# (Semi)-automated subject indexing of Swedish resources:

## Evaluating (a combination of)
## cataloguers', end users' and automated index terms
## in retrieval

Koraljka Golub

*International Conference on Dublin Core and Metadata Applications 2016 (DC-2016)*

**Linnæus University**

# Purpose

- Improve searching and browsing for information in the Swedish language, esp. in interoperable information systems like LIBRIS, SwePub, Sondera

- Subject index terms from controlled vocabularies like SAO and DDC offer:

  - Uniformity of term format

  - Provide context

  - Browsing

- But, expensive while increase of digital documents → 2 possible solutions:

  - 1) (semi)-automated solutions

  - 2) author/end-user tagging

**Linnæus University**

Växjö

SEARCH

▼ Extended search

## Your search on "växjö" resulted in 32270 hits

Sort by: Relevancy

Get it? ▼

**5905** Hits

### Archive → Show in NAD

| | |
|---|---|
| Box (3211) | Map/plan/drawing (1485) |
| Fonds (868) | Series (225) | File (90) |
| Photograph (13) | Computer file (9) |
| Item (4) | |

**KFUM Växjö**
KFUM Växjö
Fonds, 1945–1945

**Växjö landsförsamlings kyrkoarkiv**
Växjö landsförsamling
Fonds (Microfilmed) (Digitized), 1704–1939
▸ Content

**Växjö stiftsråds arkiv**
Växjö stiftsråd
Fonds, 1921–1989

**6547** Hits

### Library → Show in LIBRIS

| | |
|---|---|
| Book (4684) | Article/chapter (907) |
| Journal etc. (373) | Other (215) |
| Speech (136) | Film/video (67) |
| Map (66) | Image (36) |
| Multimedia (30) | Poster (15) |
| Musical score (12) | Music (3) |
| Manuscript (3) | |

**Välkommen till Växjö**
Book, 1985

**Välkommen till Växjö**
Book, 1989

**Växjö : gatu- och kvartersnamn / Eva Selling**
Selling, Eva, 1936–

**19818** Hits

### Sound & Image → Show in SMDB

| | |
|---|---|
| Radio (11003) | TV (8248) |
| Recorded sound (392) | Film/video (167) |
| Multimedia (11) | Texts (1) |

**[Växjö, Föreningen Öppna Kanalen Växjö, 2015-08-03--2015-08-09]**
TV, 2015-08-03--2015-08-09

**[Växjö, Föreningen Öppna Kanalen Växjö, 2015-11-02--2015-11-08]**
TV, 2015-11-02--2015-11-08

**Morgon i P4 med Peje Johansson, Anneli Koskinen och Lars-Peter Hielle [Växjö, ...**
Radio (Digitized), 2016-09-23
▸ Content

## Linnæus University
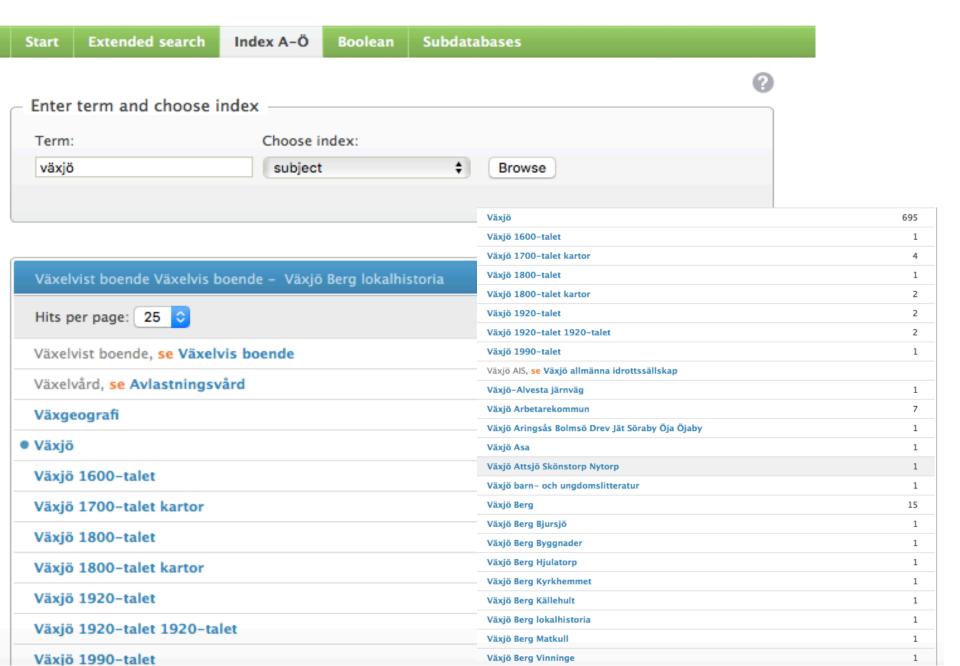
SwePub

**Refine your search**

**Availablility**
free online (944)

Sort by: relevancy ⬍   Hits per page: 10 ⬍

**Type of publication**
journal article (1129)
conference paper (1070)
book chapter (736)
reports (510)
doctoral thesis (399)
show more...

1. **Inkast : Idrottsforskning vid Växjö universitet**
   2006
   **Editorial collection** (other academic)

**Type of content**
other academic (2365)
peer-reviewed (1550)
pop. science, debate, etc.
(737)

2. Sundberg, Ann-Kari, 1962– (author)            ⤻ Read full
   **Le poids de la tradition La gestion professorale de l'altérité linguistique et culturelle en classe de FLE**
   2009
   **Doctoral thesis** (other academic)                ▸ abstract

**Author/Editor**
Khrennikov, Andrei, (130)
Hagevi, Magnus, 1965 ... (81)
Fridlund, Bengt (79)
Al-Najjar, Basim, 19 ... (70)
Milrad, Marcelo, (57)
Silander, Daniel, 19 ... (53)
show more...

3. **Blåser luft i fåglar : Kreativt skrivande i Växjö**
   2009
   **Editorial collection** (pop. science, debate, etc.)

4. Olofsson, Gunnar, 1942–, et al. (author)            ⤻ Read full
   **Akademiska yrkesutbildningar i Växjö Program, studenter och arbetsmarknad**
   2010
   **Book** (other academic)                ▸ abstract

⊞ **University**

⊞ **Language**

5. Olofsson, Tommy (author)
   **Tegnér, Växjö och den sena diktningen**
   2006
   **Book** (other academic)                ▸ abstract

**Research subject (UKÄ/SCB)**
Social Sciences (1689)
Humanities (834)
Natural sciences (771)
Engineering and Technology
(449)
Medical and Health Sciences
(263)
Agricultural Sciences (117)

6. **Policing in Scandinavia Proceedings from the Conference on Police Research in Växjö, August 2007**            ⤻ Read full
   2008
   **Editorial proceedings** (other academic)

7. Kratz, Helene, et al. (author)            ⤻ Read full
   **Lärarutbildningen vid Växjö universitet ur ett studentperspektiv**
   2004
   **Reports** (other academic)                ▸ abstract

**Linnæus University**

# LIBRIS

| Start | Extended search | **Index A–Ö** | Boolean | Subdatabases |

## Enter term and choose index

**Term:**

`växjö`

**Choose index:**

`subject` ▼

**Browse**

| | |
|---|---|
| **Växjö** | 695 |
| **Växjö 1600–talet** | 1 |
| **Växjö 1700–talet kartor** | 4 |
| **Växjö 1800–talet** | 1 |
| **Växjö 1800–talet kartor** | 2 |
| **Växjö 1920–talet** | 2 |
| **Växjö 1920–talet 1920–talet** | 2 |
| **Växjö 1990–talet** | 1 |
| Växjö AIS, *se* Växjö allmänna idrottssällskap | |
| **Växjö–Alvesta järnväg** | 1 |
| **Växjö Arbetarekommun** | 7 |
| **Växjö Aringsås Bolmsö Drev Jät Söraby Öja Öjaby** | 1 |
| **Växjö Asa** | 1 |
| Växjö Attsjö Skönstorp Nytorp | 1 |
| **Växjö barn– och ungdomslitteratur** | 1 |
| **Växjö Berg** | 15 |
| **Växjö Berg Bjursjö** | 1 |
| **Växjö Berg Byggnader** | 1 |
| **Växjö Berg Hjulatorp** | 1 |
| **Växjö Berg Kyrkhemmet** | 1 |
| **Växjö Berg Källehult** | 1 |
| **Växjö Berg lokalhistoria** | 1 |
| **Växjö Berg Matkull** | 1 |
| **Växjö Berg Vinninge** | 1 |

### Växelvist boende Växelvis boende – Växjö Berg lokalhistoria

Hits per page: `25` ▼

Växelvist boende, *se* **Växelvis boende**

Växelvård, *se* **Avlastningsvård**

**Växgeografi**

● **Växjö**

**Växjö 1600–talet**

**Växjö 1700–talet kartor**

**Växjö 1800–talet**

**Växjö 1800–talet kartor**

**Växjö 1920–talet**

**Växjö 1920–talet 1920–talet**

**Växjö 1990–talet**

## DDK:s huvudklasser ▾

| DDK-nummer | Rubrik |
|---|---|
| | DDK:s huvudklasser |
| 500 | Naturvetenskap |
| 540 | Kemi |
| **541-547** | **Kemi** |
| 541 | Fysikalisk kemi |
| 542 | Tekniker, procedurer; apparatur, utrustning, material |
| 543 | Analytisk kemi |
| 546 | Oorganisk kemi |
| 547 | Organisk kemi |

## DDK:s huvudklasser

Linnæus University

# Aims

1. Find out to what degree it is possible to apply automated subject indexing based on:

- Controlled indexing languages like the Dewey Decimal Classification (DDC) and Swedish Subject Headings (SAO)
- Derived indexing of keywords from the resource itself

2. Determine the value of automatically assigned index terms, in combination and comparison with end-user and cataloguers index terms in the process of information retrieval by end users

**Linnæus University**

# BACKGROUND

# End-user indexing

- Author keywords or social tags in Web 2.0 services also provided by library catalogues
- Cheaper, provide additional perspectives like new scientific terms
- But, no control of word forms, homonymy, polisemy or synonymy

*"…the cost savings made in the provision of low-quality indexing are cancelled out by the high costs incurred by searchers who fail either to find everything that they want (low recall) or, often more frustratingly, to avoid everything that they do not want (low precision)…"*

(Furner 2010, 1861)

**Linnæus University**

# Offering users to choose from KOS

# Offering users to choose from KOS

- The importance of controlled vocabulary suggestions for indexing and retrieval:
  - To help produce ideas of which tags to use
  - To make it easier to find focus for the tagging
  - To ensure consistency
  - To increase the number of access points in retrieval

  - However, the value and usefulness of the suggestions proved to be dependent on the quality of the suggestions, both as to conceptual relevance to the user and as to appropriateness of the terminology

  (Golub, Lykke, Tudhope 2014)

**Linnæus University**

# Automatic indexing…

3 major approaches

- Text categorization

- Document clustering

- String matching

**Linnæus University**

# …Automatic indexing…

Automatic indexing beneficial

- – Address the scale and sustainability
- – Enrich bibliographic records
- – Establish more connections across resources

Reported success of automated tools

- – Entirely replace manual indexing to machine-aided indexing (MAI)
  - • MAI example: NLM´s Medical Text Indexer

…Automatic indexing

Evaluation problem (Golub et al. 2016)

- Research comparing automatic versus manual indexing is flawed (Lancaster 2003, p. 334)

  - Out of context, laboratory conditions
  - Few reports on indexing tools in operating information systems

# Challenge A: relevance 1/2

- Purpose of indexing: making relevant documents retrievable

- Relevance
  - A complex phenomenon
    - Many possible document-query relationships
    - E.g., for children/scientists, query/information need/task…
  - Subjective
  - Multidimensional and dynamic (Borlund 2003)

# Challenge A: relevance 2/2

In practice, evaluation of IR is based on pre-existing relevance assessments

- – Initiated by Cranfield tests
- – A gold standard
    - • A test collection consisting of a set of documents
    - • A set of 'topics'
    - • A set of relevance assessments

- – *"In spite of the dynamic and multidimensional nature of relevance, in practice evaluation of information retrieval systems has been reduced to comparison against the gold standard—a set of pre-existing relevance judgments which are taken out of context. An early study on retrieval conducted by Gull in 1956 powerfully influenced the selection of a method for obtaining relevance judgments. Gull reported that two groups of judges could not agree on relevance judgments. Since then it has become common practice to not use more than a single judge or a single object for establishing a gold standard."*

(Saracevic 2008, 774)

# Challenge B: indexing

Aboutness
- Dependent on factors like interest, task, purpose, knowledge…

Exhaustivity and specificity of indexing
- Related to indexing policies at hand
- A subject correctly assigned in a high-exhaustivity system may be erroneous in a low-exhaustivity system

Terms assigned automatically but not manually might be wrong or they might be right but missed by manual indexing

→ Not good to use just the existing classes as the gold standard

**Linnæus University**

# METHODOLOGY

**Linnæus University**

# Methodology …

- Based on Golub et al. (2016)
- Data collection: a subset of SMDB, SND, SwePub

- A comparison of assigned terms against a carefully crafted 'gold standard' and in the context of actual information retrieval

  - The '**gold standard**' developed through input of
    - Professional catalogue librarians
    - End users who are experts in the subject at hand,
    - End users who are inexperienced in the subject
    - Several automated subject indexing software applications

  - **Information retrieval** will involve end users conducting actual searching on the indexed collection of resources and marking how relevant each retrieved resource is

**Linnæus University**

19

# … Methodology

- Automated subject index terms derived from several algorithms will be compared against the 'gold standard' and in the retrieval test

- The analysis will also include looking at what caused the retrieval of the document at hand: a cataloguer's term, subject expert's term, inexperienced user's term or an automated term

- Also log analysis and questionnaires to help contextualize the results

**Linnæus University**

20

# Automated indexers

- To be built/adjusted for DDC, SAO and the Swedish language

- String matching

- Machine learning

- Commercial: Data Harmony (rule-based)

**Linnæus University**

# SIGNIFICANCE

# Significance…

The value of the professional, the automated and the end-user ways of creating subject index terms will be determined

→ Allowing for informed decisions on ensuring high quality subject access points as part of the Swedish library and information infrastructure

– Cheap assignment of controlled subject terms useful at various stages of the metadata creation workflow:

• By an author creating original index terms at the time of deposit;

• By a reader annotating (for colleagues/world or for recommendation for inclusion in a collection);

• By a cataloguer

**Linnæus University**

# …Significance

– Provision of (semi)-automated solutions for assigning DDC will enable:

- Hierarchical browsing by subject

- Retrievability of Swedish resources in multilingual systems

- Integration of Swedish resources into the Semantic Web (as DDC is available as Linked Data)

– The resulting empirically tested comprehensive methodology framework will be of interest to other researchers

**Linnæus University**

# Partners

Linnaeus University
Lund University
University of South Wales

University of Aalborg
University of Buffalo
Charles Sturt University

Swedish National Library
Linnaeus University Library
Access Innovations
OCLC

**Linnæus University**

25

# References

Borlund, P. (2003). The concept of relevance in IR. *Journal of the American Society for Information Science and Technology 54(10)*, 913-925.

Golub, K., Dagobert, S., Buchanan, G., Tudhope, D., Lykke, M. & Hiom, D. (2016). A framework for evaluating automatic indexing or classification in the context of retrieval. *Journal of the Association for Information Science and Technology, 67(1),* 3-16.

Golub, K., Lykke, M. & Tudhope, Douglas (2014). Enhancing social tagging with automated keywords from the Dewey Decimal Classification. *Journal of Documentation, 70(5)*, 801-828.

Hjørland, B. (2002). Epistemology and the socio-cognitive perspective in information science. *Journal of the American Society for Information Science and Technology 53(4),* 257-270.

Lancaster, F. W. (2003). Indexing and abstracting in theory and practice. 3rd ed. Champaign: University of Illinois.

Saracevic, T. (2008). Effects of inconsistent relevance judgments on information retrieval test results: A historical perspective. *Library Trends 56(4)*, 763-783.

Soergel, D. (1994). Indexing and retrieval performance: The logical evidence. *Journal of the American Society for Information Science 45(8),* 589-599.

**Linnæus University**