Processamento  
e Recuperação  
de Informação

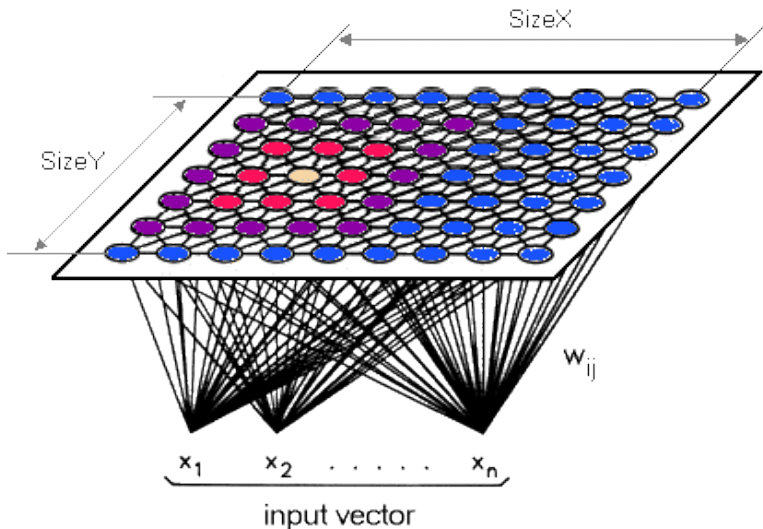
Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing





# Example

Processamento  
e Recuperação  
de Informação

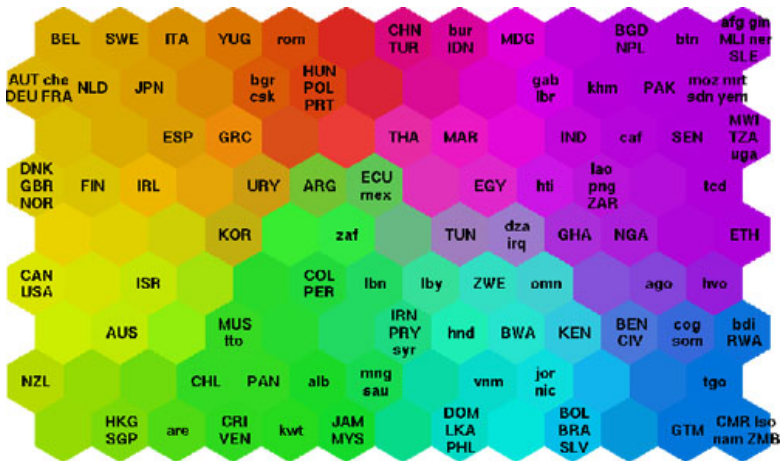
Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing





# Example

Processamento  
e Recuperação  
de Informação

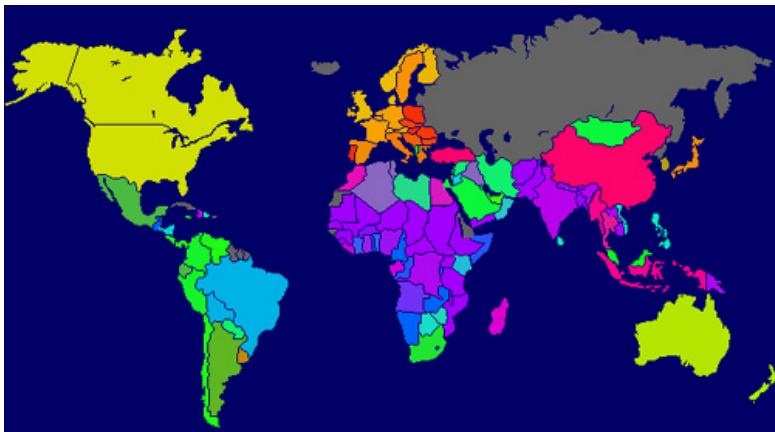
Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing



(online demo)



# Extended Similarity and LSI

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

- In the vector space model different terms have different meanings
  - car  $\neq$  automobile
- However, “car” and “automobile” are related
  - They are likely to co-occur often
- Documents having related words are related
  - Useful for search and clustering
- Two basic approaches:
  - Hand-made thesaurus (e.g., **WordNet**)
  - Co-occurrence and associations



# Latent Semantic Indexing

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

- Vector-space model
  - Distinct orthogonal direction for each term
- Not all terms are orthogonal
  - E.g., “car” and “automobile”
- We need a matrix with a lower rank than the traditional document-term matrix, where similar terms are “merged”
- This matrix can be found through **Singular Value Decomposition**





# Singular Value Decomposition

Processamento  
e Recuperação  
de Informação

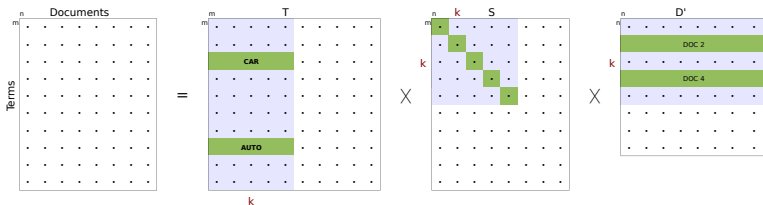
Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing





# Querying

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

- $A = T \times S \times D^T \Leftrightarrow D = A^T \times T \times S^{-1}$
- A query  $q$  is also projected to the new space

$$\hat{q} = q^T \times T \times S^{-1}$$

- Cosine similarity can now be applied to the projections
- Results are often better than standard *TF-IDF* retrieval/classification
  - SVD filters out noise and “discovers” semantic associations between terms
  - Representations in the projected vector space can be used for visualization



# Example

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

Label	Titles
B1	A Course on <u>Integral Equations</u>
B2	Attractors for Semigroups and Evolution <u>Equations</u>
B3	Automatic Differentiation of <u>Algorithms: Theory, Implementation, and Application</u>
B4	Geometrical Aspects of <u>Partial Differential Equations</u>
B5	Ideals, Varieties, and <u>Algorithms</u> An Introduction to Computational Algebraic Geometry and Commutative Algebra
B6	<u>Introduction</u> to Hamiltonian Dynamical Systems and the <u>N-Body Problem</u>
B7	Knapsack <u>Problems: Algorithms and Computer Implementations</u>
B8	<u>Methods of Solving Singular Systems of Ordinary Differential Equations</u>
B9	<u>Nonlinear Systems</u>
B10	<u>Ordinary Differential Equations</u>
B11	<u>Oscillation Theory</u> for Neutral Differential <u>Equations with Delay</u>
B12	<u>Oscillation Theory of Delay Differential Equations</u>
B13	Pseudodifferential Operators and <u>Nonlinear Partial Differential Equations</u>
B14	Sinc Methods for Quadrature and Differential <u>Equations</u>
B15	Stability of Stochastic <u>Differential Equations</u> with Respect to Semi Martingales
B16	The Boundary <u>Integral</u> Approach to Static and Dynamic Contact <u>Problems</u>
B17	The Double Mellin-Barnes Type <u>Integrals</u> and Their <u>Applications to Convolution Theory</u>



# Example

Processamento  
e Recuperação  
de Informação

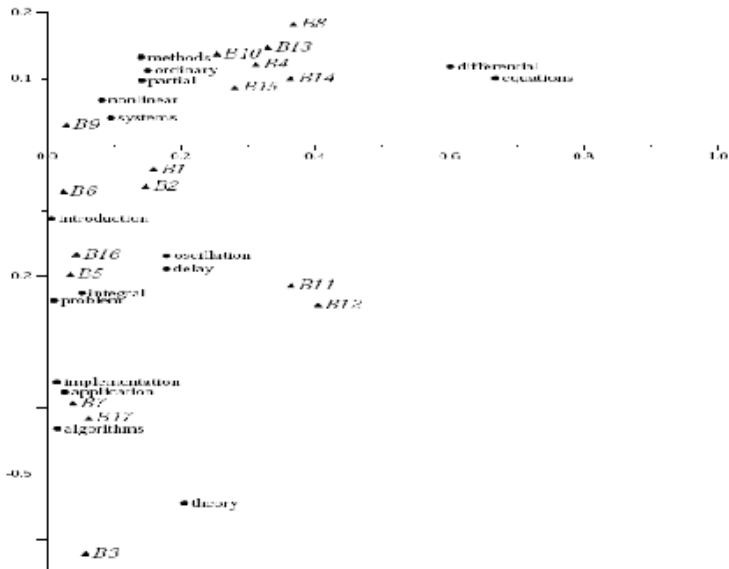
Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing





# Problems

Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

- Exact solution is computationally expensive
  - Can run in minutes to hours on a  $10^3$  to  $10^4$  collection
- Approximation algorithms frequently used in practice
  - Implementation in *scikit-learn* uses a fast randomized SVD solver [Halko , 2009]
- Most current implementations need to store the whole input matrix in memory
- Still not feasible for Web-sized collections



Processamento  
e Recuperação  
de Informação

Introduction

Partitioning  
Approaches

Dimensionality  
Reduction

Self-Organizing Maps

Latent Semantic  
Indexing

# Questions?