
Kooperative Promotion im Rahmen der TechnologieAllianzOberfranken (TAO) im Bereich *Opinion Mining*

Analyse von literarischen Texten mit Methoden aus dem Bereich des Opinion Mining

Agenda

- Einleitung und Problemstellung
- Stand der Forschung und Forschungslücke
- These der Dissertation
- Kernthema: Erzeugung einer lexikalischen Ressource
- Evaluation der lexikalischen Ressource
- Diskussion

Was ist Opinion Mining (Sentiment Analysis)?

- automatische Extraktion von Meinungen aus Texten
- verschiedene Ebenen: Dokument, Satz, **Aspekt**
- verschiedene Methoden: **Lexikon-basierter Ansatz**, Maschinelles Lernen
- verschiedene Arten von Meinungen:
regular opinions – „Das iPhone 6 ist super!“ und
comparative opinions – „Das iPhone 6 ist besser als das iPhone 5.“
- weiteres Themenfeld: *Opinion Spam Detection*
- Anwendungen: Social Media Monitoring, Analyse von literarischen Texten hinsichtlich Stimmungsverlauf (?) etc.

Was ist eine Meinung (Opinion)?

- Opinion-Quintupel (e,a,s,h,t) [1,2]
 - Entität **e**
 - Aspekt **a** der Entität **e**
 - Meinung **s** zum Aspekt **a** der Entität **e**
 - Autor **h** der Meinungsäußerung
 - Zeitpunkt **t** der Meinungsäußerung
- Beispiel

Max Mustermann (**h**) schrieb am 20.10.15 (**t**):

Das Display (**a**) des Acer Notebooks (**e**) finde ich einfach super (**s**)!

Meinungslexika (*Sentiment Lexicon*)

- Lexikon-basierter Ansatz benötigt lexikalische Ressourcen
- Liste mit meinungstragenden Wörtern (und ggf. Mehrwortphrasen) sowie „Wertung“
- verschiedene Methoden zur Erzeugung (siehe „Kernthema“)

Herausforderungen

- Güte und Vollständigkeit sowie Messung dieser Kriterien
- Umgang mit *valence shifter words* (Verstärker, Abschwächer, Negation)
- Abhängigkeit von der Sprache (Probleme bei der Erzeugung)

Agenda

- Einleitung und Problemstellung
- Stand der Forschung und Forschungslücke
- These der Dissertation
- Kernthema: Erzeugung einer lexikalischen Ressource
- Evaluation der lexikalischen Ressource
- Diskussion

Meinungslexika als lexikalische Ressource

- existieren in verschiedenen Sprachen (vor allem englische Sprache)
- in verschiedenen Detailstufen
- teilweise für unterschiedliche Domänen
- wurden erzeugt mit verschiedenen Methoden

Meinungslexika - Auswahl

- Englisch
 - SentiWordNet 3.0 [3]: 120.000 Wörter
 - Semantic Orientation of Words [4]: 90.000 Wörter
 - Subjectivity Lexicon [5]: 8.000 Wörter
- Deutsch
 - Polarity Lexicon [6]: 8.000 Wörter
 - GermanPolarityClues [7]: 10.000 Wörter
 - Sentiment Phrase List (eigene Liste) [8]: ca. 14.000 Wörter und Mehrwortphrasen
- Spanisch
 - [9]: 4.660 Wörter

Wieso eine neue lexikalische Ressource? (Forschungslücke)

- Lücken für viele Sprachen (auch für die deutsche Sprache)
- Mehrwortphrasen statt einzelner Wörter (Problem mit *valence shifter words* vermeiden)
- Sprachunabhängigkeit bei Erzeugung durch Verzicht von NLP Methoden wie POS Tagging, Lemmatisierung etc.
- Aufnahme von Redewendungen („Es ist nicht alles Gold was glänzt“)

Agenda

- Einleitung und Problemstellung
- Stand der Forschung und Forschungslücke
- These der Dissertation
- Kernthema: Erzeugung einer lexikalischen Ressource
- Evaluation der lexikalischen Ressource
- Diskussion

These:

Durch den Einsatz statistischer Verfahren, unter Berücksichtigung von Mehrwortphrasen, können Meinungslexika automatisch aus geeigneten Korpora – Rezensionen mit Titel und Bewertung – erzeugt werden.

Weitere Thesen (Diskussion)

- Sprachunabhängigkeit der Methode
- Korpora: Skalierung und Unabhängigkeit

Agenda

- Einleitung und Problemstellung
- Stand der Forschung und Forschungslücke
- These der Dissertation
- Kernthema: Erzeugung einer lexikalischen Ressource
- Evaluation der lexikalischen Ressource
- Diskussion

Methoden zur Erzeugung

- manuell
- Wörterbuch-basiert
- **Korpus-basiert**

Idee [10]

- Verwertung der Korrelation zwischen Titel und Bewertung
- Bestimmung von relevanten Wörtern und Mehrwortphrasen durch Berechnung signifikanter Kookkurrenzen sowie Häufigkeit
- Ableitung der *Sentiment Values (SV)* aus der mittleren Bewertung (Skala kontinuierlich [-1 , +1])

Beispiel

- Wörter „sehr“ und „schön“ kommen signifikant oft zusammen in positiven Amazon Rezensionen (Titel) vor
- Annahme: 50 mal in 4-Sterne Bewertungen und 50 mal in 5-Sterne Bewertungen
→ Mehrwortphrase „sehr schön“ hat damit eine durchschnittliche Bewertung von 4,5 Sternen
- Umrechnung: $SV_{\text{sehr schön}} = 0,75$

23 von 26 Kunden fanden die folgende Rezension hilfreich

★★★★★ **Sehr schön**

Von [janine](#) am 24. Juli 2015

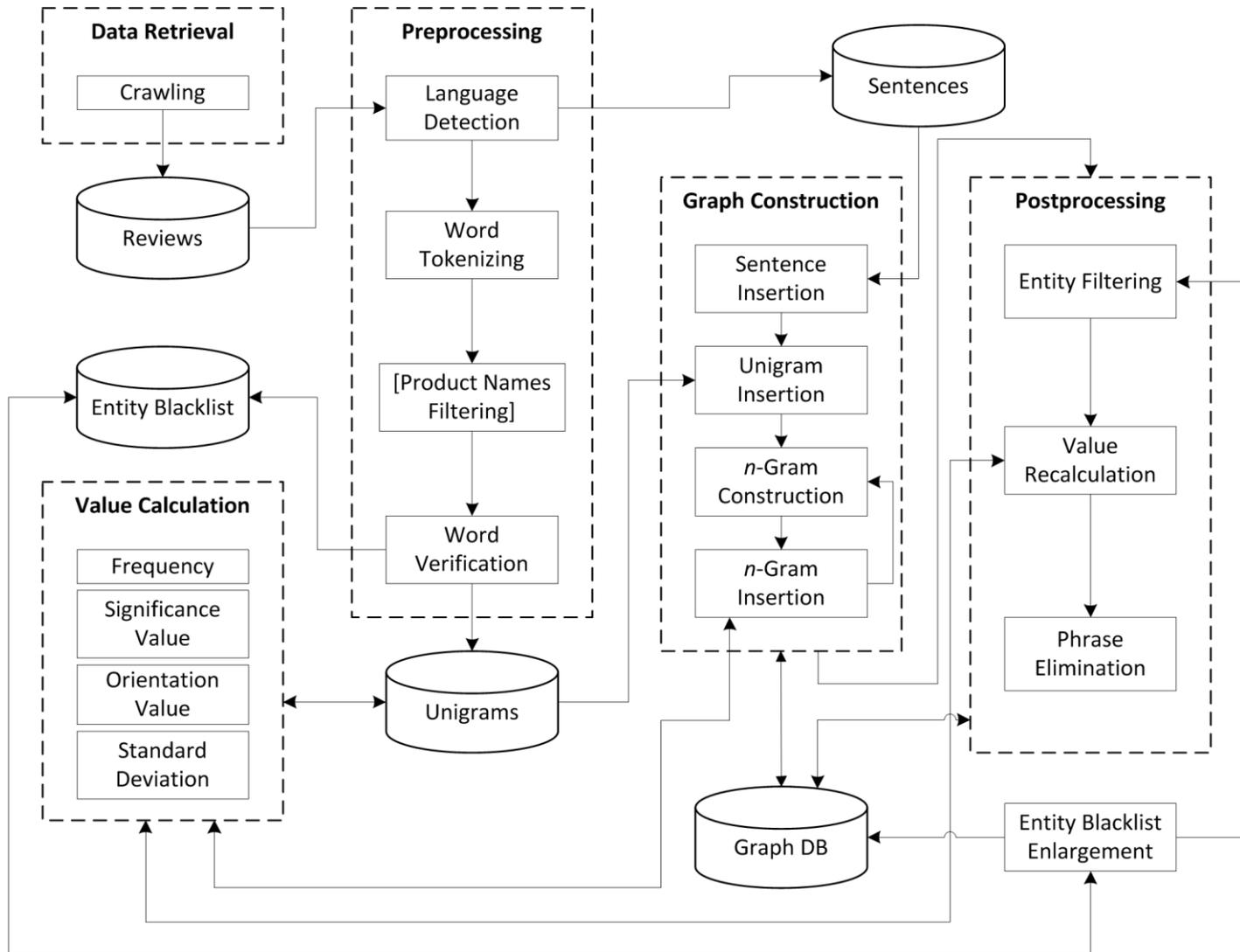
Farbe: Gold | Größe: 16 GB | **Verifizierter Kauf**

Ich habe das Iphone 6 - 16 GB in Gold bestellt. Es wurde schnell geliefert. Ich bin mit dem Gerät zufrieden, es entspricht dem was ich davor darüber gelesen habe. Bevor man ein Handy sich zulegt sollte man lesen was es bietet und kann, dann wird man auch nicht enttäuscht. Schönes Design. Ich empfehle trotz allem eine Folie für den Display Schutz und eine Schutzhülle, damit ihr auch lange von eurem Iphone etwas habt. Und wem der Akku zu schnell leer geht, es gibt Akku Bänke die bis zu 4 mal komplett aufladen bei Amazon, so eine habe ich mir lange lange vorher zugelegt. Super.

Einschub: Signifikante Kookkurrenzen

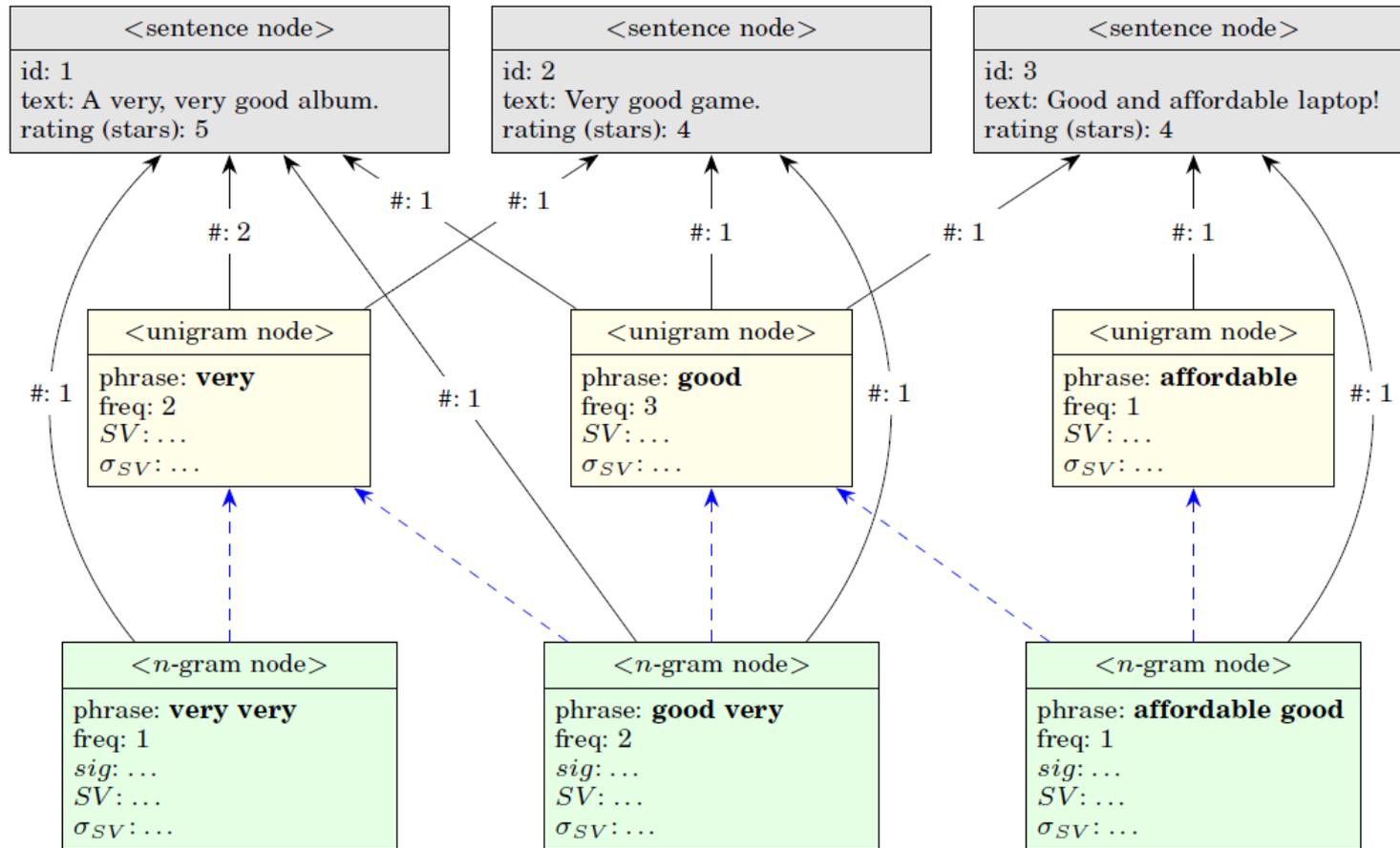
- signifikant häufiges gemeinsames Auftreten von Wortformen
- Nachbarschaftskookkurrenz und Satzkookkurrenz
- Beispiele
 - New – York
 - Harry – Potter
 - Polizei – verhaftet
 - Unfall – Krankenhaus
- Signifikanzmaße (Auswahl)
 - Log-Likelihood
 - **Signifikanzmaß von Quasthoff und Wolff [11]**

Ablauf – Erzeugung Meinungslexikon



Graphdatenbank (als Hilfsmittel)

- Aufbau der lexikalischen Ressource zeitintensiv durch wiederkehrende Berechnungen → Einsatz einer Graphdatenbank
- 3 Typen von Knoten
 - *sentence nodes*
 - *phrase nodes for unigrams*
 - *phrase nodes for n-grams*
- 2 Arten von Kanten
 - *occurrence edge*
 - *sub-phrase edge*



———→ relationship between sentence and phrase nodes (occurrence edges)
 - - - - -→ relationship between phrase nodes (sub-phrase edges)

Erste Ergebnisse

Wort / Mehrwortphrase	Sentiment Value (SV)
absolut fantastisch	1,00
sehr gut	0,98
Meisterwerk	0,94
gut	0,69
überdurchschnittlich	0,66
nicht zufriedenstellend	-0,54
einfach nur schlecht	-1,00

Erste Ergebnisse (Redewendungen)

Wort / Mehrwortphrase	Sentiment Value (SV)
klein aber oho	0,90
eierlegende Wollmilchsau	0,86
aller guten Dinge sind drei	0,72
weder Fisch noch Fleisch	-0,20
außen hui innen pfui	-0,54
Schuster bleib bei deinen Leisten	-0,77
Finger weg	-1,00

Agenda

- Einleitung und Problemstellung
- Stand der Forschung und Forschungslücke
- These der Dissertation
- Kernthema: Erzeugung einer lexikalischen Ressource
- Evaluation der lexikalischen Ressource
- Diskussion

Evaluation Meinungslexikon (Idee)

- Evaluation durch eigenes Annotationsset (**Diskussion**)
 - detaillierte Annotationsanleitung → Was soll wie annotiert werden?
 - Auswahl der Texte → Domäne, Sprache, Umfang
 - Veröffentlichung des annotierten Korpus → Probleme mit Urheberrecht etc.?
 - Interrater-Reliabilität z.B. durch Cohens Kappa prüfen
- direkte Evaluation, z.B. durch *Amazon Mechanical Turk*
 - „Turk Workers“ erhalten Liste mit meinungstragenden Wörtern und müssen diese in Reihenfolge bringen → „Ist Ausdruck A positiver als Ausdruck B?“
- Vergleich mit anderen lexikalischen Ressourcen → schwierig, da teilweise sehr unterschiedlich

Agenda

- Einleitung und Problemstellung
- Stand der Forschung und Forschungslücke
- These der Dissertation
- Kernthema: Erzeugung einer lexikalischen Ressource
- Evaluation der lexikalischen Ressource
- Diskussion

Diskussion

- Anwendungsszenarien?

Vielen Dank für Ihre Aufmerksamkeit
(Fragen und Diskussion)

Literaturverzeichnis

- [1] Bing Liu. 2012. Sentiment Analysis and Opinion Mining. Synthesis Lectures on Human Language Technologies
- [2] Bing Liu. 2015. Sentiment Analysis: Mining Opinions, Sentiments, and Emotions. Cambridge University Press
- [3] Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010. SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. In Proceedings of the 7th International Conference on Language Resources and Evaluation
- [4] Hiroya Takamura, Takashi Inui, and Manabu Okumura. 2005. Extracting Semantic Orientations of Words using Spin Model. In Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics
- [5] Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. 2005. Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis, In Proceedings of the Human Language Technology Conference
- [6] Simon Clematide and Manfred Klenner. 2010. Evaluation and Extension of a Polarity Lexicon for German. In Proceedings of the 1st Workshop on Computational Approaches to Subjectivity and Sentiment Analysis
- [7] Ulli Waltinger. 2010. GermanPolarityClues: A Lexical Resource for German Sentiment Analysis. In Proceedings of the 7th International Conference on Language Resources and Evaluation
- [8] Sven Rill, Sven Adolph, Johannes Drescher, Dirk Reinel, Jörg Scheidt, Oliver Schütz, Florian Wogenstein, Roberto V. Zicari, and Nikolaos Korfiatis. 2012. A phrase-based opinion list for the German language. In Proceedings of KONVENS 2012
- [9] Julian Brooke, Milan Tofiloski, and Maite Taboada. 2009. Cross-Linguistic Sentiment Analysis: From English to Spanish. In Proceedings of the International Conference on Recent Advances in Natural Language Processing
- [10] Sven Rill, Jörg Scheidt, Johannes Drescher, Oliver Schütz, Dirk Reinel, and Florian Wogenstein. 2012. A Generic Approach to Generate Opinion Lists of Phrases for Opinion Mining Applications. In Proceedings of the First International Workshop on Issues of Sentiment Discovery and Opinion Mining
- [11] Uwe Quasthoff and Christian Wolff. 2002. The poisson collocation measure and its applications