# Improving web search with FCA



Radim BELOHLAVEK Jan OUTRATA



State University of New York

Dept. Systems Science and Industrial Engineering Watson School of Engineering and Applied Science Binghamton University – SUNY, NY, USA

Dept. Computer Science Faculty of Science Palacky University, Olomouc, Czech Republic

# Information Retrieval $\times$ Formal Concept Analysis

web search = mining web retrieval results, part of web mining

**Information Retrieval** (IR) = retrieval of required information from textual unstructured or semistructured data (example: search by keywords, retrieval of documents), iterative and interactive process (mining):

- submitting query,
- looking at the data returned,
- submitting a refined query until appropriate data are found.

Formal Concept Analysis (FCA) = method of analysis of tabular data, extracting a hierarchically ordered collection of clusters:

- (input) tabular data = objects described by attributes,
- (output) clusters = objects having common attributes (and vice versa),
- used for data mining, knowledge discovery, preprocessing data, clustering and classification (conceptual clustering) etc.

イロト 不得下 イヨト イヨト 二日

# FCA in Information Retrieval

rationale behind using FCA in IR and document mining:

- current search engines (e.g. Google, Yahoo, etc.) provide a ranked list of retrieved documents, i.e. a "simplistic" linear view on retrieved information, without the possibility to inspect related documents at the same time,
- FCA enables structured (or categorized) view of retrieved information with contextual information,
- user is supplied with a (part of a) conceptual hierarchy of retrieved documents and he or she can browse the hierarchy to find required information more quickly,
- new type of information can be mined: most common/uncommon subjects, which subjects imply or are implied by other subjects, novel subject associations etc. → Conceptual Knowledge Processing

(日) (同) (日) (日) (日)

# Formal Concept Analysis (FCA)

FCA = method of analysis of tabular data (Wille, TU Darmstadt, 1982)

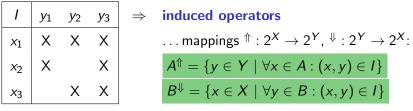
- alternatively called: concept data analysis, concept lattices, ...
- used for data mining and knowledge discovery
- input:

| 1                     | <i>y</i> 1 | <i>y</i> 2 | <i>y</i> 3 | $X = \{x_1, x_2, \dots\}$ | set of <b>objects</b>        |
|-----------------------|------------|------------|------------|---------------------------|------------------------------|
| <i>x</i> <sub>1</sub> | Х          | Х          | Х          | $Y = \{y_1, y_2, \dots\}$ | set of attributes            |
| x <sub>2</sub>        | X          |            | Х          | $I \subseteq X \times Y$  | relation to have             |
| <i>x</i> <sub>3</sub> |            | Х          | Х          | $\langle x,y angle \in I$ | object $x$ has attribute $y$ |

- output
  - concept lattice (hierarchically ordered set of clusters formal concepts)
  - attribute implications (particular attribute dependencies)

・ 同 ト ・ ヨ ト ・ ヨ ト ・ ヨ

# **FCA** basics



 $A \subseteq X \mapsto A^{\uparrow} \dots \text{attributes common to all objects from } A \\ \{x_1, x_2\}^{\uparrow} = \{y_1, y_3\} \\ B \subseteq Y \mapsto B^{\Downarrow} \dots \text{objects sharing all attributes from } B \\ \{y_1, y_2\}^{\Downarrow} = \{x_1\} \\ \text{(Birkhoff 1940s, Ore, Barbut & Monjardet, Wille 1982)}$ 

Definition (formal concept = fixed point of  $^{\uparrow}, ^{\downarrow}$ )

**Formal concept** in data is a pair  $\langle A, B \rangle$  s.t.

$$A^{\uparrow} = B$$
 and  $B^{\downarrow} = A$ .

formal concepts pprox all potentially interesting clusters in data ,

R. Belohlavek, J. Outrata (SSIE BU, CS UP)

# **FCA** basics

Definition (concept lattice = formal concepts + concept hierarchy)

**Concept lattice (Galois lattice)** of  $\langle X, Y, I \rangle$  is the set

$$\mathcal{B}(X,Y,I) = \{(A,B) \mid A^{\Uparrow} = B, B^{\Downarrow} = A\}$$

of all formal concepts PLUS concept hierarchy  $\leq$  defined by

$$(A_1, B_1) \leq (A_2, B_2)$$
 iff  $A_1 \subseteq A_2$  (iff  $B_2 \subseteq B_1$ ).

FCA ... inspired by **Port-Royal** (traditional) approach to concepts:

- **concept** (according to Port-Royal) := **extent** A + **intent** B
  - **extent** = objects covered by concept
  - $\bullet \ intent = attributes covered by concept$
- **example: DOG** (data = animals × animals' attributes)
  - extent = collection of all dogs (beagle, collie, poodle, ...)
  - intent = all dogs' attributes (barks, has four limbs, has tail,  $\dots$ )
- conceptual hierarchy  $\leq$  ... subconcept/superconcept relation
  - concept1=(extent1,intent1)  $\leq$  concept2=(extent2,intent2)
    - $\iff$  extent1  $\subseteq$  extent2 ( $\Leftrightarrow$  intent1  $\supseteq$  intent2)
  - $\bullet$  example: <code>BEAGLE  $\leq$  DOG  $\leq$  MAMMAL  $\leq$  ANIMAL</code>

# Formal concepts = maximal rectangles in data

Theorem (formal concepts = maximal rectangles)

 $\langle A,B\rangle$  is a formal concept IFF  $\langle A,B\rangle$  is a maximal rectangle.

| 1                     | <i>y</i> <sub>1</sub> | <i>y</i> <sub>2</sub> | <i>y</i> 3 | <i>y</i> 4 | 1                     | <i>y</i> <sub>1</sub> | <i>y</i> <sub>2</sub> | <i>y</i> 3 | <i>y</i> 4 | 1                     | <i>y</i> <sub>1</sub> | <i>y</i> <sub>2</sub> | <i>y</i> 3 | <i>y</i> 4 |
|-----------------------|-----------------------|-----------------------|------------|------------|-----------------------|-----------------------|-----------------------|------------|------------|-----------------------|-----------------------|-----------------------|------------|------------|
| $x_1$                 | Х                     | Х                     | Х          | Х          | <i>x</i> <sub>1</sub> | Х                     | Х                     | Х          | Х          | <i>x</i> <sub>1</sub> | Х                     | Х                     | Х          | Х          |
| <i>x</i> <sub>2</sub> | X                     |                       | Х          | Х          | <i>x</i> <sub>2</sub> | Х                     |                       | Х          | Х          | x <sub>2</sub>        | X                     |                       | Х          | X          |
| <i>x</i> 3            |                       | Х                     | Х          | X          | <i>x</i> 3            |                       | Х                     | Х          | Х          | <i>x</i> <sub>3</sub> |                       | Х                     | Х          | Х          |
| <i>x</i> 4            |                       | Х                     | Х          | X          | <i>x</i> 4            |                       | Х                     | Х          | X          | <i>x</i> <sub>4</sub> |                       | Х                     | Х          | X          |
| <i>x</i> 5            | Х                     |                       |            |            | <i>x</i> 5            | Х                     |                       |            |            | <i>x</i> 5            | Х                     |                       |            |            |

formal concepts (= maximal rectangles)

$$(A_1, B_1) = (\{x_1, x_2, x_3, x_4\}, \{y_3, y_4\})$$
$$(A_2, B_2) = (\{x_1, x_3, x_4\}, \{y_2, y_3, y_4\})$$
$$(A_3, B_3) = (\{x_1, x_2\}, \{y_1, y_3, y_4\})$$

# Literature on FCA

- books:
  - Ganter B., Wille R.: Formal Concept Analysis. Springer, 1999.
  - Carpineto C., Romano G.: Concept Data Analysis. Wiley, 2004.
- conferences: ICFCA (Int. Conf. on Formal Concept Analysis), CLA (Concept Lattices and Their Applications), ICCS (Int. Conf. on Conceptual Structures)
- web: useful resources and links at http://www.upriss.org.uk/fca/fca.html ("FCA Homepage")

### state of the art:

- Ganter B., Stumme G., Wille R. (Eds.): Formal Concept Analysis Foundations and Applications. Springer, LNCS 3626, 2005.
- theretical foundations,
- algorithms,
- increasingly popular applications (information retrieval, software engineering, ...),
- interaction with other methods of data analysis (preprocessing),
- software available.

◆□▶ ◆圖▶ ◆圖▶ ◆圖▶ ─ 圖

# Selected applications of FCA

- software engineering
- association rule mining closed frequent itemsets instead of frequent itemsets ⇒ non-redundant association rules (much less than by usual approach)
- (Boolean) factor analysis factors = selected formal concepts ... "new attributes"
- information retrieval, knowledge extraction structured view on data
- machine learning (decision making), clustering and classification preprocessing input data

• . . .

see the slides "Relational Data Analysis: Applications of Formal Concept Analysis (FCA)"

・ 同 ト ・ ヨ ト ・ ヨ ト

# FCA in Information Retrieval

pioneering work of R. Godin; C. Carpineto, G. Romano; elaborated by P. Eklund, J. Ducrou

main ideas:

- formal context = documents (objects) + index terms (attributes)
- (query/search) formal concept = (query) terms (intent) +
  retrieved documents (extent)
- query concept neighbors = minimal conjunctive refinements (specialization), enlargements (generalization) and alterations (categorization) of the query

A B M A B M

# Improving search engines with FCA

basic ideas:

- forwarding user query to a (web) search engine (Google, Yahoo etc., in a format such as SOAP), receiving ranked results (typically in XML format),
- parsing (first) results, indexing the document/snippet/title terms, optionally ranking the results,
- establishing formal context (possibly with attribute ordering = thesaurus),
- **computing** (part of the) **concept lattice** of the results, optionally ranking the results, displaying it to the user and
- enabling the user to appropriately modify the query by **navigating** through the lattice of the results (around the query concept)

more detailed treatment in Carpineto C., Romano G.: Concept Data Analysis. Wiley, 2004 (Chap. 3, 4).

イロト 不得下 イヨト イヨト 二日

# Improving search engines with FCA

indexing the document terms (studied in Information Retrieval):

- text segmentation
- word stemming using a rule-based stemmer (e.g. Porter's) or a lexical knowledge base
- stop wording
- word weighting crucial, "term frequency-inverse document frequency" (tf-idf) scheme implemented (most often) by a vector space model with a suitable weighting function, for web documents also URL, title, links etc.
- word selection removing terms with low weight
- document ranking

can be seen as a feature/attribute selection problem from data mining

超す イヨト イヨト ニヨ

# Improving search engines with FCA

document ranking (concept-lattice based ranking):

- similar to hierarchical clustering-based ranking
- **conceptual distance** between query/search concept and other document concepts in concept lattice instead of heuristic metric
- overcomes the vocabulary problem (word mismatch) seen in best-match ranking (used by current search engines)

possible difficulties:

- computational constraints → computing part of the concept lattice around the query concept = neighbor-like algorithms
- effective concept lattice visualization  $\rightarrow$  show query concept neighborhood only (focus+context techniques, tree below query concept)

existing (prototype) systems: CREDO, FooCA, SearchSleuth

くほと くほと くほと

# CREDO

- system for Conceptual REorganization of DOcuments, developed by Carpineto and Romano at Fondazione Ugo Bordoni, Italy
- displays the upper part (two levels from the top element) of the iceberg concept lattice (adding terms down the lattice), in the form of a tree
- enables "offline" navigation in concepts, narrowing the scope of the search
- Carpineto C., Romano G.: Exploiting the Potential of Concept Lattices for Information Retrieval with CREDO. J. Universal Computer Science 10(8)(2004), 985–1013
- search tool available at http://credo.fub.it
- mobile version CREDINO, http://credino.dimi.uniud.it

illustration:

- search for "dwarf" (ambiguous term), "phoenix",
- compare the results obtained by Credo vs. Google or Yahoo

# **CREDO**

| 00000  | Enter a query: formal concept analysis Search   |
|--|---|
| CREDU  | © English C Italiano <u>help terms of use</u> about   |
| <ul> <li>formal concept analysis (100)</li> <li>concepts (67) <ul> <li>fcn (23)</li> <li>using (17)</li> <li>data (12)</li> <li>mining (11)</li> <li>latices (8)</li> <li>knowledge (7)</li> <li>conceptual (7)</li> <li>introduction (5)</li> <li>structures (5)</li> <li>mathematical (5)</li> <li>view (4)</li> <li>code (4)</li> <li>code (4)</li> <li>code (4)</li> <li>code (4)</li> <li>code (4)</li> <li>design (5)</li> <li>struct (5)</li> <li>struct (2)</li> <li>ising (12)</li> <li>ising (12)</li> <li>isnowledge (9)</li> <li>tatical (9)</li> <li>conceptual (8)</li> <li>introduction (6)</li> <li>international (0)</li> </ul></li></ul> | A <u>Topological Framework for Formal Concept Analysis</u> on formal cencept analysis nives of the rich content of algebraic. topology. 1. Introduction. The idea of formal. 1Formal Concepts and     Concept Lattices     www.lates.cubk.edu.bki-cpbwong.iccs_04.pdf     Formal Concept Analysis to Learn from the Sixyphus-III Material      mixet are to idea behalf been areas Language and the state of the six of the rich content of algebraic.     Topology. 1. Introduction of the     Keeps unalgary.ex RAW KAW98 endmann     Formal Concept Analysis to Learn from the Sixyphus-III Material      Introduction of Concept Analysis areas and the site of the strongest concepts as the elements in the latice which is next to     bottom of Concept Analysis formal concept 1 and     www.cs.miron.edu/classes (sci20 fallo5 concept] pdf     Introduction to FCA     Mongest the Deares of the Dasic Concept analysis (TCA) is based on mathematical order theory and is a groups are called concepts which can     be represented     represented     www.stationa.edu/classes (sci20 fallo5 concept] pdf     Introduction to FCA     Introduction to FCA     Introduction to FCA     Mongest Analysis (the Configure Introduction to the Basic Features     In the following we try to explain the basic concept of formal concept     analysis     www.anathematic.tu-damistad.de/-bumeister ConfingIntro pdf |
|  |   |

▲ ■ ト ■ ク ۹ (~) Mar 2009 15 / 20

・ロト ・ 日 ト ・ ヨ ト ・ ヨ ト

# FooCA

- FCA + Google, developed by Bjoern Koester at Webstrategy GmbH, Darmstadt and TU Dresden, Germany
- presents search results directly in a form of formal context (documents × terms), additionaly represented by labelled Hasse diagram of the concept lattice (clicking in the table or on the diagram nodes opens a browser window with URLs)
- "online" navigation in concepts adding or removing attributes triggers new search and concept hierarchy formation
- B. Koester: FooCA Web Information Retrieval with Formal Concept Analysis. Verlag Allgemeine Wissenschaft, Mhltal, 2006. ISBN 9783-935924-06-1.

B. Koester: Conceptual Knowledge Retrieval with FooCA: Improving Web Search Engine Results with Contexts and Concept Hierarchies. Proc. ICDM 2006, Springer-Verlag, Berlin, 2006.

- search tool at http://fooca.webstrategy.de - requires
registration

### FooCA

| FooCA  |  |
|--|--|
| Search Formal Concept Analysis Yahoo 🔽 English 🔟 FooCA   |  |
| Retrieval results 10 🗘 Min. objects per attribute 2 🗘 Min. attribute length 3 🗘  |  |
| □ Stemming 17 Stopwords 17 Clarify context 17 Context refinement<br>17 Attribute ranking □ Show original results □ Show extracted attributes |  |

#### Your FooCA search for Formal Concept Analysis brought these results:



About FooCA and Terms of Use. FooCA is powered by Yahoo! Search

イロン イヨン イヨン イヨン

# SearchSleuth

- developed by Peter Eklund and Jon Ducrou within KVO (Knowledge, Visualization and Ordering), University of Wollongong, Australia, following ImageSleuth in the conceptual neighborhood paradigm
- displays the neighbors and siblings of the query/search concept (direct query generalization, specialization and categorization), in the form of text labels (links) of terms/attributes determining the concepts
- "online" navigation, multiple searches per query for neighbors of query concept, to expand the formal context
- J. Ducrou, P. Eklund: SearchSleuth: The Conceptual Neighbourhood of an Web Query. Proc. CLA 2007, LIRMM & University of Montpellier II, 2007.
- search tool available at

http://www.kvocentral.org/software/searchsleuth.html

illustration:

- search for "dwarf" (ambiguous term), "phoenix",
- compare the results obtained by SearchSleuth\_vs. Google or Yahoo one

# SearchSleuth

### -analysis

formal concept analysis ~[formal concept fca] ~[formal concept context]

+fca +data +lattice +context +based +mathematics +mining +theory +method +conceptual

### 1. Formal Concept Analysis Homepage

Formal Concept Analysis is a method of conceptual knowledge representation and data analysis. ... Christian Lindig's Concepts, (in C, older version: TkConcept? ... www.upriss.org.uk/fca/fca.html

### Formal concept analysis - Wikipedia, the free encyclopedia

... example concepts satisfy the formal definitions; the ... describing formal concept analysis for computer scientists. A Formal Concept Analysis Homepage ... en.wikipedia.org/wiki/Formal\_concept\_analysis

### 3. Formal Concept Analysis

Formal Concept Analysis is a branch of applied mathematics. ... Several books on Formal Concept Analysis have appeared, among them the first ... www.math.tu-dresden.de/~ganter/fbs.html

#### Formal Concept Analysis

Formal Concept Analysis (FCA) is a method mainly used for the analysis ... into units which are formal abstractions of concepts of human thought, allowing ... www.cs.cmu.edu/afs/cs.cmu.edu/project/jair/pub/volume24/cimiano05a-html/node3...

### 5. Linguistic Applications of Formal Concept Analysis

scribes the role that formal concept analysis can play in the automated or ... Associative and Formal Concepts. In: Priss; Corbett; Angelova (eds.), Con ... www.upriss.org.uk/papers/fcaic03.pdf

(日) (同) (日) (日) (日)

# Furher usage of the approach

(existing) usage besides web search:

- digital library search (Virtual Museum of the Pacific, requires registration),
- scientific (biology, medicine, ...) or social records mining,
- annotated multimedia archive search (ImageSleuth, DVDSleuth),
- email message search (MailSleuth),
- software documentation search,
- ... searching any other database of interest.

possible usage/improvements:

- other difficult IR tasks, e.g. natural language processing
- integration with IR techniques