

Zum State of the Art automatischer Inhaltsanalyse

Michael Scharkow, M.A.

Universität Hohenheim

Institut für Kommunikationswissenschaft (540G)

michael.scharkow@uni-hohenheim.de

Typologie der Verfahren

• deskriptive/explorative Verfahren

- keine Hypothesen
- keine Kontrolle über das Verfahren
- vor allem zur ersten Exploration der Textdaten
- vollautomatisch und effizient

• überwachte, gelenkte Verfahren

- hypothesengeleitet
- Kontrolle über die Codierung
- manuelle (Vor-)Arbeit nötig
- sinnvolle Kategorien

Ungelenkte Verfahren - Mal schauen, was da ist

Beispiele

- Zuordnung von Autoren bei den Federalist Papers
 - spezifisches Vokabular hilft bei der Identifikation
- Messung der Verständlichkeit von (politischen) Texten
 - Wortlänge und Satzlänge als Indikator von Verständlichkeit
 - Umfang des benutzten Wortschatzes als Indikator für Komplexität
- Visualisierung von Textdaten, z.B. für qualitative Interviews oder offene Antworten aus Fragebögen

Co-Occurrence und Latent Semantic Indexing

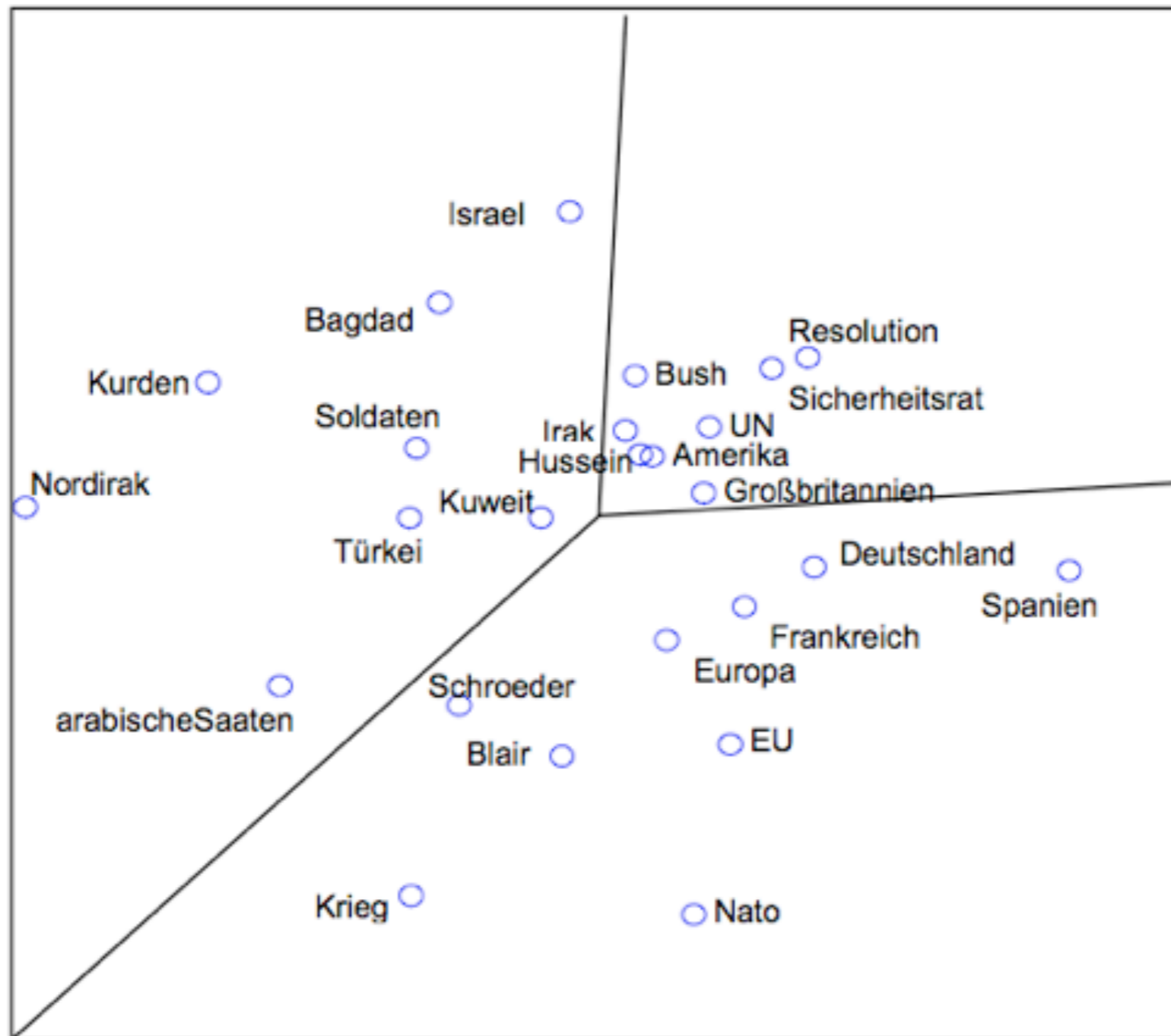
- Wörter tauchen gern **gemeinsam** in bestimmten Kontexten (Sätzen, Absätzen, Dokumenten) auf
- über alle Dokumente kann daher eine **Ähnlichkeitsmatrix** von Wörtern bzw. Wortgruppen erstellt werden (wie oft tauchen Begriffe zusammen auf)
- **Clusteranalyse** und **Multidimensionale Skalierung** arbeiten mit Ähnlichkeitsmatrizen
- LSI verdichtet zusammengehörige Wörter zu **Faktoren**, sog. Semantischen Indizes, mit denen weitergearbeitet wird (ähnlich wie bei einer klassischen **Hauptkomponentenanalyse**)

Term-Dokument-Matrix (Co-Occurrence)

	Peter	sehen	Student	Seminar	im	Telefon	Tag	am
d1	1	1	1	1	1	0	0	0
d2	0	0	0	0	0	1	1	2
d3	1	1	0	0	0	0	2	1

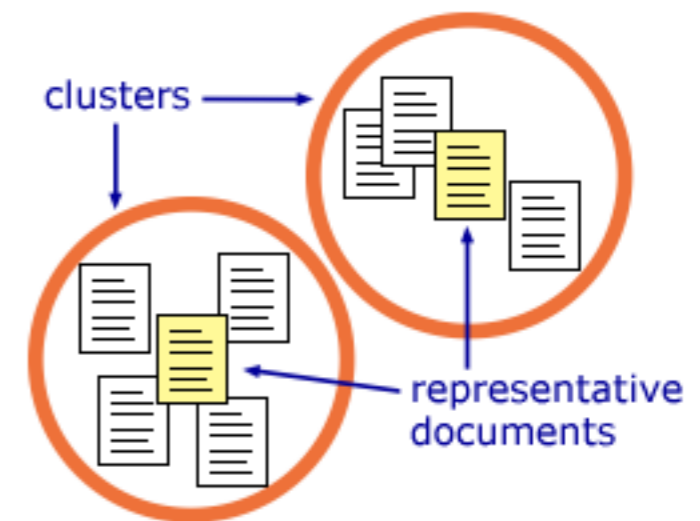
Beispiel Co-Occurrence-Analyse (Landmann/Züll)

Abbildung 7: Multidimensionale Skalierung der Texte der Frankfurter Allgemeinen Zeitung zu Kriegsbeginn



Automatische Dokumentklassifikation

- Term-Dokument-Matrix lädt zum Clustern von **Dokumenten** ein, die ein ähnliches Wortprofil haben
- **k-Means** (partitionierend) Clusteranalyse als häufigstes Verfahren
- Hierarchische Clusteranalyse (agglomerativ) skaliert sehr schlecht mit vielen Fällen und Variablen
- **Clusteranzahl** und **Interpretation** unklar, oft wenig sinnvolle Kategorien



Term-Dokument-Matrix (Document-Clustering)

	Peter	sehen	Student	Seminar	im	Telefon	Tag	am
d1	1	1	1	1	1	0	0	0
d2	0	0	0	0	0	1	1	2
d3	1	1	0	0	0	0	2	1

Beispiel Clustering von Web-Dokumenten (Google)



related:www.csee.umbc.edu/~nicholas/clustering/

Suche

[Erweiterte Suche](#)
[Einstellungen](#)

Suche: Das Web Seiten auf Deutsch Seiten aus Deutschland

Web

Ergebnisse 1 - 10 von ungefähr 18, die www.csee.umbc.edu/~nicholas/clustering/ ähn

[A Hierarchical Monothetic Document Clustering Algorithm ...](#) - [[Diese Seite übersetzen](#)]

Organizing Web search results into a hierarchy of topics and subtopics facilitates browsing the collection and locating results of interest.

citeseer.ist.psu.edu/708223.html - 25k - [Im Cache](#) - [Ähnliche Seiten](#)

von K Kummamuru - 2004 - [Zitiert durch: 84](#) - [Ähnliche Artikel](#) - [Alle 10 Versionen](#)

[A Hierarchical Document Clustering Algorithm](#) - [[Diese Seite übersetzen](#)]

Computer Science » News/Events » Project Presentations / Thesis Defences » 2004 - MSc Thesis Defence » A Hierarchical Document Clustering Algorithm ...

www.cs.dal.ca/news/def-1242.shtml - 8k - [Im Cache](#) - [Ähnliche Seiten](#)

[Introduction to Clustering Large and High-Dimensional Data ...](#) - [[Diese Seite übersetzen](#)]

Introduction to Clustering Large and High-Dimensional Data. Description · Table of contents · Excerpt · Index · Copyright · Frontmatter ...

www.cambridge.org/uk/catalogue/catalogue.asp?isbn=0521852676 - 10k -

[Im Cache](#) - [Ähnliche Seiten](#)

[\[PDF\] O-Cluster: Scalable Clustering of Large High Dimensional Data Sets](#) - [[Diese Seite übersetzen](#)]

Dateiformat: PDF/Adobe Acrobat - [HTML-Version](#)

Copyright © 2002 Oracle Corporation. 1. O-Cluster: Scalable Clustering of Large High Dimensional Data Sets. Boriana L. Milenova. Marcos M. Campos ...

www.oracle.com/technology/products/bi/odm/pdf/o_cluster_algorithm.pdf - [Ähnliche Seiten](#)

von BL Milenova - [Zitiert durch: 28](#) - [Ähnliche Artikel](#) - [Alle 9 Versionen](#)

[A new document clustering algorithm based on association rule ...](#) - [[Diese Seite übersetzen](#)]

Proceedmgs. of. the Thud Internauonal Conference on Macbme Learning and Cybemeucs , Shangha,. 26-29 August. 2004. A NEW DOCUMENT CLUSTERING ALGORITHM ...

ieeexplore.ieee.org/iel5/9459/30108/01382395.pdf - [Ähnliche Seiten](#)

von JC Song - 2004 - [Zitiert durch: 8](#) - [Ähnliche Artikel](#)

Gelenkte Verfahren - es wird ernst!

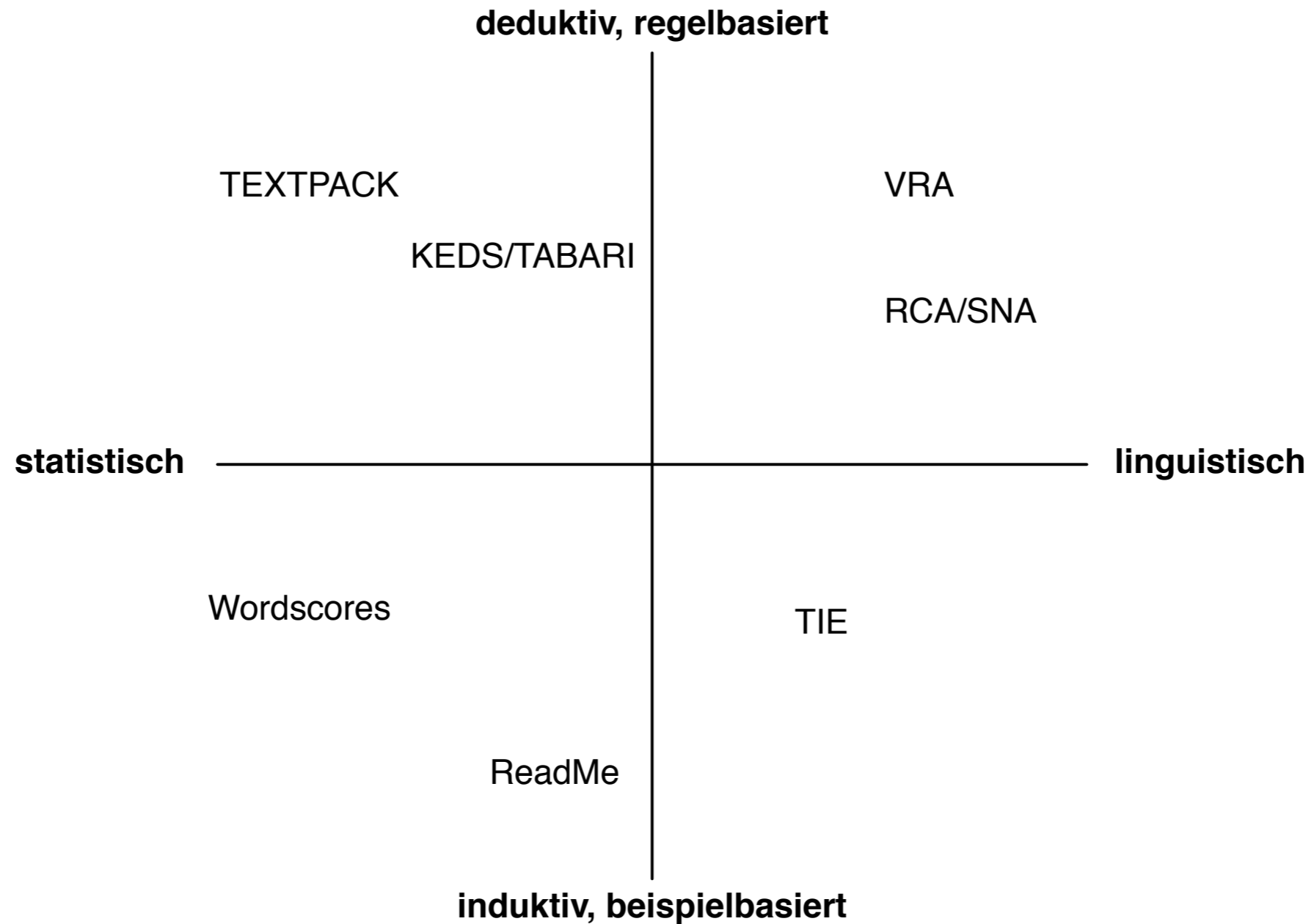
Gelenkte Verfahren - Induktion und Deduktion

- Standardfall für die Kommunikationsforschung
- **Wir haben Hypothesen, die wir überprüfen wollen!**
- Computersoftware muss gezielt gesteuert werden, aber wie?
 - durch explizite Begriffe - `^(Finanz|Banken|Kredit)krise*$`
 - durch explizite Regeln - `html.title.content.string.split[1]`
 - durch annotierte Beispiele - `doc1 = SPORT; doc2 = WIRTSCH`

Gelenkte Verfahren - Ansatz und Gegenstand

- **Bag-of-Words-Verfahren** betrachten Texte als eine (ungeordnete) Ansammlung von Wörtern bzw. Zeichen
 - rein lexikalische Ebene, Analyseeinheit ist der Text
 - klassische Inhaltsanalyse (Codierung von Themen, Nachrichtenfaktoren)
- **Syntaktisch-semantische Ansätze**
 - Texte bestehen aus Aussagen bzw. Propositionen
 - Computer soll die Bestandteile der Aussage erkennen und analysieren

Typologie gelenkter Verfahren



Diktionärsbasierte Verfahren

- **Urgestein** der automatischen Inhaltsanalyse (General Inquirer seit den 60ern)
- Forscher definieren **Wortlisten** (Diktionäre), nach denen dann Dokumente klassifiziert werden
 - „Wenn das Dokument das Wort „Fußball“ enthält, wird es als Sportmeldung codiert.“
- **Freitextrecherche** über Suchmaschinen funktioniert ähnlich, allerdings ist die Grundgesamtheit oft unbekannt, also Validität der Ergebnisse unklar.
- vollständige Reliabilität und gute Effizienz, Validität bei komplexen Konstrukten oft sehr problematisch und Diktionärsentwicklung aufwändig

Aufgabe - Entwickelt ein Diktionär für ...

- Medienresonanzanalyse bei Volkswagen (leicht!)
 - Polo, Golf, VfL Wolfsburg, Betriebsrat + Prostituierte, ...
- Nachrichtenfaktor Schaden (auch nicht schwer, aber lang)
 - sterben, Unfall, verunglücken, Tote, Verletzte, Absturz, Amok, ...
- Wirtschaftsnachrichten - ???
- Personalisierung der Wahlkampfberichterstattung - ???
- Postmaterialismus in Wahlprogrammen - ???

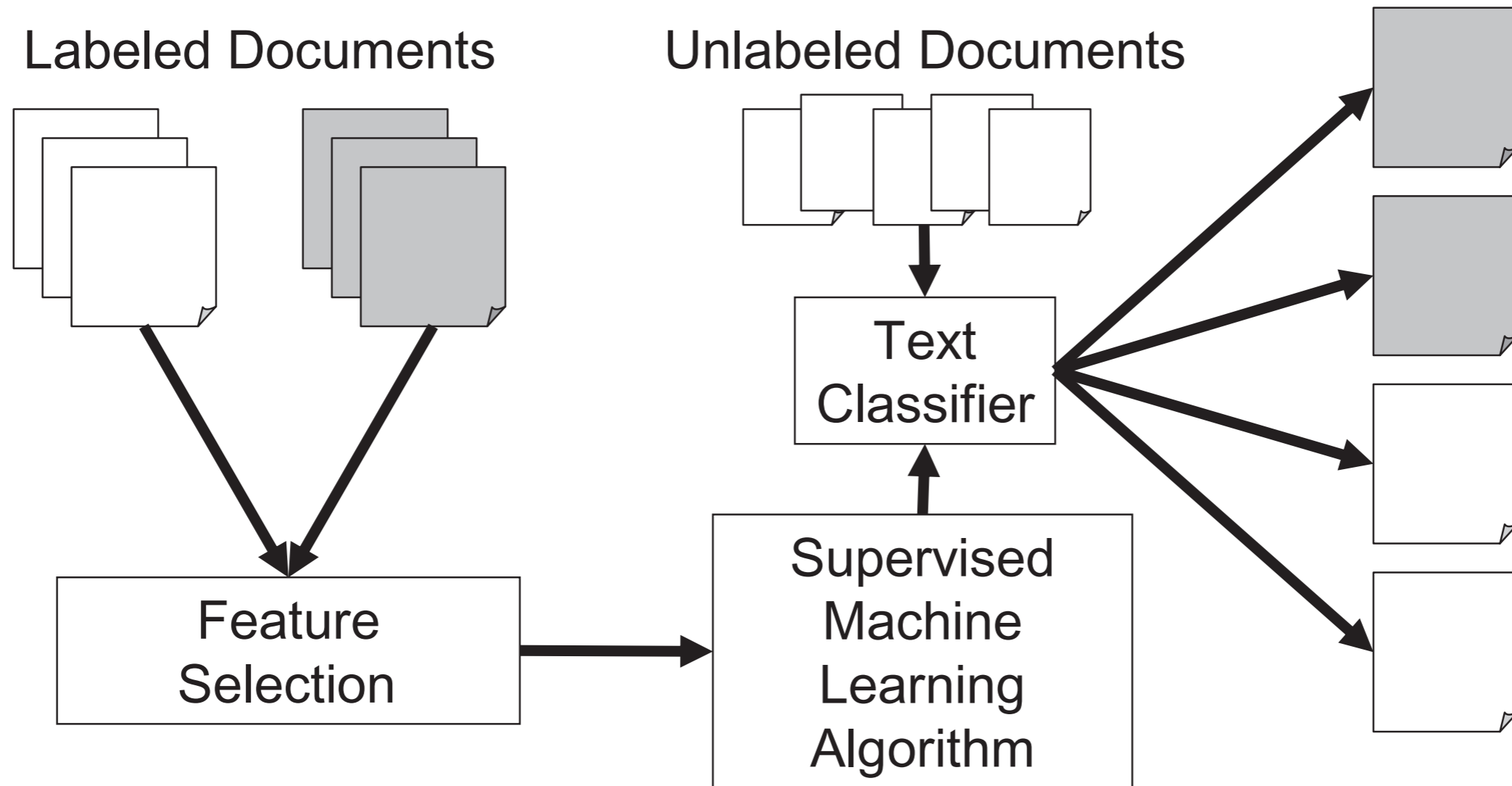
Regelbasierte linguistische Verfahren

- Umwandlung von Aussagen in analysierbare **Graphenstruktur**
- Information wird aus den Dokumenten extrahiert, d.h. offene Frage „Wer tut was mit wem?“ soll beantwortet werden
- **Syntaktische Parser** und andere Verfahren funktionieren nur bei domainspezifischen Texten, Beispiel: Extraktion von int. Ereignissen aus Tickermeldungen
- sowohl automatisches Parsing als auch die Definition von Codierregeln sind **extrem aufwändig**, meistens werden Diktionäre zur Hilfe geholt
- Verfahren eignen sich super für Online-Inhaltsanalyse, weil HTML, XML, Mails stark strukturiert sind

Induktive überwachte Textklassifikation

- da die Definition von Diktionären bzw. Regeln sehr aufwändig und wenig dem manuellen **beispielgeleiteten Codierprozess** entspricht: lernende Klassifikatoren
- Klassifikator erhält **Trainingsdaten**, die Rohtext und Kategorie enthalten, extrahiert daraus Feature-Gewichte (vgl. Regressionskoeffizienten), mit denen dann neue Dokumente klassifiziert werden
- man erhält ein **probabilistisches Diktionär**, das man nicht selbst definiert
- Beispiele: Wordscores zur Links-Rechts-Codierung von Wahlprogrammen
- viel Trainingsmaterial und gute Algorithmen verbessern die Validität

Funktionsweise induktiver Klassifikation



Induktive Informationsextraktion

- siehe linguistische, regelbasierte Codierung, aber **ohne explizite** Regeln
- der Computer soll **selbst lernen**, was ein Akteur, eine Handlung, ein Termin, ein Ort ist
- Annotierung der Aussagen und Inferenz auf **andere Texte** mit **anderen Strukturen** und **anderen Inhalten** (!)
- funktioniert bislang nur **sehr begrenzt**
- wäre aber das mächtigste Verfahren neben der vollautomatischen Textinterpretation

Fazit

- alle Verfahren haben in bestimmten Situationen noch Sinn
- Diktionäre eignen sich bei einfachen Konstrukten, Eigennamen, Marken etc.
- Lesbarkeit und Stilometrie sind bislang relativ selten in kommwiss. Forschung (außerhalb von Hohenheim!)
- Induktive Verfahren sind toll, weil sie anschlussfähig an bestehende manuelle Codierpraxis sind und trotzdem super mit vielen, vielen Dokumenten skalieren
- Linguistische Verfahren sind noch nicht weit genug für Standardfragestellungen.

**Danke für die Einladung und die
Aufmerksamkeit!**

michael.scharkow@uni-hohenheim.de