

Formal Languages and Linguistics

Pascal Amsili

Sorbonne Nouvelle, Lattice (CNRS/ENS-PSL/SN)

Cogmaster, september 2023

Overview

Formal Languages

Regular Languages

Formal Grammars

Examples

Definition

Language classes

Formal complexity of Natural Languages

Example I

$$\begin{array}{lcl} S & \rightarrow & AB \\ A & \rightarrow & aA \\ & & | \quad b \\ B & \rightarrow & bBc \\ & & | \quad \varepsilon \end{array}$$

- ▶ Rewriting system
- ▶ Auxiliary vocabulary (N for non-terminal)
- ▶ Start symbol (engendered language)
- ▶ Multiple derivations
- ▶ Syntactic tree

Example II

$$\begin{array}{l} E \rightarrow E + E \\ | E \times E \\ | (E) \\ | 0 | 1 | 2 \dots 8 | 9 \end{array}$$

- ▶ Syntactic ambiguity
- ▶ Semantic interpretation

Example III

| | | |
|-------------|---|-----------------|
| <i>NP</i> | → | <i>Det N'</i> |
| <i>N'</i> | → | <i>AdjP N'</i> |
| <i>N'</i> | → | <i>N</i> |
| <i>N'</i> | → | <i>N Cpt</i> |
| <i>AdjP</i> | → | <i>Adj AdjP</i> |
| <i>AdjP</i> | → | <i>Adj</i> |
| <i>Cpt</i> | → | <i>P NP</i> |
| <i>Det</i> | → | the my |
| <i>N</i> | → | cat friend |
| <i>Adj</i> | → | large fierce |
| <i>Prep</i> | → | of to |

- ▶ X-bar theory
- ▶ Recursive rules
- ▶ Center-embedding

Overview

Formal Languages

Regular Languages

Formal Grammars

Examples

Definition

Language classes

Formal complexity of Natural Languages

Formal grammar

Def. 12 ((Formal) Grammar)

A **formal grammar** is defined by $\langle \Sigma, N, S, P \rangle$ where

- ▶ Σ is an alphabet
- ▶ N is a disjoint alphabet (non-terminal vocabulary)
- ▶ $S \in N$ is a distinguished element of N , called the *axiom*
- ▶ P is a set of « *production rules* », namely a subset of the cartesian product $(\Sigma \cup N)^* N (\Sigma \cup N)^* \times (\Sigma \cup N)^*$.

Immediate Derivation

Def. 13 (Immediate derivation)

Let $\mathcal{G} = \langle \Sigma, N, S, P \rangle$ a grammar,

$r \in P$ a production rule, such that $r : A \rightarrow u$ with $u \in (\Sigma \cup N)^*$;

$f, g \in (\Sigma \cup N)^*$ two “(proto-)words”,

- f derives into g (immediate derivation) **with the rule r**
(noted $f \xrightarrow{r} g$) iff
 $\exists v, w$ s.t. $f = vAw$ and $g = vuw$
- f derives into g (immediate derivation) **in the grammar \mathcal{G}**
(noted $f \xrightarrow{\mathcal{G}} g$) iff
 $\exists r \in P$ s.t. $f \xrightarrow{r} g$.

Derivation

Def. 14 (Derivation)

$f \xrightarrow{\mathcal{G}^*} g$ if $f = g$ or

$\exists f_0, f_1, f_2, \dots, f_n$ s.t.

$$f_0 = f$$

$$f_n = g$$

$$\forall i \in [1, n] : f_{i-1} \xrightarrow{\mathcal{G}} f_i$$

Engendered language

Def. 15 (Language engendered by a word)

Let $f \in (\Sigma \cup N)^*$.

$$L_G(f) = \{g \in X^* / f \xrightarrow{G^*} g\}$$

Def. 16 (Language engendered by a grammar)

The *language engendered by a grammar* G is the set of words of Σ^* derived from the **axiom**.

$$L_G = L_G(S)$$

Overview

Formal Languages

Regular Languages

Formal Grammars

Examples

Definition

Language classes

Formal complexity of Natural Languages

Principle

Define language families on the basis of properties of the grammars that generate them :

1. Four classes are defined, they are included one in another
2. A language is of type k if it **can** be recognized by a type k grammar (and thus, by definition, by a type $k - 1$ grammar) ; and cannot be recognized by a grammar of type $k + 1$.

Chomsky's hierarchy

type 0 No restriction on

$$P \subset (X \cup V)^* V (X \cup V)^* \times (X \cup V)^*.$$

type 1 (*context-sensitive* grammars) All rules of P are of the shape $(u_1 S u_2, u_1 m u_2)$, where u_1 and $u_2 \in (X \cup V)^*$, $S \in V$ and $m \in (X \cup V)^+$.

type 2 (*context-free* grammar) All rules of P are of the shape (S, m) , where $S \in V$ and $m \in (X \cup V)^*$.

type 3 (*regular* grammars) All rules of P are of the shape (S, m) , where $S \in V$ and $m \in X.V \cup X \cup \{\varepsilon\}$.

Examples

type 3:

$$S \rightarrow aS \mid aB \mid bB \mid cA$$
$$B \rightarrow bB \mid b$$
$$A \rightarrow cS \mid bB$$

Examples

type 3:

$$S \rightarrow aS \mid aB \mid bB \mid cA$$

$$B \rightarrow bB \mid b$$

$$A \rightarrow cS \mid bB$$

type 2:

$$E \rightarrow E + T \mid T, T \rightarrow T \times F \mid F, F \rightarrow (E) \mid a$$

Example 1 type 0

Type 0:

$S \rightarrow SABC \quad AC \rightarrow CA \quad A \rightarrow a$

$S \rightarrow \varepsilon \quad CA \rightarrow AC \quad B \rightarrow b$

$AB \rightarrow BA \quad BC \rightarrow CB \quad C \rightarrow c$

$BA \rightarrow AB \quad CB \rightarrow BC$

generated language :

Example 1 type 0

Type 0:

$S \rightarrow SABC \quad AC \rightarrow CA \quad A \rightarrow a$

$S \rightarrow \varepsilon \quad CA \rightarrow AC \quad B \rightarrow b$

$AB \rightarrow BA \quad BC \rightarrow CB \quad C \rightarrow c$

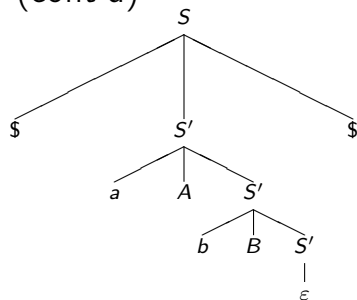
$BA \rightarrow AB \quad CB \rightarrow BC$

generated language : words with an equal number of a , b , and c .

Example 2: type 0

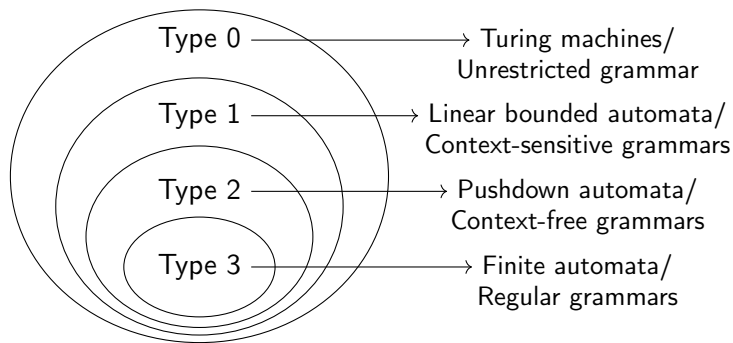
Type 0: $S \rightarrow \$S'\$$ $Aa \rightarrow aA$ $\$a \rightarrow a\$$
 $S' \rightarrow aAS'$ $Ab \rightarrow bA$ $\$b \rightarrow b\$$
 $S' \rightarrow bBS'$ $Ba \rightarrow aB$ $A\$ \rightarrow \a
 $S' \rightarrow \epsilon$ $Bb \rightarrow bB$ $B\$ \rightarrow \b
 $\$\$ \rightarrow \#$

Example 2: type 0 (cont'd)



| | | | | | |
|----|----|----|----|----|----|
| \$ | a | A | b | B | \$ |
| a | \$ | A | b | B | \$ |
| a | \$ | A | b | \$ | b |
| a | \$ | b | A | \$ | b |
| a | b | \$ | A | \$ | b |
| a | b | \$ | \$ | a | b |
| a | b | # | a | b | b |

The Chomsky-Schützenberger hierarchy



Remarks

- ▶ Type 0 (Turing-recognizable) = recursively enumerable languages
Type 1 (Turing-decidable) = recursive languages
- ▶ There are others ways to classify languages,
 - ▶ either on other properties of the grammars;
 - ▶ or on other properties of the languages
- ▶ Nested structures are preferred, but it's not necessary

The parsing problem: finding derivations

- ▶ Given a grammar G on some alphabet Σ ...

- ▶ The **parsing problem** for G :

Given some $w \in \Sigma$,
what are the derivations (if any) of w in G ?

- ▶ (Solving the parsing problem for G entails solving the recognition problem for $\mathcal{L}(G)$.)

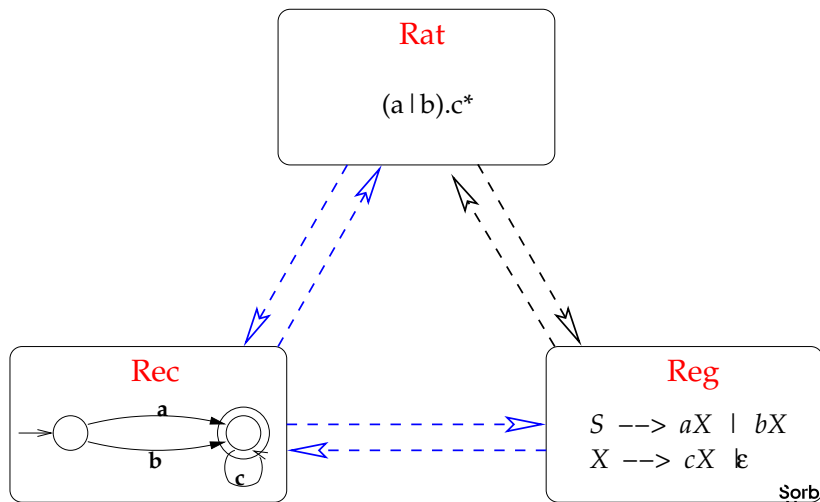
Syntactic complexity vs semantic expressivity

- ▶ Context-free grammars are commonly used to describe the syntax of many logical languages (PL, FOL), some programming languages, and parts of NL (→ Day 2).
- ▶ Untyped λ -calculus: CF syntax, Turing-complete semantics. “How is this possible?”
- ▶ → The syntactic complexity and the semantic expressivity of interpreted languages are two distinct notions.
- ▶ Jot (https://en.wikipedia.org/wiki/Iota_and_Jot) is $\{0, 1\}$, a regular language, compositionally interpreted as a Turing-complete language.

The recognition/parsing problems are very general

- ▶ Consider any binary (“yes/no”) problem P and see it as the set of inputs for which the answer is positive.
- ▶ Let str be a linearisation function for the possible inputs of P , and $L = \{str(in) \mid in \in P\}$.
- ▶ Solving P is equivalent to the recognition problem for L .
- ▶ More generally, any computable function f can be encoded as a grammar s.t. after parsing the input w , the output $f(w)$ can be read off the derivation.
- ▶ \rightarrow One can compute “syntactically”: a grammar is a program. (The parser is the machine that runs it.)
- ▶ The formalism of unrestricted grammars is a Turing-complete programming language. (syntactically regular?)

Back to regular languages



Let's play with grammars

For each of the following grammars, give the generated language, and the type they have in Chomsky's hierarchy.

$$\begin{aligned} S &\rightarrow S_1 S_2 \\ S_1 &\rightarrow a S_1 b \mid ab \\ S_2 &\rightarrow c S_2 \mid c \end{aligned}$$

$$\begin{aligned} S &\rightarrow a S B C \\ S &\rightarrow a B C \\ C B &\rightarrow B C \\ a B &\rightarrow ab \\ b B &\rightarrow bb \\ b C &\rightarrow bc \\ c C &\rightarrow cc \end{aligned}$$

Let's play with grammars (cont'd)

Give a context-free grammar that generates each of the following languages (alphabet $\Sigma = \{a, b, c\}$).

- ▶ $L_0 = \{w \in X^* / w = a^n ; n \geq 0\}$
- ▶ $L'_0 = \{w \in X^* / w = a^n b^n c a ; n \geq 0\}$
- ▶ $L_1 = \{w \in X^* / w = a^n b^n c^p ; n > 0 \text{ et } p > 0\}$
- ▶ $L_2 = \{w \in X^* / w = a^n b^n a^m b^m ; n, m \geq 1\}$
- ▶ $L'_3 = \{w \in X^* / |w|_a = |w|_b\}$
- ▶ $L_3 = \{w \in X^* / |w|_a = 2|w|_b\}$
- ▶ $L_4 = \{w \in X^* / \exists x \in X^* \text{ tq } w = x\bar{x}\}$
- ▶ $L_5 = \{w \in X^* / w = \bar{w}\}$

References I

- Bar-Hillel, Yehoshua, Perles, Micha, & Shamir, Eliahu. 1961. On formal properties of simple phrase structure grammars. *STUF-Language Typology and Universals*, 14(1-4), 143–172.
- Chomsky, Noam. 1957. *Syntactic Structures*. Den Haag: Mouton & Co.
- Chomsky, Noam. 1995. *The Minimalist Program*. Vol. 28. Cambridge, Mass.: MIT Press.
- Gazdar, Gerald, & Pullum, Geoffrey K. 1985 (May). *Computationally Relevant Properties of Natural Languages and Their Grammars*. Tech. rept. Center for the Study of Language and Information, Leland Stanford Junior University.
- Gibson, Edward, & Thomas, James. 1997. The Complexity of Nested Structures in English: Evidence for the Syntactic Prediction Locality Theory of Linguistic Complexity. *Unpublished manuscript, Massachusetts Institute of Technology*.
- Joshi, Aravind K. 1985. *Tree Adjoining Grammars: How Much Context-Sensitivity is Required to Provide Reasonable Structural Descriptions?* Tech. rept. Department of Computer and Information Science, University of Pennsylvania.
- Langendoen, D Terence, & Postal, Paul Martin. 1984. *The vastness of natural languages*. Basil Blackwell Oxford.
- Mannell, Robert. 1999. *Infinite number of sentences*. part of a set of class notes on the Internet. http://clas.mq.edu.au/speech/infinite_sentences/.
- Schieber, Stuart M. 1985. Evidence against the Context-Freeness of Natural Language. *Linguistics and Philosophy*, 8(3), 333–343.
- Stabler, Edward P. 2011. Computational perspectives on minimalism. *Oxford handbook of linguistic minimalism*, 617–643.
- Steedman, Mark, et al. . 2012 (June). *Combinatory Categorical Grammars for Robust Natural Language Processing*. Slides for NASSLLI course <http://homepages.inf.ed.ac.uk/steedman/papers/ccg/nassll112.pdf>.
- Vijay-Shanker, K., & Weir, David J. 1994. The Equivalence of Four Extensions of Context-Free Grammars. *Mathematical Systems Theory*, 27, 511–546.