

Interfaces de Visualización de *Clustering*



UCR – ECCI

CI-2414 Recuperación de Información

Prof. Kryscia Daviana Ramírez Benavides



Introducción

- Los mecanismos convencionales de una búsqueda tienen baja precisión.
- Un problema común con esto es que los usuarios deben navegar a través de muchos documentos irrelevantes antes de encontrar el tipo de documento que le interesa.
- Aún con algoritmos avanzados de *ranking*, no se puede saber de antemano que tipo de documentos prefiere el usuario.
- *Interfases Gráficas basadas en Clustering* pueden ayudar al usuario a encontrar más fácilmente lo que busca.



Nuevos Paradigmas en Visualización de Información

- *Sammon Map.*
- *Tree-Map Visualization.*
- *Radial Interactive Visualization.*



Sammon Map

- *Sammon map* genera una localización en dos dimensiones de los clusters.
- Este mapa se calcula usando una búsqueda de gradiente iterativo.
- El objetivo del algoritmo es representar puntos en un espacio n dimensiones, normalmente 2 dimensiones, mientras intenta conservar las distancia pares entre los objetos.
- Utiliza el algoritmo de *clustering* Buckshot.
- Despliega documentos en *clustering* por *keywords*.
- Al dar click a cualquier *cluster* muestra el número de documentos que contiene, y el *ranking* de las *keywords* usadas para formar el *cluster*.
- Al dar click en el número de documentos muestra las *keywords* y sus porcentajes, y un resumen de los documentos.

Information Navigator - mafia

search: mafia Find

Plain Radial Tree Map Sammon Map

Level 0 Back

Related Word	Hit Doc ...
<input type="checkbox"/> cosa	54%
<input type="checkbox"/> nostra	54%
<input type="checkbox"/> calabria	55%
<input type="checkbox"/> sicily	36%
<input type="checkbox"/> ndrangheta	33%
<input type="checkbox"/> gambino	33%
<input type="checkbox"/> bosses	30%
<input type="checkbox"/> sicilian	26%
<input type="checkbox"/> reggio	26%
<input type="checkbox"/> camorra	23%
<input type="checkbox"/> gotti	23%
<input type="checkbox"/> clans	21%
<input type="checkbox"/> campania	16%
<input type="checkbox"/> puglia	16%
<input type="checkbox"/> calabrian	16%
<input type="checkbox"/> sharking	16%
<input type="checkbox"/> riina	16%
<input type="checkbox"/> mafiosi	14%
<input type="checkbox"/> mancino	14%
<input type="checkbox"/> carabinieri	11%
<input type="checkbox"/> sacra	11%
<input type="checkbox"/> salvatore	11%
<input type="checkbox"/> paolo	8%
<input type="checkbox"/> cosenza	8%
<input type="checkbox"/> toto	8%
<input type="checkbox"/> palmi	8%
<input type="checkbox"/> catanzaro	8%
<input type="checkbox"/> catania	8%

Filter Documents

42 documents found

[Document: 345605](#)

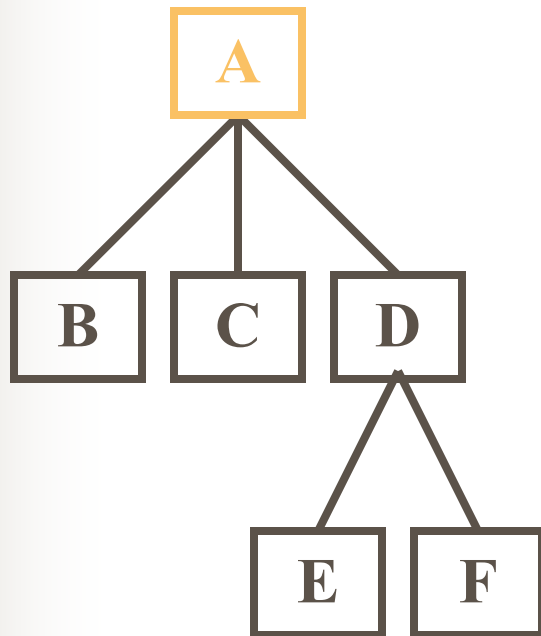
Document Type: Daily Report 26 Jan
 Document Type: Daily Report 26 Jan 1994 ITALY & VATICAN CITY Mafia's Extensio



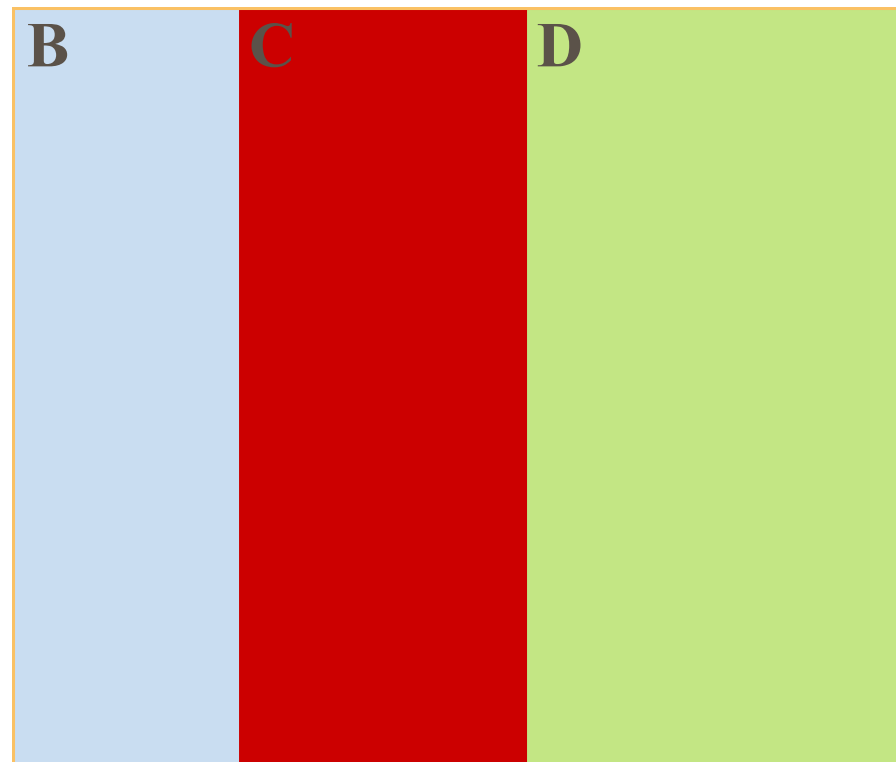
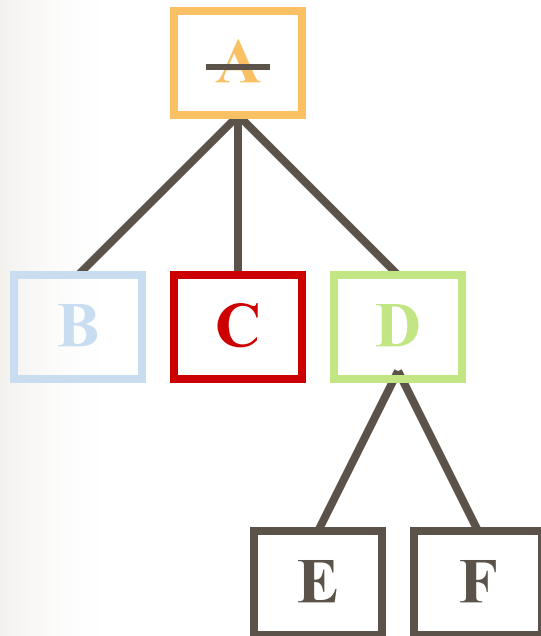
Tree-Map Visualization

- Representación jerárquica de los *clusters*.
- Los *clusters* se visualizan mediante rectángulos.
- *Clusters* similares se agrupan en Súper *Clusters*.
- Utiliza algoritmos de *clustering* jerárquicos. Aunque puede usar no jerárquicos.

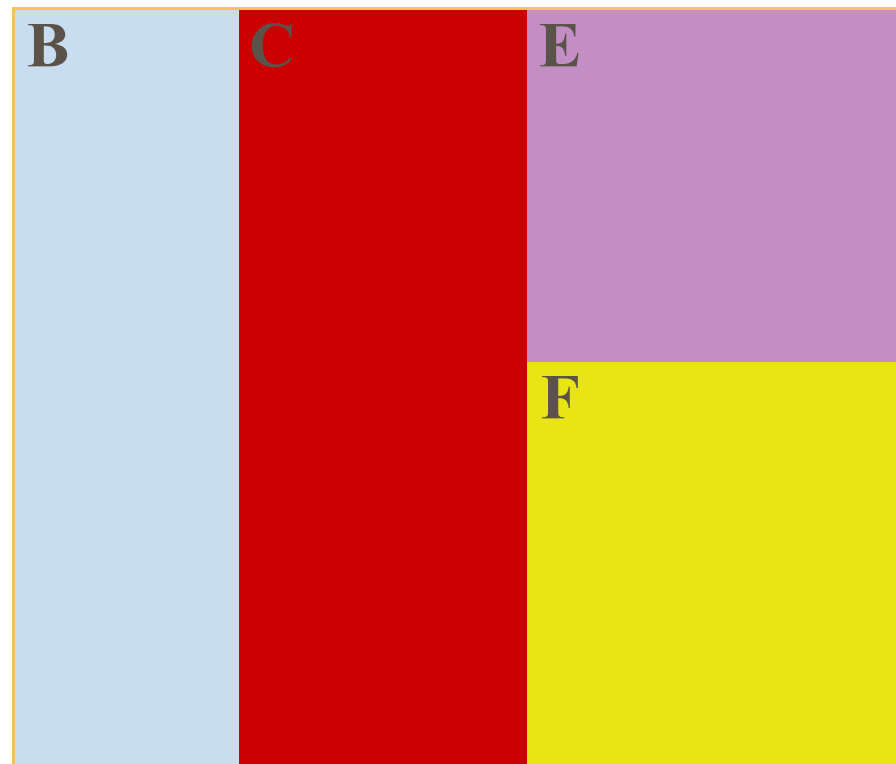
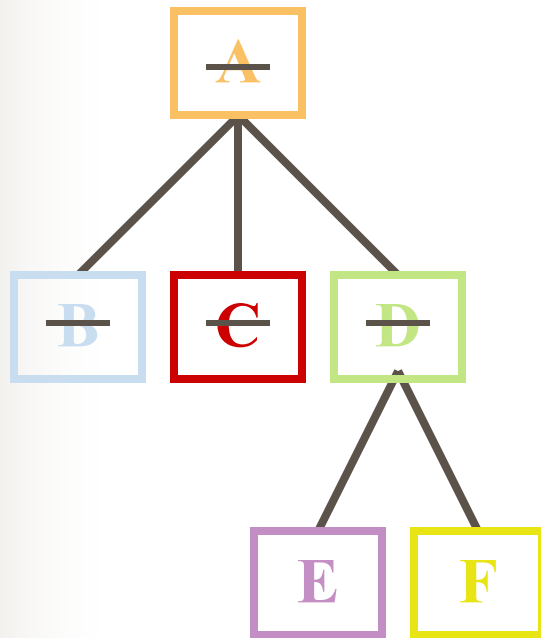
Tree-Map Visualization (cont.)



Tree-Map Visualization (cont.)



Tree-Map Visualization (cont.)



Information Navigator - mafia

search: Find

Plain Radial Tree Map Sammon Map

Related Word	Hit Doc Freq
<input type="checkbox"/> cosa	54%
<input type="checkbox"/> nostra	50%
<input type="checkbox"/> calabria	45%
<input type="checkbox"/> sicily	35%
<input type="checkbox"/> ndrangheta	33%
<input type="checkbox"/> gambino	33%
<input type="checkbox"/> bosses	30%
<input type="checkbox"/> sicilian	26%
<input type="checkbox"/> reggio	26%
<input type="checkbox"/> camorra	23%
<input type="checkbox"/> gotti	23%
<input type="checkbox"/> clans	21%
<input type="checkbox"/> campania	16%
<input type="checkbox"/> puglia	16%
<input type="checkbox"/> calabrian	16%
<input type="checkbox"/> sharking	16%
<input type="checkbox"/> riina	16%
<input type="checkbox"/> maffiosi	14%
<input type="checkbox"/> mancino	14%
<input type="checkbox"/> carabinieri	11%
<input type="checkbox"/> coors	11%

42 documents found
cosa nostra calabria

magistrates palermo sicily

<< Back Show All Documents Filter

42 documents found

[Document: 345605](#)
Document Type: Daily Report 26 Jan
 Document Type: Daily Report 26 Jan 1994 ITALY & VATICAN CITY Mafia's Extensio

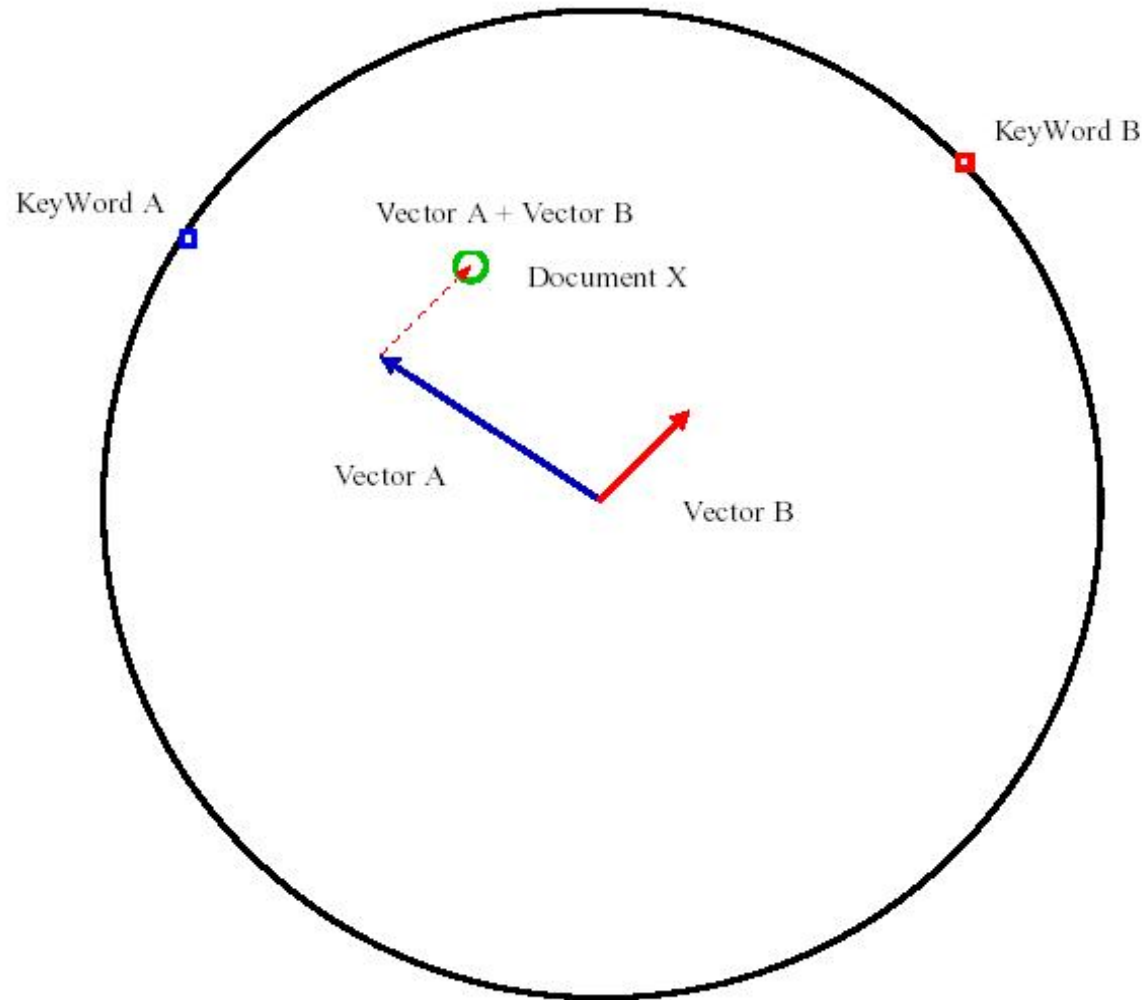
[Document: 345738](#)
Document Type: Daily Report 28 January
 Document Type: Daily Report 28 January 1994 ANNEX Italy & Vatican City Govern

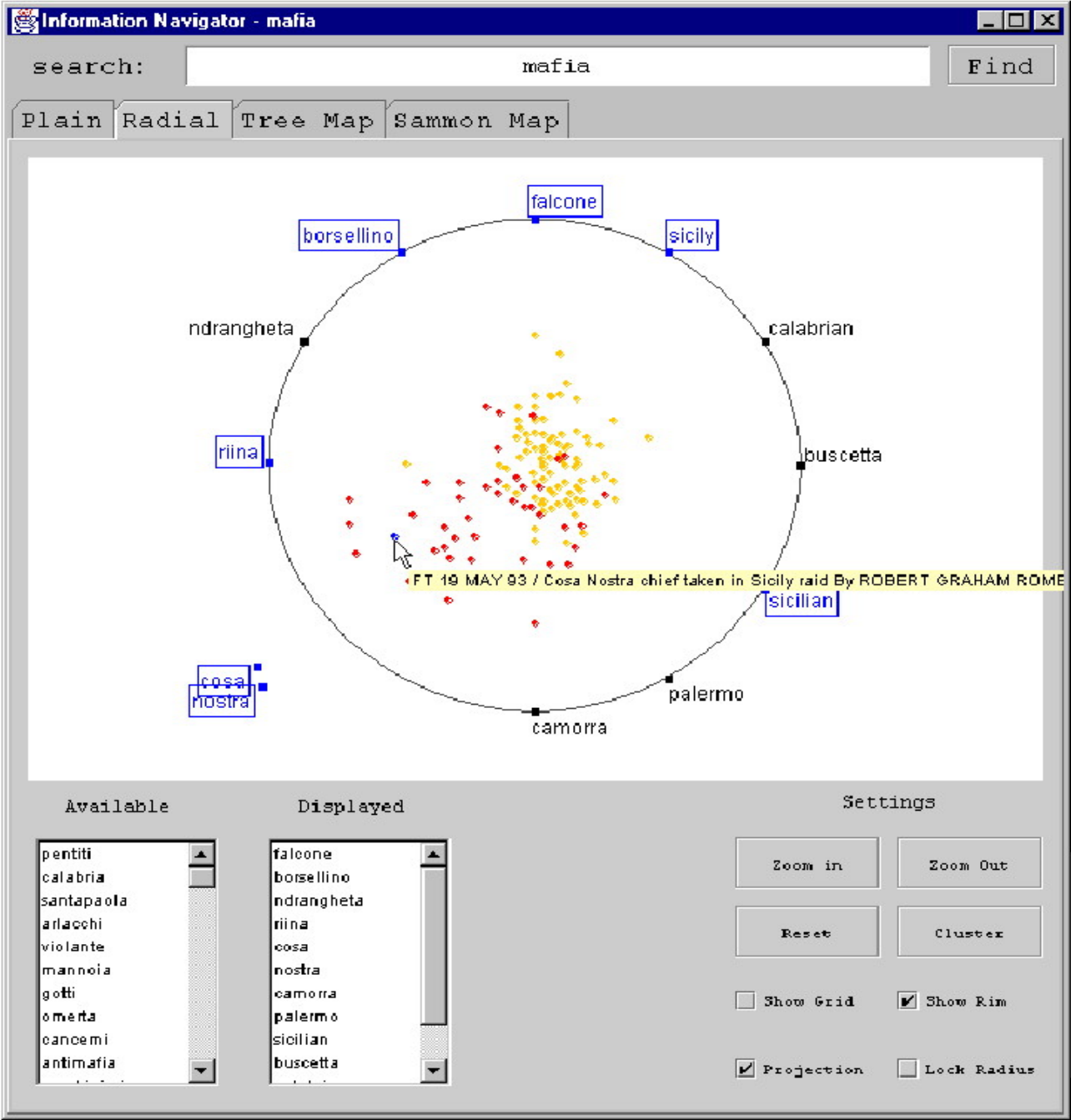
[Document: 344590](#)
Document Type: Daily Report 2 December



Radial Interactive Visualization

- Nodos con *keywords* son colocados alrededor de un círculo.
- Documentos están representados por puntos en el interior del círculo.
- Entre más relacionado un documento con un *keyword* más cercano estará de este.
- Utiliza el algoritmo de *clustering* Buckshot.
- Al dar click en una palabra resaltarán los documentos que contienen esa palabra. Al dar click en un punto mostrarán las palabras que contiene.
- Algoritmo:
 - Sea P1 la posición del *keyword* A en el círculo.
 - Sea P2 la posición del *keyword* B en el círculo.
 - Se utiliza la matriz de pesos de *keywords* para cada documento:
(P1 * wij, P2 * wij).







Características Comparativas

- Visualizaciones basadas en *clusters* dan un panorama más amplio del conjunto de resultados.
- *Sammon Map* se enfoca más en la relación que existe en un *cluster* y otro.
- *Tree Map* es más explícito en cuanto al tamaño de los *cluster* y su estructura jerárquica.
- *Radial Visualization* permite enfocarse en formar subconjuntos de *keywords*.



Facilidades de cada Interfaz

- *Sammon Map* guía en el análisis, permite reagrupar subconjuntos, y gradualmente acercar al tipo de documentos de interés.
- *Tree Map* permite enfocarse en keywords de ciertos documentos de interés, para formular una búsqueda más productiva.
- *Radial Visualization* es apropiado si el usuario está familiarizado con los *keywords* del área de su interés.



Conclusiones

- Contribuye a la visualización y navegación de un conjunto de documentos retornados mediante:
 - Identificación de *keywords* relevantes en un conjunto de documentos.
 - Desechar rápidamente *clusters* irrelevantes.
 - Operaciones de “*Drill Down*” en *clusters* relevantes.
 - Construcción de grupos personalizados.



Referencias Bibliográficas

- La información fue tomada de:
 - Libro de texto del curso.
 - Presentaciones del curso RI del estudiante Randall Mora Jiménez. Universidad de Costa Rica, 2003.
 - A Visualization Interface for Document Searching and Browsing. Matthew Carey Frank Kriwaczek Stefan M Ruger Multimedia Knowledge <http://km.doc.ic.ac.uk>.