

gesis

Leibniz Institute
for the Social Sciences



Bedeutung von Text Mining am Beispiel der Sozialwissenschaften

Philipp Mayr

*Text und Data Mining (TDM) - Trier
21.06.2018*

Text Mining

Semantic Scholar

All Fields

text mining

FAQ Contact

Data mining Sentiment analysis Text corpus Information retrieval Web search engine Social media Document classification

social **networks** text mining

social **media competitive analysis** text mining

visual text mining

text mining **tools**

market prediction text mining

examining students text mining

pizza industry text mining

study selection activity text mining

systematic reviews text mining

text mining and **bibliometrics**

text mining **techniques**

text mining **approach**

internet studies text mining

live video streaming text mining

online interaction text mining

Gooooooooooooole >

1 2 3 4 5 6 7 8 9 10

Weiter

Text Mining in den Sozialwissenschaften

- Politikwissenschaft: Analyse von Wahlkampfreden/-programmen; ...
- Wirtschaftswissenschaft: Einfluss von Social Media Sentiments auf Börsenentwicklungen; ...
- Bibliometrie/Scientometrie: Analyse von einflußreichen Themen innerhalb einer wissenschaftlichen Community; Zitationskontexte; Produktentstehung; Zitationen; ...

"Here empirical studies with real users should be performed, closely monitoring the users' actions (involving eyetracking) — see e.g. [Joachims et al. 07]."

"In contrast, the work presented in [Joachims et al. 07] uses eyetracking for observing users during Web searches, thus registering e.g. the items from the result list users were looking at."

■ ...

Statement

- Text Mining in den Sozialwissenschaften
 - ▶ Starker Anstieg an Text Mining-Aktivitäten
 - ▶ Typische autom. Analyseverfahren: Natural Language Processing, Information Retrieval, Named Entity Recognition, Machine Learning, ...
 - ▶ TM - Noch kein Standardwerkzeug der quantitativen Sozialforschung

Beispiel 2: Variablen-Erkennung in Volltexten

Mining Social Science Publications for Survey Variables

Andrea Zielinski and Peter Mutschke
 GESIS - Leibniz Institute for the Social Sciences
 Unter Sachsenhausen 6-8
 50667 Cologne, Germany
 [andrea.zielinski,peter.mutschke]@gesis.org

- Erkennung von Umfragevariablen in Volltexten
- Bisherige Analyseverfahren: manuell/intellektuell; Verlinkung von Publikation zu Variablendaten einer Studie
- Methode: Supervised Machine Learning mit Linguistische Analyse
- Ergebnisse: große Trainingsdatensätze für

Abstract

Research in Social Science is usually

vey question (e.g. *Do you believe in Heaven?*). Therefore, from the perspective of the Social Sciences, having a linkage just to the entire study

V39: Belief in life after death or in Heaven from the ISSP 1998 study

ten articles. In this paper, we present a work-in-progress study that seeks to provide a solution to the variable detection task based on supervised machine learning algorithms, using a linguistic analysis pipeline to extract a rich feature set, including terminological concepts and similarity metric scores. Further, we present preliminary results on a small dataset that has been specifically designed for this task, yielding a significant increase in performance over the random baseline.

data on the level of variables would have a number of benefits to researchers: It would enable indexing publications by survey variables and discovering publications that discuss the concept of interest (a particular variable). Moreover, it would facilitate a monitoring of the relevance of topical issues (by tracking the use of variables for research) as well as detecting research gaps (by tracking the variables not being addressed by researchers).

The problem, however, is that even though variables are usually assigned a code and a label (e.g. *V39: Belief in life after death* or *V40: Believe in Heaven* from the ISSP 1998 study) as well as the question text from the questionnaire, in practice, authors often do not adhere to citation standards, neither for study names nor for variables. Instead, authors tend to use variations of label and/or question text or combine several variables in one phrase (such as "...belief in afterlife and Heaven...") (Neporov and Nepor, 2009).

1 Introduction

In face of the growing number of scientific publications, Text Mining (TM) becomes increasingly important to make hidden knowledge explicit. A particular challenge in this regard is to identify research data citations in scholarly publications, due to their wide variety, ranging from quotations to

Diskussion

- Warum ist Text Mining noch kein Standardwerkzeug in der Sozialforschung?
 - ▶ Verfügbarkeit und Beschaffenheit von Corpora sind ein Problem!
 - ▶ Erstellung von Trainingsdatensätzen und Goldstandards z.T. sehr aufwendig!
 - ▶ Erweiterte Methodenkenntnisse notwendig!
- Weniger rechtliche Hürden als technisch-methodische Probleme!
 - ▶ Publikation von Twitter Tweets lediglich als Tweet-ID Listen

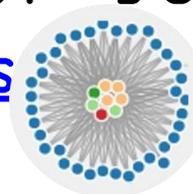
Exemplarische Projekte bei GESIS

- OpenMinTeD - Open Mining
Infrastructure for Text and Data
<http://openminded.eu/>

- InFoLiS - Integration von
Forschungsdaten und Literatur in
den Sozialwissenschaften
<http://infolis.github.io/>

- EXCITE - Extraktion von
Zitationen aus PDF-Dokumenten
<http://excite.wes.uni-koblenz.de/>

openMINTEd
Open Mining Infrastructure for Text & Data



References

- Rohrer, J. M., Brümmer, M., Schmukle, S. C., Goebel, J., & Wagner, G. G. (2017). 'What else are you worried about?' - Integrating textual responses into quantitative social science research. PLOS ONE, 12(7), e0182156.
<https://doi.org/10.1371/journal.pone.0182156>
- Zielinski, A., & Mutschke, P. (2017). Mining Social Science Publications for Survey Variables. In Proceedings of the Second Workshop on NLP and Computational Social Science. ACL 2017.
<http://aclweb.org/anthology/W17-29>

Kontakt

Dr. Philipp Mayr

Email: philipp.mayr@gesis.org

Web: <http://www.gesis.org>

@philipp_mayr

SSOAR

<http://ssoar.info>



<http://ecir2019.org/>

41st European Conference on Information Retrieval
in Cologne, Germany, on 14th April to 18th Apr