

# Automatisierte Erstellung von Pressedossiers durch Textmining

Kontextualisierung im journalistischen Umfeld

Marcel Schöneberg

[marcel.schoeneberg@haw-hamburg.de](mailto:marcel.schoeneberg@haw-hamburg.de)

Hochschule für Angewandte Wissenschaften Hamburg (HAW)  
Fakultät für Technik und Informatik  
Department Informatik

Masterseminar Präsentation

16.12.2014

- 1 Rückblick Grundprojekt
- 2 Masterarbeit Überblick
  - Einführung
  - Fragestellungen und Ziele
  - Erste Ideen
  - Related Work
  - Konkrete Schritte
  - Weitere mögliche Schritte
  - Bewertung von Ergebnissen
- 3 Aktueller Fortschritt: Hauptprojekt
- 4 Chancen und Risiken
- 5 Fragen
- 6 Literatur

- Ursprüngliche Idee: Trenderkennung in Social Media
- Ziele von Projekt 1:
  - Datenbasis schaffen
  - Vorverarbeitungen durchführen und verstehen
  - (Hierarchisches) Clustering als Experiment
    - Daten kennenlernen
    - These: 'Ausgewachsene' Trends/Themen müssten einen Cluster bilden
    - Untercluster zeigen ggf. kleinere Trends auf
    - Ggf. als Basis für spätere Zeitreihenanalysen nutzen

[Sch14]

- Vorhandene Datenbasis mit über 27000 Tweets zum Hashtag-Piraten
- Zeitraum: 11.2011 - 05.2014
- Erfahrungen mit Vorverarbeitung (z.B. Normalisierung, Stemming, NLP)
- Clustering mit verschiedenen Clustermengen und unterschiedlichen Untermengen des Korpus
- Diverse Stolpersteine kennengelernt

# Misserfolge - Begründung für den Fokuswechsel

- Weitgehend nutzlose Resultate
  - Cluster zumeist eher Sammlung aus unzusammenhängenden Begriffen (ggf. 'Zufallstreffer')
  - Auch Nutzung eines auf Nomen basierenden Korpus nicht erfolgversprechend
  - Texte grundlegend zu kurz
  - 'schlechte' Sprache (Abkürzungen usw.)
  - Clustering zu undurchsichtig
  - Trenderkennung ist ein schwammiges Ziel
- Durchaus möglich <sup>1</sup> aber u.A: komplexe **Vorverarbeitung** etc. notwendig
- Keine gute Ausgangsbasis

---

<sup>1</sup> Siehe Fazit von [Sch14]

- 1 Rückblick Grundprojekt
- 2 Masterarbeit Überblick
  - Einführung
  - Fragestellungen und Ziele
  - Erste Ideen
  - Related Work
  - Konkrete Schritte
  - Weitere mögliche Schritte
  - Bewertung von Ergebnissen
- 3 Aktueller Fortschritt: Hauptprojekt
- 4 Chancen und Risiken
- 5 Fragen
- 6 Literatur

- Ziel: Erstellung von Dossiers → Kontextualisierung eines Artikels
  - Fokus auf textbasierten Dossiers: Keine multimedialen Inhalte usw.
  - Dossiers: **Verschiedene** (Experten-) **Definitionen** vorhanden [Hä14]
    - Arbeitsthese: Zusammenstellung von 'ähnlichen' Artikeln
      - Inhaltliche / semantische Nähe
      - **Erster** Schritt: Nicht für alles geeignet → Abbildung 1
  - Ähnlichkeit als Basis → **Distanzfunktionen**
- (Zunächst) möglichst verständlich halten

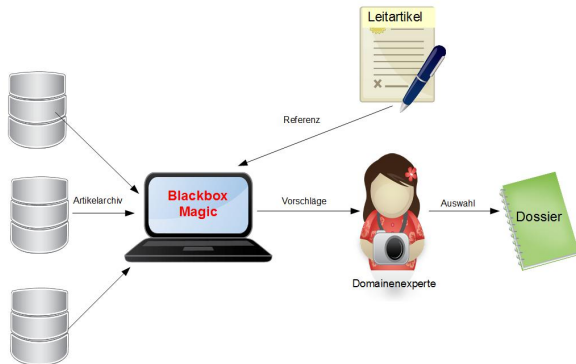


Abbildung: Realisierbarer Workflow



- Artikelarchiv: nach **Sortierung** ca. 2700 Artikel (Englisch)
  - von Journalisten
  - Basierend auf dem Eurozine Netzwerk [www.eurozine.com](http://www.eurozine.com)
  - Größere Länge als Tweets
  - Meta-Informationen vorhanden (Verlinkungen auf Inhaltsverzeichnisse, redaktionelle Focalpoints etc.) → Sinnvoll nutzbar?
  - semi-strukturiert in XML (Autor, Abstract, Überschriften usw.)

- *Beispiel*
- Artikelarchiv hat ebenfalls Probleme:
  - Nicht alle englischsprachig
  - Enthält auch Zusammenfassungen etc.
  - z.T. Gedichte (recht kurz, sprachlich besonders)
  - Invalides XML

- Automatisierbare Dossiererstellung möglich?
  - Was sind Pressedossiers bzw. wie werden diese benutzt?
  - Lassen sich fachliche Anforderungen in einen Algorithmus überführen?
  - Sind journalistische Definitionen eines Dossiers mit Informatik umsetzbar?
  - Wie kann man ein generiertes Ergebnisses bewerten?
- Versuch den Algorithmus übersichtlich zu halten

- Viele Diskussionen und Ideen
- Simple Distanzfunktionen nutzen
- Facetten / abweichende Sichtweisen finden
- Vorhandene Clustering Techniken benutzen
- Vorschlagssysteme aufgrund von Lernalgorithmen / Crowd
- Metainformationen (Autoren, Herkunft) nutzen → besserer Überblick
- Ggf. Distanzfunktionen verbessern

→ Zu einfach gedacht

- Geht nicht auf **fachliche Anforderungen** ein
- Clusteringverfahren oft sehr **undurchsichtig** (vgl. Projekt 1)
- Vorschlagssysteme brauchen (oft) eine Basis - Nicht vorhanden → **Coldstart problem**
- Viele Systeme setzen auf Personalisierung oder die Crowd → Ziel: **automatisierte und neutrale** Ergebnisse (z.B. [PL10], [KFD12], [LCLS10])
- Ideen → 'Den zweiten Schritt vor dem ersten tun'
- Erlernen von Ähnlichkeit nicht im Fokus (Neuronale Netze usw.)
  - Nicht ohne Weiteres anpassbar
  - Nicht als erster Schritt

- Forschungsprojekt der DPA
- Overview Project <http://overview.ap.org/>
- XML Ähnlichkeiten [TCY09]
- Distanzfunktionen [HRJM13]
- Content-based Recommender Systems [LGS]
- Überblick über Recommender Systems in der Breite [BOHG13]
- Article Recommender Systems [BDD<sup>+</sup>12]

- Ziel ist höchst experimentell → Schrittweise **annähern und verbessern**
- Verbesserung der Distanzfunktion zwischen Dokumenten
- Bag of Words-Ansatz
- Ausgangsbasis: **der** Leitartikel – später ggf. gewichtete Wortmenge

- **Gewichtete** Wertungen von 'wichtigeren' Teilen → Gute Werte noch zu ermitteln
  - ① Abstract
  - ② Title
  - ③ Subheadings
  - ④ Text
- Anreicherung mit themenspezifischen **Kategorien**
  - Wortmengen pro Thema (z.B. Wirtschaft, Politik usw.) einbringen
  - Gewichtung des 'Bag of words' zugunsten der (hierarchischen) Kategorien
- **Erweiterung** durch Erkenntnisse einer Fachexpertin



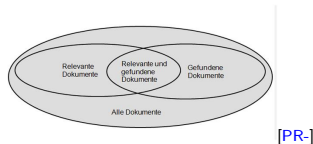
- Vieles hat wurde im Rahmen des Seminars schon angesprochen [Sip14]
  - Zunächst Nutzung einer einfachen euklidischen Distanz (auf Basis der Gewichtungen)
  - Verschiedene 'Relevanzmaße' für Termhäufigkeiten (Termvorkommen, TF-IDF)
  - Diverse Kriterien von Distanz möglich (Relevanz, Themenbreite, Neuigkeit etc.)
- Diverse Funktionen und Möglichkeiten vorhanden

- Autoren nutzen um Herkunft des Artikels zu bestimmen
- Ggf. Einbeziehung von Inhaltsverzeichnissen der (Themen-)Zeitschriften
- Linguistische Verbesserungen (z.B. Nutzung von Ontologien etc.)
- Erkennung von Facetten eines Themas (europäische Sicht auf Ursprungsartikel)

- Verifikation der Ergebnisse über Focalpoints-Zusammenfassungen
  - Testkorpus als Verifikationsbasis:
    - Focalpoint-Artikel (redaktionell eingeordnet)
    - (fast) zufällige Artikel (nicht im Focalpoint)
- Berechnete Dossierartikel sollten möglichst den Focalpointartikeln entsprechen

→ Nutzung von Eigenschaften wie:

- Recall (Trefferquote):  $\frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}$
- Precision (Genauigkeit):  $\frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$

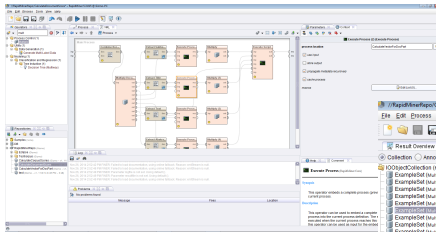


- Wissen einer Domänenexpertin

- 1 Rückblick Grundprojekt
- 2 Masterarbeit Überblick
  - Einführung
  - Fragestellungen und Ziele
  - Erste Ideen
  - Related Work
  - Konkrete Schritte
  - Weitere mögliche Schritte
  - Bewertung von Ergebnissen
- 3 Aktueller Fortschritt: Hauptprojekt
- 4 Chancen und Risiken
- 5 Fragen
- 6 Literatur

- Aufgesetzter Testkorpus (Focalpoint Artikel + beliebige Artikel als Vergleichsbasis)
- Rapidminer <sup>2</sup> Verständnis und Hürden
- KDD Preprocessing
- Erste Versuche lauffähig
- Domänenübergreifendes Projekt → Grundlegendes gemeinsames Verständnis

# Aktueller Fortschritt: Hauptprojekt II



(a) Process

The screenshot shows the 'Edit Parameter Test script' dialog in RapidMiner. It contains a Java script that processes data from a list of entries. The script uses the 'java.util.List' and 'java.util.HashMap' classes. It iterates through the entries and calculates the total frequency for each key. The script is as follows:

```
13 java.util.List<Object> subheadings = subheadingsList.getEntries();
14 java.util.List<Object> headings = headingsList.getEntries();
15 java.util.List<Object> bodies = bodiesList.getEntries();
16 java.util.List<Object> abstracts = abstractsList.getEntries();
17
18 java.util.HashMap wordmap = new HashMap();
19
20 for(Object entry : abstracts) {
21     String key = entry.getKey();
22     float value = entry.getTotalFrequency();
23     if(wordmap.containsKey(key)) {
24         wordmap.put(key, wordmap.get(key) + value * ABSTRACT_MULTIPLICATOR);
25     }
26     else {
27         wordmap.put(key, value);
28     }
29 }
30
31 for(Object entry : headings) {
32     String key = entry.getKey();
33     float value = entry.getTotalFrequency();
34     if(wordmap.containsKey(key)) {
35         wordmap.put(key, wordmap.get(key) + value * TITLE_MULTIPLICATOR);
36     }
37 }
```

(b) Script

The screenshot shows the 'Result Overview' window in RapidMiner. It displays a table with 29 rows and 4 columns: 'Row No.', 'w1', 'w2', and 'w3'. The table contains data for various categories, including 'vitalis', 'histor', 'happ', 'misp', 'subco', 'foundat', 'spoke', 'time', 'becam', 'tendenc', 'event', 'rite', 'victor', 'thesi', 'minded', 'jet', 'upm', 'hazogama', 'categor', 'shown', 'project', 'pocul', 'europen', 'parodon', 'maipolka', 'downsaw', 'distanc', 'mult', and 'conferenc'. The table is sorted by 'w1' in ascending order.

Row No.	w1	w2	w3
1	vitalis	1	2000-00-00
2	histor	2	2000-00-00
3	happ	2	2000-00-00
4	misp	1	2000-00-00
5	subco	1	2000-00-00
6	foundat	1	2000-00-00
7	spoke	1	2000-00-00
8	time	4	2000-00-00
9	becam	3	2000-00-00
10	tendenc	1	2000-00-00
11	event	3	2000-00-00
12	rite	2	2000-00-00
13	victor	1	2000-00-00
14	thesi	2	2000-00-00
15	minded	1	2000-00-00
16	jet	1	2000-00-00
17	upm	1	2000-00-00
18	hazogama	1	2000-00-00
19	categor	1	2000-00-00
20	shown	1	2000-00-00
21	project	1	2000-00-00
22	pocul	1	2000-00-00
23	europen	10	2000-00-00
24	parodon	2	2000-00-00
25	maipolka	1	2000-00-00
26	downsaw	1	2000-00-00
27	distanc	1	2000-00-00
28	mult	1	2000-00-00
29	conferenc	1	2000-00-00

(c) Resultwindow

- 1 Rückblick Grundprojekt
- 2 Masterarbeit Überblick
  - Einführung
  - Fragestellungen und Ziele
  - Erste Ideen
  - Related Work
  - Konkrete Schritte
  - Weitere mögliche Schritte
  - Bewertung von Ergebnissen
- 3 Aktueller Fortschritt: Hauptprojekt
- 4 Chancen und Risiken
- 5 Fragen
- 6 Literatur

## • Chancen

- Viel Interesse - 'Editors don't scale'
  - HAW NextMedia <sup>3</sup> gefördert durch BMBF
- DPA Forschungsprojekt

## • Risiken

- Zu abweichende Meinungen von 'Was ist ein Dossier'
- Gefahr das Experten ihr Wissen nicht weitergeben können - 'Habe ich im Gefühl' [McD83]
- Definitionen zu komplex oder zu unspezifisch um sie automatisiert umzusetzen
- Aufeinander aufbauender Prozess - wenn eins schief geht → Problem

---

<sup>3</sup> <http://nextmedia-haw.de/>



Vielen Dank für die Aufmerksamkeit!

Fragen?



BANCU, Cristian ; DAGADITA, Monica ; DASCALU, Mihai ; DOBRE, Ciprian ; TRAUSAN-MATU, Stefan ; FLOREA, Adina M.:

ARSYS – Article Recommender System.

In: **Proceedings of the 2012 14th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing.**

Washington, DC, USA : IEEE Computer Society, 2012 (SYNASC '12). – ISBN 978-0-7695-4934-7, 349–355



BOBADILLA, J. ; ORTEGA, F. ; HERNANDO, A. ; GUTIÉRREZ, A.:

Recommender Systems Survey.

In: **Know.-Based Syst.** 46 (2013), Juli, 109–132.

<http://dx.doi.org/10.1016/j.knosys.2013.03.012>. –

DOI 10.1016/j.knosys.2013.03.012. –

ISSN 0950-7051



HÄLKER, Nina:

**Dienen textminingbasierte Dossiers dem Wachsen einer gemeinsamen europäischen Erzählung?**

2014. –

Arbeitspapier



HARISPE, Sébastien ; RANWEZ, Sylvie ; JANAQI, Stefan ; MONTMAIN, Jacky:

Semantic Measures for the Comparison of Units of Language, Concepts or Entities from Text and Knowledge Base Analysis.

In: **CoRR** abs/1310.1285 (2013).

<http://arxiv.org/abs/1310.1285>



KIRSHENBAUM, Evan ; FORMAN, George ; DUGAN, Michael:

A Live Comparison of Methods for Personalized Article Recommendation at Forbes.Com.

In: **Proceedings of the 2012 European Conference on Machine Learning and Knowledge Discovery in Databases - Volume Part II.**

Berlin, Heidelberg : Springer-Verlag, 2012 (ECML PKDD'12). – ISBN 978-3-642-33485-6, 51-66



LI, Lihong ; CHU, Wei ; LANGFORD, John ; SCHAPIRE, Robert E.:

A Contextual-bandit Approach to Personalized News Article Recommendation.

In: **Proceedings of the 19th International Conference on World Wide Web.**

New York, NY, USA : ACM, 2010 (WWW '10). – ISBN 978-1-60558-799-8, 661-670



LOPS, Pasquale ; GEMMIS, Marco de ; SEMERARO, Giovanni:

In: **Recommender Systems Handbook**



MCDERMOTT, John P.:

Extracting Knowledge From Expert Systems.

In: BUNDY, Alan (Hrsg.): **IJCAI**, William Kaufmann, 1983, 100-107



PENG, Chi-Chieh ; LIU, Duen-Ren:

Combining Reputation and Content-based Filtering for Blog Article Recommendation in Social Bookmarking Websites.

In: **Proceedings of the 12th International Conference on Electronic Commerce: Roadmap for the Future of Electronic Business.**

New York, NY, USA : ACM, 2010 (ICEC '10). – ISBN 978-1-4503-1427-5, 8-14



**Precision-Recall diagram.**

[http://wikis.gm.fh-koeln.de/wiki\\_ir/uploads/InformationRetrieval/Recall/ir\\_bild.jpg](http://wikis.gm.fh-koeln.de/wiki_ir/uploads/InformationRetrieval/Recall/ir_bild.jpg)



SCHÖNEBERG, Marcel:

Erkennung von Trends in sozialen Netzwerken.

Version: 10 2014.

<http://users.informatik.haw-hamburg.de/~ubicomp/projekte/master2014-proj/schoeneberg.pdf>.  
2014. –

Projektbericht



SIPPEL, Sigurd:

Recommendations for cocktail recipes.  
(2014), 11.

<http://users.informatik.haw-hamburg.de/~ubicomp/projekte/master2014-sem/sippel/folien.pdf>



TEKLI, Joe ; CHBEIR, Richard ; YETONGNON, Kokou:

Survey: An Overview on XML Similarity: Background, Current Trends and Future Directions.

In: *Comput. Sci. Rev.* 3 (2009), August, Nr. 3, 151–173.

<http://dx.doi.org/10.1016/j.cosrev.2009.03.001>. –

DOI 10.1016/j.cosrev.2009.03.001. –

ISSN 1574–0137