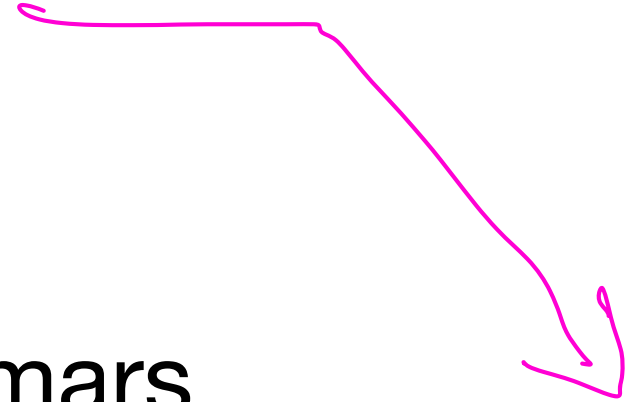


Exercise today!



Context-Free Grammars

CS 485, Fall 2023

Applications of Natural Language Processing

https://people.cs.umass.edu/~brenocon/cs485_f23/

Brendan O'Connor

College of Information and Computer Sciences

University of Massachusetts Amherst

- Syntax: how do words structurally combine to form sentences and meaning?

representations

targeted edges
→ w/o words

- Constituents

- [the big dogs] chase cats
 - [colorless green clouds] chase cats
- NP*

- Dependencies

- The **dog** ~~is~~ **chased** the cat.
- My **dog**, who's getting old, **chased** the cat.

- Idea of a *grammar* (G): global template for how sentences / utterances / phrases w are formed, via latent syntactic structure y

- Linguistics: what do G and $P(w, y | G)$ look like?
- Generation: score with, or sample from, $P(w, y | G)$
- Parsing: predict $P(y | w, G)$

↑
input

Syntax for NLP

- If we could predict syntactic structure from raw text (*parsing*), that could help with...
 - Language understanding: meaning formed from structure
 - Grammar checking
 - Preprocessing: Extract phrases and semantic relationships between words for features, viewing, etc.
- Provides a connection between the theory of *generative linguistics* and computational modeling of language
- Accurate full sentence parsing is challenging!

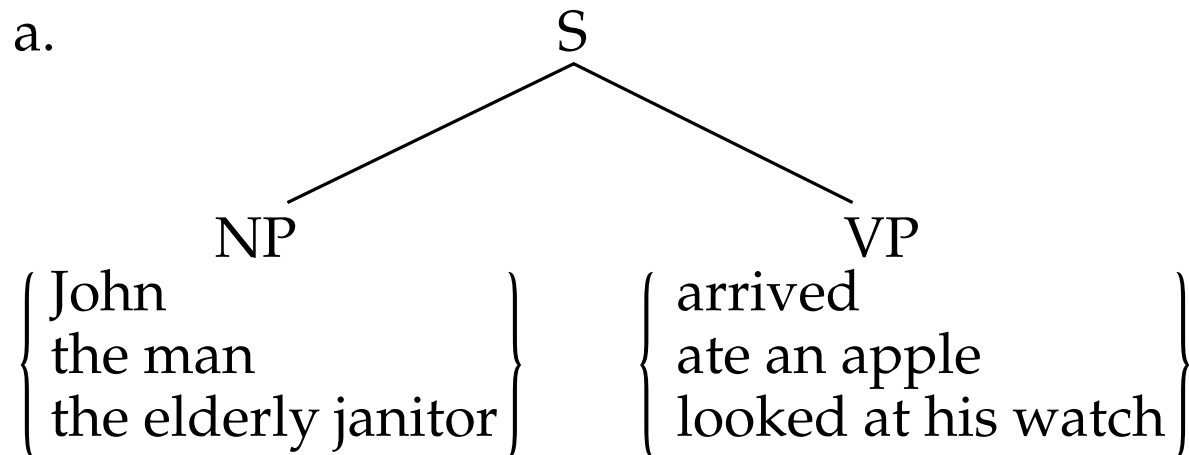


Is language context-free?

- Regular language: repetition of repeated structures
 - e.g. "base noun phrases": (Noun | Adj)* Noun
 - subset of the JK pattern
- Context-free: hierarchical recursion
- Center-embedding: classic theoretical argument for CFG vs. regular languages
 - (10.1) The cat is fat.
 - (10.2) The cat that the dog chased is fat.
 - (10.3) *The cat that the dog is fat.
 - (10.4) The cat that the dog that the monkey kissed chased is fat.
 - (10.5) *The cat that the dog that the monkey chased is fat.
- Competence vs. Performance

Hierarchical view of syntax

- “a Sentence made of Noun Phrase followed by a Verb Phrase”



b. $S \rightarrow NP VP$ (1)

Context-free grammars (CFG)

- A CFG is a 4-tuple:

N a set of non-terminals

Σ a set of terminals (distinct from N)

R a set of productions, each of the form $A \rightarrow \beta$,
where $A \in N$ and $\beta \in (\Sigma \cup N)^*$

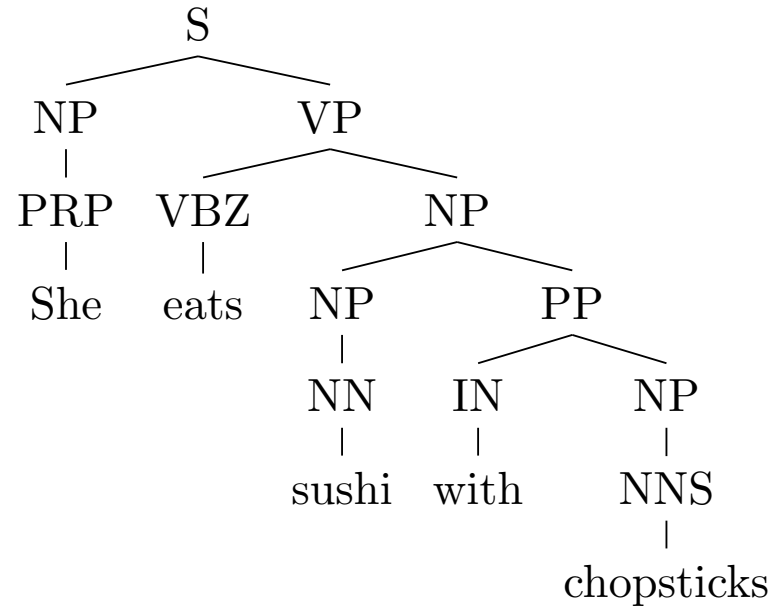
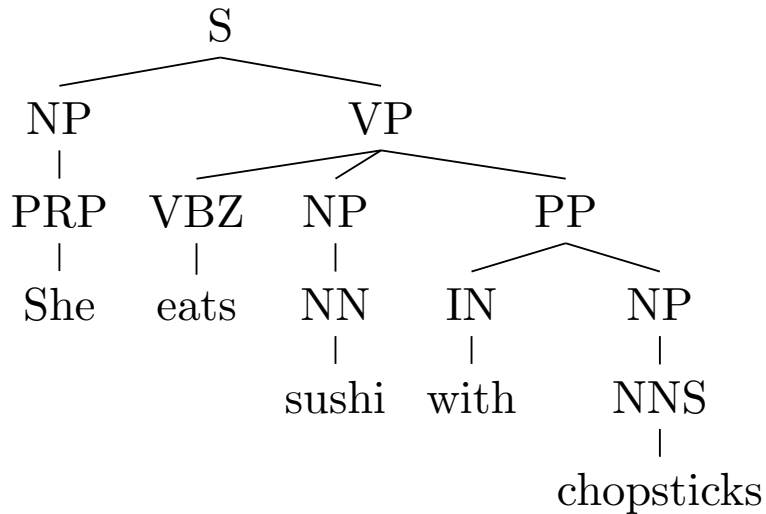
S a designated start symbol

*Example: see
handout!*

- Derivation: a sequence of rewrite steps from S to a string (sequence of terminals, i.e. words)
- Yield: the final string (sentence)
- The parse tree or constituency tree corresponds to the rewrite steps that were used to derive the string
- A CFG is a “boolean language model”
 - A grammar (4-tuple) defines to a set of strings it could generate

- Example: derivation from worksheet's grammar

Example



(S(NP(_{PRP} She))(VP(_{VBZ} eats)

(NP(_{NN} sushi))

(PP(_{IN} with)(NP(_{NNS} chopsticks))))))

(S(NP(_{PRP} She))(VP(_{VBZ} eats)

(NP(NP(_{NN} sushi))(PP(_{IN} with)(NP(_{NNS} chopsticks))))))

- All useful grammars are *ambiguous*: multiple derivations with same yield
- [Parse tree representations: Nested parens or non-terminal spans]

Constituents

- Constituent tree/parse is one representation of sentence's syntax. What should be considered a constituent, or constituents of the same category?
 - Movement tests
 - Substitution tests
 - Coordination tests
- Simple grammar of English
 - Must balance *overgeneration* versus *undergeneration*
 - Noun phrases
 - NP modification: adjectives, PPs
 - Verb phrases
 - Coordination
 - etc...
- Better coverage: machine-learned grammars, if you have a treebank (labeled dataset)

Is language context-free?

- CFGs nicely explain nesting and agreement (if you stuff grammatical features into the non-terminals)
- *The **processor** has 10 million times fewer transistors on it than today's typical micro-processors, runs much more slowly, and operates at five times the voltage...*
- - $S \rightarrow NN VP$
 - $VP \rightarrow VP3S \mid VPN3S \mid \dots$
 - $VP3S \rightarrow VP3S, VP3S, \text{ and } VP3S \mid VBZ \mid VBZ NP \mid \dots$

- **Real sentences have massively ambiguous syntax!**

Attachment ambiguity *we eat sushi with chopsticks, I shot an elephant in my pajamas.*

Modifier scope *southern food store*

Particle versus preposition *The puppy tore up the staircase.*

Complement structure *The tourists objected to the guide that they couldn't hear.*

Coordination scope *"I see," said the blind man, as he picked up the hammer and saw.*

Multiple gap constructions *The chicken is ready to eat*

Penn Treebank

```

( (S
  (NP-SBJ (NNP General) (NNP Electric) (NNP Co.) )
  (VP (VBD said)
    (SBAR (-NONE- 0)
      (S
        (NP-SBJ (PRP it) )
        (VP (VBD signed)
          (NP
            (NP (DT a) (NN contract) )
            (PP (-NONE- *ICH*-3) ))
            (PP (IN with)
              (NP
                (NP (DT the) (NNS developers) )
                (PP (IN of)
                  (NP (DT the) (NNP Ocean) (NNP State) (NNP Power) (NN project) ))))
            (PP-3 (IN for)
              (NP
                (NP (DT the) (JJ second) (NN phase) )
                (PP (IN of)
                  (NP
                    (NP (DT an) (JJ independent)
                      (ADJP
                        (QP ($ $) (CD 400) (CD million) )
                        (-NONE- *U*) )
                      (NN power) (NN plant) )
                    (, ,)
                    (SBAR
                      (WHNP-2 (WDT which) )
                      (S
                        (NP-SBJ-1 (-NONE- *T*-2) )
                        (VP (VBZ is)
                          (VP (VBG being)
                            (VP (VBN built)
                              (NP (-NONE- *-1) )
                              (PP-LOC (IN in)
                                (NP
                                  (NP (NNP Burrillville) )
                                  (, ,)
                                  (NP (NNP R.I) ))))))))))))))))
          )
        )
      )
    )
  )
)

```

Context-Free Grammar

- CFG describes a generative process for an (infinite) set of strings
 - 1. Nonterminal symbols
 - “S”: START symbol / “Sentence” symbol
 - 2. Terminal symbols: word vocabulary
 - 3. Rules (a.k.a. Productions). Practically, two types:

“Grammar”: one NT expands to ≥ 1 NT
always one NT on left side of rule

Lexicon: NT expands to a terminal

S	$\rightarrow NP VP$	I + want a morning flight
NP	\rightarrow <i>Pronoun</i>	I
	<i>Proper-Noun</i>	Los Angeles
	<i>Det Nominal</i>	a + flight
$Nominal$	\rightarrow <i>Nominal Noun</i>	morning + flight
	<i>Noun</i>	flights
VP	\rightarrow <i>Verb</i>	do
	<i>Verb NP</i>	want + a flight
	<i>Verb NP PP</i>	leave + Boston + in the morning
	<i>Verb PP</i>	leaving + on Thursday
PP	\rightarrow <i>Preposition NP</i>	from + Los Angeles

<i>Noun</i>	\rightarrow <i>flights breeze trip morning ...</i>
<i>Verb</i>	\rightarrow <i>is prefer like need want fly</i>
<i>Adjective</i>	\rightarrow <i>cheapest non-stop first latest</i> <i> other direct ...</i>
<i>Pronoun</i>	\rightarrow <i>me I you it ...</i>
<i>Proper-Noun</i>	\rightarrow <i>Alaska Baltimore Los Angeles</i> <i> Chicago United American ...</i>
<i>Determiner</i>	\rightarrow <i>the a an this these that ...</i>
<i>Preposition</i>	\rightarrow <i>from to on near ...</i>
<i>Conjunction</i>	\rightarrow <i>and or but ...</i>

Probabilistic CFGs

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

- Defines a probabilistic generative process for words in a sentence
- Can parse with a modified form of CKY
- How to learn? Fully supervised with a treebank... unsupervised learning possible too, but doesn't give great results...

PCFG as LM

Is a PCFG a *good* LM? Yes...

Is a PCFG a *good* LM? No...