

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

8. Juli 2014

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Motivation

Cohens κ

Fleiss' κ

COSMAS

Allgemeines

Zugang und Export

Einfache Suchanfragen

Wir sind hier...

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Grundlagen

Wiederholungen von Zeichen

Gruppen und Nicht-Zeichen

Liste der reservierten Zeichen

Abstände

Inter-Kodierer-Übereinstimmung

Motivation

Cohens κ

Fleiss' κ

COSMAS

Allgemeines

Zugang und Export

Einfache Suchanfragen

Abstände

Logik und Klammerung

Morphologische Annotation

Grenz-Markierungen

DWDS

Allgemeines

Zugang und Export

Einfache Suchanfragen

Abstände

Morphologische Annotation

CQP

Allgemeines

Installation und Basics

Einfache Suchanfragen

Anzeigeoptionen setzen

Suchanfragen speichern und ausgeben

Gruppieren und Sortieren

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Korpusstudien

Überlegungen vor der **quantitativen** Korpusstudie:

- Ist eine Korpusstudie überhaupt das richtige Instrument?
 - Welche **Grundgesamtheit** interessiert mich?
 - Welches Korpus repräsentiert die Grundgesamtheit? Es gibt kein „repräsentatives Korpus“ an sich!
 - Wie muss ich die Stichprobe passend zu meiner Hypothese durchführen, um die Hypothese **valid** zu testen?
 - Ist das Korpus hinreichend linguistisch annotiert?
 - Was müsste herauskommen, um meine Hypothese zu **falsifizieren**?
- Alles in allem müssen die Fragen aus Teil 2 der Statistik-VL vor der Studie lückenlos geklärt sein!

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Probleme

Wann ein Korpus evtl. nicht passt:

- Wenn klassische „**Grammatikalität**“ ein zentraler Begriff der zu testenden Theorie ist.
- Bei seltenen Phänomenen (kann zumindest Probleme machen).
- Wenn das Phänomen nicht (halb-)automatisch zu suchen ist und keine Kapazitäten für die manuelle Suche vorhanden sind.
- Wenn keine passenden Korpora existieren.

Zitieren Sie immer die Artikel der Korpus-Ersteller!

Lesen Sie sie auch!

Es könnten relevante Informationen über das Korpus drinstehen!

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Wir sind hier...

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Grundlagen

Wiederholungen von Zeichen

Gruppen und Nicht-Zeichen

Liste der reservierten Zeichen

Abstände

Inter-Kodierer-Übereinstimmung

Motivation

Cohens κ

Fleiss' κ

COSMAS

Allgemeines

Zugang und Export

Einfache Suchanfragen

Abstände

Logik und Klammerung

Morphologische Annotation

Grenz-Markierungen

DWDS

Allgemeines

Zugang und Export

Einfache Suchanfragen

Abstände

Morphologische Annotation

CQP

Allgemeines

Installation und Basics

Einfache Suchanfragen

Anzeigeoptionen setzen

Suchanfragen speichern und ausgeben

Gruppieren und Sortieren

Grober Ablauf

- Hypothesen formulieren
- Hypothesen **operationalisieren**
- Korpus wählen, Stichprobenverfahren definieren
- Stichprobe ziehen, prüfen, evtl. neu ziehen
- Stichprobe ggf. annotieren
- Hypothesen inferenzstatistisch prüfen

Alle diese Schritte sollten in Ihrem Artikel/Ihrer Arbeit **dokumentiert** und alle Entscheidungen **konzeptuell** motiviert werden.

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Zum Stichprobenverfahren

- Die Stichprobengröße kann i. d. R. nicht a priori festgelegt werden.
- Stichprobengrößen sind ein Kompromiss aus α -Niveau und Teststärke.
- Nur bei **gerichteten Hypothesen** (s. Statistik-VL) gibt es ideale Größen.
- Wahl zwischen **Quotenstichprobe** und **Zufallsstichprobe**: Ich bin vehementer Vertreter von Zufallsstichproben!
 - Quotenstichproben sind fehleranfällig,
 - sie bringen massenweise Auxiliarahypothesen mit,
 - und die Angst vor der großen Zufallsstichprobe ist meist unbegründet.
 - Die Suchanfrage(n) dürfen **nur niemals relevante Daten ausschließen!**
- Dringend empfohlen im Stichprobenprozess (und schon im Operationalisierungsprozess): **Konkordanz-Surfen!** Die besten Ideen kommen im Angesicht der Daten!

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Zur Annotation

Was sollte man annotieren?

- NUR** hypothesengebundene Variablen!
- Auf Verdacht „schnell noch Kasus (etc.) mit-annotieren“ ist Zeitverschwendung.
- Die Geschwindigkeit der AnnotatorInnen muss getestet werden.
- Lieber drei Variablen kodieren und genug Belege bekommen als acht Variablen und dann mit hundert Belegen dastehen!

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Hinweis zur Korpusqualität

Grenzen der verfügbaren linguistischen Annotation in Korpora:

- Tokenizer/Tagger/Chunker usw. sind auf Standardsprache trainiert.
- Je weniger nah an dieser das gesuchte Phänomen ist, desto **unzuverlässiger** sind die Informationen im Korpus...
- ...zusätzlich zur allgemeinen Fehlerrate.
- Prüfen Sie immer an kleinen Stichproben, ob die Suchanfrage in linguistischer Annotation korrekt läuft.
- Eventuell sind bereits konzeptuelle Entscheidungen von den Taggern anders getroffen worden, als bei Ihnen.
- Beispiele:
 - Unterscheidung von Infinitiven und Präsensformen (PI) ist nicht 100% genau.
 - Suche nach Lemmata wie „bauspüren“ führt nicht zu „baugespart“.

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

Wir sind hier...

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Grundlagen

Wiederholungen von Zeichen

Gruppen und Nicht-Zeichen

Liste der reservierten Zeichen

Abstände

Inter-Kodierer-Übereinstimmung

Motivation

Cohens κ

Fleiss' κ

COSMAS

Allgemeines

Zugang und Export

Einfache Suchanfragen

Abstände

Logik und Klammerung

Morphologische Annotation

Grenz-Markierungen

DWDS

Allgemeines

Zugang und Export

Einfache Suchanfragen

Abstände

Morphologische Annotation

CQP

Allgemeines

Installation und Basics

Einfache Suchanfragen

Anzeigeoptionen setzen

Suchanfragen speichern und ausgeben

Gruppieren und Sortieren

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer

Korpora und Korpusstudien

Korpusstudien planen

Ablauf der Studie

Kenntnisse

Reguläre Ausdrücke

Inter-Kodierer-Übereinstimmung

COSMAS

DWDS

CQP

| | | | |
|--|---|--|--|
| <h2>Rüstzeug für KorpuslinguistInnen</h2> <p>Was muss man können?</p> <ul style="list-style-type: none"> übliche Abfragesprachen von Korpus-Software Tabellenkalkulation Rohtextbearbeitung (wenigstens Notepad++ o. ä.) Statistik inkl. R, SPSS o. ä. Berechnung von Inter-Kodierer-Übereinstimmung <p>Was sollte man können?</p> <ul style="list-style-type: none"> eine Skriptsprache reguläre Ausdrücke Stochastik | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COsmAS</p> <p>Allgemeines</p> <p>Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>Allgemeines</p> <p>Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> <p>Allgemeines</p> <p>Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Platzhalter</h2> <ul style="list-style-type: none"> Oft sucht man nach Zeichenketten, die bestimmte Platzhalter enthalten. Komposita mit <i>Haus</i> als Erstglied usw. Reguläre Ausdrücke (Regexe) erlauben es, Zeichenketten allgemeiner als durch <i>literale</i> Zeichenketten anzugeben. Vereinfacht kann man sich die Regex-Syntax als <i>Platzhalter-Syntax</i> vorstellen. | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> | <h2>Zeichen</h2> <p>Man kann ein Zeichen in Regexen folgendermaßen definieren:</p> <ul style="list-style-type: none"> literal <ul style="list-style-type: none"> z. B. <code>a</code> oder <code>4</code> ⇒ genau ein Vorkommen genau dieser Zeichen Zeichenklasse in <code>[]</code>, interpretiert als Aufzählung <ul style="list-style-type: none"> <code>[a4]</code> ⇒ genau ein Vorkommen von entweder <code>a</code> oder <code>4</code> Zeichenklasse mit Bereich <ul style="list-style-type: none"> <code>[A-Z]</code> ⇒ genau ein Vorkommen eines Zeichens zwischen <code>A</code> und <code>Z</code> aufgezählte Zeichenklassen- Bereiche <ul style="list-style-type: none"> <code>[A-Za-z]</code> ⇒ genau ein Vorkommen eines Zeichens in <code>A-Z</code> oder <code>a-z</code> negierte Zeichnkasse mit <code>^</code> (Komplementbildung) <ul style="list-style-type: none"> <code>[^A-Z]</code> ⇒ genau ein Vorkommen irgendeines Zeichens außer <code>A-Z</code> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Beliebige Zeichen</h2> <ul style="list-style-type: none"> irgendein Zeichen (ob Ziffer, Buchstabe, Sonderzeichen...) Als Folge der besonderen Bedeutung einiger Zeichen wie <code>.</code> in der Regex-Syntax, müssen die zugehörigen Literale geschützt werden: <ul style="list-style-type: none"> <code>\.</code> ⇒ der Punkt <code>\\</code> ⇒ der Backslash <code>\[</code> ⇒ die eckige öffnende Klammer ... das geht mit allen Zeichen so, die eine besondere Bedeutung haben. | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COsmAS</p> <p>Allgemeines</p> <p>Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>Allgemeines</p> <p>Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> <p>Allgemeines</p> <p>Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Wiederholung von Zeichen</h2> <ul style="list-style-type: none"> Alle Literale und Zeichenklassen sowie der Platzhalter stehen für ein Zeichen in der Zeichenkette. <code>Ha.s</code> findet also <i>Hans</i>, <i>Haus</i>, <i>Ha5s</i>, aber nicht <i>Hauses</i> usw. Nach jedem Zeichen oder jeder Zeichenklasse kann daher ein Quantor stehen. Der Quantor gibt an, wie oft das Zeichen (oder irgendein Zeichen der Klasse) wiederholt werden kann oder muss. <code>*</code> ⇒ null mal oder beliebig oft <code>+</code> ⇒ mindestens einmal <code>{3}</code> ⇒ genau drei Mal <code>{3,6}</code> ⇒ drei bis sechs Mal <code>{3,}</code> ⇒ drei Mal oder öfter | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COsmAS</p> <p>Allgemeines</p> <p>Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>Allgemeines</p> <p>Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> <p>Allgemeines</p> <p>Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Grenzen von Zeichenketten</h2> <ul style="list-style-type: none"> Zwei besondere Zeichen markieren Positionen in der Zeichenkette. <code>^</code> ⇒ der Anfang der Zeichenkette <code>\$</code> ⇒ das Ende der Zeichenkette | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> | <h2>Gruppen</h2> <ul style="list-style-type: none"> Man kann Teile von Zeichenketten zusammenfassen. <code>Haus(meister)</code> ⇒ nichts anderes als <code>Hausmeister</code> Die zusammengefassten Gruppen können allerdings ein <code>Öder"</code> enthalten. <code>Haus(meister tür)</code> ⇒ <i>Hausmeister</i> oder <i>Haustür</i> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COsmAS</p> <p>DWDS</p> <p>CQP</p> |

Übungen

Schreiben Sie Regexe, die folgendes finden:

1. nur: **das** und **welches**
2. nur: **Bus**, **Bussie** und **Busfahrer**
3. **Bus**, **Bussie** und **Busfahrer**, aber nicht **Bushaltestelle**
4. **gegenderte Nomina** wie: **StudentIn**
5. **dass**, im Singular und Plural
6. alle Flexionsformen von **rot**
7. Wörter, die mit **X** oder **x** beginnen
8. Wörter, die mit **Q** beginnen und mit **ant** oder **anten** aufhören

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Wir sind hier...

Korpora und Korpusstudien
Korpusstudien planen
Ablauf der Studie
Kenntnisse
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Logik und Klammerung
Morphologische Annotation
Grenz-Markierungen
DWDS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Morphologische Annotation
CQP
Allgemeines
Installation und Basics
Einfache Suchanfragen
Anzeigoptionen setzen
Suchanfragen speichern und ausgeben
Gruppieren und Sortieren

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Reservierte Zeichen

Folgende Zeichen haben eine besondere Bedeutung in Regexen und müssen geschützt werden, wenn sie literal gemeint sind:

| Zeichen | Bedeutung |
|---------|-------------------------|
| . | beliebiges Zeichen |
| [] | Zeichenklasse |
| * | Quantor: 0–beliebig |
| + | Quantor: 1–beliebig |
| { } | numerischer Quantor |
| \ | Schutz für Literale |
| () | Gruppe |
| | Oder (Gruppe) |
| ^ | Anfang der Zeichenkette |
| \$ | Ende der Zeichenkette |

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Korpusstudien planen
Ablauf der Studie
Kenntnisse
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
Allgemeines
Zugang und Export
Einfache Suchanfragen

Wir sind hier...

Korpora und Korpusstudien
Korpusstudien planen
Ablauf der Studie
Kenntnisse
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Logik und Klammerung
Morphologische Annotation
Grenz-Markierungen
DWDS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Morphologische Annotation
CQP
Allgemeines
Installation und Basics
Einfache Suchanfragen
Anzeigoptionen setzen
Suchanfragen speichern und ausgeben
Gruppieren und Sortieren

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Merkmale und operationalisierung

- ▶ Variablen, die normalerweise von Menschen handkodiert werden sind z. B.:
 - ▶ informationsstrukturelle Kategorien wie Topic
 - ▶ anaphorische Beziehungen
 - ▶ Kasus
 - ▶ Textsorte
- ▶ Alle Variablen, die von Menschen kodiert werden, müssen möglichst so genau operationalisiert werden, dass mehrere Menschen anhand der Operationalisierung die gleichen Entscheidungen treffen würden.
- ▶ Im Idealfall lässt man deshalb zwei oder mehr Kodierer dieselben Belege kodieren und prüft dann, wie gut die Ergebnisse übereinstimmen.

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Wir sind hier...

Korpora und Korpusstudien
Korpusstudien planen
Ablauf der Studie
Kenntnisse
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Logik und Klammerung
Morphologische Annotation
Grenz-Markierungen
DWDS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Morphologische Annotation
CQP
Allgemeines
Installation und Basics
Einfache Suchanfragen
Anzeigoptionen setzen
Suchanfragen speichern und ausgeben
Gruppieren und Sortieren

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Zwei Kodierer

- ▶ Wenn mehrere Kodierer einfach irgendetwas anklicken, erwartet man trotz dem zufällige Übereinstimmung.
- ▶ Wie groß diese zufällige Übereinstimmung ist, hängt neben der Anzahl der Kodierer (hier jetzt erstmal immer 2) von der Anzahl der möglichen Ausprägungen der kodierten Variable ab.
- ▶ Cohens κ ist ein Maß dafür, um wieviel besser eine gegebene Übereinstimmung gegenüber der per Zufall erwartbaren Übereinstimmung ist.

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Kreuztabulierung der Ergebnisse

- ▶ Für das DECOW2012 haben zwei Kodierer eine Zufallsstichprobe von $n = 200$ Texten unter anderem für Zielpublikum nach den COWCat-Richtlinien kodiert.
- ▶ Die möglichen Werte waren: **allgemein** – **informiert** – **professionell**
- ▶ In der Kreuztabelle trägt man jetzt pro Fall ein, wie KodiererIn 1 (Zeilen) abhängig von KodiererIn 2 (Spalten) entschieden hat.
- ▶ Jede Abweichung von der Diagonalen markiert Nicht-Übereinstimmung.

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Kreuztabelle

| K1↓ K2→ | allgemein | informiert | professionell | Summen |
|---------------|-----------|------------|---------------|--------|
| allgemein | 151 | 14 | 4 | 169 |
| informiert | | 20 | 6 | 26 |
| professionell | | | 5 | 5 |
| Summen | 151 | 34 | 15 | 200 |

Wie man sieht tendiert K1 (salopp gesagt) dazu, das Textniveau geringer einzuschätzen und es dadurch einem weniger informierten Publikum zuzutrauen.

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Einfache Übereinstimmung

| K1↓ K2→ | allgemein | informiert | professionell | Summen |
|---------------|-----------|------------|---------------|--------|
| allgemein | 151 | 14 | 4 | 169 |
| informiert | | 20 | 6 | 26 |
| professionell | | | 5 | 5 |
| Summen | 151 | 34 | 15 | 200 |

Der einfache Anteil der Übereinstimmung ist:

$$p_0 = \frac{\text{Diagonalsumme}}{n}$$

$$\text{Bsp. } p_0 = \frac{151+20+5}{200} = \frac{176}{200} = 0.88$$

Kurzfasste Korpuslinguistik für GermanistInnen
Roland Schäfer

Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Zufällige Übereinstimmung

| K1\K2→ | allgemein | informiert | professionell | Summen |
|---------------|-----------|------------|---------------|--------|
| allgemein | 151 | 14 | 4 | 169 |
| informiert | | 20 | 6 | 26 |
| professionell | | | 5 | 5 |
| Summen | 151 | 34 | 15 | 200 |

Wenn wir die Spaltensumme in Spalte j mit SS_j und die Zeilensumme in Zeile i mit ZS_i bezeichnen, ist die zufällig erwartbare Übereinstimmung:

$$p_c = \frac{\sum_i SS_i \cdot ZS_i}{n^2}$$

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Zufällige Übereinstimmung: Beispiel

| K1\K2→ | allgemein | informiert | professionell | Summen |
|---------------|-----------|------------|---------------|--------|
| allgemein | 151 | 14 | 4 | 169 |
| informiert | | 20 | 6 | 26 |
| professionell | | | 5 | 5 |
| Summen | 151 | 34 | 15 | 200 |

$$p_c = \frac{169 \cdot 151 + 26 \cdot 34 + 5 \cdot 15}{200^2} = \frac{25519 + 884 + 75}{200^2} = \frac{26478}{40000} = 0.66$$

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Cohens κ ausrechnen

Die Formel lautet:

$$\kappa_c = \frac{p_0 - p_c}{1 - p_c}$$

Bsp.: $\kappa_c = \frac{0.88 - 0.66}{1 - 0.66} = \frac{0.22}{0.34} = 0.65$

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Interpretation von κ_c

- Wie immer ist die Interpretation eines Gütwerts keine vollständig normierte Sache.
- Ich würde den Wert prinzipiell einfach als Wert angeben.
- Die Faustregeln für die Praxis vielleicht:
 - $\kappa_c < 0.5$: Operationalisierung überprüfen und nochmal kodieren (lassen).
 - $\kappa_c > 0.7$: Ggf. im Text ausdrücklich auf die gute Übereinstimmung hinweisen.

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Wir sind hier...

Korpora und Korpusstudien
Korpusstudien planen
Ablauf der Studie
Kenntnisse
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
Allgemeines
Zugang und Export
Einfache Suchanfragen

Abstände
Logik und Klammerung
Morphologische Annotation
Grenz-Markierungen
DWDS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Morphologische Annotation
CQP
Allgemeines
Installation und Basics
Einfache Suchanfragen
Anzeigooptionen setzen
Suchanfragen speichern und ausgeben
Gruppieren und Sortieren

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Mehr als zwei Kodierer

- Wenn man mehr als zwei (also m) Kodierer hat, nimmt man Fleiss' κ_f .
- p_0 muss über m Dimensionen gerechnet werden, es gibt also mehr als nur Spalten und Zeilen.

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

In R

Für die Berechnung muss meist eine Matrix oder ein Dataframe vorliegen, der die Ergebnisse folgendermaßen tabuliert:

| Fall | K1 | K2 | ... | Km |
|------|---------------|------------|-----|-----|
| 1 | informiert | informiert | ... | ... |
| 2 | informiert | allgemein | ... | ... |
| 3 | professionell | informiert | ... | ... |
| n | ... | ... | ... | ... |

Wenn so eine Tabelle im Objekt kodierungen vorliegt:
> library(irr)
> kappa2(kodierungen)

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Andere Maße, andere Daten

- Hier wurde implizit von Nominaldaten ausgegangen.
- Es gibt andere Maße für Nominaldaten und andere Maße für andere Skalenniveaus.

Für einen Überblick die Liste aus irr anfordern:
> irr::<TAB>
Dann z.B.:
> ?kendall

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

Wir sind hier...

Korpora und Korpusstudien
Korpusstudien planen
Ablauf der Studie
Kenntnisse
Reguläre Ausdrücke
Grundlagen
Wiederholungen von Zeichen
Gruppen und Nicht-Zeichen
Liste der reservierten Zeichen
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
Allgemeines
Zugang und Export
Einfache Suchanfragen

Abstände
Logik und Klammerung
Morphologische Annotation
Grenz-Markierungen
DWDS
Allgemeines
Zugang und Export
Einfache Suchanfragen
Abstände
Morphologische Annotation
CQP
Allgemeines
Installation und Basics
Einfache Suchanfragen
Anzeigooptionen setzen
Suchanfragen speichern und ausgeben
Gruppieren und Sortieren

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

COSMAS

Vorteile:

- relativ groß
- nach Jahrgängen sortiert, immer aktuell

Nachteile:

- fast nur Zeitungssprache
- viele Regionalblätter, die die Stichproben dominieren (das berühmte *St. Galler Tageblatt*)

Momentan größter Nachteil: Obsoletheit... Das neue ganz große Ding kommt angeblich bald.

Kurzfasste Korpuslinguistik für GermanistInnen

Roland Schäfer
Korpora und Korpusstudien
Reguläre Ausdrücke
Inter-Kodierer-Übereinstimmung
Motivation
Cohens κ
Fleiss' κ
COSMAS
DWDS
CQP

| | | | |
|---|--|---|--|
| <h2>Hinweis zum Namen</h2> <p>COSMAS ist der Name der Software.</p> <p>Die IDS-Korpora pauschal als „COSMAS“ zu zitieren ist unprofessionell und unfair ggü. den Personen, die die Korpora machen, aber mit der Software COSMAS nichts zu tun haben!</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Zugang und Export</h2> <ul style="list-style-type: none"> Für das rein technische gibt es ein Videotutorial: http://www.youtube.com/watch?v=EfxSegFzPZI Wichtig: Die Konkordanzen dürfen nicht dauerhaft gespeichert werden, bitte sorgfältig die Nutzungsbedingungen lesen. Der Export kann mit Koka für eine Tabellenkalkulation vorbereitet werden: http://hpsg.fu-berlin.de/~rsling/microsites/koka/ | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Wortformen(ketten) suchen</h2> <p>Einfache Wortform suchen:</p> <ul style="list-style-type: none"> > Chuzpe <p>Kette von aufeinanderfolgenden Wortformen:</p> <p>(mit Veggelassener Verknüpfungoperator bedeutet: Wortabstand /#1)</p> <ul style="list-style-type: none"> > hat Chuzpe | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <h2>Lemma-Suche und Wortformenexpansion</h2> <p>Alle Formen einer Grundform:</p> <ul style="list-style-type: none"> > &Haus &kaufen <p>Die Expansionsliste</p> <p>Man kann die konkrete Liste der orthographischen Varianten bzw. Wortform einsehen und bearbeiten, wenn man vor dem Absetzen der eigentlichen Anfrage auf Liste zu: ... [Öffnen] klickt.</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <h2>Wortabstände</h2> <p>Beispiel: Man sucht nach einem Artikel gefolgt von maximal einem Adjektiv, dann einem Nomen.</p> <ul style="list-style-type: none"> > &ein /+w2 &Baum <p>Optimaler, weil dann das seltenere Wort zuerst gesucht wird:</p> <ul style="list-style-type: none"> > &Baum /-w2 &ein | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> |
| <h2>Satzabstand und im selben Satz suchen</h2> <p>Im selben Satz (ungerichteter Abstandoperator):</p> <ul style="list-style-type: none"> > &Haus /s0 &Hof <p>In aufeinanderfolgenden Sätzen (...wirklich?):</p> <ul style="list-style-type: none"> > &Haus /+s1 &Hof | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> | <h2>Abstandsbereiche</h2> <p>Die Abstandoperatoren können alle ungerichtet (ohne + oder -) angegeben werden, und sie können Bereiche als Argument nehmen.</p> <p>Simulation eines Nahe-Bei-Operators:</p> <ul style="list-style-type: none"> > Chuzpe /w4 Verve <p>Lösung des Problems mit aufeinanderfolgenden Sätzen – Bereiche:</p> <ul style="list-style-type: none"> > &Haus /+s1:1 &Hof <p>Also: Einzelne Zahlen n bezeichnen den Abstand bis maximal n.</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>CQP</p> |

| | | | |
|--|--|---|---|
| <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation</p> <p>CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigooptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Logik-Operatoren UND ODER NICHT</p> <p>Die Operatoren beziehen sich auf ein ganzes Dokument! Man muss dabei die Argumente des Operators als Bedingungen an das Dokument lesen („muss enthalten“).</p> <p>> Chuzpe UND Verve</p> <p>Achtung: ODER ist inklusiv. > &Kanzlerin ODER &Präsidentin</p> <p>Achtung: NICHT ist eigentlich ein logisches UND NICHT. > &Kanzlerin NICHT &Präsidentin</p> <p>(Die letzte Anfrage liefert eine Teilmenge der vorletzten.)</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> |
|--|--|---|---|

| | | | |
|---|---|---|---|
| <p>Klammern</p> <p>Mit Klammern muss man mehrere Abstandsoperatoren kombinieren. > (&können /w1 &wollen) /w4 &haben</p> <p>Vor allem ODER wird in Klammern aber eine Art echtes ODER: > (&Bad ODER &Küche) /+w1 &schrubben</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> | <p>Objektsprachliches UND ODER NICHT</p> <p>Die Operatoren müssen als normale Wörter in Anführungszeichen:</p> <ul style="list-style-type: none"> > "und" > "oder" > "nicht" | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> |
|---|---|---|---|

| | | | |
|--|---|--|---|
| <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> | <p>MORPH-Operator</p> <p>Mit dem MORPH()-Operator kann man auf POS-Tags zugreifen.</p> <p>Achtung Wie immer funktioniert es nicht so, wie die Hilfe behauptet. Vergesst STTS, nimmt einfach den MORPH-Assistenten!</p> <p>> MORPH (VRB inf v) /+w1:1 MORPH (VRB inf m) /+w1:1 MORPH (VRB fin a)</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> |
|--|---|--|---|

| | | | |
|--|---|--|---|
| <p>Beschränkungen auf Wortform und Tag</p> <p>Will man gleichzeitig Lemma/Wortform und morphologische Information einschränken, muss man den 0-Wörter-Abstand-Trick benutzen:</p> <p>> MORPH (VRB fin v) /w0 &laufen</p> <p>Man sucht also nach einer finiten Verbform und im Abstand von 0 Wörtern einer Form von <i>laufen</i>.</p> <p>Wie sucht man nach einer solchen Form im Abstand von maximal vier Wörtern nach links oder rechts zu einer nominalen Form von <i>Lauf</i>?</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> | <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation</p> <p>CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigooptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> |
|--|---|--|---|

| | | | |
|--|---|--|---|
| <p>Satz- und Absatz-Enden und -Anfänge</p> <p>Es gibt die folgenden strukturellen Attribute:</p> <ul style="list-style-type: none"> > <sa> Satzanfang > <se> Satzende > <pa> Absatzanfang > <pe> Absatzende <p>Auch sie sucht man mit dem 0-Wörter-Trick (z. B. VL-Sätze): > MORPH (VRB fin) /w0 <se></p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS CQP</p> | <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen</p> <p>DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation</p> <p>CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigooptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Logik und Klammerung Morphologische Annotation</p> <p>DWDS CQP</p> |
|--|---|--|---|

| | | | |
|--|---|---|--|
| <h2>DWDS-Kernkorpus</h2> <p>Vorteile:</p> <ul style="list-style-type: none"> ▶ balanciertes Referenzkorpus ▶ deckt das 20. Jh. gleichmäßig ab ▶ Zusatzinformationen: Kollokationen, Wortverlauf, ... <p>Nachteile:</p> <ul style="list-style-type: none"> ▶ viel zu klein ▶ einige Strata sind verzerrt (z. B. 80er Jahre Wissenschaft = Slotterdijk) ▶ deckt per Definition nur das 20. Jh. ab ▶ zu viele Treffer nicht sichtbar ▶ schlechte Exportoptionen <p>▶ Hauptnachteil: das neue Interface</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> <p>Allgemeines</p> <p>Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> |
|--|---|---|--|

| | | | |
|--|--|---|--|
| <h2>Zugang und Export</h2> <ul style="list-style-type: none"> ▶ Für das rein technische gibt es ein Videotutorial: http://www.youtube.com/watch?v=6pTW26PJeIc ▶ Der Export kann mit Koka für eine Tabellenkalkulation vorbereitet werden: http://hpsg.fu-berlin.de/~rsling/microsites/koka/ ▶ Die Anfragesprache ist einfacher, aber unterm Strich m. E. zielführender als mit COSMAS. | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> <p>Allgemeines</p> <p>Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> |
|--|--|---|--|

| | | | |
|--|--|---|--|
| <h2>Wortformen und Grundformen</h2> <p>Im Gegensatz zu COSMAS sucht DWDS immer auch nach Formvarianten!</p> <p>> Haus</p> <p>...sucht also auch <i>Hause, Häuser</i> usw.</p> <p>Man muss umgekehrt den Wortform-Operator einsetzen:</p> <p>> @Haus</p> <p>...sucht definitiv nur <i>Haus</i>.</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> | <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> <p>Allgemeines</p> <p>Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> |
|--|--|---|--|

| | | | |
|--|--|--|--|
| <h2>Wortabstand #</h2> <p>Alle Anfragen beziehen sich prinzipiell auf innerhalb eines Satzes.</p> <p>Alle Sequenzanfragen müssen in " stehen:</p> <p>> "Haus bauen"</p> <p>Der einfache Abstand muss nicht markiert werden, alles darüber hinaus mit #:</p> <p>> "der #1 Präsident"</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> | <h2>Bedingungen innerhalb eines Satzes</h2> <p>Der Satz ist sowieso die Suchdomäne. Wenn die Linearisierung egal ist, müssen logische Operatoren verwendet werden:</p> <ul style="list-style-type: none"> ▶ !Krieg (Satz enthält nicht <i>Krieg</i>) ▶ Frieden && Krieg (Satz enthält <i>Frieden</i> und <i>Krieg</i>) ▶ Frieden Krieg (Satz enthält <i>Frieden</i> und/oder <i>Krieg</i>) ▶ Frieden && !Krieg (Satz enthält <i>Frieden</i> und nicht <i>Krieg</i>) | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> |
|--|--|--|--|

| | | | |
|---|--|--|--|
| <h2>Wir sind hier...</h2> <p>Korpora und Korpusstudien</p> <p>Korpusstudien planen</p> <p>Ablauf der Studie</p> <p>Kenntnisse</p> <p>Reguläre Ausdrücke</p> <p>Grundlagen</p> <p>Wiederholungen von Zeichen</p> <p>Gruppen und Nicht-Zeichen</p> <p>Liste der reservierten Zeichen</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>Motivation</p> <p>Cohens κ</p> <p>Fleiss' κ</p> <p>COSMAS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Logik und Klammerung</p> <p>Morphologische Annotation</p> <p>Grenz-Markierungen</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> | <h2>Zugriff auf STTS-Tags</h2> <p>Suche alle infiniten Modalverben ein oder zwei Wörter nach <i>haben</i>:</p> <p>> "haben #2 \$p=VMINF"</p> <p>Kombinieren von Wortform und Tag-Information mit with...</p> <p>Suche alle finiten Formen von <i>haben</i> als Vollverb:</p> <p>> haben with \$p=VFFIN</p> <p>Kombiniert mit Abstand:</p> <p>"haben with \$p=VAFIN #4 Eis"</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>Allgemeines Zugang und Export</p> <p>Einfache Suchanfragen</p> <p>Abstände</p> <p>Morphologische Annotation</p> <p>CQP</p> |
|---|--|--|--|

| | | | |
|---|---|--|---|
| <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>CWB (bzw. CQP)</p> <p>Vorteile:</p> <ul style="list-style-type: none"> ▶ nicht an bestimmte Korpora gebunden ▶ auch für eigene Korpora ▶ sehr mächtig und flexibel ▶ eingebaute Statistikfunktionen <p>Nachteile:</p> <ul style="list-style-type: none"> ▶ Rechenzeit- und Speichergierig ▶ Beschaffung der Korpora | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>Von uns verfügbare Korpora</p> <ul style="list-style-type: none"> ▶ COW2012 Token-Größen: <ul style="list-style-type: none"> ▶ DE: 9 Mrd. ▶ ES: 1,5 Mrd. ▶ NL: 2,3 Mrd. ▶ SE: 2,2 Mrd. ▶ UK: 1 Mrd. ▶ COW2014: Alle über 10 Mrd, Deutsch wahrscheinlich über 20 Mrd. ▶ http://corporafromtheweb.org/ | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>Installation und Basics</p> <ul style="list-style-type: none"> ▶ Video-Tutorial zur Installation (Win): http://www.youtube.com/watch?v=5xR3QikalcS ▶ Video-Tutorial zur Korpusinstallation: http://www.youtube.com/watch?v=Ku5Qvd6pax8&hd=1 ▶ Video-Tutorial zur Icon-Installation/Grafikanpassung (Win): http://www.youtube.com/watch?v=eNL_KI3DAkI&hd=1 | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Das schwarze Fenster mit dem Text</p> <ul style="list-style-type: none"> ▶ CQP funktioniert nach dem Textkonsolen/Shell-Prinzip. ▶ Man schreibt einen Befehl, drückt RETURN... ▶ ... und das System regiert mit einer Ausgabe. ▶ Pfeil-hoch und Pfeil-runter navigieren durch bereits getippte Kommandos. | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>Hinweis zu rcqp</p> <p>Eine sehr praktische Art, auf CQP zuzugreifen, ist die Schnittstelle rcqp für R.</p> <p>Von den hier beschriebenen CQP-Methoden ist nur die Syntax der Suchanfragen in rcqp gleich. Korpusauswahl, KWIC-Anzeige usw. gehen grundlegend anders.</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>Korpuswahl und Basisbefehle</p> <p>Verfügbare Korpora anzeigen: > show corpora;</p> <p>Informationen über ein Korpus abrufen: > info DECOW2012-C00X1M;</p> <p>Shell verlassen: > exit;</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Suchanfragen</p> <ul style="list-style-type: none"> ▶ Neben einigen Steuerbefehlen ist die typische Eingabe in CQP eine Abfrage in dem Korpus. ▶ Man kann dabei nur abfragen: <ul style="list-style-type: none"> ▶ Ketten von Tokens und ▶ innerhalb bestimmter struktureller Attribute (Sätze, Absätze, Dokumente usw.). ▶ Die Reaktion der Shell ist die Ausgabe der Suchanfrage als KWIC. | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |

| | | | |
|---|--|---|--|
| <h2>Wortformen direkt oder mit regulären Ausdrücken</h2> <p>Einfache Wortform(enkette) abfragen, jedes Wort in ":</p> <ul style="list-style-type: none"> > "zu"; > "zu" "versuchen"; <p>Mit regulären Ausdrücken:</p> <ul style="list-style-type: none"> > ".+lein"; > ".+lein."; | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <h2>Volle reguläre Ausdruckssyntax</h2> <p>Character classes:</p> <ul style="list-style-type: none"> > "[MN].+lein" <p>Alternativen:</p> <ul style="list-style-type: none"> > "Mann" "und oder" "Frau" <p>Quantoren:</p> <ul style="list-style-type: none"> > "[a-z]{2,3}" | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |
| <h2>Übungen – bisher nicht unbedingt perfekt lösbar</h2> <ol style="list-style-type: none"> Alle Flexionsformen von <i>rot</i>. Alle Indikativ-Formen von <i>kaufen</i>. Alle Indikativformen von <i>laufen</i>. Dative mit und ohne Monoflexion der Adjektive. | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <h2>Voller Zugriff auf Wortinformationen</h2> <ul style="list-style-type: none"> Wenn mehr als nur die Wortform eingeschränkt werden soll, muss jedes Such-Token in [] gesetzt werden. In den eckigen Klammern werden dann mit & die Suchkriterien für ein Token verbunden. Meist: <ul style="list-style-type: none"> word= Wortform (wie bei der einfachen Suchanfrage). lemma= Lemma. pos= Wortart. Die einfache Suchanfrage mit ist eigtl. nur eine Abkürzung für [word=]. | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |
| <h2>Kombinationen, Beispiele</h2> <p>Finite Formen von <i>haben</i> als Auxiliar:</p> <ul style="list-style-type: none"> > [lemma="haben"& pos="VAFIN"]; <p>ne als Artikel:</p> <ul style="list-style-type: none"> > [word="ne"& pos="ART"]; <p>-lein-Diminutive:</p> <ul style="list-style-type: none"> > [pos="NN"& lemma=".+lein"]; | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <h2>Abstände, endlich konsistent</h2> <p>[] ist das <i>matchall</i>-Token und kann daher für die Festlegung von Abständen benutzt werden:</p> <ul style="list-style-type: none"> > [pos="ART"] []{0,4} [pos="NN"]; | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |
| <h2>Übungen</h2> <ul style="list-style-type: none"> Jetzt nochmal die Monoflexion bitte. Massangaben im Stil von <i>eine Tasse heißer/heißen Kaffee/s</i>. | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <h2>Logische Operatoren</h2> <ul style="list-style-type: none"> In der [] können außer & auch andere logische Operatoren verwendet werden: <ul style="list-style-type: none"> oder nicht ! <p>Finite Formen von <i>gehen</i> außer Präteritum:</p> <ul style="list-style-type: none"> > [lemma="gehen"& pos="VFIN"& !word="ging.*"]; | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |
| <h2>Strukturelle Attribute</h2> <ul style="list-style-type: none"> Auf strukturelle Attribute kann man entweder in <> direkt zugreifen... oder mit <i>within</i>. <p>Satzinitiale Adverbien:</p> <ul style="list-style-type: none"> > <s> [pos="ADV"]; <p>Satzfinale Adverbien:</p> <ul style="list-style-type: none"> > [pos="ADV"] </s>; <p>Innerhalb eines Satzes:</p> <ul style="list-style-type: none"> > [lemma="Fortuna"] [lemma="Düsseldorf"] [] * [lemma="gewinnen"]; | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <h2>Übungen</h2> <ul style="list-style-type: none"> Satzinitiale Adverbien alleine im Vorfeld. Sätze, die auf einen Nebensatz enden. W-Exklamativa. | <p>Kurzfasste Korpuslinguistik für GermanistInnen</p> <p>Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |

| | | | |
|---|---|--|---|
| <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kontext</p> <p>Kontext ist genau 1 Satz: > set Context 1 s;</p> <p>Kontext sind 20 Zeichen links und rechts: > set Context 20;</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>PrintStructures</p> <p>Dokument-URL anschalten: set PrintStructures doc_url; (Nicht definiert in COW-X-Korpora!)</p> <p>Suchen in Dokumentnamen/Verfassernamen etc.:</p> <pre>> [lemma="Raubkopie"] :: match.doc_url = "http.+heise\\.de.+";</pre> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Cohens κ Fleiss' κ COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>Suchanfragen benennen</p> <p>Eigentlich sollte man Suchanfragen immer speichern und dann erst mit cat ausgeben: > myQuery = [lemma="Chuzpe"]; cat myQuery;</p> <p>... macht äußerlich dasselbe wie: > [lemma="Chuzpe"];</p> <p>Suchanfrage persisten machen: > save myQuery;</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Last</p> <ul style="list-style-type: none"> Wenn die letzte Suchanfrage nicht explizit gespeichert wurde, ist sie als Last verfügbar. Man kann sie speichern in eine eigens benannte Suchanfrage: myNewQuery = Last; | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>Suchanfrage exportieren</p> <pre>> myQuery > set Context 1 s; > set PrintStructures doc_url; > myQuery = "Chuzpe"; > cat myQuery > "dateiname.txt";</pre> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Übungen</p> <ul style="list-style-type: none"> Zwei der gemachten Übungen nochmal machen, aber mit Speichern, catten und dann exportieren. Export durch Koka jagen, Konkordanz in LibreOffice formatieren. | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |
| <p>Suchanfragen reduzieren</p> <p>Wenn man sehr allgemeine Suchanfragen abfragt: > [pos="NN"]; kann es sein, dass man nur eine Zufallsauswahl möchte.</p> <p>Trefferanzahl abfragen: > size myQuery;</p> <p>Auf eine bestimmte Zufallstrefferzahl beschränken: > reduce Last to 100; oder > reduce Last to 5%;</p> | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Übungen</p> <p>Bitte bis zur Konkordanz bringen:</p> <ul style="list-style-type: none"> Monoflexion als zwei 100er-Sample (Mono/nicht Mono). 10% aller <i>lein</i>-Diminutive im Korpus. | <p>Kurzgefaste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien Reguläre Ausdrücke Inter-Kodierer-Übereinstimmung COSMAS DWDS CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> |

| | | | |
|--|--|---|--|
| <p>Wir sind hier...</p> <p>Korpora und Korpusstudien Korpusstudien planen Ablauf der Studie Kenntnisse Reguläre Ausdrücke Grundlagen Wiederholungen von Zeichen Gruppen und Nicht-Zeichen Liste der reservierten Zeichen Inter-Kodierer-Übereinstimmung Motivation Colons & Flies' & COSMAS Allgemeines Zugang und Export Einfache Suchanfragen</p> <p>Abstände Logik und Klammerung Morphologische Annotation Grenz-Markierungen DWDS Allgemeines Zugang und Export Einfache Suchanfragen Abstände Morphologische Annotation CQP Allgemeines Installation und Basics Einfache Suchanfragen Anzeigeoptionen setzen Suchanfragen speichern und ausgeben Gruppieren und Sortieren</p> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>match und matchend</p> <ul style="list-style-type: none"> ▶ Jeder Treffer wird durch den Match charakterisiert: Anfang und Ende des Satzstücks, das der Suchanfrage entsprach. ▶ Das erste Token davon heißt <code>match [0]</code>, das letzte <code>matchend [0]</code>. ▶ Man kann auf Wörter davor oder danach zurückgreifen, indem man von 0 entsprechend Werte abzieht. ▶ Nach diesem Tokens kann man sortieren. | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |
| <p>Sortieren</p> <p>Nach einzelmem Token:</p> <pre>> Bla = "mit" [pos="ADJA"& word=".+em"] [pos="ADJA"] [pos="NN"]; > sort Bla by word on match [2];</pre> <p>Nach einem Bereich:</p> <pre>> sort Bla by word on match [1] .. match [2];</pre> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Gruppieren/Wortverteilungen</p> <p>Mit <code>group</code> kann man Resultate gruppieren und die Zählraten ausgeben lassen.</p> <pre>> mySchnell = [lemma="schnell"] [pos="VVFIN"]; > group mySchnell match [1] lemma;</pre> <p>Man kann nach allem gruppen:</p> <pre>> group mySchnell match [-1] pos;</pre> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |
| <p>Übungen</p> <ul style="list-style-type: none"> ▶ Vergleiche die Verteilung der finiten satzfinalen Verb-Arten (voll, aux, modal) mit der der nicht-satzfinalen. ▶ Später: Rechne einen χ^2-Test dazu! ▶ Finde heraus, mit welchen Wörtern (möglichst Satzzeichen und Zahlen und kreativen Sprachgebrauch ausschließen) <i>statt</i> in welcher Häufigkeit rechts-kombiniert (Lemmata und POS). | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | <p>Mehrdimensionale Gruppierung</p> <p>Man kann an die erste Grupierung eine zweite anhängen:</p> <pre>> myAdjKonj = [pos="ADJA"] "und" [pos="ADJA"]; > group myAdjKonj match lemma by matchend lemma;</pre> | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> |
| <p>Übungen</p> <ul style="list-style-type: none"> ▶ Finde die häufigsten Adjektiv-Nomen-Bigramme. ▶ Finde die Häufigsten Bigramme aus infinitem Vollverb und infinitem Modalverb. | <p>Kurzfasste Korpuslinguistik für GermanistInnen Roland Schäfer</p> <p>Korpora und Korpusstudien</p> <p>Reguläre Ausdrücke</p> <p>Inter-Kodierer-Übereinstimmung</p> <p>COSMAS</p> <p>DWDS</p> <p>CQP</p> <p>Allgemeines Installation und Basics</p> <p>Einfache Suchanfragen</p> <p>Anzeigeoptionen setzen</p> <p>Suchanfragen speichern und ausgeben</p> <p>Gruppieren und Sortieren</p> | | |