# Context-free languages and grammars

September 19, 2019

CS/ECE 374 A

Ian Ludden

# Reminders

- Midterm 1: Monday, Sep 30, 7-9 p.m.
  - DRES: reserve ASAP
  - Review session(s) next week
  - This is the last material that may be covered by the exam

- Homework 3 due next Tuesday

# Learning Objectives

By the end of this lecture, you will be able to:

• Recall the definition of a context-free grammar/language (CFG/CFL).

• Give examples of CFGs/CFLs.

• Derive strings generated by CFGs using parse trees.

• Determine the CFL generated by a CFG.

• Compare/contrast CFLs with regular languages.

• Identify CFGs in Chomsky normal form.

# Context-free = regular + recursion

Regular languages

- Sequencing (A · B)

- Branching (A + B)
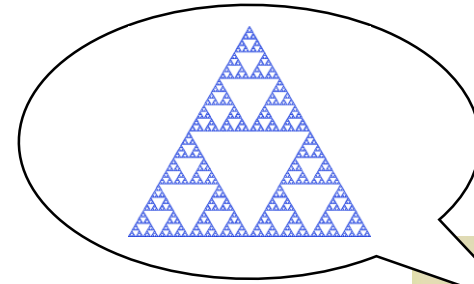
- Repetition (A*)

# Context-free = regular + recursion

~~Regular~~ <span style="color:red">Context-free</span> languages

- Sequencing (A · B)

- Branching (A + B)

- Repetition (A*)

- <span style="color:red">Recursion</span>

All regular languages are context-free.
(Proof in Section 5.5 of lecture notes.)

# Motivation

- Not all languages are regular
  - $L = \{0^n1^n \mid n \geq 0\}$

- Recursive languages occur in nature
  - Gentner, T. Q.; Fenn, K. M.; Margoliash, D.; Nusbaum, H. C. (2006). "Recursive syntactic pattern learning by songbirds." *Nature.* **440** (7088): 1204-1207.

- Natural Language Processing (NLP)
  - Charniak, E. (1997). "Statistical Parsing with a Context-Free Grammar and Word Statistics." *AAAI/IAAI.*

- Probabilistic modeling of RNA structures
  - Sakakibara Y.; Brown M.; Hughey R.; Mian I. S.; et al. (1994). "Stochastic context-free grammars for tRNA modelling." *Nucleic Acids Research.* **22** (23): 5112–5120.

# Formal Definition

- A ***context-free grammar*** is a structure defined by:
  - A finite set Σ of *symbols* or **terminals**

  - A finite set Γ of ***non-terminals*** (disjoint from Σ)

  - A finite set $R$ of ***production rules*** of the form *A -> w*, where A is a non-terminal and w is a string of symbols and non-terminals

  - A starting non-terminal, typically *S*

$G = (Σ, Γ, R, S)$

# Example

Context-free grammar for (a subset of) English sentences
Symbols are words, strings are sentences

⟨sentence⟩ → ⟨noun phrase⟩⟨verb phrase⟩⟨noun phrase⟩
⟨noun phrase⟩ → ⟨adjective phrase⟩⟨noun⟩
⟨adj. phrase⟩ → ⟨article⟩ | ⟨possessive⟩ | ⟨adjective phrase⟩⟨adjective⟩
⟨verb phrase⟩ → ⟨verb⟩ | ⟨adverb⟩⟨verb phrase⟩

⟨noun⟩ → dog | trousers | daughter | nose | homework | time lord | pony | · · ·
⟨article⟩ → the | a | some | every | that | · · ·
⟨possessive⟩ → ⟨noun phrase⟩'s | my | your | his | her | · · ·
⟨adjective⟩ → friendly | furious | moist | green | severed | timey-wimey | little | · · ·
⟨verb⟩ → ate | found | wrote | killed | mangled | saved | invented | broke | · · ·
⟨adverb⟩ → squarely | incompetently | barely | sort of | awkwardly | totally | · · ·

# Example

$\Sigma = \{\textcolor{red}{0}, \textcolor{red}{1}\}$ — Terminals

$\Gamma = \{S, A, B, C\}$ — Non-terminals

$S \to A \mid B$

$A \to \textcolor{red}{0}A \mid \textcolor{red}{0}C$

$B \to B\textcolor{red}{1} \mid C\textcolor{red}{1}$

$C \to \varepsilon \mid \textcolor{red}{0}C\textcolor{red}{1}$

Production rules
'|' means 'or'

$xAy \rightsquigarrow xwy$
(produces immediately)

$S \rightsquigarrow^* w$
(produces eventually)

# Example

Σ = {0, 1}

Γ = {S, A, B, C}

S → A | B

A → 0A | 0C

B → B1 | C1

C → ε | 0C1

S → A

→ 0A

→ 00C

→ 000C1

→ 0000C11

→ 0000ε11

→ 000011

Surely there's a more descriptive way to write this derivation...

# Parse trees visualize string derivations.

Σ = {0, 1}
Γ = {S, A, B, C}

S → A | B
A → 0A | 0C
B → B1 | C1
C → ε | 0C1

⟨sentence⟩
├─ ⟨noun phrase⟩
│  ├─ ⟨adj. phrase⟩
│  │  ├─ ⟨adj. phrase⟩
│  │  │  ├─ ⟨adj. phrase⟩
│  │  │  │  └─ ⟨posessive⟩ — your
│  │  │  └─ ⟨adjective⟩ — furious
│  │  └─ ⟨adjective⟩ — green
│  └─ ⟨noun⟩ — time lord
├─ ⟨verb phrase⟩
│  ├─ ⟨adverb⟩ — barely
│  └─ ⟨verb phrase⟩
│     └─ ⟨verb⟩ — mangled
└─ ⟨noun phrase⟩
   ├─ ⟨adj. phrase⟩
   │  └─ ⟨posessive⟩
   │     ├─ ⟨noun phrase⟩
   │     │  ├─ ⟨adj. phrase⟩
   │     │  │  └─ ⟨possessive⟩ — my
   │     │  └─ ⟨noun⟩ — dog
   │     └─ 's
   └─ ⟨noun⟩ — trousers

# Exercise: Parse trees

$\Sigma = \{1, 2, +, x\}$

$\Gamma = \{S, A, M, C\}$

$S \rightarrow A \mid M \mid 1 \mid 2$

$A \rightarrow S + S$

$M \rightarrow S \; x \; S$

**Activity (2 min.)**

1. Derive $2 + 1 \; x \; 1$ from this grammar using a parse tree.
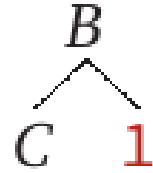
2. Compare with neighbor(s).

# Exercise: Parse trees

$\Sigma = \{1, 2, +, x\}$
$\Gamma = \{S, A, M, C\}$

$S \rightarrow A \mid M \mid 1 \mid 2$
$A \rightarrow S + S$
$M \rightarrow S \times S$



2 + (1 x 1) = 4

(2 + 1) x 1 = 3

# Ambiguity

- A string w is ***ambiguous*** with respect to a grammar if there is more than one parse tree for w.

- A grammar G is ***ambiguous*** if some string is ambiguous with respect to G.

- A context-free language L is ***inherently ambiguous*** if every context-free grammar that generates L is ambiguous.
  (Contrived examples)

# Disambiguating

$\Sigma = \{1, 2, +, x\}$
$\Gamma = \{S, A, M, C\}$

$S \rightarrow A \mid M \mid 1 \mid 2$
$A \rightarrow S + S$
$M \rightarrow S \times S$

$\Longrightarrow$
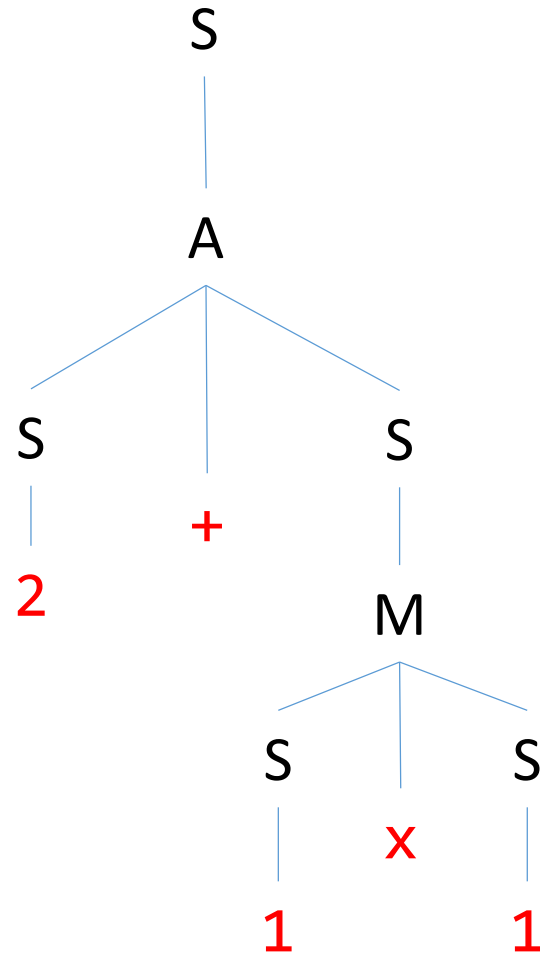
$\Sigma = \{1, 2, +, x, (, )\}$
$\Gamma = \{S, A, M, C\}$

$S \rightarrow A \mid M \mid 1 \mid 2$
$A \rightarrow (S + S)$
$M \rightarrow (S \times S)$

No longer ambiguous!

# Arithmetic Expressions

- Arithmetic expressions, possibly with redundant parentheses, over the variables X and Y:

$$E \rightarrow E+T \mid T \qquad \text{(expressions)}$$
$$T \rightarrow T\times F \mid F \qquad \text{(terms)}$$
$$F \rightarrow (E) \mid X \mid Y \qquad \text{(factors)}$$

Every $E$expression is a sum of $T$erms, every $T$erm is a product of $F$actors, and every $F$actor is either a variable or a parenthesized $E$expression.

# Regular Expressions

- Regular expressions over the alphabet $\{0, 1\}$ *without* redundant parentheses

$$S \to T \mid T+S \qquad \text{(Regular expressions)}$$

$$T \to F \mid FT \qquad \text{(Terms = summable expressions)}$$

$$F \to \emptyset \mid W \mid (T+S) \mid X\star \mid (Y)\star \qquad \text{(Factors = concatenable expressions)}$$

$$X \to \emptyset \mid \varepsilon \mid 0 \mid 1 \qquad \text{(Directly starrable expressions)}$$

$$Y \to T+S \mid F \bullet T \mid X\star \mid (Y)\star \mid ZZ \qquad \text{(Starrable expressions needing parens)}$$

$$W \to \varepsilon \mid Z \qquad \text{(Words = strings)}$$

$$Z \to 0 \mid 1 \mid ZZ \qquad \text{(Non-empty strings)}$$

# From grammars to languages

- For non-terminal A, L(A) is the set of all strings generated by A.

- Given context-free grammar $G = (\Sigma, \Gamma, R, S)$, L($G$) = L($S$).

- A **context-free language** is the language generated by a context-free grammar.

- Often easiest to examine "later" rules first when determining L($G$).

# What language do you speak?

$\Sigma = \{0, 1\}$

$\Gamma = \{S, A, B, C\}$

$S \rightarrow A \mid B$

$A \rightarrow 0A \mid 0C$

$B \rightarrow B1 \mid C1$

$C \rightarrow \varepsilon \mid 0C1$

<u>Lemma</u>: $L(C) = \{0^m1^n \mid m = n \geq 0\}$.

<u>Proof</u> ($\supseteq$):

Let n be an arbitrary non-negative integer.

Assume $C \rightsquigarrow^* 0^m1^m$ for all $m < n$.

Two cases:

- $n = 0$, $0^n1^n = \varepsilon$, $C \rightarrow \varepsilon$, done.
- $n > 1$, $C \rightarrow 0C1$. By I.H.,

$0C1 \rightsquigarrow^* 0(0^{n-1}1^{n-1})1 = 0^n1^n$, done.

Thus $L(C) \supseteq \{0^m1^n \mid m = n \geq 0\}$.

# What language do you speak?

$\Sigma = \{0, 1\}$

$\Gamma = \{S, A, B, C\}$

$S \rightarrow A \mid B$

$A \rightarrow 0A \mid 0C$

$B \rightarrow B1 \mid C1$

$C \rightarrow \varepsilon \mid 0C1$

<u>Lemma</u>: $L(C) = \{0^m1^n \mid m = n \geq 0\}$.

<u>Proof</u> ($\subseteq$):

Fix w in L(C).

Assume for all x in L(C) with $|x| < |w|$,

$x = 0^m1^m$ for some $m \geq 0$.

Two cases (first production):

- $C \rightarrow \varepsilon$, $w = \varepsilon = 0^01^0$, done.

- $C \rightarrow 0C1$ So w = 0x1 for some x in L(C). By I.H., x $= 0^m1^m$ for some $m \geq 0$.

So $w = 0(0^m1^m)1 = 0^{m+1}1^{m+1}$, done.

Thus $L(C) \subseteq \{0^m1^n \mid m = n \geq 0\}$.

# What language do you speak?

$\Sigma = \{0, 1\}$

$\Gamma = \{S, A, B, C\}$

$S \to A \mid B$

$A \to 0A \mid 0C$

$B \to B1 \mid C1$

$C \to \varepsilon \mid 0C1$

$L(C) = \{0^m 1^n \mid m = n \geq 0\}$.

$L(B) = ?$

$L(A) = ?$

$L(S) = ?$

# What language do you speak?

$\Sigma = \{0, 1\}$

$\Gamma = \{S, A, B, C\}$

$S \rightarrow A \mid B$

$A \rightarrow 0A \mid 0C$

$B \rightarrow B1 \mid C1$

$C \rightarrow \varepsilon \mid 0C1$

$L(C) = \{0^m 1^n \mid m = n \geq 0\}.$

$L(B) = \{0^m 1^n \mid m < n \geq 0\}.$

$L(A) = ?$

$L(S) = ?$

# What language do you speak?

$\Sigma = \{0, 1\}$
$\Gamma = \{S, A, B, C\}$

$S \to A \mid B$
$A \to 0A \mid 0C$
$B \to B1 \mid C1$
$C \to \varepsilon \mid 0C1$

$L(C) = \{0^m1^n \mid m = n \geq 0\}$.

$L(B) = \{0^m1^n \mid m < n \geq 0\}$.

$L(A) = \{0^m1^n \mid m > n \geq 0\}$.

$L(S) = ?$

# What language do you speak?

$\Sigma = \{0, 1\}$

$\Gamma = \{S, A, B, C\}$

$S \rightarrow A \mid B$

$A \rightarrow 0A \mid 0C$

$B \rightarrow B1 \mid C1$

$C \rightarrow \varepsilon \mid 0C1$

$L(C) = \{0^m1^n \mid m = n \geq 0\}.$

$L(B) = \{0^m1^n \mid m < n \geq 0\}.$

$L(A) = \{0^m1^n \mid m > n \geq 0\}.$

$L(S) = \{0^m1^n \mid m \neq n \geq 0\}.$

# The grammar that generates a CFL is not unique.

$\Sigma$ = {0, 1}

How to generate 0*1*?

S → ε | 0S | S1        vs.        S → AB
                                   A → ε | 0A
                                   B → ε | 1B

# More fun CFLs/CFGs

- Binary palindromes:

$S \rightarrow$ 0S0 | 1S1 | 0 | 1 | ε

- Binary strings with same number of 0s and 1s:

$S \rightarrow$ 0S1 | 1S0 | SS | ε         (This is HW 0.3)

or... $S \rightarrow$ 0S1S | 1S0S | ε

- Balanced strings of parentheses:

$S \rightarrow$ (S) | SS | ε                 or        $S \rightarrow$ (S)S | ε

# Are all languages context-free?

- **No**. Canonical example: $L = \{0^n1^n0^n \mid n \geq 0\}$ is **not** context-free.

(To get a feel for why, try to create a context-free grammar that generates L.)

- Counting argument: The set of possible CFGs over $\Sigma$ is *countably* infinite, but the set of all languages over $\Sigma$ is *uncountably* infinite.

- There are also techniques for proving a specific language is not context-free. If curious, search "pumping lemma."

# Chomsky Normal Form (CNF)

- Developed by Noam Chomsky in 1959

- A context-free grammar is in CNF if every rule is one of:

$A \rightarrow BC$  (A can be S, but neither B nor C can be S)

$A \rightarrow$ <span style="color:red">a</span>

$S \rightarrow \varepsilon$     (only if $\varepsilon$ is in L(G))

# Why are CNF grammars nice/useful?

- Full binary parse trees
  - Easy brute-force check to see if a given string can be generated

- Simple structure makes proofs easier

- Assumed by popular parsing algorithms

# Every CFG has a CNF equivalent.

- By "equivalent," we mean "defines the same language."

- Lecture notes: 5.9 CNF Conversion Algorithm

- No more than quadratic size increase

# Recap: Learning Objectives

By the end of this lecture, you will be able to:

- Recall the definition of a context-free grammar/language (CFG/CFL).
- Give examples of CFGs/CFLs.
- Derive strings generated by CFGs using parse trees.
- Determine the CFL generated by a CFG.
- Compare/contrast CFLs with regular languages.
- Identify CFGs in Chomsky normal form.

# Recap: Learning Objectives

By the end of this lecture, you will be able to:

• Recall the definition of a context-free grammar/language (CFG/CFL).

• Give examples of CFGs/CFLs.

• Derive strings generated by CFGs using parse trees.

• Determine the CFL generated by a CFG.

• Compare/contrast CFLs with regular languages.

• Identify CFGs in Chomsky normal form.

By the end of tomorrow's lab, you will be able to:

• Construct and describe CFGs that generate given languages.