



OTTO-VON-GUERICKE UNIVERSITY
MAGDEBURG

FACULTY OF COMPUTER SCIENCE

BACHELOR'S THESIS

Interactive Visualization of Large Concept Lattices for Exploratory Search

Author:

Johannes FILTER

Advisors:

Prof. Dr. Andreas NÜRNBERGER

Otto-von-Guericke University Magdeburg

Prof. Dr. Ana GARCÍA-SERRANO

Universidad Nacional de Educación a Distancia

July 6, 2015

Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Inhaltsangabe

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Acknowledgements

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Contents

1	Introduction	11
2	Background	13
2.1	Formal Concept Analysis	13
2.1.1	Definition	13
2.1.2	Formal Concept Analysis for Information Retrieval . .	14
2.2	Information Seeking Mantra	15
2.2.1	The Mantra	16
3	Related Work	17
3.1	Local View	17
3.2	Transform to Tree	18
3.3	Pruning Nodes	19
3.4	Conclusions	20
4	Fancy FCA 1.0	22
4.1	My Idea	22
4.2	Implementation	22
5	Evaluation of the User Interface	23
5.1	Fundamentals	23
5.1.1	Categories of Usability Studies	24
5.1.2	Data Collections	25
5.1.3	Participants	26
5.1.4	Tasks	26
5.2	Methods	26
5.2.1	Hypothesis	26
5.2.2	Evaluation Design	26
5.3	Results	26

5.4 Conluciosn	26
6 Fancy FCA 2.0	27
7 Second User Evaluations	28
8 Conclusions	29

1. Introduction

The digital revolution is affecting every part of our life. Also the humanities scholars stand before a big change in their ways when huge analog collections are digitized. They have to apply computer science methods to organize and analyze huge amount of data. The term "Digital Humanities" evolved during the last 10 years which can be defined as an "intersection between the humanities and information technology" [31].

The Information Retrieval Department of the Universidad Nacional de Educación a Distancia (UNED) in Madrid (Spain) cooperates with human scholars to conduct research in the Digital Humanities. In this project, there are historical maps which have been digitized and annotated. To extract knowledge from the collection the research group advocates for the use of Formal Concept Analysis (FCA) for topic organization [8, 9].

They successfully implemented a FCA algorithm but lack an interactive web-based user interface which will be developed in this thesis.

While FCA is a mathematically well-funded principle, the resulting traditional visualization of large concept lattices are a problem. Large concept lattices arose when you apply FCA to large amount of entities (Details will be explained in Chapter 2). When applying FCA to a document collection, you are likely to occur huge amount of entities. That is why alternative visualization techniques are important to get the insights of FCA, even if the lattice is large.

This thesis does not focus on the visualization of FCA itself. It focus on the visualization of FCA for information retrieval - especially exploratory search. Exploratory search is vaguely defined as the need to find something you did not know about before.

This user interface will be running in the browser. Because of the fast-changing environment of the web, it is important to keep with the latest technologies and techniques to not fall apart. Besides others, the software utilizes the frameworks d3.js and Bootstrap to create a pleasant user interface. The website is fast-responding because it reduces the communication between browser and web server to a minimum. In most cases, instead of reloading the page, the interface only changes DOM elements.

The remainder of this thesis is structured as follows: The background of Formal Concept Analysis will be presented in Chapter 2 and the background of User Search Interfaces in Chapter 3. In Chapter 4 I will present my (first) approach and the implementation, which will be evaluated in Chapter 5. Built on the Evaluation, I will adjust my work and present an updated version of my work in Chapter 6. I conclude in Chapter 7 and give ideas for future work.

2. Background

This chapter describes Formal Concept Analysis in the first part and the 'information seeking mantra' a guideline for user interface design in the second.

2.1 Formal Concept Analysis

2.1.1 Definition

Formal Concept Analysis (FCA) [17] is a mathematically well-funded technique to find relationships among objects. Formally, a *formal context* is defined as a tripple $K = (G, M, I)$ where G is a set of objects, M is a set of attributes and I is a binary relation $I \subseteq G \times M$. I specifies whether an object has an attribute or not. Table 2.1 illustrates an example (from David Eppstein [16]) where G comprises the integers from 1 to 10 and M comprises the attributes composite, even, odd, prime and square.

A *formal concept* is a pair of (A, B) where $A \subseteq G$, a set of objects called *extent*, and where $B \subseteq M$, a set of attributes called *intent*. From the example in 2.1, we can derive several formal concepts. For example:

- $C_1 = (A_1, B_1)$, with $A_1 = \{4, 6, 8, 10\}$ and $B_1 = \{composite, even\}$
- $C_2 = (A_2, B_2)$, with $A_2 = \{2, 4, 6, 8, 10\}$ and $B_2 = \{even\}$
- $C_3 = (A_3, B_3)$, with $A_3 = \{9\}$ and $B_3 = \{composite, odd, square\}$

A partial order among formal concepts is defined as follows:

$$(A_i, B_i) \leq (A_j, B_j) \iff A_i \subseteq A_j \quad (2.1)$$

It can be seen, that $C_1 \leq C_2$ and that C_3 is unrelated to C_1 , and that C_3 is unrelated to C_2 .

Table 2.1: Formal context, integers 1 to 10 as objects and attributes

	composite	even	odd	prime	square
1			×		×
2		×		×	
3			×	×	
4	×	×			×
5			×	×	
6	×	×			
7			×	×	
8	×	×			
9	×		×		×
10	×	×			

The set of formal concepts with the partial order are called *concept lattice*. A concept lattices can be visually represented in a *Hasse diagram*. The vertices represent the formal concepts and the edges are there if they are directly related in the partial order. The most general formal concepts are in the top and the most specific ones are in the bottom. There have been added two special formal concepts. One for the most specific on the top containing all possible obojects, and one most general containing no objects in the bottom. The Figure 2.1 is the corresponding Hasse diagram to the formal context presented in Table 2.1.

We will describe in the next section how we can apply FCA to Information Retrieval.

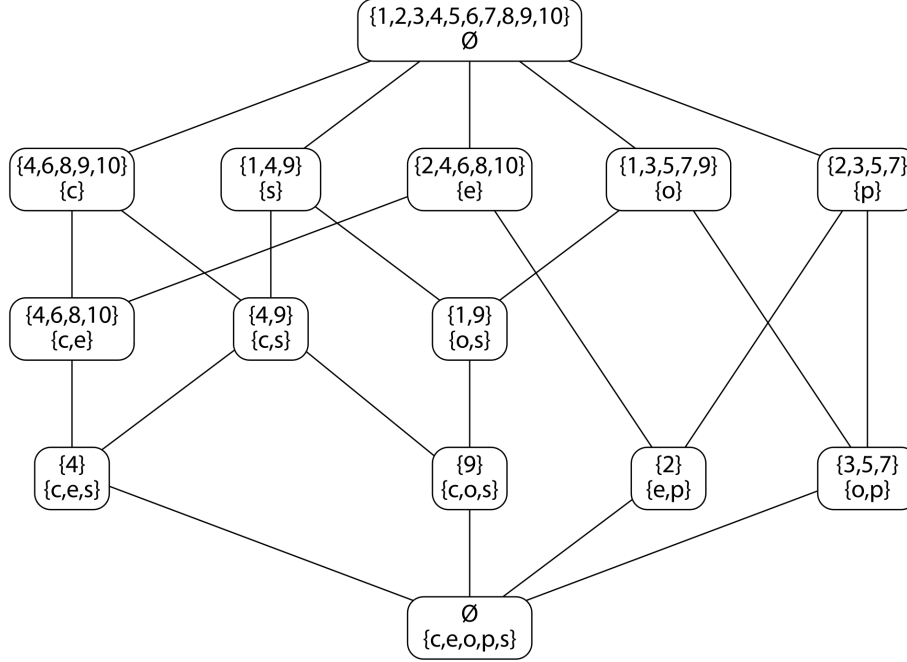
2.1.2 Formal Concept Analysis for Information Retrieval

Carpineto et al.[7] describe the start of FCA in information retrieval:

In the 80's, basic ideas were put forth - essentially that a concept can be seen as a query (the intent) with a set of retrieved documents (the extent).

This essentially means the appliance of the Standard Boolean Retrieval Model, which, as Manning et al. [26] describe, "is a model for information retrieval in which we can pose any query which is in the form of a Boolean

Figure 2.1: Hasse diagram, with the integers 1 to 10 as objects and attributes square (s), prime (p), composite (c), even (e), and odd (o)



expression of terms [...]. The model views each document as just a set of words.” . In a simple case, a system only supports conjunction of terms.

According to Poelmans et al. has been FCA ”applied in many disciplines such as software engineering, knowledge discovery and information retrieval” [28] and they did two comprehensive surveys on the application of FCA [28, 29]. This thesis focusses on the visualization of a concept lattice and not on the act of creating of one. The interested reader can read literature from Carpineto et al. [5, 7] to further investigate this area.

2.2 Information Seeking Mantra

The Information Seeking Mantra (The mantra) was introduced by Ben S. and summeraized user interface design principles. Albeit it was intended to be a description, he wrtoe ””, it is well respected in the scientific communtiy. Some paper question it an make clear, that there is no evidence for its

success.

2.2.1 The Mantra

3. Related Work

It has been shown, that FCA can be applied to information retrieval, but for now we did not consider the visualization of the concept lattice. While some user studies proclaim that non-experts can read Hasse diagram[14], the study has been conducted on relatively small lattices. On the field of information retrieval the objects can easily outreach a few dozens. Kuznetsov et al. [23] describe this resulting visualization.

Representing concept lattices constructed from large contexts often results in heavy, complex diagrams that can be impractical to handle and, eventually, to make sense of.

Especially enormous edge crossing can hinder the visual representation. Take a look at the appendix for the first results of the research group.

The visual representation of Hasse diagrams can be improved by fine-tuning visual components like labels, edges etc.. Or some ideas like a Fish-Eye Views XXX-QUOTE have to be applied to FCA. But this action does not scale well and won't help us with large concept lattices. To cope with large lattices, three reduction techniques exist which will be presented in the following: One where you visualize only a part of the lattice, one to transform it into a tree and one where you remove nodes from the concept lattice which means, that you modify the structure of the lattice.

3.1 Local View

Instead of showing the whole Hasse diagram to the user, only a small part of the lattice is visualized. The focus lies on one concept and its neighborhood. There exist several names and small variations of this idea. Eklund et al. name this idea *conceptual neighborhood* [13, 15]. The user can query the system or navigate through the lattice by going up (removing terms) or going down (adding terms). Only adjacent nodes are displayed in this

model. The user can incrementally browse the whole lattice. Eklund et al. applied this approach to a broad range of topic: for the 'Virtual Museum of the Pacific' [13, 15], image browsing [12, 11] and search engines [10].

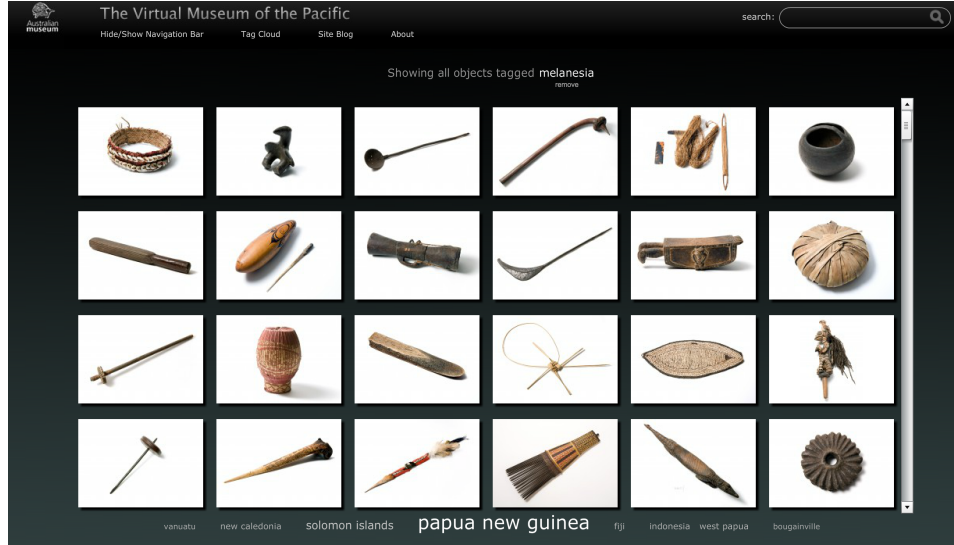


Figure 3.1: Screenshot of 'Virtual Museum of the Pacific', focus on concept 'melanesia'

Carpineto and Romano developed a search engines ULYSSES [3, 4], which visualizes the results in a similar way. But it visualizes a small sublattice - more than just the directly adjacent nodes. The size of the sublattices varies and can be fine-tuned by parameters. For instance, you can the degree of children or parents to visualize, which is the minimal distance between two nodes. You can see a screenshot of the software in Figure??.

In their following work, CREDO [6], Carpineto and Romano they restricted the system to only show directly neighboring nodes which a folding mechanism.

3.2 Transform to Tree

While transforming the lattice into a tree sounds promising, because you could apply sophisticated tree visualizing techniques to reduce edge crossing, it comes with several drawbacks. One naiv approach is described by Carpineto and Romano [5]: If a node hast more than one parent, remove

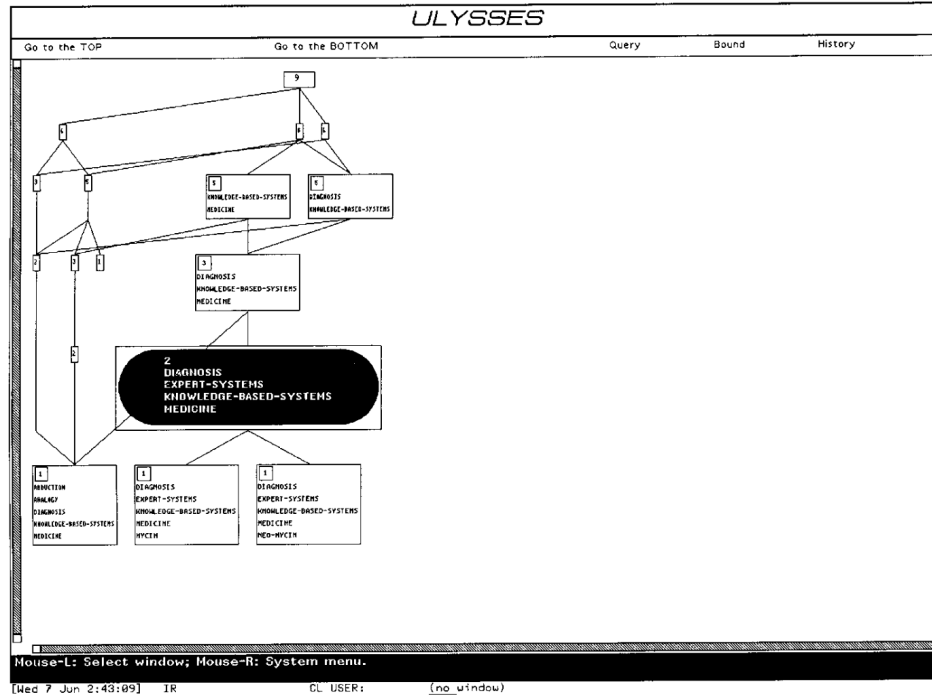


Figure 3.2: Display screen of ULYSSES, focusing on the black node [4]

that parent an insert a copy of then node and attach it to that parent. This is problematic, because we dramatically increase the number of nodes.

There exist another approach [27]: select the 'best' parent and hide edges to all other parents. While this technically not breaks the concept, the visual representation does not correspond to the underlying model.

3.3 Pruning Nodes

A prominent approach to prune lattices called "iceberg lattice" [30]. A variation from the frequent item-set mining which specific min-support and min-confidence[1]. It creates a top of the lattice but has some drawbacks because "One should be careful not to overlook small but interesting groups, for example, "exotic" or "emergent" groups not yet represented by a large number of objects, or, groups that contain objects who are not members of any other group." [23] The iceberg lattice just focuses on the concepts that contain a lot of documents. That is why an other approaches exist: Stability[24]: "A concept is stable if its intent does not depend much on

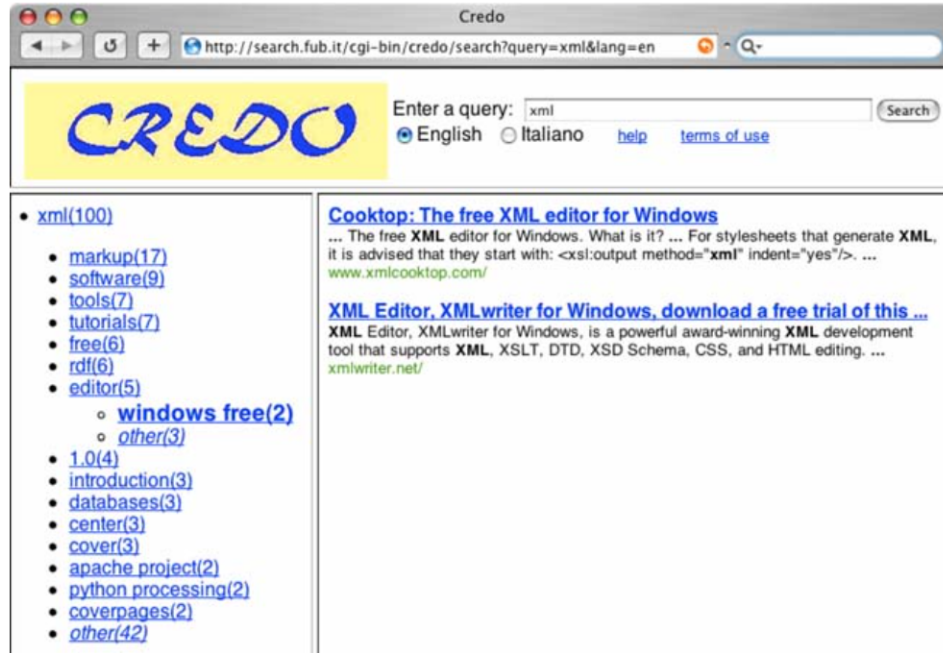


Figure 3.3: Screenshot of CREDO, after query 'xml' and browsing after 'editor(5)' and 'windows free(2)' [6]

each particular object of the extent.” [23]

It is also possible to apply traditional cluster techniques to FCA.[2]

3.4 Conclusions

We showed that there exist different possibilities to reduce the complexity of large lattices. But in our case, we cannot change the structure of the lattice. We apply FCA to explore the data and get insights about it. When pruning the nodes, you are losing many data relationships, many formal concepts and, consequently, the "power" of FCA as exploratory technique is significantly reduced. The application of the local view technique in combination with a query interface is best for our needs. Especially the use of a query interface is familiar to all user with the rise of web search engines. The 'transform to tree' approach without a query interface seems cumbersome and all together, it is really similar to the local view approach.

Before I critically analyze the current applications of local views in chapter XX, I want to give background information in (search) interface design

principles in the next chapter.

4. Fancy FCA 1.0

4.1 My Idea

We just want to find out if our approach, treating FCA as internal data structure and offering an user search interfaces + guidance (type ahead, showing neighboring concepts), is helpful for human scholars. This can be answered with a usability study.

4.2 Implementation

I will describe important parts of the implementation here.

5. Evaluation of the User Interface

The design of an interface is highly subjective. To judge if the user like it, there exist methods to evaluate the work. A small introduction into this field from the computer science perspective gives Hearst in his book on User Search Interfaces in Chapter 2. [18]. A comprehensive guide into the "Methods for Evaluating Interactive Information Retrieval Systems with User" gives Kelly [21]. The interested reader is advised read those papers because we will only scratch the surface.

Computer scientists are not experts in human studies and Zobel [32] proclaims: "Far too many human studies in computer science are amateurish and invalid." Nevertheless, I tried to be as scientific as possible to conduct human studies even with strict resource limitations.

I will describe the idea of the experiment first and refer to literature to illustrate my choices, I will describe the process of the experiment, and after that explain and evaluate the outcome of the results.

5.1 Fundamentals

If we walk about an evaluation, it is important to make clear what what aspects are evaluated. Hearst writes [18]:

Search interfaces are usually evaluated in terms of three main aspects of usability: effectiveness, efficiency, and satisfaction, which are defined by ISO 9241-11, 1998 [20] as:

- Effectiveness: Accuracy and completeness with which users achieve specified goals.

- Efficiency: Resources expended in relation to the accuracy and completeness with which users achieve goals.
- Satisfaction: Freedom from discomfort, and positive attitudes towards the use of the product.

These are the criteria that ideally should be measured when evaluating a search user interface.”

It also nice to know the difference between an experiment and evaluation.

They different ways how to evaluate the interfaces are described now.

5.1.1 Categories of Usability Studies

The different studies Hearst [18] propose can roughly be categorized as follows:

Informal Usability Testing

Hearst [18] describes the process shortly as ”Showing designs to participants and recording their responses”. It is often used in short iterative cycles to quickly evaluate a design. Thinking out loud.

Formal Studies and Controlled Experiments

Hearst [18] says that it is a ”form of controlled experiments aim to advance the field’s understanding of how people use interfaces, to determine which design concepts work well under what circumstances, and why.” In contrast to informal studies, it is important to isolate factors and not treat the whole system as a black box. Using eye tracker in a laboratory with 2-way windows is one example.

Longitudinal Studies

In contrast to the studies above, this focus more on a longer time period and observes user real behavior. Drawing conclusions from the analysis of Google Search Queries is an example.

Log Analysis

In contrast to the studies above, this focus more on a longer time period and observes user real behavior. Drawing conclusions from the analysis of Google Search Queries is an example.

Bucket testing (A/B Testing)

The traffic to a particular website is split and alternative view. It is evaluated how the users of the alternative website reacts to the new site. For example: Amazon changes its search filters and evaluates if the users buy more.

As you can see, this is only a small categorization and Kelly describes in her work the different approaches more in detail. Because a complete coverage of this topic would exceed this thesis, only some parts are covered here.

5.1.2 Data Collections

Kelly described different ways to collect data.

Questionnaires

A questionnaire comprises set of questions and is cheap and fast way to gather information from people. Kelly et al. [22] describe two types of questions as follows:

Questionnaires can be comprised of closed questions, open questions or a mixture of both. *Closed questions* are questions that provide a fixed set of responses with which subjects must respond. It is common practice for usability questionnaires to include closed questions in the form of statements such as, the system was easy to learn to use. Subjects are typically provided with 5–7-point Likert-type scales for responding, where one scale end-point represents strong agreement and the other represents strong disagreement. [...] *Open questions*, on the other hand, do not provide a response set and subjects are able to provide any type of response they feel is appropriate.

In study Kelly et al. conducted, they conducted

Hornbæk and Law did a well respected meta-analysis of usability studies and as one of their conclusions they "recommend that standard questionnaires be used when possible, given their higher reliability, and that the more complex effectiveness measures be used when feasible (as they are more likely to give information that cannot be obtained by measures in the other categories)." [19] We stick to their advice and use the USE questionnaire [25] which was (partly) used in the investigation from Kelly et al. [22].

Online Pen-and-Paper Interview

Thinking Alloud

Search Logs

5.1.3 Participants

It is important to reduce a structural bias of an experiment. Zobel [32] mentions that "the sample of human subjects should be representative (a class of computer science students may not be typical of users of mobile devices)". We tried to vary the users or at least focus on human scholars because that is the user group that is important for our stuff.

5.1.4 Tasks

To let them

5.2 Methods

5.2.1 Hypothesis

The interface is useful for the human scholars and they like the interface. It offers rich possibilities for the users to navigate along the different documents. Because of the similarity to popular search engines, they know how to use it. But there some stuff that is not implemented yet that they would like to see.

5.2.2 Evaluation Design

Because of our restricted resources, an informal usability study is the most attractive choice for us. The logs of the system are recored and can evaluated in future, but because of the long-term duration, they cannot be evaluated in this thesis.

5.3 Results

5.4 Conluciosn

6. Fancy FCA 2.0

7. Second User Evaluations

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

8. Conclusions

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Bibliography

- [1] Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. Mining association rules between sets of items in large databases. *ACM SIGMOD Record*, 22(2):207–216, 1993.
- [2] Ch Aswani Kumar and S. Srinivas. Concept lattice reduction using fuzzy K-Means clustering. *Expert Systems with Applications*, 37(3):2696–2704, 2010.
- [3] Claudio Carpineto and Giovanni Romano. ULYSSES: a lattice-based multiple interaction strategy retrieval interface. *Human-Computer Interaction*, pages 91–104, 1995.
- [4] Claudio Carpineto and Giovanni Romano. Information retrieval through hybrid navigation of lattice representations. *International Journal of Human-Computer Studies*, 45(5):553–578, 1996.
- [5] Claudio Carpineto and Giovanni Romano. *Concept data analysis: Theory and applications*. John Wiley & Sons, 2004.
- [6] Claudio Carpineto and Giovanni Romano. Exploiting the potential of concept lattices for information retrieval with CREDO. *Journal of Universal Computer Science*, 10(8):985 – 1013, 2004.
- [7] Claudio Carpineto and Giovanni Romano. Using Concept Lattices for Text Retrieval and Mining. *Formal Concept Analysis*, pages 161–179, 2005.
- [8] Ángel Castellanos, Ana García-serrano, Etsi Informática Uned, Juan Cigarrán, and Etsi Informática Uned. Concept-based Organization for semi-automatic Knowledge Inference in Digital Humanities : Modelling and Visualization.

- [9] J Cigarrán and A Castellanos. A step forward in Topic Detection Algorithms : An Approach based on Formal Concept Analysis . 16.
- [10] Frithjof Dau, Jon Ducrou, and Peter Eklund. Concept similarity and related categories in SearchSleuth. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 5113 LNAI(October 2007):255–268, 2008.
- [11] Jon Ducrou and Peter Eklund. an Intelligent User Interface for Browsing and Searching Mpeg-7 Images Using Concept Lattices. *International Journal of Foundations of Computer Science*, 19(02):359–381, 2008.
- [12] Jon Ducrou, Björn Vormbrock, and Peter Eklund. FCA-based browsing and searching of a collection of images. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 4068 LNAI:203–214, 2006.
- [13] P. Eklund, P. Goodall, T. Wray, B. Bunt, a. Lawson, L. Christidis, V. Daniel, and M. Van Olffen. Designing the digital ecosystem of the Virtual Museum of the Pacific. *2009 3rd IEEE International Conference on Digital Ecosystems and Technologies, DEST '09*, (June):377–383, 2009.
- [14] Peter Eklund, Jon Ducrou, and Peter Brawn. Concept lattices for information visualization: Can novices read line-diagrams? *Concept Lattices. Lecture Notes in Computer Science*, 2961:57–73, 2004.
- [15] Peter Eklund, Tim Wray, Peter Goodall, and Amanda Lawson. Design, information organisation and the evaluation of the Virtual Museum of the Pacific digital ecosystem. *Journal of Ambient Intelligence and Humanized Computing*, 3(4):265–280, 2012.
- [16] David Eppstein. Hass diagramm, with the integers 1 to 10 as objects and attributes square, prime, composite, even, and odd. https://commons.wikimedia.org/wiki/File:Concept_lattice.svg, 2006. Accessed on 16. June 2015.
- [17] Bernhard Ganter and Rudolf Wille. *Formal concept analysis: mathematical foundations*. Springer Science & Business Media, 2012.
- [18] Marti a Hearst. Search User Interfaces. *Search User Interfaces*, 54(Ch 1):404, 2009.

- [19] Hornb, Aelig, Kasper K, and Effie Lai-Chong Law. Meta-analysis of correlations among usability measures. *Proceedings of ACM CHI 2007 Conference on Human Factors in Computing Systems*, 1:617–626, 2007.
- [20] ISO. Ergonomic requirements for office work with visual display terminals (vdts) - part 11 : Guidance on usability. ISO 9241-11, International Organization for Standardization, Geneva, Switzerland, 1998.
- [21] Diane Kelly. Methods for Evaluating Interactive Information Retrieval Systems with Users. *Foundations and Trends® in Information Retrieval*, 3(1—2):1–224, 2007.
- [22] Diane Kelly, David J. Harper, and Brian Landau. Questionnaire mode effects in interactive information retrieval experiments. *Information Processing and Management*, 44(1):122–141, 2008.
- [23] Sergei Kuznetsov, Sergei Obiedkov, and Camille Roth. Reducing the Representation Complexity of Lattice-Based Taxonomies. pages 241–254, 2007.
- [24] Sergei Kuznetsov, Sergei Obiedkov, and Camille Roth. Reducing the Representation Complexity of Lattice-Based Taxonomies. pages 241–254, 2007.
- [25] Arnold M Lund. Measuring usability with the use questionnaire. *Usability interface*, 8(2):3–6, 2001.
- [26] Christopher D. Manning and Prabhakar Raghavan. An Introduction to Information Retrieval, 2009.
- [27] Cassio Melo, Bénédicte Le-Grand, Marie Aude Aufaure, and Anastasia Bezerianos. Extracting and visualising tree-like structures from concept lattices. *Proceedings of the International Conference on Information Visualisation*, 6:261–266, 2011.
- [28] J Poelmans and So Kuznetsov. Formal concept analysis in knowledge processing: a survey on models and techniques. *Expert systems with ...*, (2003):1–40, 2013.
- [29] J Poelmans and So Kuznetsov. Formal concept analysis in knowledge processing: a survey on models and techniques. *Expert systems with ...*, (2003):1–40, 2013.

- [30] Gerd Stumme, Rafik Taouil, Yves Bastide, Nicolas Pasquier, and Lotfi Lakhal. Computing iceberg concept lattices with TITANIC. *Data and Knowledge Engineering*, 42(2):189–222, 2002.
- [31] Patrik Svensson. The Landscape of Digital Humanities. *DHQ: Digital Humanities Quarterly*, 4(1):1–34, 2010.
- [32] Justin Zobel. *Writing for Computer Science*. 2004.