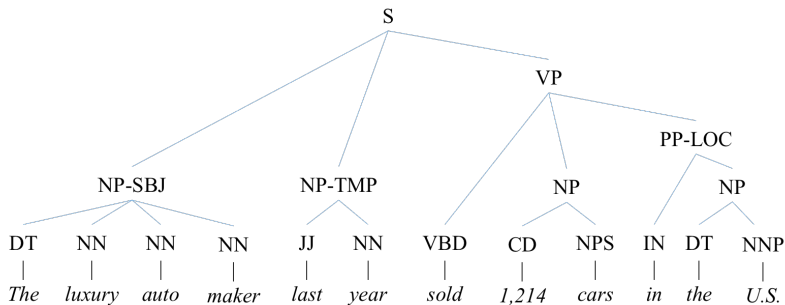


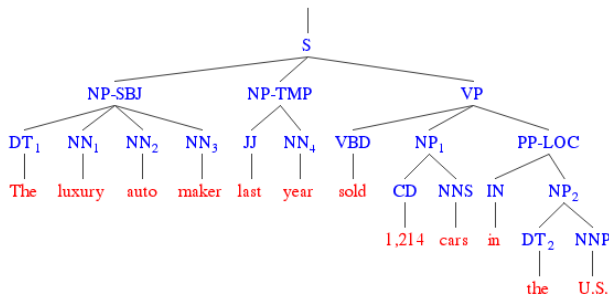
LING439/539 - Statistical NLP
Probabilistic context-free grammars

Monday, September 29 2016

Probabilistic context-free grammars



Tree drawing...



[[S

[NP-SBJ [DT The] [NN luxury] [NN auto] [NN maker]]

[NP-TMP [JJ last] [NN year]]

[VP [VBD sold]

[NP [CD 1,214] [NNS cars]]

[PP-LOC [IN in]

[NP [DT the] [NNP U.S.]]]]]]

<http://ironcreek.net/phpsyntaxtree>

Treebanks as grammars

S	→	NP-SBJ NP-TMP VP
NP-SBJ	→	DT NN NN NN
NP-TMP	→	JJ NN
VP	→	VBD NP PP-LOC
NP	→	CD NNS
PP-LOC	→	IN NP
NP	→	DT NNP
...		

Treebanks as grammars

S	→	NP-SBJ NP-TMP VP
NP-SBJ	→	DT NN NN NN
NP-TMP	→	JJ NN
VP	→	VBD NP PP-LOC
NP	→	CD NNS
PP-LOC	→	IN NP
NP	→	DT NNP
DT	→	<i>The</i>
NN	→	<i>luxury</i>
NN	→	<i>auto</i>
NN	→	<i>maker</i>
...		

Context-free grammar

$$A \rightarrow \gamma$$

where $A \in V$, and $\gamma \in (V \cup \Sigma)^*$

Context-free grammar: $S \rightarrow aSa$
 $S \rightarrow bSb$
 $S \rightarrow \epsilon$

Context-free language: $S \Rightarrow aSa \Rightarrow aaSaa \Rightarrow aabSbaa \Rightarrow aabbbaa$

$$L(G) = \{ww^R : w \in \{a, b\}^*\}$$

Chomsky normal form (CNF)

A grammar where every production is either of the form

$$A \rightarrow BC$$

or where $A, B, C \in V$ and $c \in \Sigma$.

$$A \rightarrow c$$

Converting a CFG to the Chomsky normal form

1. START: Eliminate the start symbol from right-hand sides
 - ▶ Introduce a new start symbol S_0 , and a new rule $S_0 \rightarrow S$
 - ▶ where S is the previous start symbol. This doesn't change the grammar's produced language, and S_0 won't occur on any rule's right-hand side.

2. TERM: Eliminate rules with nonsolitary terminals

- ▶ To eliminate each rule $A \rightarrow X_1 \dots a \dots X_n$ with a terminal symbol a being not the only symbol on the right-hand side,
- ▶ introduce a new nonterminal symbol N_a for every such terminal, $A \rightarrow X_1 \dots N_a \dots X_n$
- ▶ and a new rule $N_a \rightarrow a$

3. BIN: Eliminate right-hand sides with more than 2 nonterminals
- ▶ Replace each rule $A \rightarrow X_1X_2...X_n$ with more than 2 nonterminals $X_1, ..., X_n$ by rules
 - ▶ $A \rightarrow X_1A_1$,
 - ▶ $A_1 \rightarrow X_2A_2$,
 - ▶ \dots ,
 - ▶ $A_{n-2} \rightarrow X_{n-1}X_n$ where A_i are new nonterminal symbols.

4. UNIT: Eliminate unit rules

- ▶ A unit rule is a rule of the form $A \rightarrow B$ where A, B are nonterminal symbols, and
- ▶ $B \rightarrow X_1 \dots X_n$ where $X_1 \dots X_n$ is a string of nonterminals and terminals,
- ▶ remove them and add rule $A \rightarrow X_1 \dots X_n$ unless this is a unit rule which has already been removed.

5. DEL: Eliminate ϵ -rules

- ▶ An ϵ -rule is a rule of the form $A \rightarrow \epsilon$, where A is not the grammar's start symbol.
- ▶ To eliminate all rules of this form, first determine the set of all nonterminals that derive ϵ (nonterminals nullable).
 - ▶ If a rule $A \rightarrow \epsilon$ exists, then A is nullable.
 - ▶ If a rule $A \rightarrow X_1 \dots X_n$ exists, and each X_i is nullable, then A is nullable, too.
- ▶ Obtain an intermediate grammar by replacing each rule $A \rightarrow X_1 \dots X_n$ by all versions with some nullable X_i omitted. By deleting in this grammar each ϵ -rule, unless its left-hand side is the start symbol, the transformed grammar is obtained.

For example, in the following grammar, with start symbol S_0 ,

- ▶ $S_0 \rightarrow AbB|C$
- ▶ $B \rightarrow AA|AC$
- ▶ $C \rightarrow b|c$
- ▶ $A \rightarrow a|\epsilon$

where the nonterminal A , and hence also B , is nullable, while neither C nor S_0 is.

- ▶ $S_0 \rightarrow AbB|C$
 - ▶ $S_0 \rightarrow AbB|\cancel{AbB}|\cancel{AbB}|\cancel{AbB}|C$
- ▶ $B \rightarrow AA|AC$
 - ▶ $B \rightarrow AA|\cancel{AA}|\cancel{AA}|\cancel{AA}|AC|\cancel{AC}$
- ▶ $C \rightarrow b|c$
- ▶ $A \rightarrow a|\epsilon$
 - ▶ $A \rightarrow a|\epsilon$

where the nonterminal A , and hence also B , is nullable, while neither C nor S_0 is.

- ▶ $S_0 \rightarrow AbB|Ab|bB|b|C$
- ▶ $B \rightarrow AA|A|AC|C$
- ▶ $C \rightarrow b|c$
- ▶ $A \rightarrow a$

Exercise

Consider the CFG:

- ▶ $S \rightarrow aXbX$
- ▶ $X \rightarrow aY \mid bY \mid \epsilon$
- ▶ $Y \rightarrow X \mid c$

Which one is nullable?

DEL:

- ▶ $S \rightarrow aXbX$
 - ▶ $S \rightarrow aXbX | a\cancel{X}bX | aXb\cancel{X} | a\cancel{X}b\cancel{X}$
- ▶ $X \rightarrow aY | bY | \epsilon$
 - ▶ $X \rightarrow aY | a\cancel{Y} | bY | b\cancel{Y} | \epsilon$
- ▶ $Y \rightarrow X | c$

Remove Unit $Y \rightarrow X$ (UNIT)

- ▶ $S \rightarrow aXbX|abX|aXb|ab$
- ▶ $X \rightarrow aY|a|bY|b$
- ▶ $Y \rightarrow X|c$
 - ▶ $Y \rightarrow aY|a|bY|b|c$

Add $A \rightarrow a$, $B \rightarrow b$, $C \rightarrow c$ (TERM)

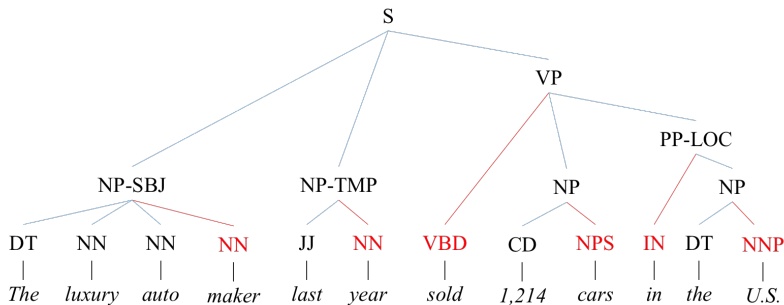
- ▶ $S \rightarrow aXbX|abX|aXb|ab$
 - ▶ $S \rightarrow AXBX|ABX|AXB|AB$
- ▶ $X \rightarrow aY|a|bY|b$
 - ▶ $X \rightarrow AY|A|BY|B$
- ▶ $Y \rightarrow aY|a|bY|b|c$
 - ▶ $Y \rightarrow AY|A|BY|B|C$
- ▶ $A \rightarrow a$
- ▶ $B \rightarrow b$
- ▶ $C \rightarrow c$

BIN:

- ▶ $S \rightarrow \cancel{A}X\cancel{M}B\cancel{X}N | \cancel{A}B\cancel{X}N | \cancel{A}X\cancel{M}B | AB$
 - ▶ $M \rightarrow AX$
 - ▶ $N \rightarrow BX$
- ▶ $X \rightarrow AY | A | BY | B$
- ▶ $Y \rightarrow AY | A | BY | B | C$
- ▶ $A \rightarrow a$
- ▶ $B \rightarrow b$
- ▶ $C \rightarrow c$

- ▶ $S \rightarrow MN|AN|MB|AB$
- ▶ $M \rightarrow AX$
- ▶ $N \rightarrow BX$
- ▶ $X \rightarrow AY|A|BY|B$
- ▶ $Y \rightarrow AY|A|BY|B|C$
- ▶ $A \rightarrow a$
- ▶ $B \rightarrow b$
- ▶ $C \rightarrow c$

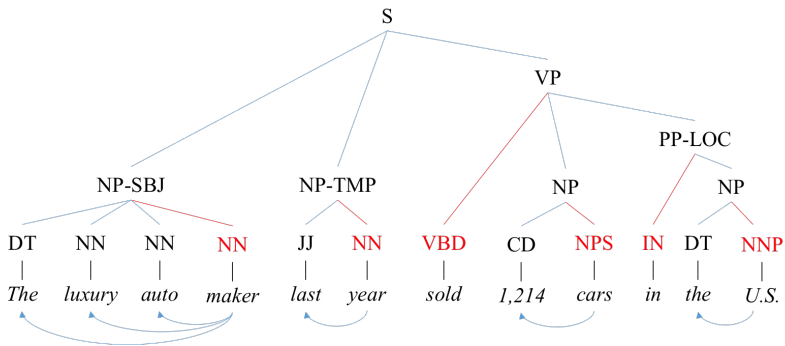
Head and dependencies

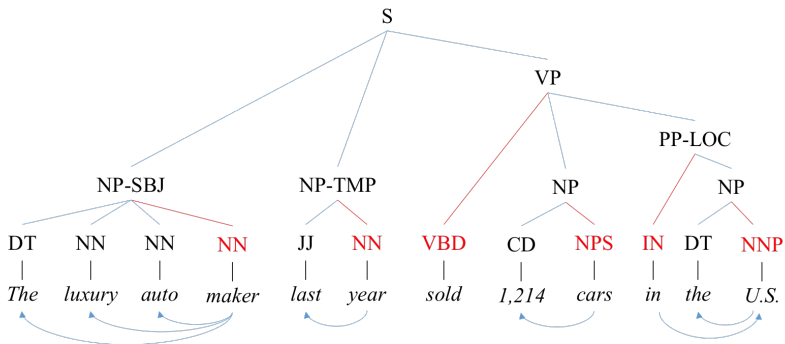


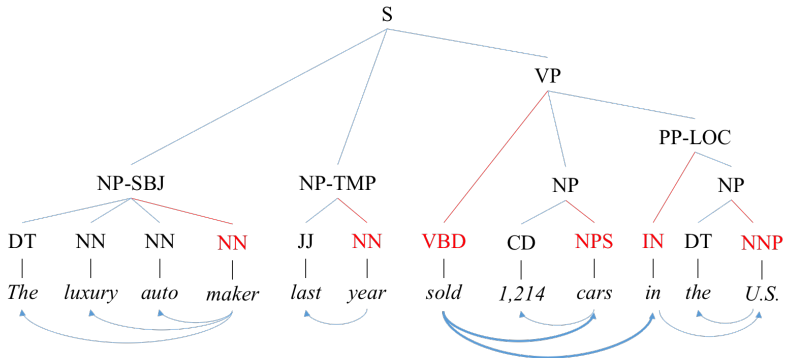
Head percolation

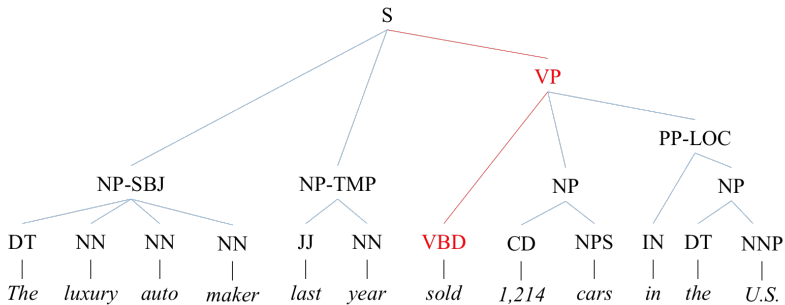
Michael Collins. *Head-Driven Statistical Models for Natural Language Parsing*. PhD

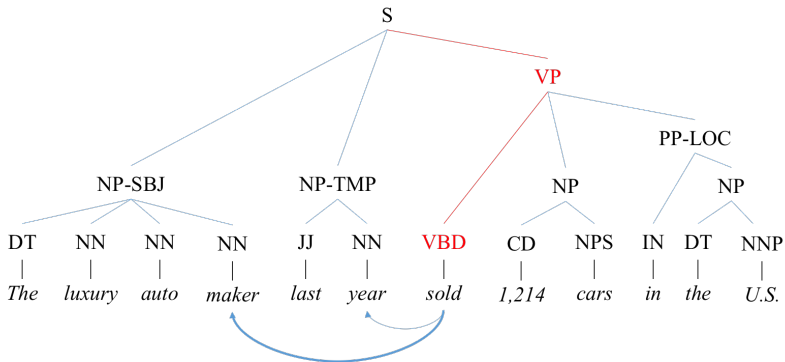
Dissertation, University of Pennsylvania, 1999.



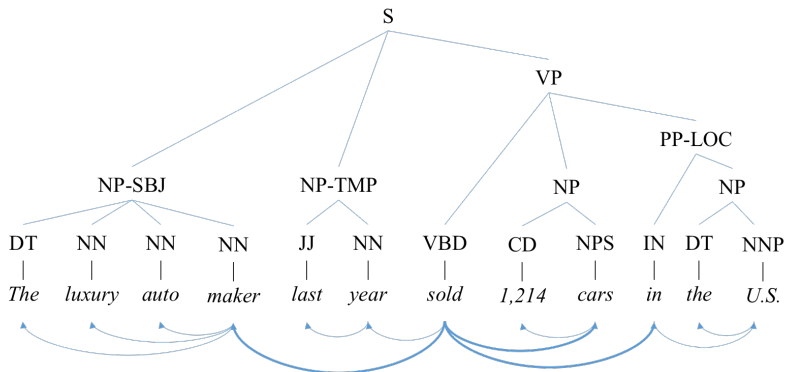




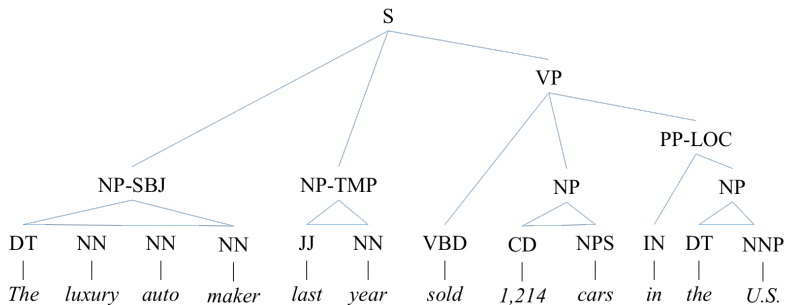




Dependencies



Chunking: Shallow parsing



Abney, Steven (1991), Parsing By Chunks. *Principle-Based Parsing* Kluwer Academic Publishers, pp. 257–278.