

Specialized Grammars

Linear Grammars and Chomsky Normal Form

Monday, November 5, 2007

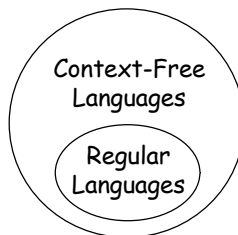
Reading: Stoughton 4.4, 4.8-4.9; Sipser 2.1; Kozen 21

CS235 Languages and Automata

Department of Computer Science
Wellesley College

Overview

1. Introduce **right-linear** and **left-linear** grammars and use these to prove that regular languages are a proper subset of context-free languages.



2. Introduce **Chomsky Normal Form**, and show how to convert any context-free grammar to Chomsky Normal form.

Linear Grammars

A CFG is **linear** iff every production has at most one variable in its RHS.

A CFG is **right-linear** iff every production has one of these two forms, where V and W are any variables and x and y are any string of terminals.

$$\begin{array}{l} V \rightarrow x \\ V \rightarrow yW \end{array}$$

A CFG is **left-linear** iff every production has one of these two forms, where V and W are any variables and x and y are any string of terminals.

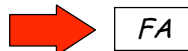
$$\begin{array}{l} V \rightarrow x \\ V \rightarrow Wy \end{array}$$

Specialized Grammars 26-3

Right-Linear CFGs Generate All Regular Languages

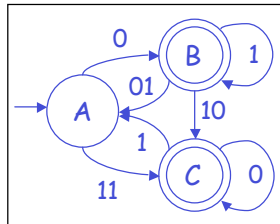
- (1) We'll show that every right-linear grammar can be converted to an FA, and so generates a regular language.

$$\begin{array}{l} S \rightarrow \epsilon \mid 0S \mid 00P \mid 11Q \\ P \rightarrow 01P \mid 0Q \mid 1R \\ Q \rightarrow \epsilon \mid 0S \mid 1Q \mid 00R \\ R \rightarrow \epsilon \mid 10R \end{array}$$

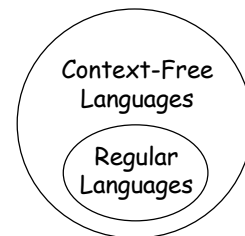


FA

- (2) We'll show that every FA (and thus every regular language) can be converted to a right-linear grammar.



right-linear grammar



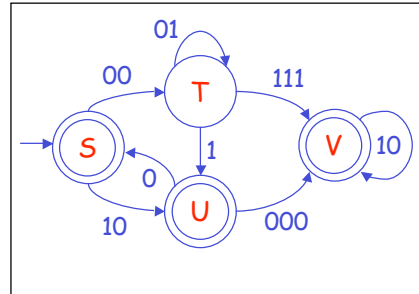
Together, (1) and (2) imply that regular languages are a subset of the context-free languages.

Non-regular CFLs (like 0^n1^n) show that the subset relation is proper.

Specialized Grammars 26-4

Converting Right-Linear Grammars to FAs

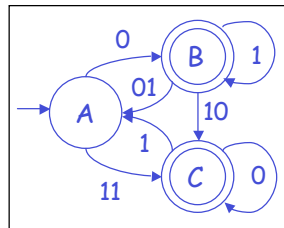
$S \rightarrow \% \mid 00T \mid 10U$
 $T \rightarrow 01T \mid 1U \mid 111V$
 $U \rightarrow \% \mid 0S \mid 000V$
 $V \rightarrow \% \mid 11V$



- (1) The **states** of the FA are named by the **variables** of the CFG.
- (2) The **start state** is the state named by the **start variable**.
- (3) A state **Q** is an **accepting state** iff there is a production $Q \rightarrow \%$.
- (4) There is a **transition** (P, x, Q) for each production $P \rightarrow xQ$.

Specialized Grammars 26-5

Converting FAs to Right-Linear Grammars



$A \rightarrow 0B \mid 11C$
 $B \rightarrow \% \mid 01A \mid 1B \mid 10C$
 $C \rightarrow \% \mid 1A \mid 0C$

- (1) The **variables** of the CFG are the states of the CFG.
- (2) The **nonterminals** of the CFG are the transition symbols in the FA.
- (3) The **start variable** of the CFG is the FA's start state.
- (4a) There is a **production** $Q \rightarrow \%$ for each accepting state **Q**.
- (4b) There is a **production** $P \rightarrow xQ$ for each transition (P, x, Q) .

Specialized Grammars 26-6

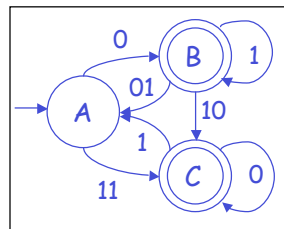
What About Left-Linear Grammars?

We have seen that a language is regular iff it can be expressed via a right-linear grammar.

It turns out that the same holds for left-linear grammars:

A language is regular iff it can be expressed via a left-linear grammar.

(You will work out the details on PS5.)



$$\begin{aligned} A &\rightarrow A10 \mid B1 \mid C0 \\ B &\rightarrow \% \mid A11 \mid B01 \mid C101 \\ C &\rightarrow \% \mid B010 \mid C00 \end{aligned}$$

left-linear grammar



Careful! Regularity is not guaranteed if right-linear and left-linear productions are mixed! E.g.:

$$\begin{aligned} A &\rightarrow \% \mid 0B \\ B &\rightarrow A1 \end{aligned}$$

Specialized Grammars 26-7

Chomsky Normal Form

Sometimes it's helpful to require a CFG to be in a standard form. E.g., we'll see this soon in regards to a *pumping lemma for CFLs*.

One such form is **Chomsky Normal Form (CNF)**, in which all productions must have one of the following two forms:

$$\begin{aligned} V &\rightarrow UW && \text{Variable rewrites to two variables; } U, V \text{ can't be start variable} \\ V &\rightarrow t && \text{Variable rewrites to a single terminal} \end{aligned}$$

Chomsky Normal Form can generate any CFL not containing $\%$. In order to allow languages with $\%$, we also allow the production:

$$S \rightarrow \% \quad \% \text{ production allowed only for start variable } S$$

Intuitions:

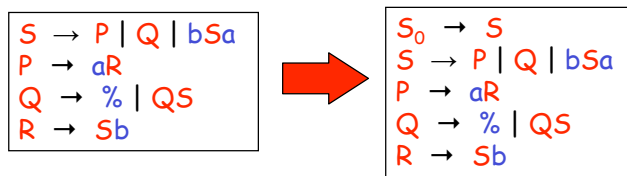
- Parse trees for CNF have variables arranged in binary trees.
- Every step in a derivation from a CNF grammar makes nontrivial progress toward the terminal string. Can't have subtrees yielding $\%$ or long sequences of unit productions: $A \rightarrow B \rightarrow C \rightarrow \dots$

Specialized Grammars 26-8

CFG to CNF, Step 1: Add New Start Variable*

Introduce a new start variable that rewrites to the given one. (This guarantees that the new start variable does not occur in a RHS.)

Our running example (what language does it generate?):



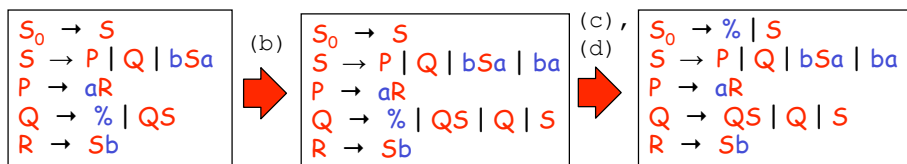
* Sipser does this, but Stoughton and Kozen do not (because they don't handle languages containing %).

Specialized Grammars 26-9

CFG to CNF, Step 2: Remove Nullable Variables

- (a) Find all **nullable variables** = those variables that can yield %.
- (b) For each nullable V and production $W \rightarrow \alpha V \beta$, (where at least one of α or β is nonempty) add a production $W \rightarrow \alpha \beta$.
- (c) If S_0 is nullable, add the production $S_0 \rightarrow \%$.
- (d) Remove all productions of the form $V \rightarrow \%$, where $V \neq S_0$.

Our running example:



(a) Nullable variables
= $\{S_0, S, Q\}$

Specialized Grammars 26-10

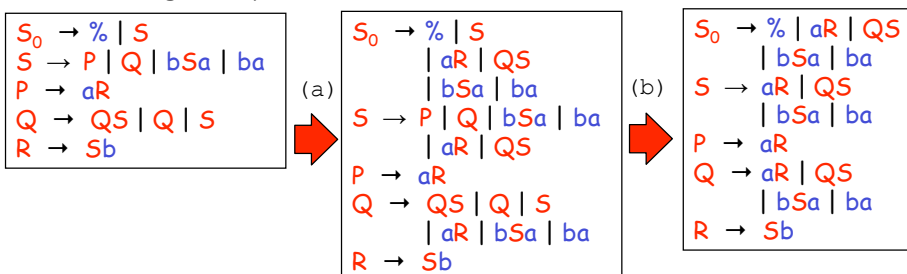
CFG to CNF, Step 3: Remove Unit Productions

A **unit production** is one of the form $V \rightarrow W$.

(a) For each rewrite sequence $V_1 \Rightarrow^* V_n \Rightarrow \alpha$, where α is not a single variable, add a production $V_1 \rightarrow \alpha$ (if it doesn't already exist).

(b) Remove all unit productions.

Our running example:



Specialized Grammars 26-11

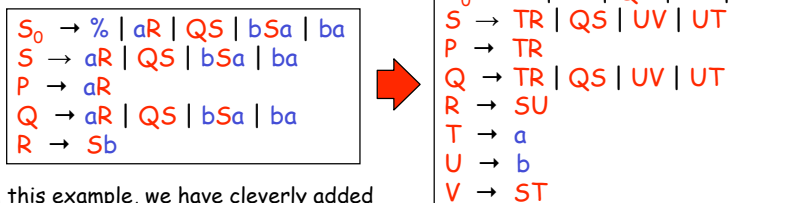
CFG to CNF, Step 4: Completing CNF

Every production now has one of the following three forms:

- (1) $V \rightarrow t$
- (2) $V \rightarrow \alpha$, where α has at least two variables and/or terminals.
- (3) $S_0 \rightarrow \%$

Introduce new variables and productions to replace all productions of form (2) not in CNF by CNF productions.

Our running example:



- In this example, we have cleverly added only three new variables, but this is not required and a straightforward implementation would add **many** more.

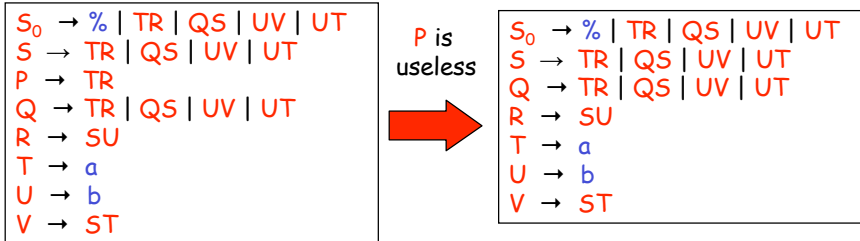
Specialized Grammars 26-12

Simplification

A variable in a CFG is **useless** if it can never appear in a parse tree rooted at the start symbol.

It is always safe to simplify a CFG by removing all productions that mention a useless variable.

Our running example:



See Stoughton 4.4 for details on simplification.

Specialized Grammars 26-13

Greibach Normal Form (GNF)

Another normal form for CFGs is **Greibach Normal Form (GNF)**, in which every production (except for $S_0 \rightarrow \%$) has the form:

$$V_0 \rightarrow \dagger V_1 V_2 \dots V_n, n \geq 0$$

GNF has the nice property that every rewrite (except $S_0 \rightarrow \%$) makes progress toward the final string by adding a terminal.

Idea: To obtain GNF, start with CNF, and create a GNF production of the above form for every collection of CNF productions with the following form:

$V_0 \rightarrow W_n V_n$
$W_n \rightarrow W_{n-1} V_{n-1}$
...
$W_2 \rightarrow W_1 V_1$
$W_1 \rightarrow \dagger$

See Kozen Lecture 21 for details.

Specialized Grammars 26-14