

Introduction to Computational Linguistics

Frank Richter

fr@sfs.uni-tuebingen.de.

**Seminar für Sprachwissenschaft
Eberhard Karls Universität Tübingen
Germany**

Sentence Segmentation

Task:

Determining how a text should be divided into sentences for further processing.

Terminology:

- sentence boundary detection
- sentence boundary disambiguation
- sentence boundary recognition

Sentences I

O. Dittrich:

Ein Satz ist eine modulatorisch abgeschlossene Lautung, wodurch der Hörende veranlasst wird, eine vom Sprechenden als richtig anerkennbare, relativ abgeschlossene apperzeptive (beziehende) Gliederung eines Bedeutungstatbestandes zu versuchen.

D. Jespersen:

eine (relativ) vollständige und unabhängige menschliche Äußerung, deren Vollständigkeit und Unabhängigkeit sich in ihrem Alleinstehen zeigt, d.h. darin, daß sie für sich allein geäußert wird.

Sentences II

A. Meillet:

eine Gemeinsamkeit von Artikulationen, die untereinander durch gewisse grammatische Beziehungen verbunden sind, grammatisch von keiner anderen Gesamtheit abhängen und sich selbst genügen.

W. Meyer-Lübke:

ein Wort oder eine Gruppe von Wörtern, die in der gesprochenen Sprache als Ganzes erscheinen, die sich als eine Mitteilung eines Sprechenden an einen anderen darstellen.

Sentences III

A. Nehring:

der sprachliche Ausdruck für eine vom Sprechenden jeweils hergestellte Ordnung einer gegebenen Mannigfaltigkeit von Sachverhalten.

W. Porzig:

ein Bedeutungsgefüge von derjenigen Form, durch die (in der betreffenden Sprache) Sachverhalte als abgeschlossene gemeint werden.

A. Stöhr:

eine mehrfache Benennung desselben Geschehnisses durch logisch gleichwertige Satzglieder.

Sentences IV

H. Paul:

der sprachliche Ausdruck, das Symbol dafür, daß sich die Verbindung mehrerer Vorstellungen oder Vorstellungsgruppen in der Seele des Sprechenden vollzogen hat, und das Mittel dazu, die nämliche Verbindung der nämlichen Vorstellungen in der Seele des Hörenden zu erzeugen. Jede engere Definition des Begriffes Satz muß als unzulänglich zurückgewiesen werden.

L. Bloomfield

an independent linguistic form, not included by virtue of any grammatical construction in any larger linguistic form.

Sentences in Real Life 1

There was nothing so VERY remarkable in that; nor did Alice think it so VERY much out of the way to hear the Rabbit say to itself, ‘Oh dear! Oh dear! I shall be late!’ (when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually TOOK A WATCH OUT OF ITS WAISTCOAT-POCKET, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket, or a watch to take out of it, and burning with curiosity, she ran across the field after it, and fortunately was just in time to see it pop down a large rabbit-hole under the hedge.

Sentences in Real Life 2

The holes certainly were rough—"Just right for a lot of vagabonds like us," said Bigwig—but the exhausted and those who wander in strange country are not particular about their quarters.

Sentences in Real Life 2

The holes certainly were rough—"Just right for a lot of vagabonds like us," said Bigwig—but the exhausted and those who wander in strange country are not particular about their quarters.

1. Two high-ranking positions were filled Friday by Penn St. University President Graham Spanier.
2. Two high-ranking positions were filled Friday by Penn St. University President Graham Spanier announced the appointments.

Problems with Sentence Segmentation

- Strict punctuation rules might exist, adherence varies
- Different punctuation marks or characters
- Task: disambiguate all punctuation marks that denote sentence boundaries:
periods, question marks, exclamation point,
semicolons, dashes, commas
- Use can vary with text types

Contextual Factors

- Case distinctions
- Part of speech
- Word length
- Lexical endings (to exclude abbreviations)
- Prefixes and suffixes before and after the punctuation mark
- Abbreviation classes

Incremental Linguistic Analysis

- tokenization
- morphological analysis (lemmatization)
- part-of-speech tagging
- named-entity recognition
- partial chunk parsing
- full syntactic parsing
- semantic and discourse processing

Potential Tasks

- Tokenize arbitrary text
- Subtask: Recognize date expressions
- Assign correct suffixes respecting vowel harmony
- Given an inflected verb: Find a base form of verbs and their agreement features
- Given a base form of verbs and their agreement features: find the appropriate inflected form
- Morphology: derivation: English verbs + suffix *-able* (yields an adjective: desirable, printable, readable, etc.)
- Assign syntactic categories to tokens in preprocessed text
- Bracketing of syntactic chunks in arbitrary text

Formal Languages & Computation

The language perspective

1. Type 3: regular expression languages
2. Type 2: context free languages
3. Type 1: context sensitive languages
4. Type 0: recursively enumerable languages

Formal Languages & Computation

The language perspective

1. Type 3: regular expression languages
2. Type 2: context free languages
3. Type 1: context sensitive languages
4. Type 0: recursively enumerable languages

The automata perspective

1. Finite automata
2. Pushdown automata
3. Linear automata (Turing machines with finite tapes)
4. Turing machines

Form of Grammars of Type 0–3

For $i \in \{0, 1, 2, 3\}$, a grammar $\langle N, T, \Pi, s \rangle$ of Type i , with N the set of non-terminal symbols, T the set of terminal symbols (N and T disjoint, $\Sigma = N \cup T$), Π the set of productions, and s the start symbol ($s \in N$), obeys the following restrictions:

Type 3: Every production in Π is of the form $A \rightarrow aB$ or $A \rightarrow \epsilon$, with $B, A \in N$, $a \in T$.

Type 2: Every production in Π is of the form $A \rightarrow x$, with $A \in N$ and $x \in \Sigma^*$.

Type 1: Every production in Π is of the form $x_1 A x_2 \rightarrow x_1 y x_2$, with $x_1, x_2 \in \Sigma^*$, $y \in \Sigma^+$, $A \in N$ and the possible exception of $C \rightarrow \epsilon$ in case C does not occur on the righthand side of a rule in Π .

Type 0: No restrictions.

An Example of a Type 2 Grammar

Let $\langle N, T, \Pi, S \rangle$ be a grammar with N, T and Π as given below:

- $N = \{S, NP, VP, V\}$
- $T = \{\text{John}, \text{walks}\}$
- $\Pi = \{S \rightarrow NP\ VP, NP \rightarrow \text{John}, VP \rightarrow V, V \rightarrow \text{walks}\}$