

Chomsky normal form EXAMPLE

$$G: \begin{array}{l} S \rightarrow YZ \mid ZZZ \\ X \rightarrow \epsilon \\ Y \rightarrow XY \\ Z \rightarrow \epsilon \mid aZ \end{array}$$

① Eliminate useless symbols

ⓐ Eliminate symbols that cannot generate a terminal string since Y is not generating

$$G': \begin{array}{l} S \rightarrow ZZZ \\ X \rightarrow \epsilon \\ Z \rightarrow \epsilon \mid aZ \end{array}$$

ⓑ Eliminate symbols that are not reachable (i.e., X)

$$G'': \begin{array}{l} S \rightarrow ZZZ \\ Z \rightarrow \epsilon \mid aZ \end{array}$$

ⓒ Eliminate ϵ productions ($Z \rightarrow \epsilon, S \xrightarrow{*} \epsilon$)

$$G''': \begin{array}{l} S \rightarrow ZZZ \mid ZZ \mid Z \\ Z \rightarrow aZ \mid a \end{array} \quad \text{Note } L(G''') = L(G) \setminus \{\epsilon\}$$

ⓓ Eliminate unit productions ($S \rightarrow Z$)

$$G''': \begin{array}{l} S \rightarrow ZZZ \mid ZZ \mid aZ \mid a \\ Z \rightarrow aZ \mid a \end{array}$$

ⓔ Convert to CNF by splitting larger productions

$$G'''': \begin{array}{l} S \rightarrow BZ \mid ZZ \mid AZ \mid a \\ B \rightarrow ZZ \\ Z \rightarrow AZ \mid a \\ A \rightarrow a \end{array}$$

This is the CNF-form grammar that we obtain mechanically.

Looking closer, it can be simplified to

$$S \rightarrow a \mid SS$$

with $L(S) = L(a^+)$, compared to $L(G) = L(a^*) = L(a^+ + \epsilon)$.

Decision problems.

① Emptiness:

Given: CFL L in form of a CFG G

Question: Is $L(G) = \emptyset$?

Idea: If $L(G) = \emptyset$ then S is not generating
 $\Rightarrow S$ is useless.

Easy to check?

② Membership

Given CFL L , string w

Q: Is $w \in L$?

Idea (Much harder than for regular languages) Assume $L = L(G)$ for CFG G

1. If $w = \epsilon$, check if S is nullable

2. Construct a Chomsky normal form grammar G' for $L \setminus \{\epsilon\}$.

> HW #4, problem 1

If $|w| = n$ then $S' \xRightarrow{*} w$ in exactly $2n-1$ steps.

4. Try all possible derivations of length $2n-1$

$\in |P'|^{2n-1}$ possibilities, $|P'| = \#$ productions in G' .

Better approach:

Dynamic Programming $V_{ij} = \{X \in V \mid X \xRightarrow{*} w_i w_{i+1} \dots w_j\}$

③ Equality Given CFGs G_1, G_2

Question: Is $L(G_1) = L(G_2)$?

Claim: There is no decision procedure? (show later)

Remark:

1. Easy to test equality for regular languages

$$L_1 = L_2 \iff (L_1 \cap L_2) \cup (L_1^c \cap L_2^c) = \Sigma^*$$

2. [Séniérgues, 2001]:

Many other undecidable properties of CFGs, but

algorithm for deciding if $L(M) = L(M')$

for deterministic PDAs M, M'

170 pages, Gödel prize

Deciding if $L(M) \subseteq L(M')$ still impossible. [Gödel]

or $L = M$ also OKAY [Ginsburg, Greibach 1966]

\uparrow CFL \uparrow REG

Dynamic programming alg. to decide membership in a CFL

Take a CFG in CN form.

Goal: For each substring $x(i,j) := x_i x_{i+1} \dots x_j$
compute $V_{ij} := \{X \in V \mid X \xRightarrow{*} x(i,j)\}$

Induction in $j-i$:

Base case $j=i$: Easy,

Induction:

$$X \xRightarrow{*} x(i,j)$$

must start with $X \Rightarrow BC$

where for some k , $i \leq k \leq j$,

$$B \xRightarrow{*} x(i,k), \quad C \xRightarrow{*} x(k+1,j) \quad \checkmark$$

\rightarrow
 $O(n)$ time

where $n = |x|$

$\Rightarrow O(n^3)$ time total, if CFG size is constant.

(then also conversion to CNF is $O(|x|)$)

Pumping Lemma for Context-Free Languages: [H-M-U 7.2]

Question: Of $L_1 = \{a^n b^n \mid n \geq 1\}$
 $L_2 = \{a^n b^{n+m} c^m \mid n, m \geq 1\}$
 $L_3 = \{a^n b^{2n} c^n \mid n \geq 1\}$,
 which are CFLs?

Answer: L_1, L_2 are CFLs, L_3 is not a CFL.

Intuitively, because the number of b's depends on its context, both to its left and right.

Pumping Lemma

Let L be a CFL.

Then there exists a constant N such that

For all $z \in L$ with $|z| > N$

can write $z = uvwxy$ with

- a) $|vwx| \leq N$
- b) $|vx| > 0$
- c) $\forall j \geq 0, uv^jwx^jy \in L$

Pumping Lemma for regular languages

Let L be regular.

Then $\exists N$

$\forall z \in L, |z| > N$

$z = uvw$ with

- a) $|uv| \leq N$
- b) $|v| > 0$
- c) $\forall j \geq 0, uv^jw \in L$

Application: Claim: $L_3 = \{a^n b^{2n} c^n \mid n \geq 1\}$ is not CFL.

Proof: ① Assume L_3 is CFL and apply P.L.

② Get constant $N > 0$

③ Choose $z = a^N b^{2N} c^N \in L_3$

④ Get u, v, w, x, y such that

$$z = uvwxy$$

$$\text{and } |vwx| \leq N, |vx| \geq 1$$

⑤ We choose $j=0$ and claim

$$z' = uv^0wx^0y = uwy \notin L_3. \quad (\text{contradiction})$$

Why? Observe: $|vwx| \leq N$ implies that either vwx has no c's or it has no a's

Case I [vwx has no c's]:

$$|z'| = |z| - |vx| < |z| = 4N$$

but $n_c(z') = N$, so more than $\frac{1}{4}$ of its symbols are c's.

Every string in L_3 has exactly $\frac{1}{4}$ c's, so $z' \notin L_3$. ✓

Case II [vwx has no a's]: By a similar argument, $z' \notin L_3$.

Corollary: Unlike regular languages, CFLs are not necessarily closed under intersection.

$$L_3 = L_2 \cap \{a^n b^m c^n \mid m \geq 0, n \geq 1\}$$

Proof of the CFL Pumping Lemma:

Given CFL L

Let $G = (V, T, P, S)$ be a Chomsky normal form grammar for $L \subseteq \Sigma^*$.

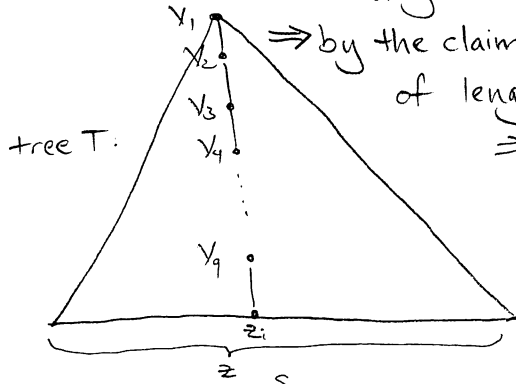
Claim: If in a parse tree for G , all root-leaf paths have length $\leq l$, then the yield of the tree has $\leq 2^{l-1}$ terminals.

Proof: By induction in l . Intuitively obvious since the tree is binary.

Worst cases: $l=1$: $\begin{matrix} A \\ | \\ a \end{matrix}$ $l=2$: $\begin{matrix} A \\ / \quad \backslash \\ B \quad C \\ | \quad | \\ b \quad c \end{matrix}$ $l=3$: $\begin{matrix} A \\ / \quad \backslash \\ B \quad C \\ / \quad \backslash \quad / \quad \backslash \\ D \quad E \quad F \quad G \\ | \quad | \quad | \quad | \\ d \quad e \quad f \quad g \end{matrix}$

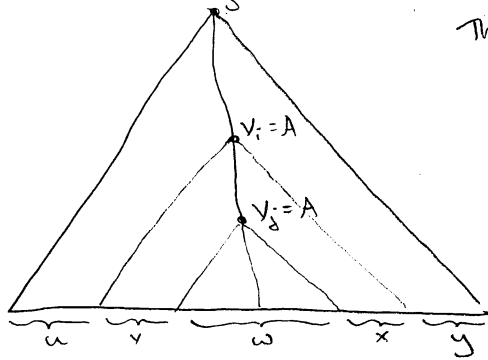
Now define $k = |V| = \# \text{variables}$, $N = 2^k$.

Consider any $z \in L(G)$ with $|z| > N$.



\Rightarrow By the pigeon-hole principle, some variable repeats along the path Q . In fact the last $k+1$ variables on Q must contain a repetition.

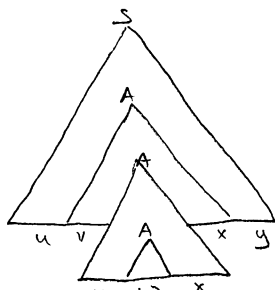
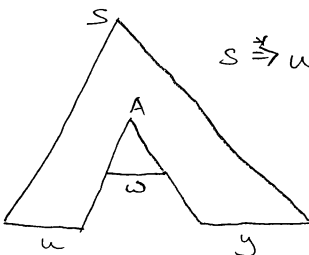
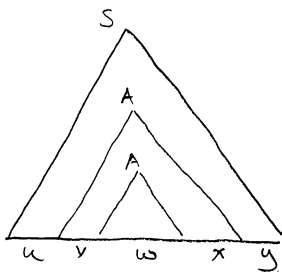
Let i, j be such that $q-k \leq i < j \leq q$, $V_i = V_j = A$.



and so $j \geq 0$,

$S \xRightarrow{*} uAy \xRightarrow{*} uv^jwx^jy$,

implying $uv^jwx^jy \in L(G)$. \checkmark



Also: $|vx| > 0$ since CNF G has no ϵ or unit productions.

Finally: $|vw| \leq N = 2^k$, since $V_i = A$ has height $\leq k+1$ and so by the claim its yield vw has length at most $2^{(k+1)-1} = 2^k$ terminals. \square