

Language and Computation

week 8, Thursday, March 06, 2014

Tamás Biró

Yale University

tamas.biro@yale.edu

<http://www.birot.hu/courses/2014-LC/>



Practical matters

- **Post-reading:** Chapters 12, 16 and 14.1
- **Pre-reading:** Sections 13.1-3
- **Homework 3** returned after the break.
- **Midterm:** will be posted tomorrow.
- Proof of HW 2, part 3 posted.
- To come: Viterbi and Forward-Backward – an example



Today

- Formal definition of Formal Grammars
- Chomsky Hierarchy
- The Pumping Lemma
- Beyond regular languages: Context-Free Grammars
- Probabilistic Context Free Grammars

(Parsing to come after the break)



Formal Grammars: an example



Formal Grammars

A toy grammar for English:

Phrase Structure Rules: **Lexical Insertion Rules:**

$S \rightarrow NP \ VP$

$VP \rightarrow V$

$VP \rightarrow V \ NP$

$NP \rightarrow N$

$NP \rightarrow Det \ N$

$V \rightarrow \{ \text{eat, love, walk, sleep} \}$

$V \rightarrow \{ \text{eats, loves, walks sleeps} \}$

$N \rightarrow \{ \text{John, Marry. . .} \}$

$N \rightarrow \{ \text{apple, pear. . .} \}$

$N \rightarrow \{ \text{apples, pears. . .} \}$

$Det \rightarrow \{ \text{the, a, an, } \emptyset \}$



Formal Grammars

A toy grammar for English:

$S \Rightarrow NP VP \Rightarrow Det N VP \Rightarrow$

$\Rightarrow Det N V NP \Rightarrow Det N V N \Rightarrow$

$\Rightarrow The N V N \Rightarrow The John V N \Rightarrow$

$\Rightarrow The John sleep N \Rightarrow The John sleep apple$

This is a sentence derived from this grammar.



Formal Grammars

A toy grammar for English – lessons:

- Introduce additional categories:
 $V_{\text{transitive}}$ vs. $V_{\text{intransitive}}$.
- Proper names as NPs .
- Agreement

→ more general formalism needed (feature structures: Ch. 15)

Formal Grammars: an example

$V \rightarrow V$ and... what?

Subcategorization frames for a set of example verbs:

Frame	Verb	Example
\emptyset	eat, sleep	I ate
NP	prefer, find, leave	Find [NP the flight from Pittsburgh to Boston]
$NP NP$	show, give	Show [NP me] [NP airlines with flights from Pittsburgh]
$PP_{\text{from}} PP_{\text{to}}$	fly, travel	I would like to fly [PP from Boston] [PP to Philadelphia]
$NP PP_{\text{with}}$	help, load	Can you help [NP me] [PP with a flight]
VP_{to}	prefer, want, need	I would prefer [VP_{to} to go by United airlines]
VP_{brst}	can, would, might	I can [VP_{brst} go from Boston]
S	mean	Does this mean [S AA has a hub in Boston]

Formal Grammars



Formal Grammars

N a set of **non-terminal symbols** (or **variables**)

Σ a set of **terminal symbols** (disjoint from N)

R a set of **rules** or productions, each of the form $A \rightarrow \beta$,
where A is a non-terminal,

β is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$

S a designated **start symbol**

Capital letters like A , B , and S

S

Lower-case Greek letters like α , β , and γ

Lower-case Roman letters like u , v , and w

Non-terminals

The start symbol

Strings drawn from $(\Sigma \cup N)^*$

Strings of terminals

Formal Grammars

Given formal grammar $G = (N, \Sigma, R, S)$:

Def: Given strings a and $b \in (\Sigma \cup N)^*$, $a \Rightarrow_G b$ iff there exist $p, q, r, s \in (\Sigma \cup N)^*$ such that

- $a = p + q + s$,
- $b = p + r + s$, and
- $(q \rightarrow r) \in R$

Def: A string $a \in \Sigma^*$ is grammatical in grammar G iff $S \Rightarrow^* a$.

The Chomsky Hierarchy



Generative power of a formalism

What is the set of languages generated by a formalism?

- **Overgeneration:** too powerful a formalism, also generating languages that we don't want.
- **Undergeneration:** too weak a formalism, not generating the languages we would like to.

Generative power of a formalism

What is the set of languages generated by a formalism?

Goal: generate exactly the attested human languages.

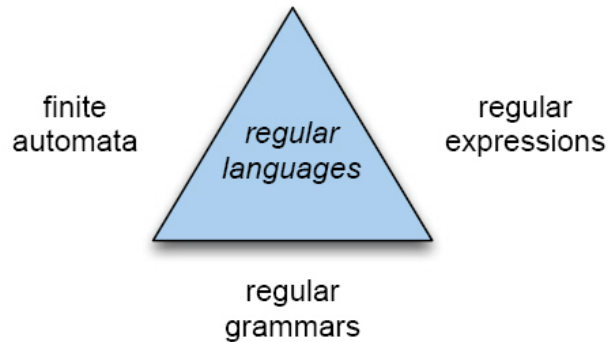
If reached: our formalism *accounts* for human languages.

Making happy

- the theoretical linguist wishing to characterize the possible languages of the world, who is now offered a mathematical tool to do so.
- the cognitive scientist wishing to decipher the “mental software” run by our brain.



Regular languages



But this is too weak a formalism for natural languages!

What can we do with formal grammars?

Chomsky hierarchy

Type	Common Name	Rule Skeleton	Linguistic Example
0	Turing Equivalent	$\alpha \rightarrow \beta$, s.t. $\alpha \neq \epsilon$	HPSG, LFG, Minimalism
1	Context Sensitive	$\alpha A \beta \rightarrow \alpha \gamma \beta$, s.t. $\gamma \neq \epsilon$	
–	Mildly Context Sensitive		TAG, CCG
2	Context Free	$A \rightarrow \gamma$	Phrase-Structure Grammars
3	Regular	$A \rightarrow xB$ or $A \rightarrow x$	Finite-State Automata

NB:

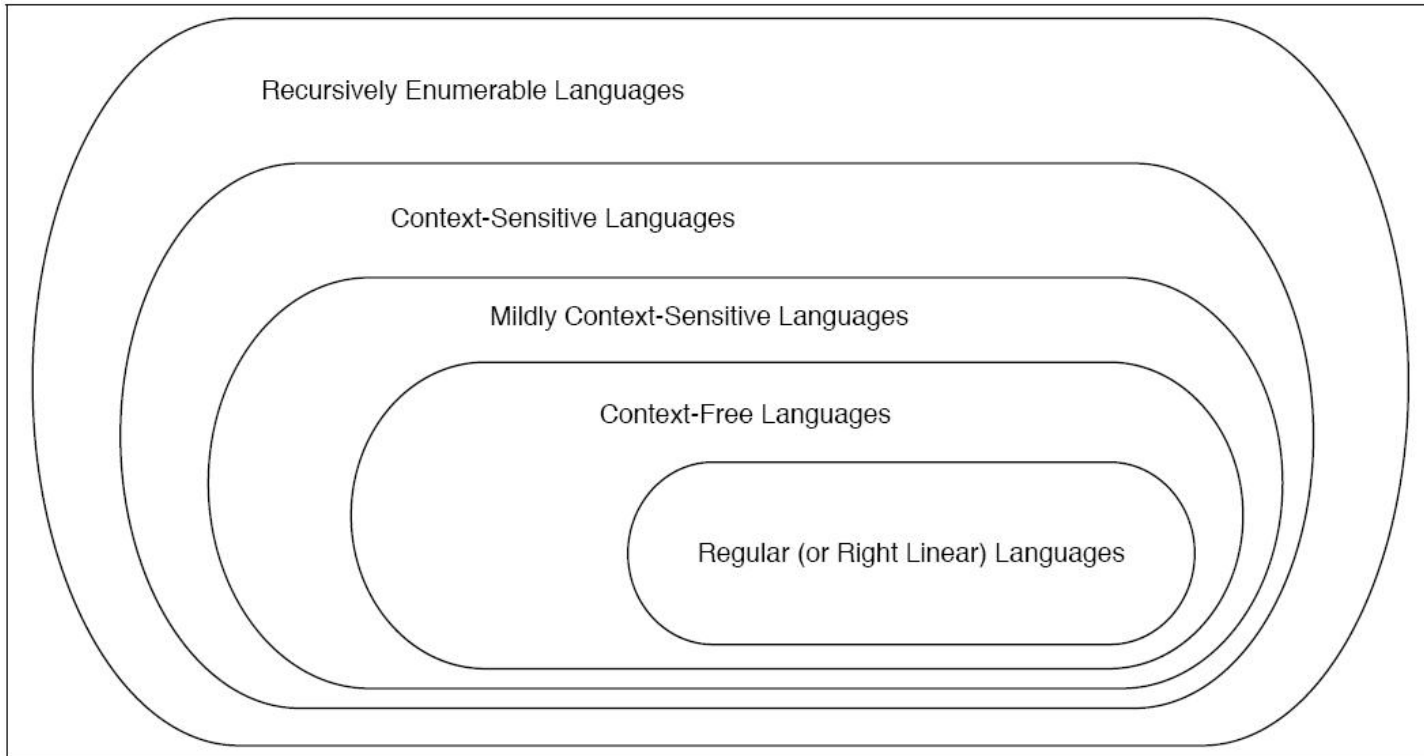
0: Turing machine

1: Linear bounded automaton

2: Non-deterministic push-down automaton

3: Finite-state automaton

The Chomsky Hierarchy



Weak and strong equivalence

$$\{a^n b^m \mid n, m \in \mathbb{N}^+\}$$

- Regular expression: $/a^+ b^+ /$
- Finite State Automaton: initial state q_0 , state q_1 , end state q_2 , arc $q_0 \rightarrow q_0$ with label a , arc $q_1 \rightarrow q_1$ with label b , arc loop $q_0 \rightarrow q_1$ with label a , arc loop $q_1 \rightarrow q_2$ with label b .
- Regular grammar: $S \rightarrow a S$, $S \rightarrow a A$, $A \rightarrow b A$, $A \rightarrow b$
- Context Free Grammar: $S \rightarrow A B$, $A \rightarrow A A$, $B \rightarrow B B$, $A \rightarrow a$, $B \rightarrow b$



The Pumping Lemma

For all L infinite regular languages,
there are strings x , y and z such that

$y \neq \epsilon$ and

$xy^Nz \in L$ for all $N \geq 0$.

Example: $\{a^n b^n\}$ is not regular.

The Pumping Lemma

Example: $L = \{xx^{rev} \mid x \in \{a, b\}^*\}$ is not regular.

where x^{rev} is the string x reversed. The strings in L are symmetrical.

Proof:

1. Intersect L with regular language aa^*bbaa^* .

If L were regular, then intersection would also be regular, because regular languages are closed for intersection (J&M 2.3).

2. Resulting language is $a^n b^2 a^n$, which is not regular, due to the pumping lemma. Therefore L cannot be regular, either.

Probabilistic Context Free Grammars



Probabilistic Context Free Grammars

N a set of **non-terminal symbols** (or **variables**)

Σ a set of **terminal symbols** (disjoint from N)

R a set of **rules** or productions, each of the form $A \rightarrow \beta [p]$,
where A is a non-terminal,

β is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$,
and p is a number between 0 and 1 expressing $P(\beta|A)$

S a designated **start symbol**

$$\sum_{\beta \in (N \cup \Sigma)^*} P(A \rightarrow \beta) = 1$$

Probabilistic Context Free Grammars

Probability of tree T (which yields sentence S):

$$P(T, S) = \prod_{i=1}^n P(RHS_i | LHS_i)$$

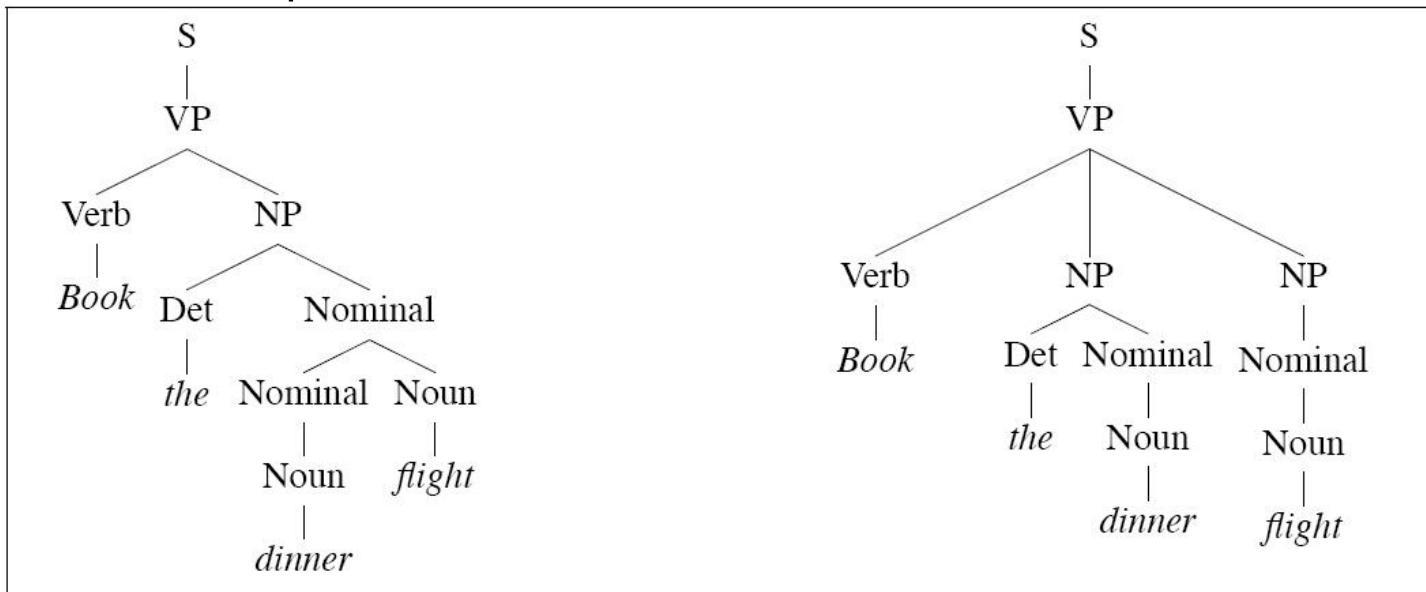
the product of the probabilities of the n rules used to expand each of the n non-terminal nodes in parse tree T (J&M 14.1.1).

Probabilistic CFG: an example

Grammar		Lexicon
$S \rightarrow NP VP$	[.80]	<i>Det</i> \rightarrow <i>that</i> [.10] <i>a</i> [.30] <i>the</i> [.60]
$S \rightarrow Aux NP VP$	[.15]	<i>Noun</i> \rightarrow <i>book</i> [.10] <i>flight</i> [.30]
$S \rightarrow VP$	[.05]	<i>meal</i> [.15] <i>money</i> [.05]
$NP \rightarrow Pronoun$	[.35]	<i>flights</i> [.40] <i>dinner</i> [.10]
$NP \rightarrow Proper-Noun$	[.30]	<i>Verb</i> \rightarrow <i>book</i> [.30] <i>include</i> [.30]
$NP \rightarrow Det Nominal$	[.20]	<i>prefer</i> ; [.40]
$NP \rightarrow Nominal$	[.15]	<i>Pronoun</i> \rightarrow <i>I</i> [.40] <i>she</i> [.05]
$Nominal \rightarrow Noun$	[.75]	<i>me</i> [.15] <i>you</i> [.40]
$Nominal \rightarrow Nominal Noun$	[.20]	<i>Proper-Noun</i> \rightarrow <i>Houston</i> [.60]
$Nominal \rightarrow Nominal PP$	[.05]	<i>NWA</i> [.40]
$VP \rightarrow Verb$	[.35]	<i>Aux</i> \rightarrow <i>does</i> [.60] <i>can</i> [.40]
$VP \rightarrow Verb NP$	[.20]	<i>Preposition</i> \rightarrow <i>from</i> [.30] <i>to</i> [.30]
$VP \rightarrow Verb NP PP$	[.10]	<i>on</i> [.20] <i>near</i> [.15]
$VP \rightarrow Verb PP$	[.15]	<i>through</i> [.05]
$VP \rightarrow Verb NP NP$	[.05]	
$VP \rightarrow VP PP$	[.15]	
$PP \rightarrow Preposition NP$	[1.0]	

Probabilistic Context Free Grammars

An example:



(booking a flight serving dinner vs. booking a flight on behalf of 'dinner'.)

Parsing and grammar learning

- **Parsing:** Given (probabilistic) CFG G , given sentence s , find (possible/most probable) parse tree for s in G .
- **Learning:** Given set of sentences, build a (probabilistic) context free grammar.

Have a nice break, and
see you after the spring recess!

