# Context-Free Grammars

CS 485, Spring 2024
Applications of Natural Language Processing
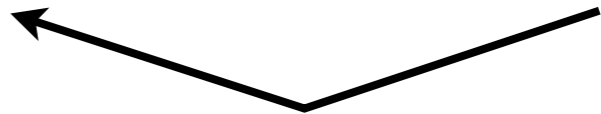https://people.cs.umass.edu/~brenocon/cs485_s24/

Brendan O'Connor

College of Information and Computer Sciences
University of Massachusetts Amherst

- *Do we have to make notes for every single text? Or just have one note encompassing everything? I am a little confused on that part.*

- Just makes notes on any of the texts where it makes sense to do so, in the "annotator notes" column referred to in 1.2 of the HW2 document. There are no requirements for the number of notes.
- The more informative your notes, the better a job you can do in Phase 2 when you (and separately, your groupmates) analyze the differences between annotators.

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

- Constituents
  - [the big dogs] chase cats
  - [colorless green clouds] chase cats

- Dependencies
  - The **dog** ← **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)

# 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

- Practically, accurate full sentence parsing is challenging....
  - ... but the same challenges exist for all NLP tasks/models/systems

# Is language context-free?

[Examples from Eisenstein (2017)]

# Is language context-free?

- Regular language: repetition of repeated structures

[Examples from Eisenstein (2017)]

# Is language context-free?

- Regular language: repetition of repeated structures
  - e.g. "base noun phrases":  (Noun | Adj)* Noun

[Examples from Eisenstein (2017)]

# Is language context-free?

- Regular language: repetition of repeated structures
  - e.g. "base noun phrases":  (Noun | Adj)* Noun
    - subset of the JK pattern

# 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

[Examples from Eisenstein (2017)]

# 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

[Examples from Eisenstein (2017)]

# 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.

# 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.

[Examples from Eisenstein (2017)]

# 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.

[Examples from Eisenstein (2017)]

# 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.

[Examples from Eisenstein (2017)]

# 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.
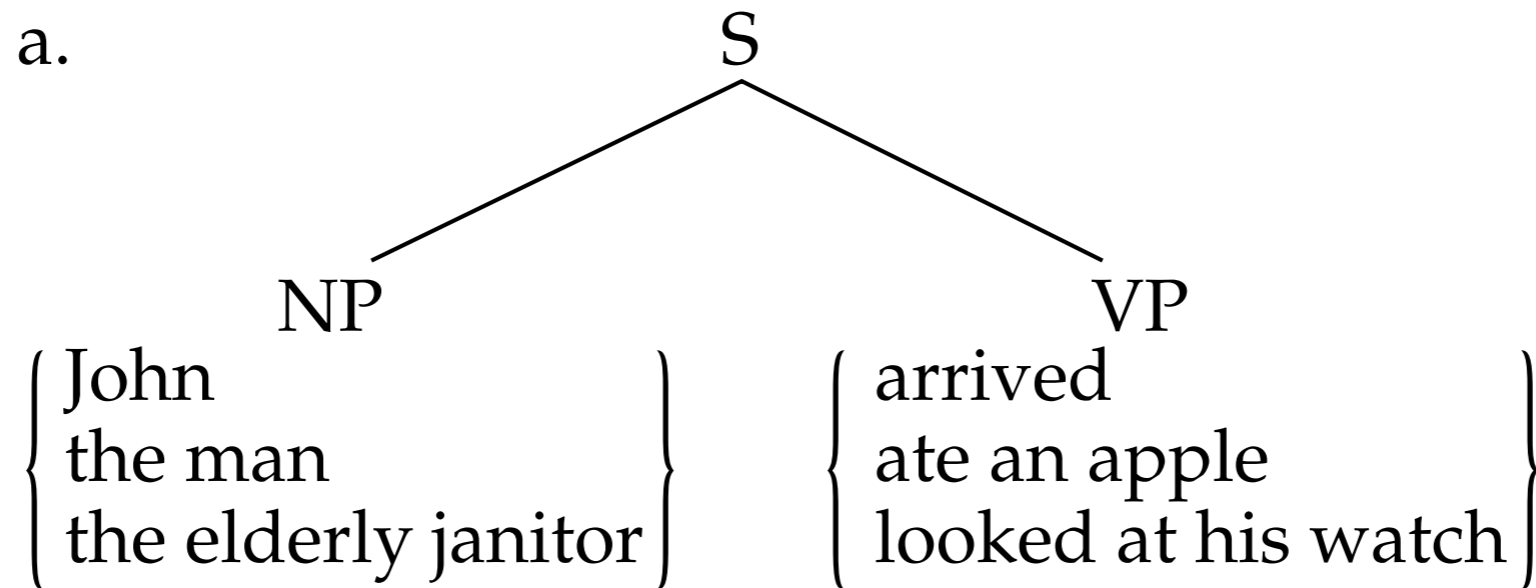
[Examples from Eisenstein (2017)]

# 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

[Examples from Eisenstein (2017)]

# Hierarchical view of syntax

- "a Sentence made of Noun Phrase followed by a Verb Phrase"

a.

$$S \rightarrow \begin{array}{cc} NP & VP \\ \left\{ \begin{array}{l} \text{John} \\ \text{the man} \\ \text{the elderly janitor} \end{array} \right\} & \left\{ \begin{array}{l} \text{arrived} \\ \text{ate an apple} \\ \text{looked at his watch} \end{array} \right\} \end{array}$$

b. $\quad$ S $\quad \rightarrow \quad$ NP VP $\hspace{4cm}$ (1)

# Context-free grammars (CFG)

- A CFG is a 4-tuple:

  $N$   a set of non-terminals

  $\Sigma$   a set of terminals (distinct from $N$)

  $R$   a set of productions, each of the form $A \to \beta$,
     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

# Context-free grammars (CFG)

*R:* production rules typically split into two groups

<u>Core grammar:</u> 1 NT expands to >=1 NT

|  |  |  |
|---|---|---|
| $S \rightarrow$ | $NP\ VP$ | I + want a morning flight |
| $NP \rightarrow$ | $Pronoun$ | I |
| $\mid$ | $Proper\text{-}Noun$ | Los Angeles |
| $\mid$ | $Det\ Nominal$ | a + flight |
| $Nominal \rightarrow$ | $Nominal\ Noun$ | morning + flight |
| $\mid$ | $Noun$ | flights |
| $VP \rightarrow$ | $Verb$ | do |
| $\mid$ | $Verb\ NP$ | want + a flight |
| $\mid$ | $Verb\ NP\ PP$ | leave + Boston + in the morning |
| $\mid$ | $Verb\ PP$ | leaving + on Thursday |
| $PP \rightarrow$ | $Preposition\ NP$ | from + Los Angeles |

<u>Lexicon:</u> NT expands to a terminal

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

- Example: derivation from worksheet's grammar

- Why not?

  S -> ADVP S

# Ambiguity

```
          S
        /   \
      NP      VP
       |     / | \
      PRP  VBZ NP      PP
       |    |   |    /    \
      She  eats NN  IN     NP
                |   |       |
             sushi with    NNS
                            |
                        chopsticks
```

$(_S(_{NP}(_{PRP}$ *She*$)(_{VP}(_{VBZ}$ *eats*$)$

$\qquad (_{NP}(_{NN}$ *sushi*$))$

$\qquad (_{PP}(_{IN}$*with*$)(_{NP}(_{NNS}$ *chopsticks*$))))))$

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

[Examples from Eisenstein (2017)]

# Ambiguity



$(_S(_{NP}(_{PRP}$ *She*$)(_{VP}(_{VBZ}$ *eats*$)$
$\quad (_{NP}(_{NN}$ *sushi*$))$
$\quad (_{PP}(_{IN}$*with*$)(_{NP}(_{NNS}$ *chopsticks*$))))))$

$(_S(_{NP}(_{PRP}$ *She*$)(_{VP}(_{VBZ}$ *eats*$)$
$\quad (_{NP}(_{NP}(_{NN}$ *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]

[Examples from Eisenstein (2017)]

# 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 todays typical micro-processors, **runs** much more slowly, and **operates** at five times the voltage...*

  - $S \rightarrow NN\ VP$
    $VP \rightarrow VP3S\ |\ VPN3S\ |\ \dots$
    $VP3S \rightarrow VP3S, VP3S, and\ VP3S\ |\ VBZ\ |\ VBZ\ NP\ |\ \dots$

[Examples from Eisenstein (2017)]

- ## 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) )))))))))))))))))
```

- stopped here 10/17