**ICMC USP**
SÃO CARLOS

**Instituto de Ciências Matemáticas e de Computação**

| Universidade de São Paulo |

# DATA VISUALIZATION

Multidimensional Projections and
Similarity Trees/
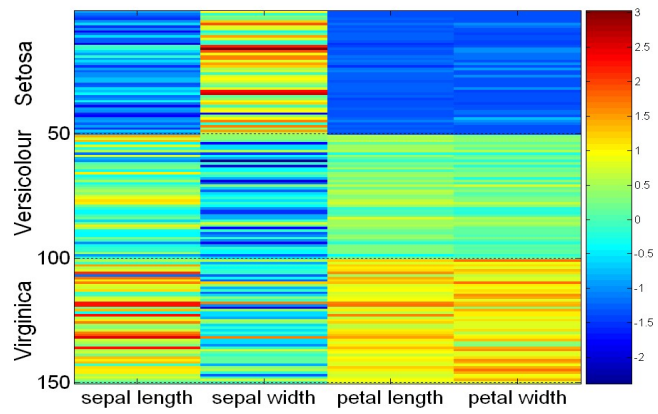Text / other applications

Rosane Minghim
2018-2

---

2

## Multidimensional Visualization
## Projections/Multidimensional Projections
## Document Collections
## Image Collections

- Visualization
- Visual Mining and Visual Analysis
- Projections
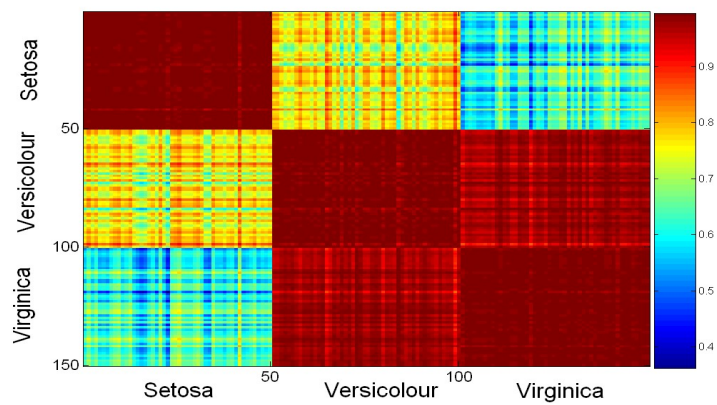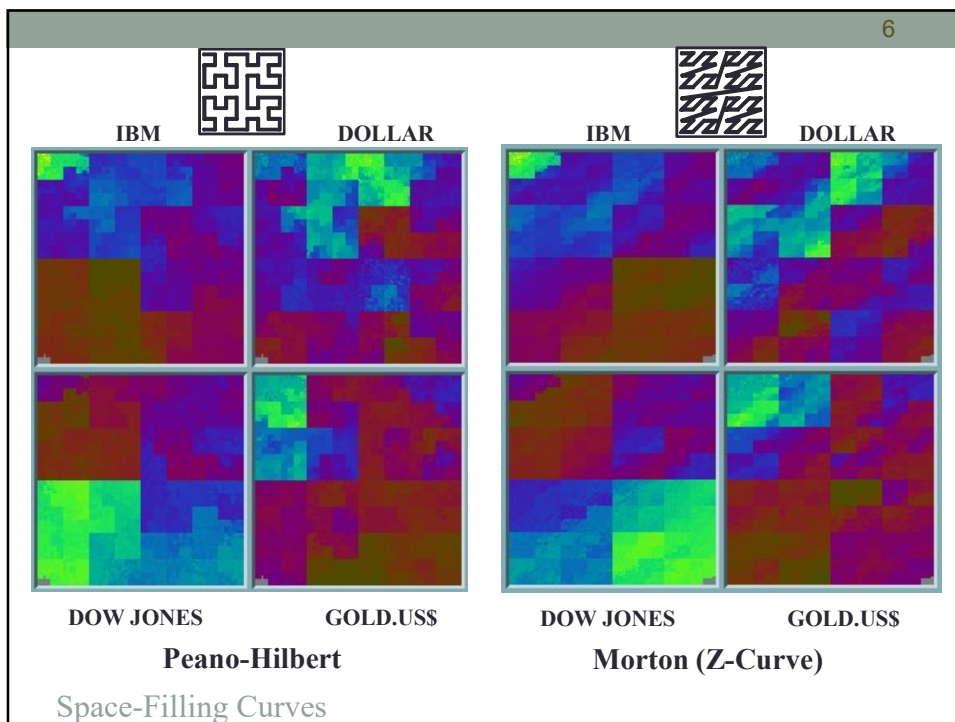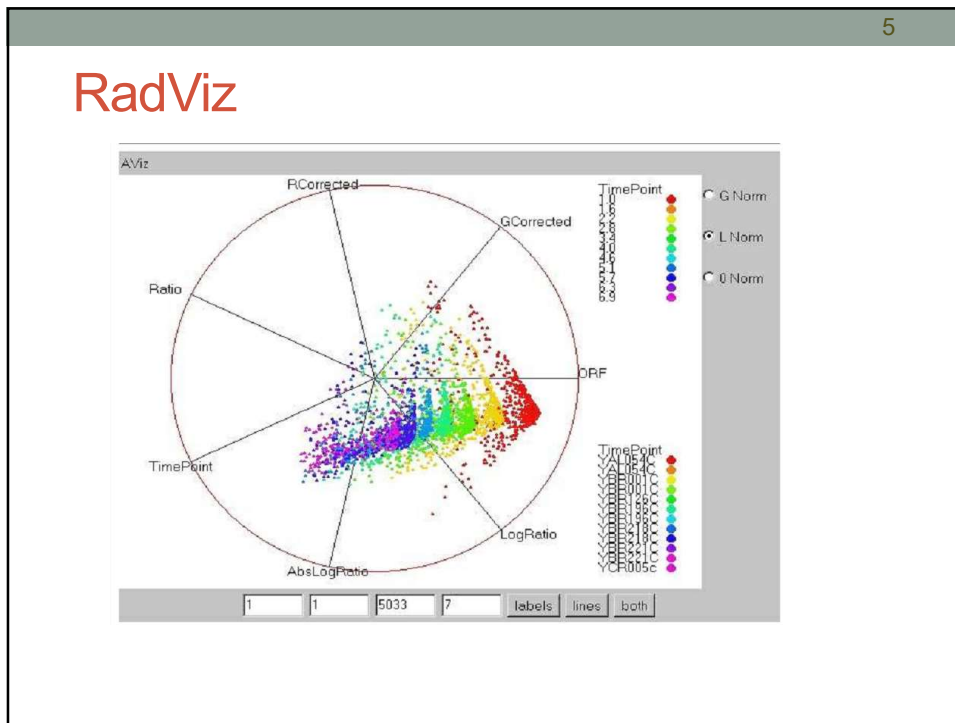- Examples:
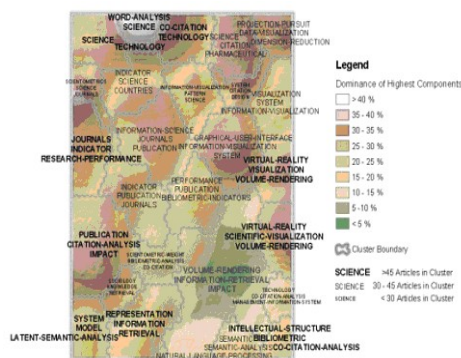  - Text and Images

# Data Matrix



# Correlation Matrix

RadViz



Peano-Hilbert      Morton (Z-Curve)

Space-Filling Curves

# SOM based

- **Self-Organization Maps (SOMs) cartográficos (ex. Skurpin 2002)**



---

# Clustered Force Layout

- http://bl.ocks.org/mbostock/1747543

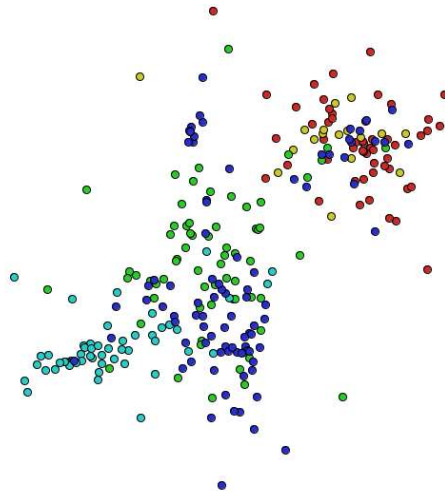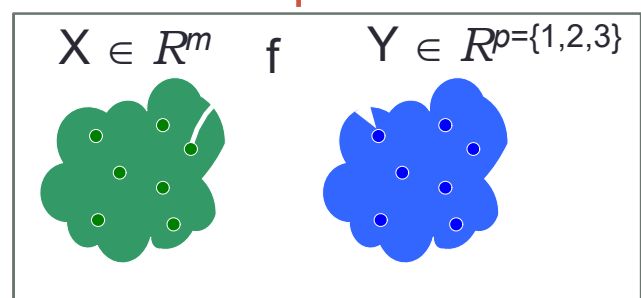Mapeamento para o plano permitindo a exploração.
Ex: Patents surgery, drugs, molecular bio

# Projection Techniques

$$X \in R^m \qquad f \qquad Y \in R^{p=\{1,2,3\}}$$

- $\delta: x_i, x_j \to R, x_i, x_j \in X$
- $d: y_i, y_j \to R, y_i, y_j \in Y$
- $f: X \to Y, |\delta(x_i, x_j) - d(f(x_i), f(x_j))| \approx 0, \forall x_i, x_j \in X$

# Problems PCA

390 dimensions

# Problems PCA

# Problems PCA

## Ex: Sammon Mapping

- Let **X** be the points in the original space $R^n$, we apply a distance measure $d_{ij}^*$ between Xi an Xj., and find **Y, the projected point**, ex. $R^2$ and $d_{ij}$ the Euclidean distance between them.

- Sammon's method applies an error function to measure the target.

## Force Based Point Placement

17

# Force Scheme [Tejada et al., 2003]



18

# Force Scheme [Tejada et al., 2003]

# Force Scheme [Tejada et al., 2003]

# Force Scheme [Tejada et al., 2003]

# Force Scheme [Tejada et al., 2003]

1. Map each point X to the plane (fastmap, nnp, etc.)
2. For each projected point x
   1. For each projected point q' ≠ x'
      1. Compute the vector **v** of <x' to q'>
      2. Move q' in direction of **v**, one fraction of Δ

$$\Delta = \frac{\delta(x,q) - \delta_{min}}{\delta_{max} - \delta_{min}} - d(x',q')$$

3. Normalize the coordinates between [0,1]

---

22

# LSP [Paulovich et al., 2006/2008]

- Least-Square Projection (LSP)
- Core idea: project a sub-set of points and interpolate the rest.
- Interpolation seeks to preserve the neighborhood between points.
- Each point is mapped within the convex hull of its neighbors.

# LSP [Paulovich et al., 2006/2008]

- Three main steps:
  1. Select a subset of points(control points) and Project these in $R^p$
  2. Determine the neighborhood of points
  3. Create a linear system whose answers are the Cartesian coordinates of points $p_i$ in $R^p$

# LSP: Laplacian Matrix

- Let $V_i = \{p_{i1},\ldots,p_{iki}\}$ be the neighborhood of a point $p_i$ and $c_i$ the coodinates of $p_i$ in $R^p$

$$c_i - \frac{1}{ki} \sum_{p_j \in V_i} c_j = 0$$



- Each $p_i$ will be the centroid of points in $V_i$

# LSP: Laplacian Matrix

$L\mathbf{x_1}=0, L\mathbf{x_2}=0, …, L\mathbf{x_p}=0$

where $x_1, x_2,…, x_p$ are vectors containing the Cartesian coordinates of the points

and L is the matrix defined by:

$$Lij = \begin{cases} 1 & i = j \\ -\dfrac{1}{ki} & p_j \in V_i \\ 0 & otherwise \end{cases}$$

$$L \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ ... \\ 0 \\ 0 \end{pmatrix}$$

# LSP: Matriz Laplaciana

$L\mathbf{x_1}=0, L\mathbf{x_2}=0, …, L\mathbf{x_p}=0$

onde $x_1, x_2,…, x_p$ são vetores contendo as coordenadas cartesianas dos pontos e $L$

• é a matriz dada por

$$Lij = \begin{cases} 1 & i = j \\ -\dfrac{1}{ki} & p_j \in V_i \\ 0 & caso\ contrário \end{cases}$$

$$L \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ ... \\ 0 \\ 0 \end{pmatrix}$$
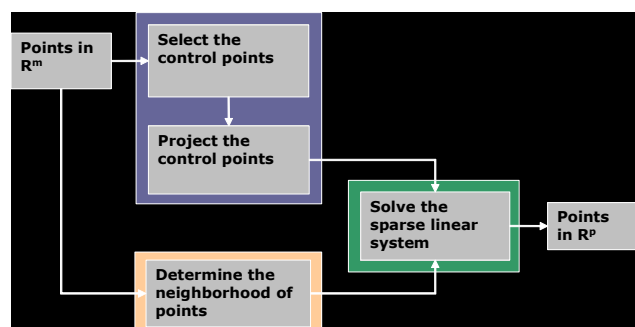
## LSP: Adding control points

$$A = \begin{pmatrix} L \\ C \end{pmatrix} \qquad Cij = \begin{cases} 1 & p_j \text{ is a control point} \\ 0 & \text{otherwise} \end{cases}$$

$$b_i = \begin{cases} 0 & i \leq n \\ x_{p_{c_i}} & n < i \leq n+nc \end{cases}$$

## LSP: Overview

## LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$

$v_2 = \{p_5 p_4 p_6\}$

$v_3 = \{p_1 p_5 p_6\}$

$v_4 = \{p_1 p_6\}$

$v_5 = \{p_3 p_2 p_6\}$

$v_6 = \{p_1 p_2 p_4 p_5\}$

$$L = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \end{bmatrix}$$

## LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$

$v_2 = \{p_5 p_4 p_6\}$

$v_3 = \{p_1 p_5 p_6\}$

$v_4 = \{p_1 p_6\}$

$v_5 = \{p_3 p_2 p_6\}$

$v_6 = \{p_1 p_2 p_4 p_5\}$

$$L = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \end{bmatrix}$$

31

# LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$
$v_2 = \{p_5 p_4 p_6\}$
$v_3 = \{p_1 p_5 p_6\}$
$v_4 = \{p_1 p_6\}$
$v_5 = \{p_3 p_2 p_6\}$
$v_6 = \{p_1 p_2 p_4 p_5\}$

$$L = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \end{bmatrix}$$

---

32

# LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$
$v_2 = \{p_5 p_4 p_6\}$
$v_3 = \{p_1 p_5 p_6\}$
$v_4 = \{p_1 p_6\}$
$v_5 = \{p_3 p_2 p_6\}$
$v_6 = \{p_1 p_2 p_4 p_5\}$

$$A = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \, L$$

$pc = \{p_3 p_6\}$

---

## LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$

$v_2 = \{p_5 p_4 p_6\}$

$v_3 = \{p_1 p_5 p_6\}$

$v_4 = \{p_1 p_6\}$

$v_5 = \{p_3 p_2 p_6\}$

$v_6 = \{p_1 p_2 p_4 p_5\}$

$$A = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \\ \hline 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

L

$pc = \{p_3 p_6\}$

C

---

## LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$

$v_2 = \{p_5 p_4 p_6\}$

$v_3 = \{p_1 p_5 p_6\}$

$v_4 = \{p_1 p_6\}$

$v_5 = \{p_3 p_2 p_6\}$

$v_6 = \{p_1 p_2 p_4 p_5\}$

$$A = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ c_{x_3} \\ c_{x_6} \end{bmatrix}$$

$pc = \{p_3 p_6\}$

---

## LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$

$v_2 = \{p_5 p_4 p_6\}$

$v_3 = \{p_1 p_5 p_6\}$

$v_4 = \{p_1 p_6\}$

$v_5 = \{p_3 p_2 p_6\}$

$v_6 = \{p_1 p_2 p_4 p_5\}$

$pc = \{p_3 p_6\}$

$$A = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ c_{x_3} \\ c_{x_6} \end{bmatrix}$$

## LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$

$v_2 = \{p_5 p_4 p_6\}$

$v_3 = \{p_1 p_5 p_6\}$

$v_4 = \{p_1 p_6\}$

$v_5 = \{p_3 p_2 p_6\}$

$v_6 = \{p_1 p_2 p_4 p_5\}$

$pc = \{p_3 p_6\}$

$$A = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ c_{y_3} \\ c_{y_6} \end{bmatrix}$$

# LSP: Exemplo de Sistema

$v_1 = \{p_3 p_4 p_6\}$
$v_2 = \{p_5 p_4 p_6\}$
$v_3 = \{p_1 p_5 p_6\}$
$v_4 = \{p_1 p_6\}$
$v_5 = \{p_3 p_2 p_6\}$
$v_6 = \{p_1 p_2 p_4 p_5\}$

$pc = \{p_3 p_6\}$

$$A = \begin{bmatrix} 1 & 0 & -1/3 & -1/3 & 0 & -1/3 \\ 0 & 1 & 0 & -1/3 & -1/3 & -1/3 \\ -1/3 & 0 & 1 & 0 & -1/3 & -1/3 \\ -1/2 & 0 & 0 & 1 & 0 & -1/2 \\ 0 & -1/3 & -1/3 & 0 & 1 & -1/3 \\ -1/4 & -1/4 & 0 & -1/4 & -1/4 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ c_{y_3} \\ c_{y_6} \end{bmatrix}$$

---

# LSP: Solving the system

- It is necessary to solve $A\mathbf{x} = \mathbf{b}$
- The system is solved by using least squares

$$\|Ax - b\|^2$$

- The analytical solution is

$$A^T A\mathbf{x} = A^T \mathbf{b} \quad \Rightarrow \quad \mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$$

- $A^T A$ is symmetric and sparse and can be solved using the factorization of Cholesky

# LSP: Resolvendo o Sistema

- É necessário resolver $A\mathbf{x} = \mathbf{b}$

- Este sistema é resolvido usando mínimos quadrados

$$\|Ax - b\|^2$$

- A única solução analítica será

$$A^T A\mathbf{x} = A^T \mathbf{b} \quad \Rightarrow \quad \mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$$

- $A^T A$ é simétrica e esparsa e pode ser resolvida usando a fatoração de *Cholesky*
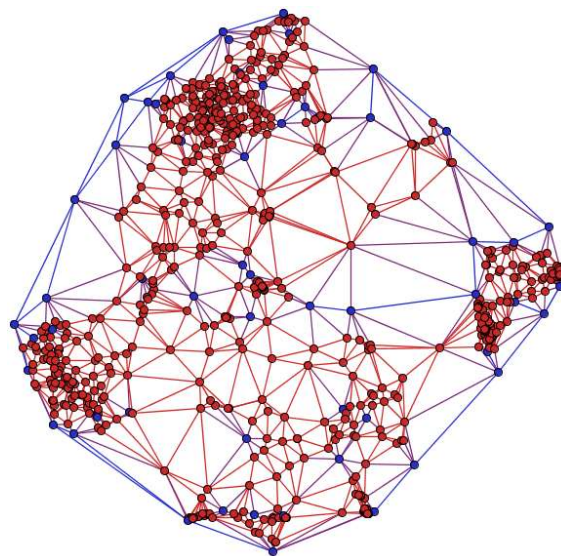
# Choosing the Control Points

- In order to select the control points
  - the space $R^m$ is split into *nc* clusters using k-medoids.
  - the control points are the medoids of each cluster

# Choosing the Control Points

- Once the control points are chosen, these points are projected onto $R^d$ through a fast dimensionality reduction method
  - Fast Projection (Fastmap or NNP)
  - Force Placement

Control points in blue

# Content – based by Projections
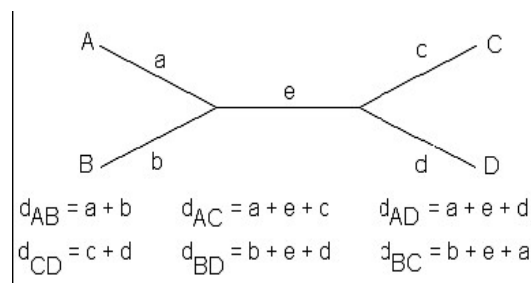
# Projection

# Projection: HDI

# Projection: Voting

## Point Placement by Phylogenetic Tree Construction Algorithms (N-J Trees)

## Point Placement by Phylogenetic Tree Construction Algorithms (N-J Trees)

$$d_{AB} + d_{CD} \leq \max\left(d_{AC} + d_{BD},\ d_{AD} + d_{BC}\right)$$



$d_{AB} = a + b \qquad d_{AC} = a + e + c \qquad d_{AD} = a + e + d$

$d_{CD} = c + d \qquad d_{BD} = b + e + d \qquad d_{BC} = b + e + a$

## Algorithm Neighbor-joining

Input: distance matrix
1. Criate a star tree for n objects.
2. Iteration
   1. Select a node pair (i,j) with smaller Sij (branch size)

$$S_{ij} = \frac{1}{2(n-2)} \sum_{k=3}^{N}(D_{ik} + D_{jk}) + \frac{1}{2}D_{ij} + \frac{1}{n-2} \sum_{3 \le m < n} D_{ij}$$

   2. Combine nodes i and j in a new node and calculate the branch size of the new node.

$$L_{ix} = \frac{D_{ij} + D_{iz} - D_{jz}}{2} \qquad\qquad L_{jx} = \frac{D_{ij} + D_{jz} - D_{iz}}{2}$$

## Algorithm Neighbor-joining

3. Calculate new distance matrix, computing the new distances from the new node to the remaining nodes.

$$D_{(i-j),k} = \frac{(D_{ik} + D_{jk})}{2} \qquad\qquad (3 \le k \le N)$$

4. Eliminate previous nodes i and j
5. If n>2 then iterate again.

51

- Alternate view (N-J Tree)

## Projections & Trees

55

# Application to Visual Data Classification

- Sample selection
- Classification
- Evolution of models
- Cooperation IC/UNICAMP
  - Helio Pedrini
  - William Schwartz (now UFMG)
- *Aplications: GPS data, Systems biology Data, data on quality of text*
  - *Cooperations with UNESP / Presidente Prudente, LNBio/CNPEM, NILC/ICMC*

56

# Visualization by Projections

57

# The case of document collections

- Applications
  - Teaching/Research
  - Search
  - Investigation

  - Patents
  - Medical reports
  - News

58

# O caso de coleções de documentos

- Aplicações
  - Ensino/Pesquisa
  - Busca
  - Investigação

  - Patentes
  - Laudos médicos
  - Notícias

- Maps of text Collections
  - Based on Relationships (Borner & Chen)
    - Co-authorship, co-citation
  - Based on Content
    - Similarity and Grouping
    - Common underlying subject
    - ➔Topics

# Relationships :
# Topic Busts and co-word



(Mane and Borner)
2004

# Relationships :
# Citation and Co-citation



(Borner)
(2003)

# Content-based Text Mapping

- Approach 1: Dimension reduction
  - ex. MSD, SVD, PCA
- Approach 2: Point Placement (PP)
- Approach 3: Clustering
- Approach 4: Projections
  - ex. FASPMAP, NNP, LSP

## Content - based

63



(Skupin)
(2002)
(abstracts)
SOM

## SOM based

- **Self-Organization Maps (SOMs) cartográficos (ex. Skurpin 2002)**

# Content - based



(Dimensional
Reduction)
News flash
IN-SPIRE
(PNL)

# VxInsight

- Sandia National Laboratories, mountain metaphor (Boyack et al., 2002).

# Content – based by Projections

# Exemplo de Projeção

# Exemplo de Projeção: IDH

# Exemplo de Projeção: Votação

# Basic Concepts

- Text Preprocessing

- Data and text mining

- Projection techniques

- Point Placement Strategies

# Text Preprocessing

1. Stopwords elimination
2. Extraction of words radicals (stemming)
3. Creation of n-grams
4. Frequency count and Luhn's lower cut (n-grams appearing less then x times are ignored)
5. Weighting process (*term-frequency inverse document-frequency - (tfidf)*)

# Result is a Vector Model

- Attributes: terms (n-grams)
- Value: term weight
- Table Data

# Vector Representation – term weighting

- tf – term frequency
- tfidf – tf x idf = tf x inverse document frequency

$$w_{ik} = tf_{ik} \times log\left(\frac{N}{n_k}\right)$$

# Vector Representation

|  | term$_1$ | term$_2$ | term$_3$ | term$_4$ | ... | term$_m$ |
|---|---|---|---|---|---|---|
| Doc$_1$ | 0.92 | 0.62 | 0.92 | 0.10 | ... | 0.67 |
| Doc$_2$ | 0.13 | 0.11 | 1.00 | 0.34 | ... | 0.33 |
| Doc$_3$ | 0.52 | 0.00 | 0.00 | 0.44 | ... | 0.77 |
| ... | ... | ... | ... | ... | ... | ... |
| Doc$_n$ | 0.02 | 0.12 | 0.22 | 0.92 | ... | 0.00 |

# Vector Representation – Similarity calculatin

| EUCLIDEAN |
|---|
| $sim_{i,j} = \sqrt{(w_{i,1} - w_{j,1})^2 + \ldots + (w_{i,k} - w_{j,k})^2}$ |
| MANHATAN |
| $sim_{i,j} = |w_{i,1} - w_{j,1}| + \ldots + |w_{i,1} - w_{j,1}|$ |
| COSINE |
| $sim_{i,j} = \dfrac{(w_{i,1} \times w_{j,1}) + \ldots + (w_{i,k} \times w_{j,k})}{(w_{i,1}^2 + \ldots + w_{i,k}^2) \times (w_{j,1}^2 + \ldots + w_{j,k}^2)}$ |
|  |

77

## Vector Representation – distance calculation

$$dis(doc_i, doc_j) = \sqrt{2*(1 - sim(doc_i, doc_j))}$$

$$sim(doc_i, doc_j) = \frac{doc_i \times doc_j}{\|doc_i\| * \|doc_j\|}$$

78

## Visualization

- Attribute Reduction
  - Co-clustering
  - PCA
  - SVD

followed by

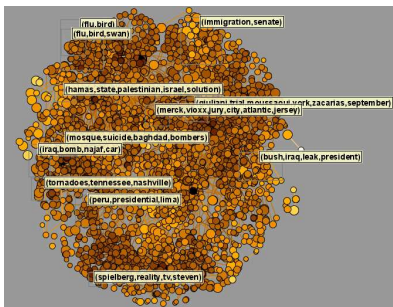- Projection by Dimension Reduction

# Alternative to Vector Representation

- Similarity Calculation text against text
  - Ex: NCD Normalized Compression Distance
    - Approximation of Kolmogorov Complexity
    - Ver: G. P. Telles, R. Minghim, and F. V. Paulovich. 2007. Visual Analytics: Normalized compression distance for visual analysis of document collections. *Comput. Graph.* 31, 3 (June 2007), 327-337. DOI=http://dx.doi.org/10.1016/j.cag.2007.01.024
  - Editing distance
    - Dice's coefficient
    - Matching's coefficient
    - Overlap's coefficient
    - Qgram Distance
    - Ver: Frizzi San Roman Salazar. Um estudo sobre o papel de medidas de similaridade na visualização de coleções de documentos. 2012. Dissertação (Mestrado em Ciências da Computação e Matemática Computacional) - Instituto de Ciências Matemáticas e de Computação, Fundação de Amparo à Pesquisa do Estado de São Paulo. Orientador: Maria Cristina Ferreira de Oliveira.

# Problems with Stress measurements

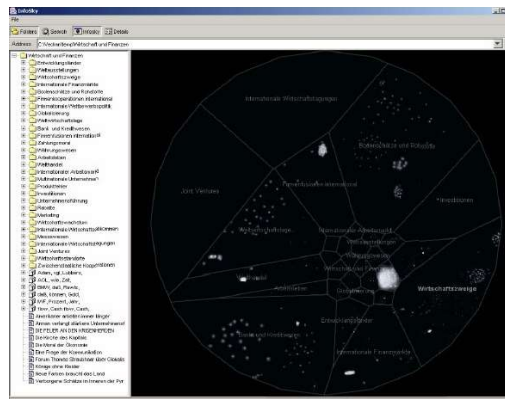Visual representations: graphs, surfaces, volumes, triangulations

## Exploration



## IN-SPIRE

- Spatial Paradigm for Information Retrieval - Pacific Nortwest National Laboratories

- Two Visualization Metaphors:
  - Galaxies – dimensional reduction
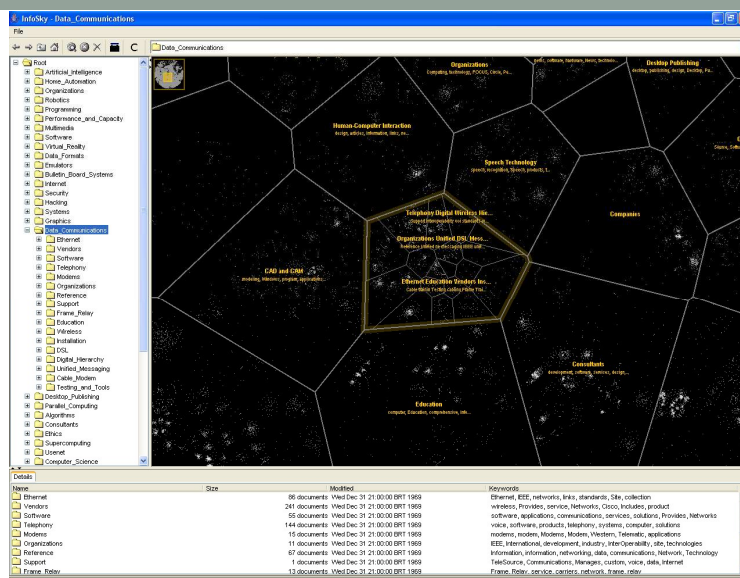  - Themescape

# InfoSky

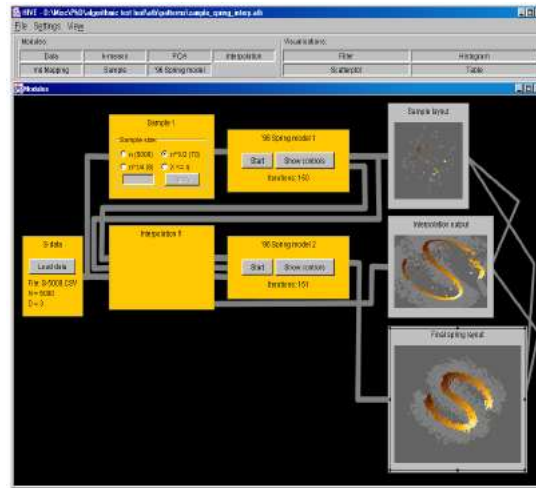Granitzer (Granitzer et al., 2004) also employs galaxy metaphor



86



http://en.know-center.at/forschung/wissenserschliessung/downloads_demos/infosky_demo
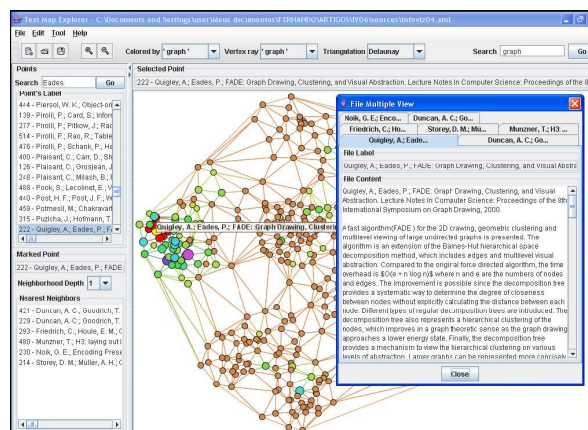
87

# HIVE (Ross and Chalmers 2003)

- Interconnected components:
  - Import
  - Transform
  - Render multi-dim data



88

# Projection Explorer (PEx)

- Projection and Point placement

- Precision

- Graphs and surfaces (Super Spider)

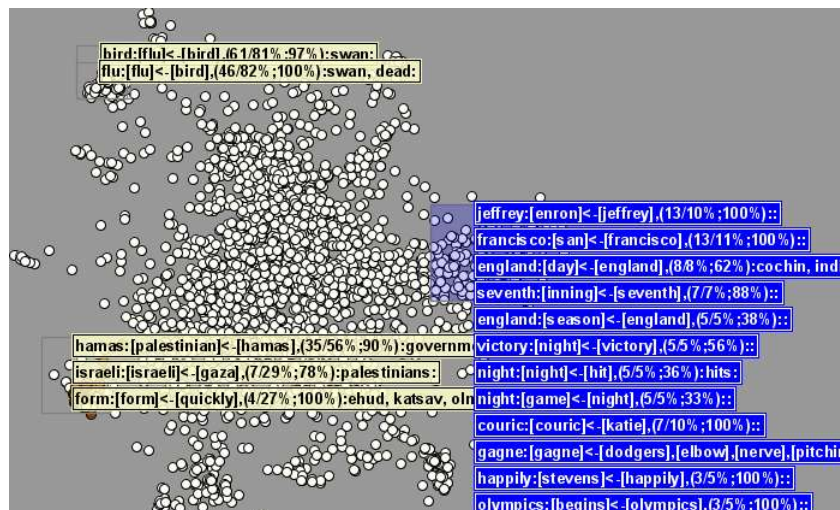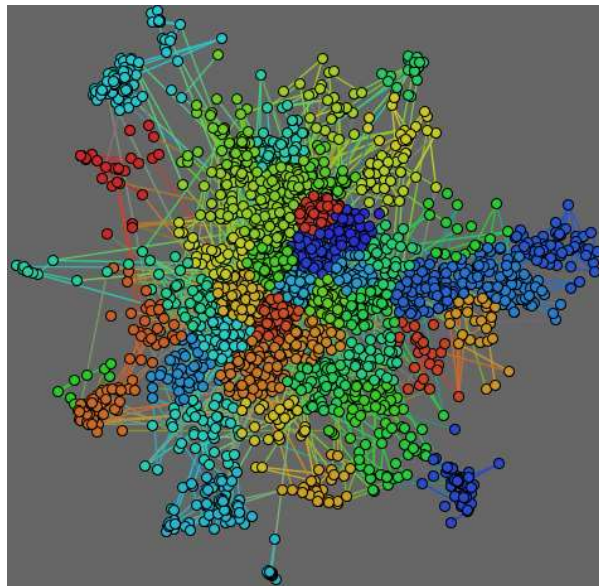## Mapping Text Collections via Projections and Point Placement

• Positioning and labeling

91

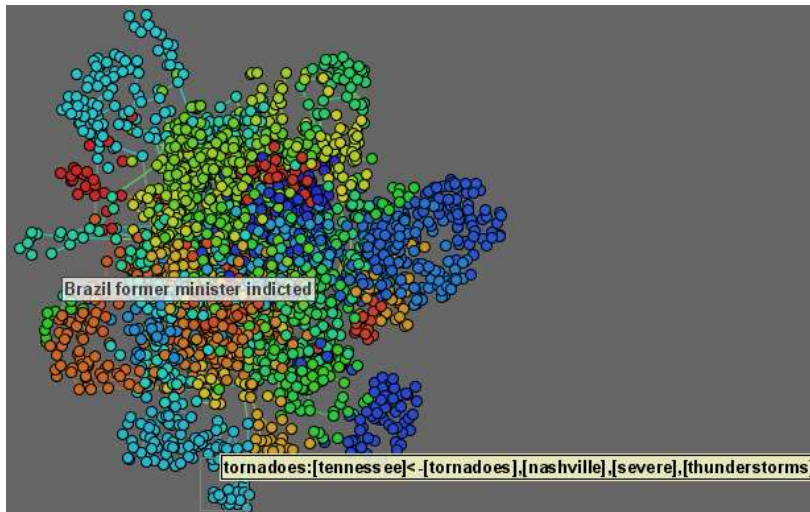- Detailing topics

92

• Finding Relationships

$\Longrightarrow$

95

- Untangling

$\Longrightarrow$

96



Brazil former minister indicted

tornadoes:[tennessee]< -[tornadoes],[nashville],[severe],[thunderstorms],

- Building a mesh

$\Longrightarrow$

• Coloring by degree of proximity

→

- Coordinating

• Building a Surface

# Explorando



- Case-Base Reasoning
- Information Retrieval
- Inductive Logic Programming

# Exemplos de Mapas
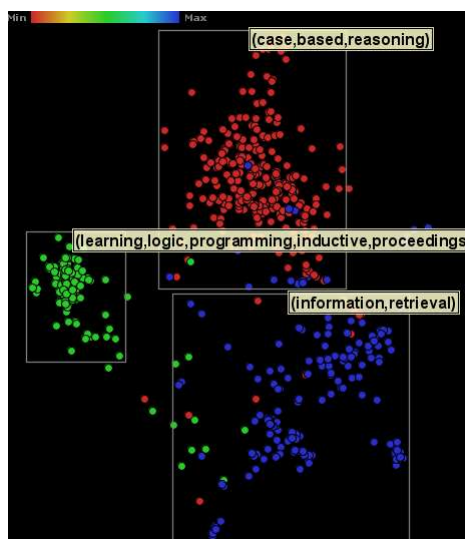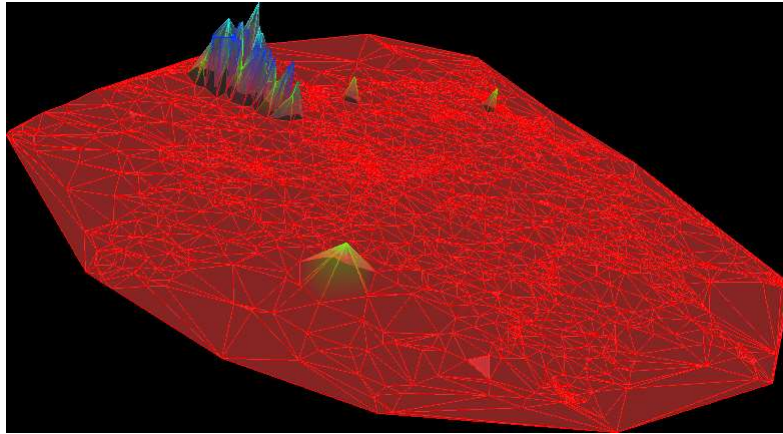
# Exemplos de Mapas

# RSS News Flash



Bird and Flu

Palestinian

Bush                              Iraq

Bush                              Iraq

## Curvas de Nível

# Time Series - Flow in Hydro. P. P.



Figure 2. Power plants of the basin Paraná

# Further Example - patents

# Further Example

- Cattle performance data
  - Translated to text from categorical information, e.g.,
    - Ranges of weight to words such as:

      {weight_below_fifty_percent;
      weight_between_fifty_seventy_five; etc..}
  - 9135 individuals

# Cattle performance data

## Cattle performance data

Colored
by word
'top'

## Cattle performance data

Colored
by
female

# Cattle performance data



Colored
by farm

# Images

# Images

# Cattle performance data



Colored
by word
'top'

## Scalability
## The Visual Super Tree

Cooperation:

Guilherme Pimentel Telles
IC/UNICAMP

## Application: Gene Expression and Systems biology data

- Cooperation ICMC/UNICAMP/LNBio (Campinas)/ Embrapa (Campinas)
- Complex Networks for Biological Data.
- Venn diagrams for Biological Data:
  - interactivenn.net, interactivenn.org

## Application: comparison of sets

- Cooperation ICMC/UNICAMP/LNBio (Campinas)/ Embrapa (Campinas)
- Fig.: Comparison of gene lists from different species
- www.interactiVenn.net



## Context: Visual Data Mining

- Definition [Ankerst 2000]
  - step in process of knowledge discovery / extraction (KDD)
  - utilizes visualization as communication channel between computer and user
  - to support identification of new and interpretable patterns

# Homework 2

- Explore the data sets left in tidia-ae using:
  - Vispipeline
  - Any other tool available to you

- For the news data set:
  - Mention 5 headlines of importance
  - Describe generally what happened regarding each one.

- Create or obtain a new text or image data set.
  - Format using .data or .dmat (and .zip, if text) for Vispipeline
  - Explore using both projections and trees.
  - Write and illustrate your findings in two pages.

# References

- Cuadros, A. M, Paulovich, F. V., Minghim, R., Telles, G. P - Point Placement by Phylogenetic Trees and its Application to Visual Analysis of Document Collections IEEE VAST 2007, Sacramento, CA, USA, IEEE CS Press, pp.99-106.
- Paulovih, F. V., Oliveira, M.C.F., Minghim, R. - The Projection Explorer: A Flexible Tool for Projection-based Multidimensional Visualization, IEEE Sibgrapi 2007, IEEE CS Press, Belo Horizonte, Brazil,pp. 27-34.
- Lopes, A. A., Minghim, R., Melo, V., Paulovich, F.V.; Mapping texts through dimensionality reduction and visualization techniques for interactive exploration of document collections, **SPIE Conference on Visualization and Data Analysis**, San Jose, CA, USA Jan. 2006, 6060T-11.
- Minghim, R., Paulovich, F.V., Lopes, A. A.; Content-based text mapping using multidimensional projections for exploration of document collections, **SPIE Conference on Visualization and Data Analysis,** San Jose, CA, USA Jan. 2006, 6060T-11.

# References

- Pinho, R. D. ; Oliveira, M. C. F. ; Minghim, R. ; Andrade, M. G. . Voromap: A Voronoi-based Tool for Visual Exploration of Multidimensional Data. In: **10th International Conference on Information Visualization**, 2006, Londres. Proceedings of Information Visualisation 2006, 2006. v. 1. p. 39-44
- Paulovich, F. V. ; Minghim, R. . Text Map Explorer: a Tool to Create and Explore Document Maps. In: Information Visualisation 2006 (IV06) **10th International Conference on Information Visualisation**, 2006, Londres. Proceedings of Information Visualisation 2006, 2006. v. 1. p. 245-251.
- Paulovich, F. V. ; Nonato, L. G. ; MINGHIM, R. ; Levkowitz, H. . Least Square Projection: a fast high precision multidimensional projection technique and its application to document mapping. IEEE Transactions on Visualization and Computer Graphics, 2008.
- Minghim, R. ; Levkowitz, H. ; Nonato, L. G. ; Watanabe, L. S. ; Salvador, V. C. L. ; Lopes, H. ; Pesco, S. ; Tavares, G. . Spider Cursor: A simple versatile interaction tool for data visualization and exploration. In: **ACM GRAPHITE'05** - 3rd International Conference on Computer Graphics and Interactive Techniques in Australasia and Southeast Asia, 2005, Dunedin. Proceedings of Graphite 2005, 2005. p. 307-314.
- Heberle, H.; Meirelles, G. V.; da Silva, F. R.; Telles, G. P.; Minghim, R. *InteractiVenn: a web-based tool for the analysis of sets through Venn diagrams*. BMC Bioinformatics 16:169 (2015).

# 5.3 Topic Extraction and Visualization

- Topic Definition by Covariance

- Topic Extraction by Seeded Generation of Association Rules  (pruning by relevant terms)

- Labeling and Viewing

# Topic Extraction and Visualization

Topic Definition by Covariance

- Pair of words with highest covariance
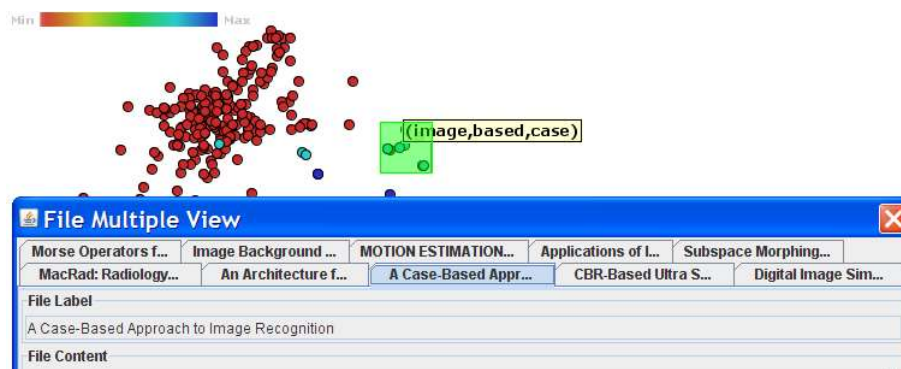
$$cov(t_i, t_j) = \frac{1}{n-1} \sum_{k=1}^{n} (t_{ki} - \bar{t}_i)(t_{kj} - \bar{t}_j)$$

- For all the other words, highest mean covariance compared to first two.
- Add to label if above threshold.

# Topic Definition by Covariance

# Topic Extraction and Visualization

Topic Extraction using Association Rules

- Use relevant words as seeds

- Prune the case by rule weighting

## Topic Extraction using Association Rules

| Transactions | Items |
|--------------|-------|
| 1 | Trousers, t-shirt, snickers |
| 2 | T-shirt, snickers |
| 3 | shorts, snickers |
| 4 | Trousers, sandals |

| Frequent Itemsets | Support |
|-------------------|---------|
| {snickers} | 75% |
| {Trousers} | 50% |
| {T-shirt} | 50% |
| {T-shirt, snickers} | 50% |

Min. support = 50% (2 transactions).

Min. confidence = 50%.

## Topic Extraction using Association Rules

$$\text{tenis} \longrightarrow \text{t-shirt}$$

$$support = support(\{snickers, t\text{-}shirt\}) = 50\%$$

$$confidence = \frac{support(\{snickers, t - shirt\})}{support(\{t - shirt\})} = \frac{50}{50} = 100\%$$

$$\text{T-shirt} \longrightarrow \text{snickers}$$

$$support = support(\{snickers, \ t - shirt\}) = 50\%$$

$$confidence = \frac{support(\{snickers, t - shirt\})}{support(\{snic\,ker\,s\})} = \frac{50}{75} = 66,6\%$$

## Topic Extraction using Association Rules (example)

## Topic Extraction using Association Rules (example)

image:[recognition]<-[image],[medical],[segmentation],[images],(5/71%;100%):vision, digital, symposium:

visualization:[sound]<-[visualization],[sonification],[visual],[scientific],(5/50%;100%):displays, martin, tool, minghim,

reality:[environments]<-[reality],(4/40%;100%)::

visually:[blind]<-[ch],[roth],[visually],[hci],(4/67%;100%):speech, visual, audio, nov, device:

## Topic Extraction using Association Rules

- Topics using AR
  - Term co-occurance in documents <=> subject
  - Transaction => Document
  - Item => term

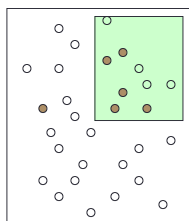## Topic Extraction using Association Rules

- Issues
  - Discovered rules amount
  - Term relevance (items)
  - Rule relevance measure (filtering)

  - High Sup. & conf. => few interesting rules
  - Low Sup. & conf. => huge amount of rules

## Locally weighted and seeded AR

- Weighting Terms and Rules

$$w_{i,j} = \frac{\sum_{j=1}^{k} Tf_{i,j}}{\sum_{j=1}^{k} Tf_i}$$

$wi,s = 5/6 = \textbf{0.83333}$

## Steps

1. S: set of user selected documents
2. Picked 10 most relevant terms

$$W_{t_j S_k} = \frac{\sum T f_{t_j S_k}}{\sum T f_{t_j C}}$$

## Steps

1. Initial item sets: Tr x T
   - Relevant Terms x All Terms
2. Items Sets discovered by Apriori altorighm
3. Sorted by weight:

$$\sum W_{t_j S_k}$$

# Steps

6. Highest weight item set selected
7. Covered documents removed from S
8. Further item sets are selected if there is support over residual S ( repeats 6 e 7 )

9. If all items sets are considered and |S residual| >0, repeats whole process with residual S.

# Sequential covering with Multiple restart

- Variance and Coverage
- Partitioning Strategies
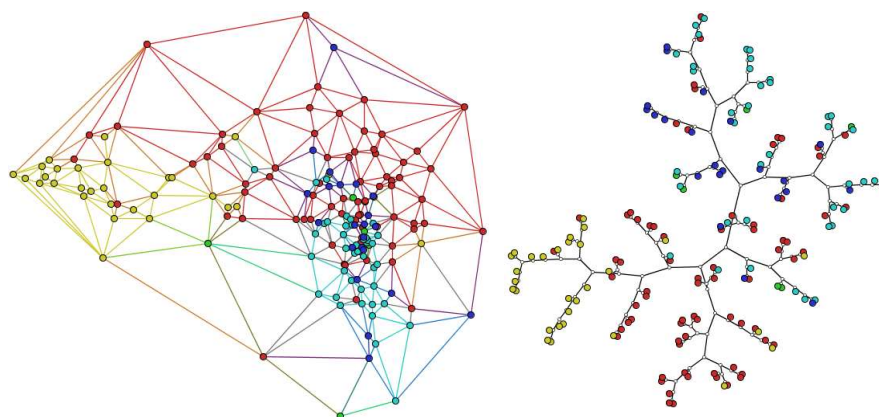- Grid
  - Resize
  - Slide
- Cluster
  - Cluster number

# 5.4 Further Examples

- RSS Patent Data, recovered from the Web
  http://www.freepatentsonline.com/
- Case 1:
  - 170 files
  - Graphics processing, printer, database, document, ai
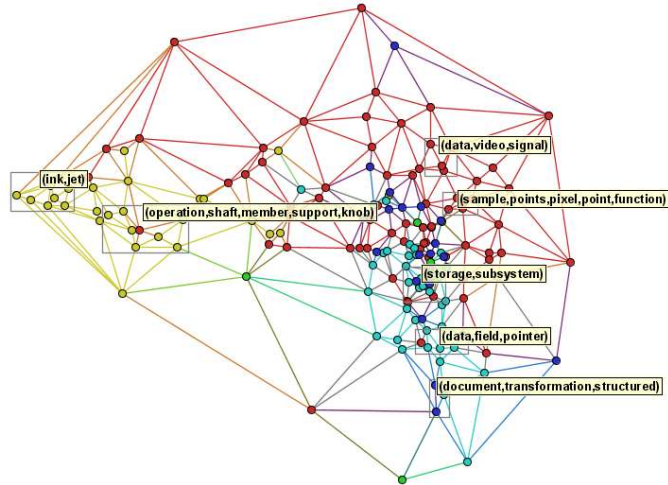
# Further Examples

# Further Examples

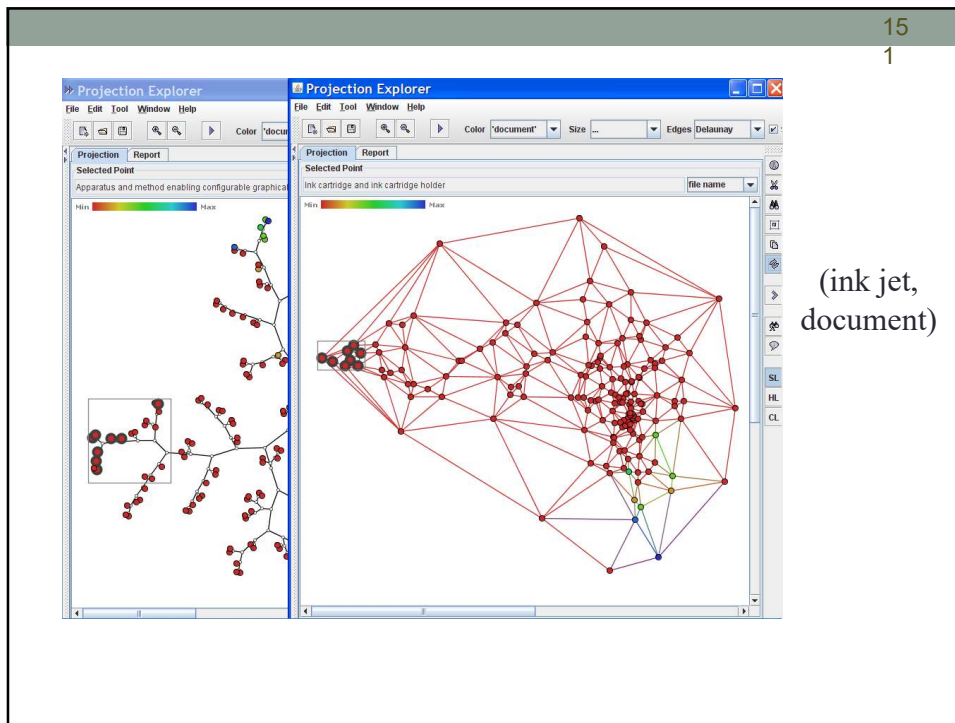# Further Examples

(ink jet,
document)

153

---

# Patents – case 2

- http://www.freepatentsonline.com/
- 172 files
- surgery (2), drugs(2), molecular biology

Patents surgery, drugs, molecular bio



Patents surgery, drugs, molecular bio
stopwords selection

Patents surgery, drugs, molecular bio topics



Patents surgery, drugs, molecular bio

Patents surgery, drugs, molecular bio



Projection Explorer (PEx)

http://infoserver.lcad.icmc.usp.br/