

Korpuslinguistik

Jürgen Spitzmüller

Vorlesungsreihe »Methoden der Linguistik«

(WS 2004/05)

Universität Zürich, 7.12.2004

1. Teil: Theorie – Grundlegende theoretische Fragestellungen:

- ▶ Was sind überhaupt Korpora?
- ▶ Wozu Korpora?
- ▶ Was sollen Korpora leisten/was können sie *nicht* leisten?
- ▶ Welche Korpusarten gibt es?
- ▶ Wie müssen linguistische Korpora aufbereitet sein?

2. Teil: Praxis – Konkrete Arbeit mit einem vorgefertigten Korpus

- ▶ Vorstellung der Korpora des IDS und des COSMAS-Systems
- ▶ Versuch, einige kleine lexikologische, semantische und syntaktische Fragestellungen korpusanalytisch zu lösen

Linguistische Datenerhebung

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

1. Introspektion

Linguistische Datenerhebung

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

1. Introspektion
2. Experiment

Linguistische Datenerhebung

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

1. Introspektion
2. Experiment
3. Datenerhebung

Linguistische Datenerhebung

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpusarten

Beispiele

1. Introspektion
2. Experiment
3. Datenerhebung
4. Korpora

»Corpus linguists« vs. »armchair linguists«

»Armchair linguistics«: Bilden die Daten selbst (Introspektion) und beurteilen ihre »Richtigkeit« mit Hilfe eigener Sprachkompetenz

»Empirical linguistics«: Sammeln die Daten im konkreten »Sprechalltag« und beurteilen sie auf der Grundlage ihres Vorkommens

»Corpus linguists« vs. »armchair linguists«

»Armchair linguistics«: Bilden die Daten selbst (Introspektion) und beurteilen ihre »Richtigkeit« mit Hilfe eigener Sprachkompetenz

- ▶ Daten entsprechen nicht dem tatsächlich Geäußerten; Linguist verlässt sich zu stark auf sein eigenes »Sprachgefühl« und auf die Varietäten, die er verwendet.
- ▶ Wichtige Daten werden übersehen, andere überschätzt

»Empirical linguistics«: Sammeln die Daten im konkreten »Sprechalltag« und beurteilen sie auf der Grundlage ihres Vorkommens

- ▶ Die Performanz ist unvollständig
- ▶ Die Performanz enthält Fehler
- ▶ Die Performanz ist irrelevant (gg. der Kompetenz)

»Corpus linguists« vs. »armchair linguists«

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpusarten

Beispiele

»These two don't speak to each other often, but when they do, the corpus linguist says to the armchair linguist, ›Why should I think that what you tell me is true?‹ and the armchair linguist says to the corpus linguist, ›Why should I think that what you tell me is interesting?‹«

Charles Fillmore (1992): »Corpus linguistics« or »Computer-aided armchair linguistics«. In: Jan Svartvik (Hg.): Directions in Corpus Linguistics. Proceedings of the Nobel Symposium 82, Stockholm, 4.–8. August 1991, Berlin/New York, S. 35–60.

Korpora: Sinn und Unsinn

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

Was sollen/können Korpora leisten?

Korpora: Sinn und Unsinn

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

Was sollen/können Korpora leisten?

- ▶ Empirische Überprüfung linguistischer Fragen

Was sollen/können Korpora leisten?

- ▶ Empirische Überprüfung linguistischer Fragen
- ▶ Möglichst realitätsnaher »Ausschnitt« aus einer bestimmten Sprechpraxis

Was sollen/können Korpora leisten?

- ▶ Empirische Überprüfung linguistischer Fragen
- ▶ Möglichst realitätsnaher »Ausschnitt« aus einer bestimmten Sprechpraxis
- ▶ Vermeidung der Fehl- und Überinterpretation eigener Sprechweisen des Forschers und künstlich zu Stande gekommener Daten

Was sollen/können Korpora leisten?

- ▶ Empirische Überprüfung linguistischer Fragen
- ▶ Möglichst realitätsnaher »Ausschnitt« aus einer bestimmten Sprechpraxis
- ▶ Vermeidung der Fehl- und Überinterpretation eigener Sprechweisen des Forschers und künstlich zu Stande gekommener Daten
- ▶ Quantitative (z. B. statistische) Aussagen sind (mit Einschränkungen) möglich

Korpora: Sinn und Unsinn

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

Was können Korpora nicht leisten?

Was können Korpora nicht leisten?

- ▶ Niemals »vollständiges« Abbild einer Sprache oder Varietät

Was können Korpora nicht leisten?

- ▶ Niemals »vollständiges« Abbild einer Sprache oder Varietät
- ▶ Auf eine Auswahl von Textsorten beschränkt

Was können Korpora nicht leisten?

- ▶ Niemals »vollständiges« Abbild einer Sprache oder Varietät
- ▶ Auf eine Auswahl von Textsorten beschränkt
- ▶ Aussagen beziehen sich *immer* zunächst nur auf das Korpus selbst

Was können Korpora nicht leisten?

- ▶ Niemals »vollständiges« Abbild einer Sprache oder Varietät
- ▶ Auf eine Auswahl von Textsorten beschränkt
- ▶ Aussagen beziehen sich *immer* zunächst nur auf das Korpus selbst
- ▶ Quantität ist nicht alles!

Korpustypen

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

Korpustypen

Grundsätzlich: Korpuswahl abhängig vom
Forschungsziel

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

Korpustypen

Grundsätzlich: Korpuswahl abhängig vom Forschungsziel

1. Zusammensetzung

- ▶ Homogen: thematisch, textsortenspezifisch, medial eingeschränkt

Grundsätzlich: Korpuswahl abhängig vom Forschungsziel

1. Zusammensetzung

- ▶ Homogen: thematisch, textsortenspezifisch, medial eingeschränkt
- ▶ Heterogen: möglichst breite Auswahl verschiedenartigster Texte

Grundsätzlich: Korpuswahl abhängig vom Forschungsziel

1. Zusammensetzung

- ▶ Homogen: thematisch, textsortenspezifisch, medial eingeschränkt
- ▶ Heterogen: möglichst breite Auswahl verschiedenartigster Texte

2. Metasprachliche Informationen

Grundsätzlich: Korpuswahl abhängig vom Forschungsziel

1. Zusammensetzung

- ▶ Homogen: thematisch, textsortenspezifisch, medial eingeschränkt
- ▶ Heterogen: möglichst breite Auswahl verschiedenartigster Texte

2. Metasprachliche Informationen

3. Genese

- ▶ Vorgefertigtes Korpus

Grundsätzlich: Korpuswahl abhängig vom Forschungsziel

1. Zusammensetzung

- ▶ Homogen: thematisch, textsortenspezifisch, medial eingeschränkt
- ▶ Heterogen: möglichst breite Auswahl verschiedenartigster Texte

2. Metasprachliche Informationen

3. Genese

- ▶ Vorgefertigtes Korpus
- ▶ Selbst erstelltes Korpus

Grundsätzlich: Korpuswahl abhängig vom Forschungsziel

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpusarten

Beispiele

1. Zusammensetzung

- ▶ Homogen: thematisch, textsortenspezifisch, medial eingeschränkt
- ▶ Heterogen: möglichst breite Auswahl verschiedenartigster Texte

2. Metasprachliche Informationen

3. Genese

- ▶ Vorgefertigtes Korpus
- ▶ Selbst erstelltes Korpus

4. Frage der »Repräsentativität«

Grundsätzlich: Korpuswahl abhängig vom Forschungsziel

1. Zusammensetzung

- ▶ Homogen: thematisch, textsortenspezifisch, medial eingeschränkt
- ▶ Heterogen: möglichst breite Auswahl verschiedenartigster Texte

2. Metasprachliche Informationen

3. Genese

- ▶ Vorgefertigtes Korpus
- ▶ Selbst erstelltes Korpus

4. Frage der »Repräsentativität«

5. Analyseverfahren

- ▶ quantitativ (computergestützt)
- ▶ qualitativ (hermeneutisch)

Anforderungen an ling. Korpora

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

Anforderungen an ling. Korpora

Annotationen (metasprachliche Informationen und Informationen zu den Dokumenten):

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpustypen

Beispiele

Anforderungen an ling. Korpora

Annotationen (metasprachliche Informationen und Informationen zu den Dokumenten):

- ▶ Meta-Informationen: Textsorte, Verfasser, Quelle, Entstehungszeit etc.

Anforderungen an ling. Korpora

Korpuslinguistik

Jürgen
Spitzmüller

Annotationen (metasprachliche Informationen und Informationen zu den Dokumenten):

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpusarten

Beispiele

- ▶ Meta-Informationen: Textsorte, Verfasser, Quelle, Entstehungszeit etc.
- ▶ Morphosyntaktische Annotationen: Wortart, Flexionsform, Stellung im Satz, Stellung des Satzes im Text, Phrasenstruktur, satzübergreifende Verweisen etc.

Anforderungen an ling. Korpora

Annotationen (metasprachliche Informationen und Informationen zu den Dokumenten):

- ▶ Meta-Informationen: Textsorte, Verfasser, Quelle, Entstehungszeit etc.
- ▶ Morphosyntaktische Annotationen: Wortart, Flexionsform, Stellung im Satz, Stellung des Satzes im Text, Phrasenstruktur, satzübergreifende Verweisen etc.
- ▶ Semantische Annotationen: Bedeutung, zum Wortfeld etc.

Anforderungen an ling. Korpora

Annotationen (metasprachliche Informationen und Informationen zu den Dokumenten):

- ▶ Meta-Informationen: Textsorte, Verfasser, Quelle, Entstehungszeit etc.
- ▶ Morphosyntaktische Annotationen: Wortart, Flexionsform, Stellung im Satz, Stellung des Satzes im Text, Phrasenstruktur, satzübergreifende Verweisen etc.
- ▶ Semantische Annotationen: Bedeutung, zum Wortfeld etc.
- ▶ Kontext- und Kotextannotationen: Wortumfeld, Verwendungsspezifik etc.

Annotationen (metasprachliche Informationen und Informationen zu den Dokumenten):

- ▶ Meta-Informationen: Textsorte, Verfasser, Quelle, Entstehungszeit etc.
- ▶ Morphosyntaktische Annotationen: Wortart, Flexionsform, Stellung im Satz, Stellung des Satzes im Text, Phrasenstruktur, satzübergreifende Verweisen etc.
- ▶ Semantische Annotationen: Bedeutung, zum Wortfeld etc.
- ▶ Kontext- und Kotextannotationen: Wortumfeld, Verwendungsspezifik etc.
- ▶ *Treebanks*: Möglichkeit, Strukturbäume zu generieren

Auszeichnungssprachen (Markup Languages):

- ▶ HTML (Hypertext Markup Language)
- ▶ SGML (Standard Generalized Markup Language)
- ▶ XML (Extensible Markup Language)

Methode (bei allen): Setzen von »Tags« (Marken), die ein Programm (»Parser«) lesen und verarbeiten kann, die aber im Text unsichtbar bleiben.

Annotationsverfahren: Beispiele

Korpuslinguistik

Jürgen
Spitzmüller

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpusarten

Beispiele

- (1) Chomskys Bücher sind gerade wieder
<ADJ>modern</ADJ>.

Annotationsverfahren: Beispiele

- (1) Chomskys Bücher sind gerade wieder
<ADJ>modern</ADJ>.
- (2) Chomskys Bücher <VRB>modern</VRB> in der
Bibliothek vor sich hin.

Annotationsverfahren: Beispiele

- (1) Chomskys Bücher sind gerade wieder
<ADJ>modern</ADJ>.
- (2) Chomskys Bücher <VRB>modern</VRB> in der
Bibliothek vor sich hin.
- (3) <NOU CAS=GEN>Chomskys</NOU>Bücher
<VRB NUM=Plu PERS=3 MOD=Ind GV=akt>modern</VRB>

Das COSMAS II-System

Ort: IDS Mannheim

(<http://www.ids-mannheim.de/cosmas2/>)

Größe: ca. 1,52 Milliarden Wortformen

Davon annotiert: 26 Mio. Textwörter

Struktur: Modular (156 Teilkorpora)

Zugangsmöglichkeit: Online (passwortgeschützt) und vor Ort, jeweils mit Hilfe eines „Client-Programms“ (derzeit nur Windows)

Textsorten: Zeitungstexte, Belletristik, Fachtexte, transkribierte Gespräche, auch historische Texte

Besonderheiten: Sehr mächtige Recherchemöglichkeiten; Recherche in Gesprächstranskripten; sehr umfangreich; modulare Zusammenstellung; Teilkorpora sehr gut dokumentiert; Belege jederzeit verifizierbar

Einschränkungen: Texte werden (aus urheberrechtlichen Gründen) nur auszugsweise angezeigt (nicht geeignet für diskursanalytische Studien)

Das Leipziger Wortschatz-Projekt

Korpuslinguistik

Jürgen
Spitzmüller

Ort: Uni Leipzig

(<http://wortschatz.uni-leipzig.de>)

Zugang: Online einsehbar (keine Anmeldung erforderlich), Recherche mit Hilfe des Browsers (Eingabeformular)

Überblick

Datenerhebung

Korpora: Sinn
und Unsinn

Korpusarten

Beispiele

Größe: ca. 518 Mio. Wortformen

Textsorten: v. a. Zeitungstexte, z. T. auch spez. Wortlisten und Fachtexte (Lexika) sowie Belletristik (Projekt Gutenberg)

Besonderheiten: Kollokationen (grafisch), Suche nach Anagrammen

Einschränkungen: Keine detaillierte Auskunft über die Korpuszusammensetzung, Ausgabe nur einiger ausgewählter Textstellen

Das TIGER-Korpus

Ort: Institut für Maschinelle
Sprachverarbeitung Stuttgart
([http://www.ims.uni-stuttgart.de/
projekte/TIGER/](http://www.ims.uni-stuttgart.de/projekte/TIGER/))

Zugang: Muss komplett heruntergeladen und mit
einem speziellen Programm
(<http://www.tigersearch.de>) eingesehen
werden. Anmeldung (für Programm und
Korpus jeweils separat) erforderlich.

Größe: 700.000 Wortformen

Textsorten: Zeitungstexte (*Frankfurter Rundschau*)

Besonderheiten: *Treebanks* (Generierung von
Strukturbäumen)

Einschränkungen: relativ klein, relativ homogen.

Strukturbäume mit dem TIGER

The screenshot shows the TIGERGraphViewer interface. The main window displays a syntax tree for the sentence "Direktor des in Frankfurt ansässigen Genossenschaftsverbandes". The root node is NP, which branches into NK (Direktor) and AG. AG branches into NP, which further branches into NK (des), NK (in), AP (ansässigen), and NK (Genossenschaftsverbandes). The AP node branches into MO (in) and HQ (ansässigen). The MO node branches into PP, which branches into AC (in) and NK (Frankfurt). The AC node branches into AC (in). The HQ node branches into HQ (ansässigen). The NK node branches into NK (Genossenschaftsverbandes). The AC node branches into AC (in). The HQ node branches into HQ (ansässigen). The NK node branches into NK (Genossenschaftsverbandes).

Below the tree, the words and their grammatical categories are listed:

Direktor	des	in	Frankfurt	ansässigen	Genossenschaftsverbandes
NN	ART	APPR	NE	ADJA	NN
Masc.Nom.Sg	Masc.Gen.Sg	Dat	Neut.Dat.Sg	Pos.Masc.Gen.Sg	Masc.Gen.Sg
Direktor	der	in	Frankfurt	ansässig	Genossenschaftsverband

At the bottom of the window, there are navigation controls and statistics:

- Graphs: 2
- Subgraphs: 3
- Navigation: Previous, Next, First, Last
- Subgraph: 1 / 1

s26878: Und Walter Weinkauf , Direktor des in Frankfurt ansässigen Genossenschaftsverbandes , ergänzt : `` Auf jeden freigesetzten Arbeitsplatz kommen im Mittelstand 1,5 Neueinstellungen , während in Großunternehmen zehn Prozent mehr Arbeitsplätze wegrationalisiert als geschaffen werden . ``