

Session 17

Chomsky Normal Form

Chomsky Normal Form (CNF)

- For any CFG $G = (V, \Sigma, S, P)$ there is a CFG $G' = (V', \Sigma, S, P')$ in CNF so that $L(G') = L(G) - \{\Lambda\}$
- A CFG is in CNF if every production is of one of these two type:

$$A \rightarrow BC$$

$$A \rightarrow a$$

where A, B and C are variables, and a is a terminal symbol

- If a grammar is unambiguous (i.e. it is already unambiguous or there is an unambiguous grammar generating the same language) its corresponding grammar in CNF is also unambiguous!

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Chomsky's Normal Form

- Consider $G = (V, \Sigma, S, P)$ and $S \Rightarrow^* x$ where $x \in \Sigma^*$ and $|x| = k$
 - Let l be the length of a string
 - Let t be the number of terminal symbols
 - For $S : l + t = 1 + 0 = 1$
 - For $x : l + t = k + k = 2k$
- If there are no Λ -productions (i.e. of form $A \rightarrow \Lambda$) and Unit productions (i.e. of form $T \rightarrow F$), for any derivation $\alpha \Rightarrow \beta$:
 - $l + t$ of $\beta > l + t$ of α
 - β has either more variables, increasing l , or more terminals, increasing t , or both

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Chomsky's Normal Form

- An interesting property:
 - If there are no Λ -productions (i.e. of form $A \rightarrow \Lambda$) and Unit productions (i.e. of form $T \rightarrow F$), for any derivation $\alpha \Rightarrow \beta$: the value of $l + t$ is increased by rewriting a variable by a production of form:

$$A \rightarrow \gamma \quad \text{where } \gamma \in (V \cup \Sigma)^*$$
 - In particular, $l + t$ increases by one if the productions have the form:

$$A \rightarrow BC \quad (\text{i.e. } l \text{ is increased in one})$$

$$A \rightarrow a \quad (\text{i.e. } t \text{ is increased in one})$$
 - So, a derivation $S \Rightarrow^* x$ (from $l + t = 1$ to $l + t = 2k$) has at most $2k - 1$ productions!

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Removing Λ -productions:

- A nullable variable in a CFG $G = (V, \Sigma, S, P)$ is defined as follows:
 - If there is a production of form $A \rightarrow \Lambda$ in P then A is nullable
 - If P contains the production $A \rightarrow B_1B_2\dots B_n$ where $B_1B_2\dots B_n$ are nullable then A is nullable
 - No other variables are nullable

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Removing Λ -productions:

- Given CFG $G = (V, \Sigma, S, P)$ construct a CFG $G_1 = (V, \Sigma, S, P_1)$ with no Λ -productions as follows:
 - Let $P_1 = P$
 - Find all nullable variables in V
 - For every production $A \rightarrow \alpha$ in P , augment P_1 with every production that can be obtained from $A \rightarrow \alpha$ by deleting one or more occurrences of nullable variables in α
 - Delete all Λ -productions from P_1 , duplications of a production and productions of form $A \rightarrow A$

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Removing Λ -productions:

- Let the CFG $G = (V, \Sigma, S, P)$ where P are the productions:

$$\begin{aligned} S &\rightarrow AACD \\ A &\rightarrow aAb \mid \Lambda \\ C &\rightarrow aC \mid a \\ D &\rightarrow aDa \mid bDb \mid \Lambda \end{aligned}$$
- Eliminating Λ -transitions: nullable variables are A and D :

$$\begin{aligned} S &\rightarrow AACD \mid ACD \mid AAC \mid CD \mid AC \mid C \\ A &\rightarrow aAb \mid ab \\ C &\rightarrow aC \mid a \\ D &\rightarrow aDa \mid bDb \mid aa \mid bb \end{aligned}$$
- Eliminating nullvariables in a CFG is like eliminating Λ -transitions in a NFA- Λ

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Removing unit-productions:

- Let the CFG $G = (V, \Sigma, S, P)$ where P has no Λ -productions:

$$\begin{aligned} S &\rightarrow AACD \mid ACD \mid AAC \mid CD \mid AC \mid C \\ A &\rightarrow aAb \mid ab \\ C &\rightarrow aC \mid a \\ D &\rightarrow aDa \mid bDb \mid aa \mid bb \end{aligned}$$
- Eliminating unit-productions: $S \rightarrow C$

$$\begin{aligned} S &\rightarrow AACD \mid ACD \mid AAC \mid CD \mid AC \mid aC \mid a \\ A &\rightarrow aAb \mid ab \\ C &\rightarrow aC \mid a \\ D &\rightarrow aDa \mid bDb \mid aa \mid bb \end{aligned}$$

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Normalizing form of productions

- Right sides only variables or only terminals:

$$\begin{aligned} S &\rightarrow AACD \mid ACD \mid AAC \mid CD \mid AC \mid aC \mid a \\ A &\rightarrow aAb \mid ab \\ C &\rightarrow aC \mid a \\ D &\rightarrow aDa \mid bDb \mid aa \mid bb \end{aligned}$$
- Replace $S \rightarrow aC$ by $S \rightarrow X_aC$ and $X_a \rightarrow a$:

$$\begin{aligned} S &\rightarrow AACD \mid ACD \mid AAC \mid CD \mid AC \mid X_aC \mid a \\ A &\rightarrow X_aX_b \mid X_aX_b \\ C &\rightarrow X_aC \mid a \\ D &\rightarrow X_aDX_a \mid X_bDX_b \mid X_aX_a \mid X_bX_b \\ X_a &\rightarrow a \\ X_b &\rightarrow b \end{aligned}$$

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Obtain Chomsky normal form

- Replace $S \rightarrow ABC\alpha$ by $S \rightarrow AT$ and $T \rightarrow BC\alpha$

$$\begin{aligned} S &\rightarrow AACD \mid ACD \mid AAC \mid CD \mid AC \mid X_aC \mid a \\ A &\rightarrow X_aX_b \mid X_aX_b \\ C &\rightarrow X_aC \mid a \\ D &\rightarrow X_aDX_a \mid X_bDX_b \mid X_aX_a \mid X_bX_b \\ X_a &\rightarrow a \\ X_b &\rightarrow b \end{aligned}$$
- Obtain Chomsky Normal Form:

$$\begin{aligned} S &\rightarrow AT_1 & T_1 &\rightarrow AT_2 & T_2 &\rightarrow CD \\ S &\rightarrow AU_1 & U_1 &\rightarrow CD \\ S &\rightarrow AV_1 & V_1 &\rightarrow AC \end{aligned}$$

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

The Grammar in CNF:

$$\begin{array}{lll} S \rightarrow AT_1 & T_1 \rightarrow AT_2 & T_2 \rightarrow CD \\ S \rightarrow AU_1 & U_1 \rightarrow CD \\ S \rightarrow AV_1 & V_1 \rightarrow AC \\ S \rightarrow CD \mid AC \mid X_aC \mid a \\ A \rightarrow X_aW_1 & W_1 \rightarrow AX_b & A \rightarrow X_aX_b \\ C \rightarrow X_aC \mid a \\ D \rightarrow X_aY_1 & Y_1 \rightarrow DX_a \\ D \rightarrow X_bZ_1 & Z_1 \rightarrow DX_b \\ D \rightarrow X_aX_a \mid X_bX_b \\ X_a \rightarrow a & & \\ X_b \rightarrow b & & \end{array}$$

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

An analogy!

- “Ambiguous” CFG grammar \iff RE or NFA- Λ
- Nullable variables \iff Λ -closure
- Removing Λ -productions \iff NFA
e.g. $C \rightarrow aC \mid a$
- CNF \iff DFA

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Ambiguity as a device to express abstraction!

“Ambiguous” CFG grammar \longleftrightarrow RE or NFA- Λ

Nullable variables \longleftrightarrow Λ -closure

Removing Λ -productions \longleftrightarrow NFA
 $C \rightarrow aC \mid \Lambda \Rightarrow C \rightarrow aC \mid a$

The implementation: CNF \longleftrightarrow DFA

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Chomsky's Normal Form

- Chomsky normal form (CNF):

$$A \rightarrow BC$$

$$A \rightarrow a$$

- Regular grammar:

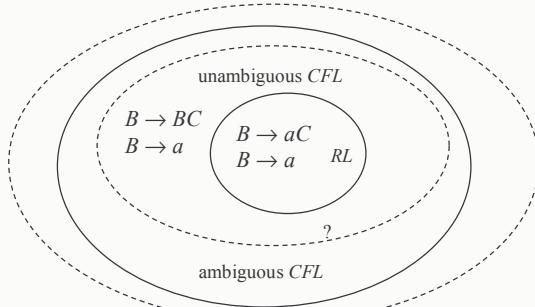
$$B \rightarrow aC$$

$$B \rightarrow a$$

- From regular languages to unambiguous CFL (almost!)

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Is there a class of ambiguous CFL



- There is no algorithm to tell whether a grammar is ambiguous
- There is no way to tell when a language is inherently ambiguous!

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003