

Mathematische Grundlagen III

Korpora und Sprachressourcen

Garance PARIS

26. April 2011

Überblick

- Korpora (Singular Korpus, neutrum!):
Große Mengen von Text
- Wörterbücher, Lexika, Thesauri,
manche Enzyklopädien
- Ontologien, semantische Netze und
sonstige Formen von Wissensrepräsentation

Definition

- *Ein Korpus (n.!) ist eine endliche Sammlung von konkreten sprachlichen Äußerungen, die als Grundlage für sprachwissenschaftliche Untersuchungen dienen (Lexikon der Sprachwissenschaft)*
- *Eine idealerweise repräsentative, möglicherweise auf einem Bereich eingeschränkte Sammlung von Texten einer gegebenen Sprache, die zum Zwecke linguistischer Analyse zusammengestellt wurde (Francis, 1964)*

(Francis and Kučera: Ersteller des Brown Corpus, eins der ersten Korpora fürs Englische)

Definition

- Die Daten sollten aus möglichst natürlichen Gesprächssituationen stammen
- Meist ist es wünschenswert, sich auf Muttersprachler zu begrenzen
- Sie sollten repräsentativ sein, wobei sich die Frage stellt, was denn repräsentativ ist!
- Heutzutage sind Korpora nur noch elektronisch vorstellbar
- Sie werden manchmal für einen bestimmten Zweck erhoben (z. B. Fehleranalyse bei Sprachlernern im CALL-Bereich)

- Die größten Korpora sind rohe Korpora (heute: das Internet selbst)
- Einsatz in der Lexikographie:
Manuelle Sichtung Beispiele (Konkordanz),
um Wortbedeutungen zu bestimmen,
sowie Neologismen und Kollokationen zu entdecken

| |
|--|
| hängen , Packpferde mit Brennholz ; Frauen backen Brot , Kinder hüten Ziegen . Von Zeit zu Zeit |
| unmusikalisch . Aber sie kann Pfannkuchen Brot , backen Nun folgt die konkrete Utopie (oder was m |
| Bei 170 Grad , Gas : Stufe 3 etwa 1 1/4 Std. backen . Vor dem Herausnehmen erkalten lassen . R |
| Leute an . Laßt uns anfangen , ich muß Brot " backen , meinte er unwirsch und genehmigt sich un |
| kann doch nicht jeder seine eigenen Brötchen . backen , mahnte Scherf . Dann wieder Fragen : Ob |
| e , und die zieht er formvollendet durch : Wir backen einen guten Kurzfilm . An der Idee blieb auc |
| ssen . In heißem Backfett kleine Pfannkuchen backen und mit saurer Sahne und Kaviar servieren . |
| zu besticken , Kaffee zu kochen und Kekse zu backen , um so ihrer Verpflichtung gegenüber dem |
| . Im Moment aber muß er ganz kleine Brötchen backen . Der Grüne sieht sich einer erdrückenden sc |
| , 1/2 Stunde ziehen lassen , dann goldbraun backen . Mit Erdbeeren garnieren . Alle Rezepte aus |
| Halloween höhlen sie einen Kürbis aus und " backen Pumpkin-Pie . Die Prices sind eine durchsch |
| ade oder Quark . Schwaben südlich der Donau backen Brot , wie die riesigen Knauzawecka , noch |
| schwimmen gehen , nachtwandern , Stockbrot backen , die Bauern besuchen , basteln , spielen . Bei |

- Alte Korpora: Ad-hoc Format
- Interlinear format (hier: Wort_PoS_Lemma):
John_PN_john left_VBP_leave ._PUNC_period
- Spalten (Susanne, 1. Spalte: Satz- und Wort-Id)

| | | | |
|----------|------|--------|------|
| A12:0210 | John | john | PN |
| A12:0211 | left | leave | VBN |
| A12:0212 | . | Period | PUNC |
- SGML Mark-up (veraltet, Vorgänger von XML)
- Heute: meist XML (Vorteil: Allgemeine Tools)

```
<s n=0001>
<w NN1>INTRODUCTION
</head>
<p>
<s n=0002>
<w AT0>The <w AJ0>extensive <w NN1>upland <w NN2-VVZ>landscapes <w PRF>of
<w AT0>the <w NPO>UK<c PUN>, <w CJC>and <w AT0>the <w AJ0>varied <w CJC>and
<w AJ0>rich <w NN1>wildlife <w PNP>they <w VVB>support<c PUN>, <w VBB>are <w AT0>the
<w NN1>product <w PRF>of <w NN2>centuries <w PRF>of <w AV0>predominantly
<w AJ0>pastoral <w AJ0-NN1>agricultural <w NN1>activity<c PUN>.
<s n=0003>
<w PRP>In <w AT0>the <w AJ0-NN1>past<c PUN>, <w AT0>the <w NN1>use <w PRF>of
<w DT0>these <w NN2>uplands <w PRP>for <w NNO>sheep <w CJC>and <w NN1>beef
<w NN2>cattle <w NN1-VVG>rearing <w VHZ>has <w XX0>not <w VVN>conflicted
<w AV0>significantly <w PRP>with <w AT0>the <w NN1>need <w TOO>to <w VVI>retain
<w NN2>habitats <w PRP>such as <w NN2>moorlands<c PUN>, <w NN1>hill
<w NN2>grasslands<c PUN>, <w AJ0>high <w NN1>altitude <w AJ0>montane
<w NN1>vegetation<c PUN>, <w AJ0-VVD>enclosed <w NN2>pastures <w CJC>and
<w NN1-VVB>hay <w NN2>meadows<c PUN>, <w NN2>wetlands <w CJC>and <w AJ0>native
<w NN2>woodlands<c PUN>, <w DTQ>which <w VVB>form <w AT0>the <w NN1>basis <w PRF>of
<w AT0>the <w NN1>nature <w NN1>conservation <w NN1>interest <w PRF>of <w AT0>the
<w CRD>9.68 <w CRD>million <w NN2>hectares <w PRF>of <w NN1>upland <w PRP>in
<w AT0>the <w NPO>UK<c PUN>.
```

- Roher Text genügt oft nicht, daher wird es ergänzt um Satzgrenzen, Wortkategorien,...
- Korpus-Annotation macht enthaltene linguistische Information explizit, kann aber falsch sein

Prinzipien für Korpus-Annotation (Leech, '93)

- Sowohl Annotation als originales (rohes) Korpus sollte von einander trennbar sein
- Die Annotation sollte Theorieunabhängig und neutral sein
- Die Annotationsmethode (manuell, maschinell, oder Kombination davon) sollte bekannt sein
- Die Annotationsrichtlinien sollten mit allen Details verfügbar sein

Annotation: Hinzufügen linguistischer Information

Probleme:

- Welcher Tagset, welcher (Grammatik-) Formalismus, ...?
- Interpretation (HPSG/LFG vs. funktionale Grammatik)
- Wegen Ambiguität ist Annotation nicht einfach
- Es ist aufwendig, da Handarbeit
 - Annotationsaufwand für ein Wort: 30 Sekunden
 - 1M Worte: 500 000 Minuten = 5 Jahre
 - Plus Aufwand fuer Qualitätssicherung
- Außerdem ist es auch fehlerbehaftet und möglicherweise inkonsistent
- Aber automatische Annotation ist nicht 100 % zuverlässig und macht systematische Fehler

- Annotationsmöglichkeiten gering halten (z. B. kleines Tagset), um schwierige Entscheidungen aus dem Weg zu gehen
- Bei Unsicherheit mehrere Tags zuweisen (dem User wissen lassen, dass es Unsicherheit gab)
z. B.: “Ambiguity Tags” im BNC
AJ0-AV0 (Adjectiv oder Adverb), mit Präferenz für AJ0
- Automatische Annotation mit der Überprüfung durch menschliche Annotatoren kombinieren
- Bootstrapping

Merkmale von Korpora

- Sprache: monolingual vs. bilingual vs. multilingual; vergleichbar vs. parallel, aligniert
- Textart, Inhalt, Genre, Domäne:
 - Spontansprache: Usenet, Wizard-of-Oz Experimente
 - Editiert: Zeitungsartikel, Romane, Fachtexte, Lyrik,...
 - Ausgewogenheit:
homogen vs. heterogen, unbalanciert vs. balanciert
- Geschriebene Sprache vs. gesprochene Sprache
- Umfang (Tokens, Types), Zeitraum
- Format (s. oben), Text oder Binär (indexiert)
- Medium (Text, Audio, Transkripte, Video, usw.)
- Aufbereitung und Annotation
- Urheber- und Nutzungsrechte, Preis
- Standard-Referenz: Allgemeine Verfügbarkeit

Korpora mit Wortarten

- Standardkorpora:
 - British National Corpus (BNC), 100M Worte
 - American National Corpus (ANC), 22M Worte
 - Huge German Corpus (HGC), 200M Worte
- Einsatz: Training von Taggern

Annotation: PoS-Tagging

- Manuell oder automatisch?
- Tag Sets sind unterschiedlich groß; sie variieren in sowohl innerhalb als auch unter Kategorien in ihrer Granularität

| | Brown | Penn | Claws 1-8 | STTS |
|-------|--------|------|-----------|------|
| Größe | 77/177 | 45 | 60-160 | 54 |

- Sie sind sprachspezifisch
- Manchmal enthalten sie Seltsamkeiten
 - Brown:
VBG für Present Participles und für Gerunde
John is purchasing apples
The Fulton County purchasing department
 - Penn:
TO sowohl für Präpositionen als auch vor Infinitiven
(I want to go to the store)

| | | | | | |
|------|-------------------------------|--------|---------------------------------------|-----|------------------------------------|
| - | dash | EX | existential there | | |
| , | comma | FW | foreign word | | |
| : | colon | HV | have | | |
| . | sentence closer (. ; ? *) | HVD | had (past tense) | | |
| (| left paren | HVG | having | QL | qualifier (very, fairly) |
|) | right paren | HVN | had (past participle) | QLP | post-qualifier (enough, indeed) |
| * | not, n't | HVZ | have, pres., 3rd p. sg. | | |
| ABL | pre-qualifier (quite, rather) | IN | preposition | RB | adverb |
| ABN | pre-quantifier (half, all) | JJ | adjective | RBR | comparative adverb |
| ABX | pre-quantifier (both) | JJR | comparative adjective | RBT | superlative adverb |
| AP | post-determiner | JJS | semantic superl. adj. (chief, top) | RN | nominal adverb (here, indoors) |
| AT | article (a, the, no) | | | RP | particle (about, off, up) |
| BE | be | JJT | superlative adjective | TO | to (before infinitive) |
| BED | were | MD | modal auxiliary | UH | interjection |
| BEDZ | was | NC | cited word | VB | verb, base form |
| BEG | being | NN | singular or mass noun | VBD | verb, past tense |
| BEM | am | NNS | plural noun | VBG | pres. part./gerund |
| BEN | been | NP | proper noun | VBN | verb, past part. |
| BER | are, art | NPS | plural proper noun | VBZ | verb, 3rd p. sg. pres. |
| BEZ | is | NR | adverbial noun | WDT | wh- determiner |
| CC | coordinating conjunction | OD | ordinal numeral | WPO | wh- pronoun, object |
| CD | cardinal numeral | PN | nominal pronoun | WPS | wh- pronoun, nom. |
| CS | subordinating conjunction | PP\$ | determiner, possessive | WQL | wh- qualifier (how) |
| DO | do | PP\$\$ | pronoun, possessive | WRB | wh- adverb |
| DOD | did | PPL | sg. reflexive pers. pron. | | |
| DOZ | does | PPLS | pl. reflexive pers. pron. | | |
| DT | sg. determiner (this, that) | PPO | personal pronoun | | |
| DTI | sg. or pl. det. (some, any) | PPS | 3rd p. sg. nom. pron. | | |
| DTS | pl. determiner (these, those) | PPSS | other nominative pers. pron. | | |
| DTX | double conjunction (either) | | | | |

| | | | |
|-----|-------------------------------------|-----|--------------------------------------|
| AJ0 | Adjective | TO0 | Infinitive marker TO |
| AJC | Comparative adjective | UNC | Foreign words |
| AJS | Superlative adjective | VBB | The present tense of BE, except is |
| AT0 | Article | VBD | The past tense of BE |
| AV0 | Adverb | VBG | The -ing form of BE |
| AVP | Adverb particle (e.g. up, off, out) | VBI | The infinitive of BE |
| AVQ | Wh-adverb | VBN | The past participle of BE |
| CJC | Coordinating conjunction | VBZ | IS, 'S |
| CJS | Subordinating conjunction | VDB | The finite base form of DO |
| CJT | that | VDD | The past tense of DO |
| CRD | Cardinal number | VDG | The -ing form of DO |
| ORD | Ordinal numeral | VDI | The infinitive of DO |
| DPS | Possessive determiner or pronoun | VDN | The past participle of DO |
| DT0 | General determiner-pronoun | VDZ | The -s form of DO |
| DTQ | Wh-determiner-pronoun | VHB | The finite base form of HAVE |
| EX0 | Existential there | VHD | The past tense of HAVE |
| NN0 | Common noun, neutral for number | VHG | The -ing form of HAVE |
| NN1 | Singular common noun | VHI | The infinitive of HAVE |
| NN2 | Plural common noun | VHN | The past participle of HAVE |
| NP0 | Proper noun | VHZ | The -s form of HAVE |
| PN1 | Indefinite pronoun | VM0 | Modal auxiliary verb |
| PNP | Personal pronoun | VVB | The finite base of lexical verbs |
| PNQ | Wh-pronoun | VVD | The past tense of lexical verbs |
| PNX | Reflexive pronoun | VVG | The -ing form of lexical verbs |
| POS | The genitive marker 'S or ' | VVI | The infinitive of lexical verbs |
| PRF | The preposition OF | VVN | The past participle of lexical verbs |
| PRP | Preposition, except OF | VVZ | The -s form of lexical verbs |
| PUL | Punctuation: left bracket | XX0 | NOT or N'T |
| PUN | Punctuation: general | ITJ | Interjection |
| PUQ | Punctuation: quotation mark | ZZ0 | Alphabetical symbols |
| PUR | Punctuation: right bracket | | |

Syntax-Korpora ("Baumbanken")

- Penn Treebank: 1M Worte aus dem Wall Street Journal
- Deutsch:
 - NEGRA
(20.000 Sätze Frankfurter Rundschau, 400K Worte)
 - TIGER
(50.000 Sätze Frankfurter Rundschau = 1M Worte)
- Prague Dependency Treebank (Czech)
- Neuerdings auch für viele andere Sprachen:
Chinesisch, Französische, usw.

NEGRA

- Als SQL Datenbank gespeichert; kann man in Bäume umwandeln
- Annotation:
 - PoS-tagged
 - Morphologische Annotation (60K)
 - Grammatische Funktionen
- Vorgehen:
 - Kombination aus automatischer Analyse und menschlicher Arbeit
 - Abfrage mit speziell dafür entwickelte Tools

Semantik-Korpora

| | |
|------------|-----------|
| [Peter] | Agent |
| gibt | |
| [Maria] | Recipient |
| [ein Buch] | Theme |

- Satzteilen werden semantische Rollen zugeordnet
- Einsatz: Training semantischer Parsern
- Korpora:
 - Englisch: PropBank, auf Grundlage der Penn Treebank
 - Deutsch: SALSA, auf Grundlage von TIGER

Diskurs-Korpora

[Peter ist müde]. Grund
Deshalb [schläft er]. Folge

- Ordne Paaren von Sätzen Diskursrelationen zu:
z. B. Begründung (weil), Zweck (damit),...
- Training von “Diskurs-Parsern”
- Korpora:
DiscourseBank, auf Grundlage der Penn Treebank

Bilinguale Korpora

- Vergleichbare Daten:
Crater corpora (English, French, Spanish)
- Parallel: Hansard Corpus, EUROPARL

Keine Korpora verfügbar

- Pragmatik:
Intention der Sprecher, “was wirklich gemeint ist”
- Viele andere Sprachen, besonders für höhere Ebenen