

# Active Subtopic Detection in Multitopic Data

Benjamin Bergner   Georg Krempf

Knowledge Management & Discovery

Faculty of Computer Science

Otto-von-Guericke University

Magdeburg, Germany

## Outline

- ▶ Introduction
- ▶ Related Work: Clustering by Intent (CBI)
- ▶ Enhancements: Multitopic CBI
- ▶ Evaluation
  - ▶ Aims
  - ▶ Setup
  - ▶ Results
- ▶ Summary & Outlook

# App: Vocabulary Creator for Language Learning

- ▶ User reads foreign language texts
- ▶ Words of interested topic are highlighted
- ▶ How to build comprehensive topic vocabulary sets?



Kategorien

Konjugieren

Impressum

Technology tecnologia.com



## Google poderá comprar Twitter

O mercado das 'dotcom' ressurge com uma efervescência tremenda nas últimas semanas. A Microsoft saiu à rua e comprou a rede social profissional LinkedIn por 26 mil milhões de dólares, estando os assessores financeiros da Google a pressionar a empresa a comprar o Twitter já que esta rede social está prestes a ver a sua cotação fortemente valorizada. Este rumor foi apontado pelo site Market Watch, cabendo a operação financeira à cúpula de gestão da companhia, a Alphabet. O negócio deverá estar concluído até final do ano passando para o lado da Google um lote de 300 milhões de utilizadores. Para este negócio pairam ainda no ar outros nomes de grandes empresas do setor como o Facebook, ou até mesmo a Apple. Como consequência da entrada da Microsoft no LinkedIn, o valor do Twitter aumentou em cerca de 7%, mas ainda assim está 58% abaixo dos seus máximos históricos. Fonte: USA Today Deixe o seu comentário comentários

Verstecke Voca-Bubbles

telephone internet

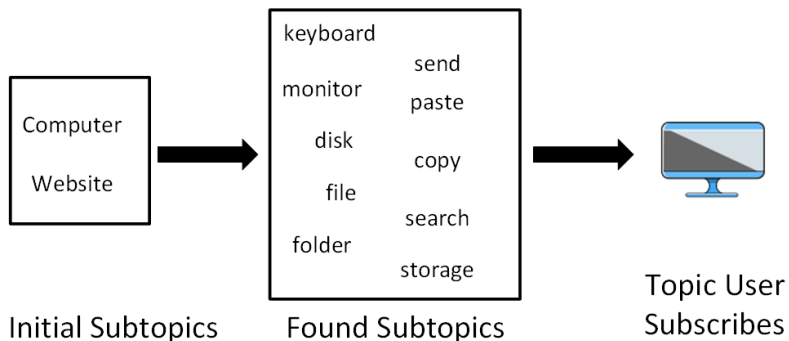


## Beispielsatz-Generator

Klicke jetzt auf ein Voca-Bub deiner Muttersprache zu generi

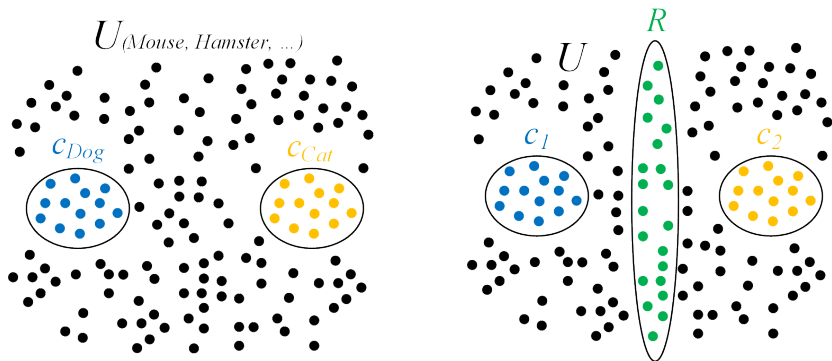
# App: Vocabulary Creator for Language Learning

- ▶ Input:
  - ▶ set of documents
  - ▶ an **intention**, e.g. two words from same topic → they act as subtopics
- ▶ Desirable Output: More subtopics of same topic (and documents they occur in)
- ▶ How: Probabilistic Active Incremental Clustering



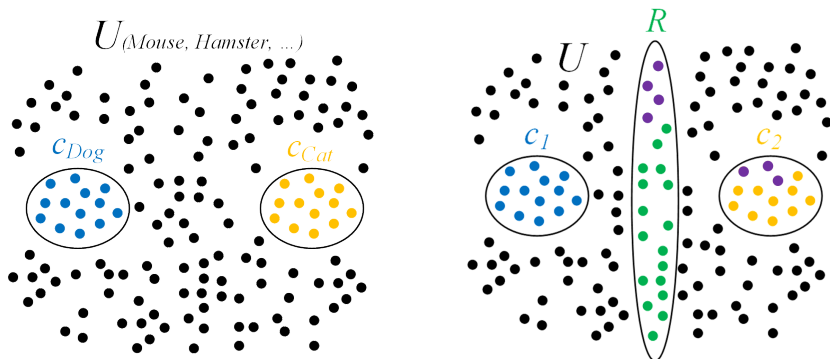
## Related Work: Clustering by Intent (Forman et. al 2015)

- ▶ Dots: represent documents
- ▶  $C_1, C_2$ : documents that contain predefined subtopics
- ▶ Residual Set  $R$ : documents most unsure how to cluster
- ▶ Get feedback for words in  $R$  that best discriminate  $R$  from  $L$
- ▶ Positive Feedback creates new subtopics and increases  $|L|$ , repeat
- ▶ Forman: Hewlett Packard analyzes support logs, customer surveys with CBI to find meaningful groups based on subtopics



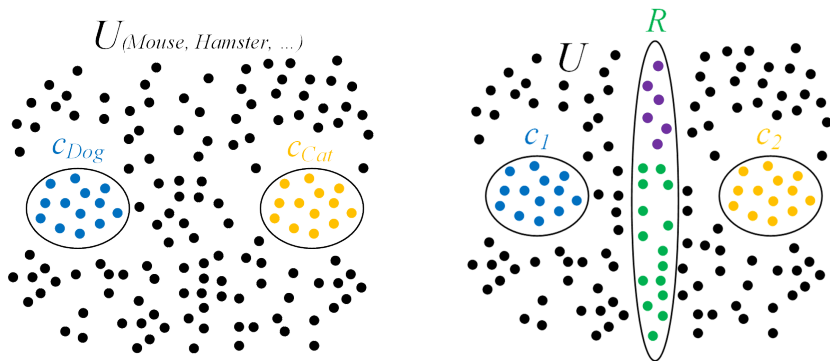
## Related Work: Clustering by Intent (Forman et. al 2015)

- ▶ Dots: represent documents
- ▶  $C_1, C_2$ : documents that contain predefined subtopics
- ▶ Residual Set  $R$ : documents most unsure how to cluster
- ▶ Get feedback for words in  $R$  that best discriminate  $R$  from  $L$
- ▶ Positive Feedback creates new subtopics and increases  $|L|$ , repeat
- ▶ Forman: Hewlett Packard analyzes support logs, customer surveys with CBI to find meaningful groups based on subtopics



## Related Work: Clustering by Intent (Forman et. al 2015)

- ▶ Dots: represent documents
- ▶  $C_1, C_2$ : documents that contain predefined subtopics
- ▶ Residual Set  $R$ : documents most unsure how to cluster
- ▶ Get feedback for words in  $R$  that best discriminate  $R$  from  $L$
- ▶ Positive Feedback creates new subtopics and increases  $|L|$ , repeat
- ▶ Forman: Hewlett Packard analyzes support logs, customer surveys with CBI to find meaningful groups based on subtopics



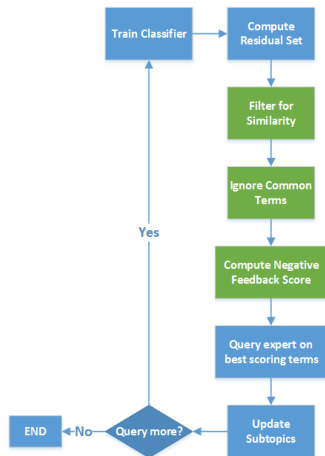
# Multitopic Clustering by Intent (MCBI)

## Problem and Objective

- ▶ In multitopic datasets, also unrelated documents will be considered as residual
- ▶ This work aims to extend CBI to make it usable for multitopic datasets: **Multitopic Clustering by Intent (MCBI)**

## Contributions

- ▶ Similarity Set
- ▶ Improve usage of negative feedback
- ▶ Ignore common words

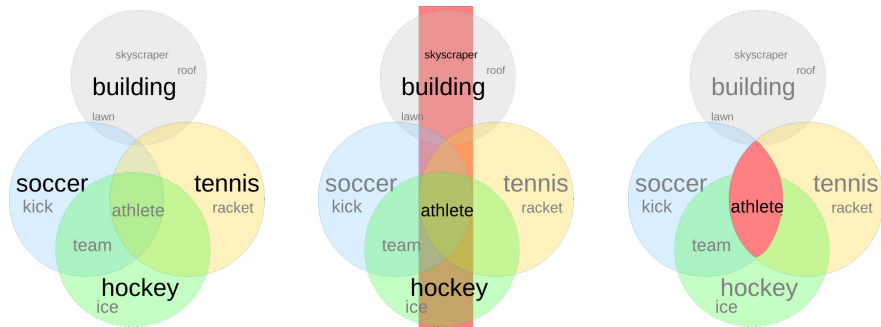


Remark: Concepts are explained per example, for formulae: see paper



# MCBI Contributions: From Residual to Similarity Set

- ▶ Subtopic's word bags with **intent: sports**
- ▶ initially given: soccer, tennis
- ▶ aim: find hockey, multi-topic environment: building documents
- ▶ residual set contains hockey AND building related documents
- ▶ word bags of soccer, tennis and hockey overlap when considering similarity



## Negative Feedback

- ▶ Do not allow negative feedback terms to occur further times
- ▶ Candidates that have many relative cooccurrences with negative feedback terms will be punished
- ▶ they are likely to be estimated negatively, too
- ▶ Ex.: Russia/hockey are competing given previously rejected candidate Finland, Russia and Finland occur more often together (also non-sports-relationships) than hockey and Finland (only sports-relationship)

## Ignore common words

- ▶ Candidate Terms occurring often in  $L$  are topic specific, not subtopic specific
- ▶ e.g. athlete or stadium occur often in  $L$
- ▶ We ignore them to
  - ▶ save annotation time
  - ▶ prevent them from being added to negative word list which would punish real subtopics

# Evaluation: Aim, Setup

## Objective

- ▶ Optimize interaction with expert
- ▶ i.e. find many actively chosen subtopics in few tries

## General Setting

- ▶ Compare against **random** and **CBI** baseline  
(for comparison between CBI and clustering see Forman et. al 2015)
- ▶ Consider average number of positive feedback over iterations

## Dataset

- ▶ Over **4000 wikipedia articles** drawn from most common nouns → heterogenous/multi-topic dataset
- ▶ Stop word removal, lemmatization, tf-idf → ca. **6500 unique words**
- ▶ Focusing on one closed category for testing: **countries**
- ▶ build a list of all countries, languages, denonyms for **auto-evaluation**

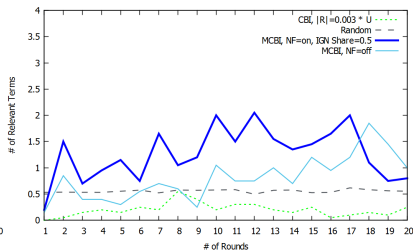
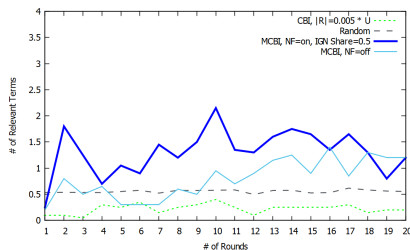
## Evaluation Results

Iteration Setting # subtopics, NF, IGN	Average subtopics found		
	Count.	Lang. & Denonym	Total
<b>Random</b>	<b>5.52</b>	<b>5.60</b>	<b>11.12</b>
<b>2, CBI, R = 0.005 ·   U  </b>	<b>0.9</b>	<b>3.9</b>	<b>4.8</b>
<b>4, CBI, R = 0.003 ·   U  </b>	<b>1.05</b>	<b>2.95</b>	<b>4.0</b>
2, <i>off</i> , –	7.75	8.20	15.95
2, <i>on</i> , 0	15.40	8.60	24.00
2, <i>on</i> , 0.1	15.55	8.9	24.45
2, <i>on</i> , 0.3	15.35	8.9	24.25
<b>2, on, 0.5</b>	<b>15.90</b>	<b>10.30</b>	<b>26.20</b>
2, <i>on</i> , 0.7	15.60	10.15	25.75
2, <i>on</i> , 1.0	10.15	7.05	17.20
4, <i>off</i> , –	9.00	7.10	16.10
4, <i>on</i> , 0	15.40	7.90	23.30
4, <i>on</i> , 0.1	15.90	8.00	23.90
4, <i>on</i> , 0.3	15.95	8.45	24.40
<b>4, on, 0.5</b>	<b>16.30</b>	<b>9.05</b>	<b>25.35</b>
4, <i>on</i> , 0.7	15.90	9.35	25.25
4, <i>on</i> , 1.0	9.00	6.90	15.90

Table: Detailed evaluation results for random, CBI and MCBI.

# Evaluation Results

- ▶ Average number of positive feedback over
- ▶ 20 iterations with differing starting values, 20 rounds per iteration, 20 queries per round → 400 queries per iteration
- ▶ *MCBI* with *NF = on* and *IGNShare = 0.5* performs almost always better than *MCBI* without those settings



## Summary

- ▶ MCBI extends CBI to multi-topic environments with residual sets within same topic → **similarity sets**
- ▶ Active feedback querying from user on candidates
- ▶ Incorporation of positive and **negative** user feedback
- ▶ **Ignoring** common words
- ▶ Evaluation on a wikipedia corpus for most common nouns
- ▶ Promising results compared to random and CBI

## Outlook

- ▶ More extensive experimental evaluation needed
- ▶ Parameter tuning
- ▶ Tests for applications such recommender systems, information retrieval

## Bibliography

- ▶ Forman, G.; Nachlieli, H. & Keshet, R. Bifet, A.; May, M.; Zadrozny, B.; Gavalda, R.; Pedreschi, D.; Bonchi, F.; Cardoso, J. & Spiliopoulou, M. (Eds.) Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2015, Porto, Portugal, September 7-11, 2015, Proceedings, Part III Clustering by Intent: A Semi-Supervised Method to Discover Relevant Clusters Incrementally Springer International Publishing, 2015, 20-36
- ▶ Jurafsky, Daniel, and James H. Martin. 2009. Speech and Language Processing: An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics. 2nd edition. Prentice-Hall.
- ▶ Settles, Burr. 2010. Active Learning Literature Survey. Computer Sciences Technical Report 1648. University of Wisconsin–Madison

Thank you

Thank you for your attention

Interested? Ask Questions



## Content-based Recommender Systems

- ▶ Key terms in already visited pages act as subtopics
- ▶ Find new subtopics
- ▶ Check terms from unseen items
- ▶ Recommend items that have maximum of found subtopics
- ▶ **Advantage:** Find surprising results

## Information Retrieval

- ▶ Recognize relations between documents
- ▶ Relate those that share many of found subtopics to same category

$$p_{s|u} = \Pr(s|\vec{u}) \propto \log \Pr(s) + \sum_{i=1}^{|\mathcal{V}_u|} \log \Pr(u_i|s) \quad (1)$$

$$\Pr(s) = \frac{|\mathcal{L}_s|}{|\mathcal{L}|} \quad (2)$$

$$\Pr(u_i|s) = \frac{\sum_{u_j \in \mathcal{L}_s} u_j + 1}{\sum_{v \in \mathcal{L}_s} v_j + |\mathcal{V}_L|} \quad (3)$$

$$\text{uncertainty}_u = p_{s|u} - p_{s'|u} \quad (4)$$

$$\text{similarity}_u = p_{s|u} + p_{s'|u} \quad (5)$$

$$\text{rejscore}_i \leftarrow \max_{n \in \mathcal{V}_{REJ}} \left( \frac{|\text{getDocumentsWithWords}(U \cup L, \{w, n\})|^2}{|\text{getDocumentsWithWords}(U \cup L, n)|} \right) \quad (6)$$

$$\text{discscore}_i \leftarrow |\text{getDocsWithWords}(U, w)| - |\text{getDocsWithWords}(L, w)| \quad (7)$$