
Mildly Context-Sensitive Grammar

Formalisms:

Natural Languages are not Context-Free

Laura Kallmeyer
Heinrich-Heine-Universität Düsseldorf
Sommersemester 2011

Grammar Formalisms 1 Natural Languages and CFG

Kallmeyer Sommersemester 2011

Overview

1. CFG and Natural Languages
2. Cross-serial Dependencies
3. Swiss German is not Context-Free
4. LCFRS and Cross-serial Dependencies

Grammar Formalisms 2 Natural Languages and CFG

CFG and Natural Languages

- For a long time there has been a debate about whether CFGs are sufficiently powerful to describe natural languages. Several approaches have used CFGs, oftentimes enriched with some additional mechanism of transformation (Chomsky, 1956) or with features (Gazdar et al., 1985) for natural languages.
- In the 80's Stuart Shieber was able to prove in (Shieber, 1985) that there are natural languages that cannot be generated by a CFG. Before that, (Bresnan et al., 1982) made already a similar argument but their proof is based on the tree structures obtained with CFGs while Shieber argues on the basis of weak generative capacity, i.e., of the string languages.
- The phenomena considered in both papers are *cross-serial dependencies*.

Grammar Formalisms 3 Natural Languages and CFG

Kallmeyer Sommersemester 2011

Cross-serial Dependencies (1)

Cross-serial dependencies in Dutch (Bresnan et al., 1982):

- (1) ... dat Jan de kinderen zag zwemmen
... that Jan the children saw swim
'... that Jan saw the children swim'

The colours mark the dependencies between the two verbs and the two NPs: *the children* is an argument of *swim* while *Jan* is an argument of *saw*. The dependency links are in a crossing configuration.

Grammar Formalisms 4 Natural Languages and CFG

Cross-serial Dependencies (2)

This phenomenon can be iterated:

- (2) ... dat Jan Piet de kinderen zag helpen zwemmen
... that Jan Piet the children saw help swim
'... that Jan saw Piet help the children swim'
- (3) ... dat Jan Piet Marie de kinderen zag helpen leren zwemmen
... that Jan Piet Marie the children saw help teach swim
'... that Jan saw Piet help Marie teach the children to swim'

Cross-serial Dependencies (3)

- In principle, an unbounded number of crossed dependencies is possible.
- However, except for the first and last verb any permutation of the NPs and the verbs is grammatical as well (even though with a completely different dependency structure since the dependencies are always cross-serial).
- Therefore, the dependencies are not visible on the strings and the string language of Dutch cross-serial dependencies amounts roughly to $\{n^k v^k \mid k > 0\}$ which is a context-free language.

This is different for Swiss German because Swiss German has case marking.

Cross-serial Dependencies (4)

Cross-serial dependencies in Swiss German (Shieber, 1985):

- (4) ... das mer em Hans es huus hälfed aastriiche
... that we Hans_{Dat} house_{Acc} helped paint
'... that we helped Hans paint the house'
- (5) ... das mer d'chind em Hans es huus lönd hälfe aastriiche
... that we the children_{Acc} Hans_{Dat} house_{Acc} let help paint
'... that we let the children help Hans paint the house'

Swiss German

- uses case marking
- and displays cross-serial dependencies.

Swiss German is not Context-Free (1)

Proposition 1 *The language L of Swiss German is not context-free (Shieber, 1985).*

The argumentation of the proof goes as follows:

- We assume that L is context-free.
- Then the intersection of a regular language with the image of L under a homomorphism must be context-free as well.
- We find a particular homomorphism and a regular language such that the result obtained in this way is a non context-free language.
- This is a contradiction to our assumption and, consequently, the assumption does not hold.

Swiss German is not Context-Free (6)

If L is context-free, then L' must be context-free as well.

- Then the image of L' under a homomorphism f' with $f'(w) = f'(x) = f'(y) = \varepsilon$, $f'(a) = a$, $f'(b) = b$, $f'(c) = c$, $f'(d) = d$ is also context-free. This image is

$$f'(L') = L'' = \{a^i b^j c^i d^j \mid i, j \geq 0\}$$

- Consequently, L'' satisfies the pumping lemma for context-free languages. Inspecting the word $a^k b^k c^k d^k$ where k is the constant from the pumping lemma, this can be shown to lead to a contradiction.

Consequently, L'' is not context-free, and neither are L' and L .

LCFRS and Cross-serial Dependencies

LCFRS for Dutch cross-serial dependencies:

- $S(X Y zag Z) \rightarrow NP(X) VP(Y, Z)$
- $VP(X Y, leren Z) \rightarrow NP(X) VP(Y, Z)$
- $VP(X Y, helpen Z) \rightarrow NP(X) VP(Y, Z)$
- $VP(X, zwemmen) \rightarrow NP(X)$
- $NP(Jan) \rightarrow \varepsilon$
- $NP(Marie) \rightarrow \varepsilon$
- $NP(Piet) \rightarrow \varepsilon$
- $NP(de kinderen) \rightarrow \varepsilon$

References

- Bresnan, Joan, Ronald M. Kaplan, Stanley Peters, and Annie Zaenen. 1982. Cross-serial dependencies in Dutch. *Linguistic Inquiry*, 13(4):613–635. Reprinted in (Savitch et al., 1987).
- Chomsky, Noam. 1956. Three models for the description of language. *IRE Transactions on Information Theory*, 2:113–124.
- Gazdar, Gerald, Ewan Klein, Geoffrey Pullman, and Ivan Sag. 1985. *Generalized Phrase Structure Grammar*. Harvard University Press, Cambridge, Massachusetts.
- Savitch, Walter J., Emmon Bach, William Marxh, and Gila Safran-Naveh, editors. 1987. *The Formal Complexity of Natural Language*. Studies in Linguistics and Philosophy. Reidel, Dordrecht, Holland.
- Shieber, Stuart M. 1985. Evidence against the context-freeness of natural language. *Linguistics and Philosophy*, 8:333–343. Reprinted

in (Savitch et al., 1987).