

Text Mining als Methode der qualitativen Sozialforschung

VO Qualitative Methoden der empirischen
Sozialforschung, 3.6.2015

Nikolaus Pöchhacker, Institut für Höhere Studien

Was ist Text Mining?

- ▶ Analyse von unstrukturiertem Text
- ▶ Basiert auf Algorithmen
- ▶ Suche nach (Sinn-)Strukturen im Text
- ▶ Oft ein „uninformiertes“ Verfahren



Warum machen wir das?

- ▶ Empirische Sozialforschung basiert sehr stark auf Textanalyse
- ▶ Oft große Datenmengen
- ▶ Generiert neue Perspektiven auf die Texte
- ▶ Geschwindigkeit!



Strukturierte vs. Unstrukturierte Daten

▶ Strukturierte Daten

- ▶ Haben eine vordefinierte Form
- ▶ Lassen sich formal beschreiben
- ▶ Beispiele: Excel Spreadsheets, Graphen, etc.

▶ Unstrukturierte Daten

- ▶ Zeigen keine offensichtlichen Muster
- ▶ Formale Beschreibung schwer
- ▶ Beispiele: Diese Folien, ihre Bachelorarbeit, etc.



Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern

Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern
- ▶ Document Clustering: Gruppieren von Dokumenten (oder Teile davon) anhand von Ähnlichkeiten

Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern
- ▶ Document Clustering: Gruppieren von Dokumenten (oder Teile davon) anhand von Ähnlichkeiten
- ▶ Document Classification: Tagging von Dokumenten aufgrund von vorher definierten Modellen



Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern
- ▶ Document Clustering: Gruppieren von Dokumenten (oder Teile davon) anhand von Ähnlichkeiten
- ▶ Document Classification: Tagging von Dokumenten aufgrund von vorher definierten Modellen
- ▶ Web Mining: Text Mining von Internettextran (Twitter, Blogs, Websites, etc.). Aufgrund von Links können Netzwerke erstellt werden.



Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern
- ▶ Document Clustering: Gruppieren von Dokumenten (oder Teile davon) anhand von Ähnlichkeiten
- ▶ Document Classification: Tagging von Dokumenten aufgrund von vorher definierten Modellen
- ▶ Web Mining: Text Mining von Internettextran (Twitter, Blogs, Websites, etc.). Aufgrund von Links können Netzwerke erstellt werden.
- ▶ Information Extraction: Verwandeln von unstrukturiertem Text in strukturierte Daten. Muster und Beziehungen werden hier aufgezeigt.

Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern
- ▶ Document Clustering: Gruppieren von Dokumenten (oder Teile davon) anhand von Ähnlichkeiten
- ▶ Document Classification: Tagging von Dokumenten aufgrund von vorher definierten Modellen
- ▶ Web Mining: Text Mining von Internettextran (Twitter, Blogs, Websites, etc.). Aufgrund von Links können Netzwerke erstellt werden.
- ▶ Information Extraction: Verwandeln von unstrukturiertem Text in strukturierte Daten. Muster und Beziehungen werden hier aufgezeigt.
- ▶ Natural Language Processing (NLP): „Verstehen“ von Texten. Erkennen von Phrasen und einfachen Aussagen.

Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern
- ▶ Document Clustering: Gruppieren von Dokumenten (oder Teile davon) anhand von Ähnlichkeiten
- ▶ Document Classification: Tagging von Dokumenten aufgrund von vorher definierten Modellen
- ▶ Web Mining: Text Mining von Internettexten (Twitter, Blogs, Websites, etc.). Aufgrund von Links können Netzwerke erstellt werden.
- ▶ Information Extraction: Verwandeln von unstrukturiertem Text in strukturierte Daten. Muster und Beziehungen werden hier aufgezeigt.
- ▶ Natural Language Processing (NLP): „Verstehen“ von Texten. Erkennen von Phrasen und einfachen Aussagen.
- ▶ Concept Extraction: Gruppieren von Wörter und Phrasen in semantisch ähnliche Gruppen.

Arten des Text Minings

- ▶ Search and Information Retrieval (IR): Abspeichern und wiederauffinden von Dokumenten anhand von Schlüsselwörtern
- ▶ Document Clustering: Gruppieren von Dokumenten (oder Teile davon) anhand von Ähnlichkeiten
- ▶ Document Classification: Tagging von Dokumenten aufgrund von vorher definierten Modellen
- ▶ Web Mining: Text Mining von Internettexten (Twitter, Blogs, Websites, etc.). Aufgrund von Links können Netzwerke erstellt werden.
- ▶ **Information Extraction: Verwandeln von unstrukturiertem Text in strukturierte Daten. Muster und Beziehungen werden hier aufgezeigt.**
- ▶ Natural Language Processing (NLP): „Verstehen“ von Texten. Erkennen von Phrasen und einfachen Aussagen.
- ▶ **Concept Extraction: Gruppieren von Wörter und Phrasen in semantisch ähnliche Gruppen.**

Beispiel: Das Projekt

Evaluation des Förderinstruments „Diskursprojekte zu ethischen, rechtlichen und sozialen Fragen in den modernen Lebenswissenschaften“ des BMBF

Mein Arbeitspaket: Analyse der Webseiten



INSTITUT FÜR HÖHERE STUDIEN
INSTITUTE FOR ADVANCED STUDIES
Vienna



Fraunhofer

ISI



Beispiel: Inhaltsanalyse

- ▶ Mittels Skript die Textinhalte herunterladen
- ▶ Stopp-Words entfernen
 - ▶ Häufig genutzte Wörter wie: *und, oder, mit, ...*
 - ▶ Wörterbücher online verfügbar
- ▶ Word Cloud erzeugen
- ▶ Concept Extraction: Gibt Aufschluss über die thematische Ausrichtung



Probleme?

- ▶ Nur ca. 10% aller Wörter ausgewertet



Probleme?

- ▶ Nur ca. 10% aller Wörter ausgewertet
- ▶ Kategorisierung muss trainiert werden
 - ▶ Zeitaufwendig
 - ▶ Kategorien je nach Analyse unterschiedlich



Probleme?

- ▶ Nur ca. 10% aller Wörter ausgewertet
- ▶ Kategorisierung muss trainiert werden
 - ▶ Zeitaufwendig
 - ▶ Kategorien je nach Analyse unterschiedlich
- ▶ Mehrdeutigkeit der Wörter
 - ▶ „Ungenauigkeit“ der Sprache



Mehrdeutigkeit der Sprache

- ▶ **Kontext ist wichtig**

- ▶ Ich gehe zur Prüfung

vs.

- ▶ Ich gehe nur dann zur Prüfung, **wenn die Hölle zufriert**

- ▶ **Wörter sind mehrdeutig**

- ▶ Ich sitze auf der **Bank**

vs.

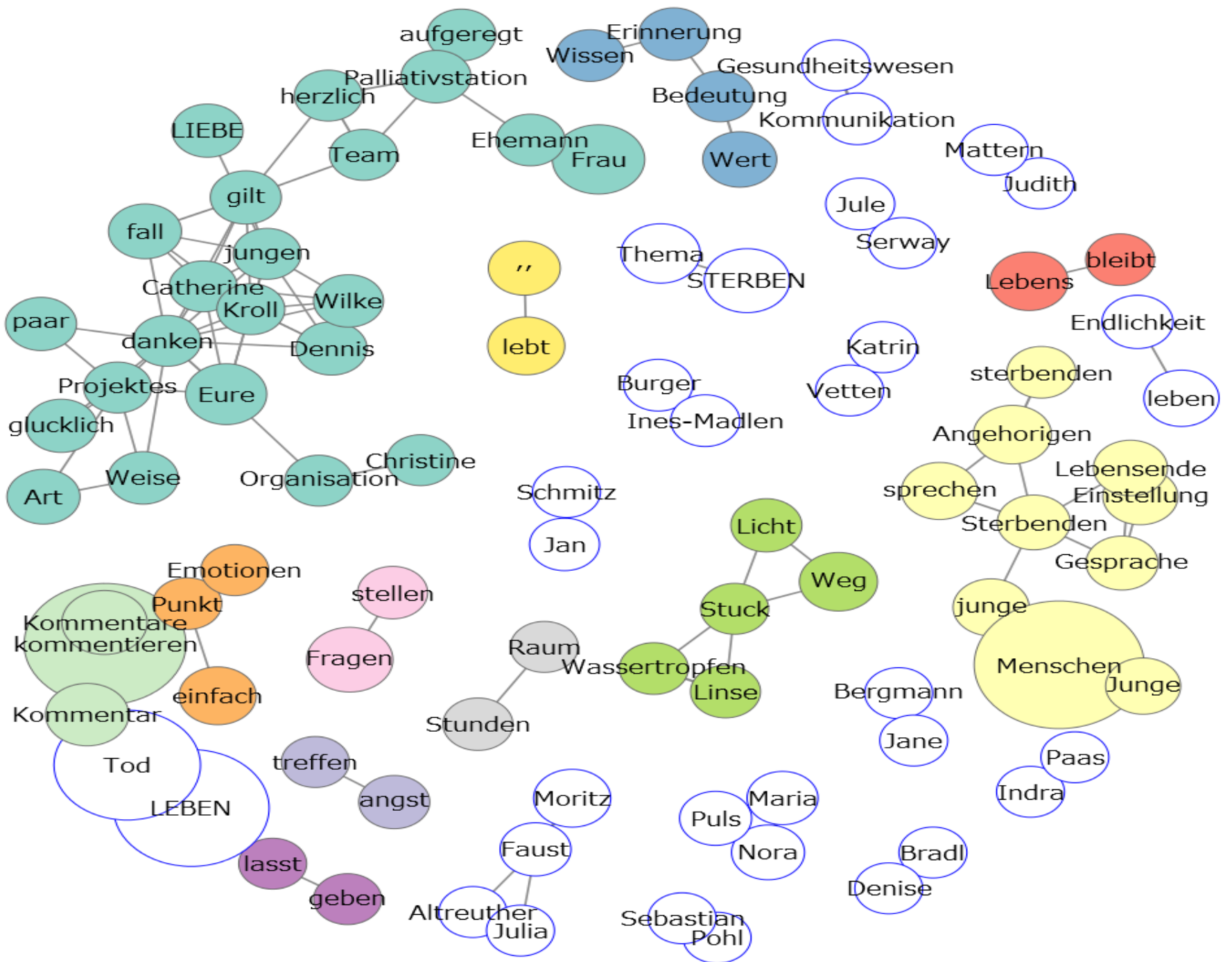
- ▶ Ich habe auf der **Bank** ein Konto

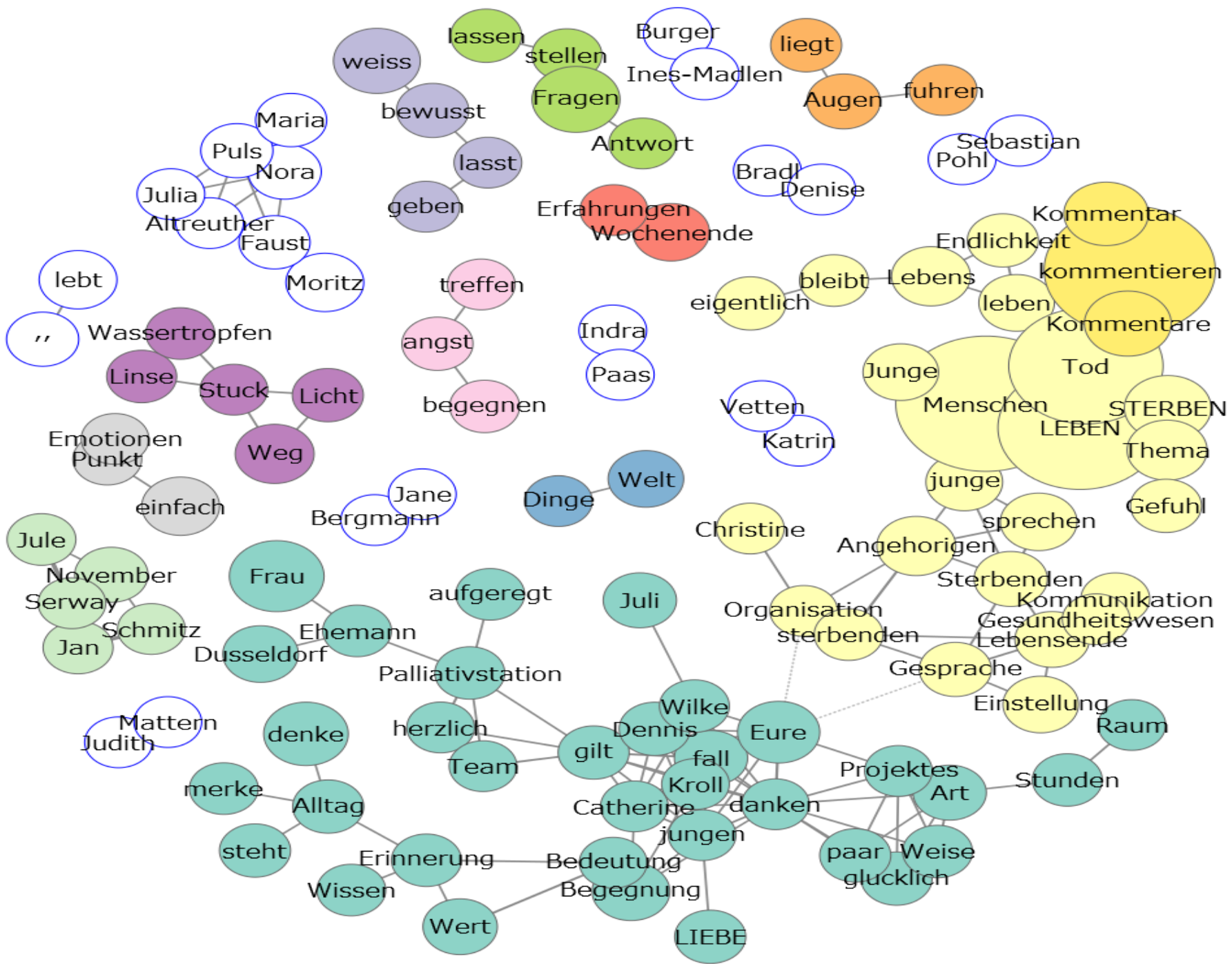


Beispiel: Beziehungen

- ▶ Welche Eigenschaften sind wichtig?
- ▶ Welchen Eindruck erhalten sie davon?
- ▶ Wo liegen ihrer Ansicht nach die Probleme?







Conclusio

- ▶ **Aufbrechen von Texten**

- ▶ Überraschen lassen
- ▶ Beziehungen finden, die man sonst übersehen hätte

- ▶ **Überblick über große Datenmengen erhalten**

- ▶ Stimmungslage
- ▶ Grobe Einschätzungen möglich

- ▶ **Fundiertes Sampling für genauere Analysen**

- ▶ Abseits vom Mainstream analysieren

- ▶ **Was es nicht ist:**

- ▶ Ersatz für tiefergehende Analysen
- ▶ Ausgelagerte Intuition

