

①

CONTEXT-FREE GRAMMARS:

A more powerful way to specify languages.

- A CFG ~~has~~ consists of a (1) A set of symbols called terminals (Σ) (2) A set of symbols called non-terminals or variables V (3) A set of rules for generating sequences: productions. (4) A special variable called the start symbol: S .

Example $\Sigma = \{a, b\}$ $V = \{S\}$
 $S \xrightarrow{1} \epsilon$ $S \xrightarrow{2} a S b$

How does this produce a string in Σ^* ?

$S \xrightarrow{2} a S b \xrightarrow{2} a a S b b \xrightarrow{2} a a a S b b b \xrightarrow{1} a a a b b b$.

When you produce a string without variables you stop. The sequence is called a derivation.

This grammar produces the language $\{a^n b^n / n \geq 0\}$ which is not regular.

Used in linguistics to model sentence generation:

$\langle \text{SENTENCE} \rangle \rightarrow \langle \text{N.P.} \rangle \langle \text{V.P.} \rangle$

$\langle \text{N.P.} \rangle \rightarrow \langle \text{C.N.} \rangle \mid \langle \text{C.N.} \rangle \langle \text{PREP-PHASE} \rangle$

$\langle \text{V.P.} \rangle \rightarrow \langle \text{C.V.} \rangle \mid \langle \text{C.V.} \rangle \langle \text{PREP-P.} \rangle$

$\langle \text{P. PHASE} \rangle \rightarrow \langle \text{PREP} \rangle \langle \text{C.N.} \rangle$

ARITHMETIC EXPRESSIONS: $\Sigma = \{0, 1, 2, \dots, 9, +, \times, (,)\}$.

$V = \{ \langle \text{EXP} \rangle, \langle \text{NUM} \rangle, \langle \text{NZ} \rangle, \langle \text{N} \rangle \}$

$\langle \text{EXP} \rangle \rightarrow \langle \text{EXP} \rangle + \langle \text{EXP} \rangle \mid \langle \text{EXP} \rangle * \langle \text{EXP} \rangle \mid (\langle \text{EXP} \rangle) \mid \langle \text{NUM} \rangle$

$\langle \text{NUM} \rangle \rightarrow 0 \mid \langle \text{NZ} \rangle$

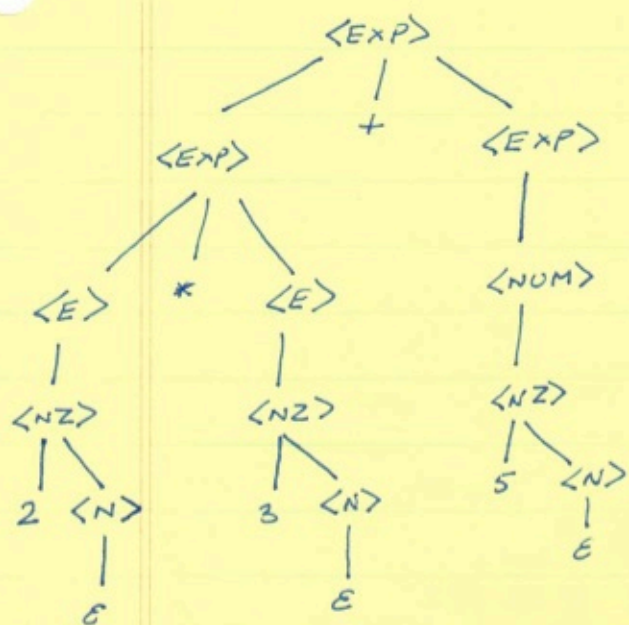
$\langle \text{NZ} \rangle \rightarrow 1 \mid 2 \mid \dots \mid 9 \mid \langle \text{N} \rangle$

$\langle \text{N} \rangle \rightarrow 0 \mid 1 \mid \dots \mid 9 \mid \epsilon$.

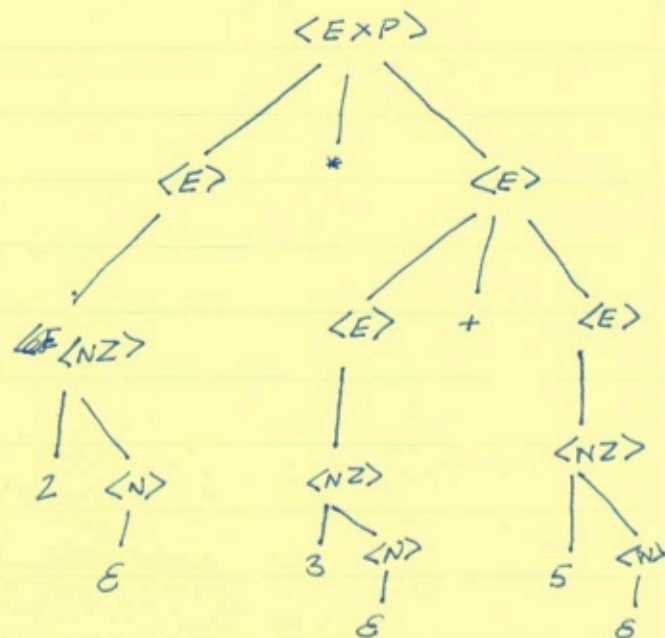
To display derivations it is far better to use a tree:

parse tree of a string. CFG's capture tree structure

(2)



2 * 3 + 5



2 * 3 + 5

AMBIGUOUS!

Def A grammar is said to be ambiguous if there are two distinct parse trees for the same string.

NOTE This is a property of the grammar not of the language.

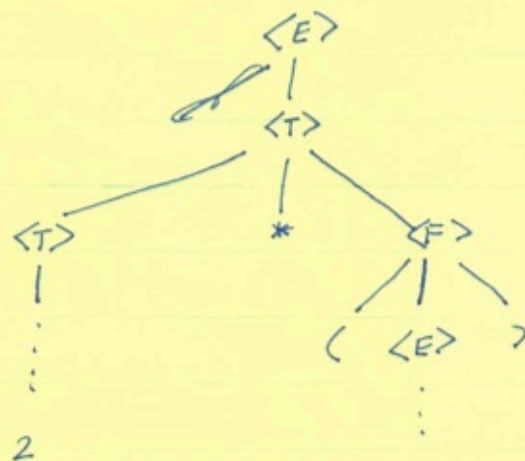
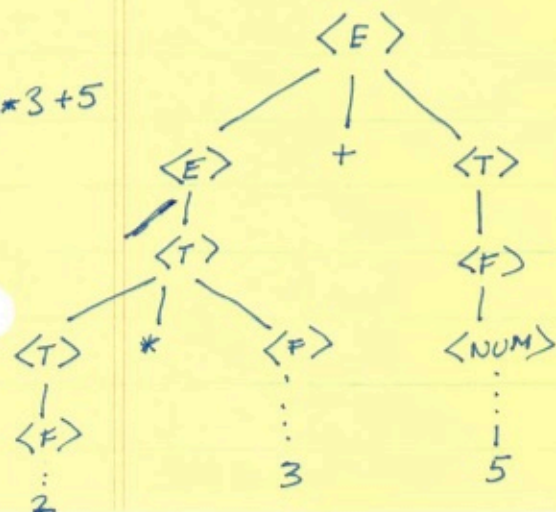
$V = \{ \langle E \rangle, \langle T \rangle, \langle F \rangle \}$

$\langle E \rangle \rightarrow \langle E \rangle + \langle T \rangle \mid \langle T \rangle$

$\langle T \rangle \rightarrow \langle T \rangle * \langle F \rangle \mid \langle F \rangle$

$\langle F \rangle \rightarrow (\langle E \rangle) \mid \langle \text{NUM} \rangle$

2 * 3 + 5



2 * (3 + 5) PARENS ARE FORCED

(3)

Can we do everything? No! $\{a^n b^n c^n \mid n \geq 0\}$.

Closure properties:

A language is context-free if there is a CFG for it.

$\{a^n b^n c^n \mid n \geq 0\}$ is not a CFL.

Closure properties: If L_1, L_2 are CFLs so is $L_1 \cup L_2$

$S_1 \dots L_1 \quad S_2 \dots L_2$

New grammar $V = V_1 \cup V_2 \cup \{S\}$

$S \rightarrow S_1 \mid S_2$

$L_1 = a^n b^n c^n$ is CFL, $L_2 = a^n b^n c^n$ is a CFL
 $L_1 \cup L_2$ is NOT.

Any regular language is a CFL:

$M = (Q, q_0, \delta, F)$

$G = (V = Q, S = q_0,$

if $\delta(q, a) = q'$ add the rule

$\langle q \rangle = a \langle q' \rangle$

If $q \in F$ add the rule $\langle q \rangle \rightarrow \epsilon$.

A grammar in which every rule has the form

$A \rightarrow aB$ or $A \rightarrow \epsilon$ is called a regular grammar.

Chomsky Normal Form

A CFG is in CNF if every rule has the form

$A \rightarrow BC$ or $A \rightarrow a$

B, C are not allowed to be the start variable

& we do allow $S \rightarrow \epsilon$.

Then

Every CFL is generated by a CFG in CNF.

Proof.

First add a new start S_0 & add the rule $S_0 \rightarrow S$. This ensures that S_0 does not occur on the RHS of a rule. However we will have to deal with this ~~new~~ new rule later.

We remove all rules of the form $A \rightarrow \epsilon$: For any rule with A on the RHS we add a new rule with each occurrence of A removed.

e.g. $X \rightarrow uAvAw$ will be replaced ~~by~~ have added $X \rightarrow uvw / uAvw / uVAw$.

If there is a rule $X \rightarrow A$ we add $X \rightarrow \epsilon$ unless we had previously removed $X \rightarrow \epsilon$. We keep going until we remove all rules of the form $A \rightarrow \epsilon$.

Third we remove rules like $A \rightarrow B$. If we have $B \rightarrow u$ we add $A \rightarrow u$ unless this is a unit rule previously removed.

Fourth : Now we have no ϵ -rules or unit rules we need to get rid of "long" R.H sides.

$$A \rightarrow u_1 u_2 \dots u_k \text{ where } u_i \in \Sigma \cup V.$$

We replace this with

$$A \rightarrow u_1 A_1 \quad A_1 \rightarrow u_2 A_2 \dots A_{k-2} \rightarrow u_{k-1} u_k.$$

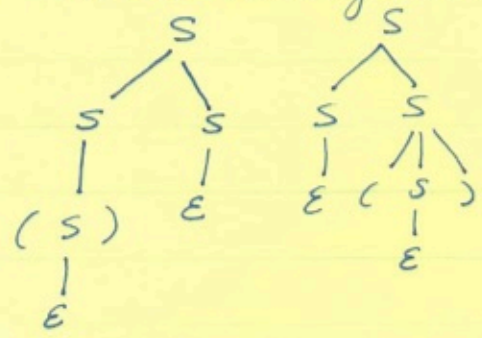
The A_i 's are new variables. Now any ~~non-term~~ terminal u_i is replaced by a new variable V_i & we add $V_i \rightarrow u_i$.

CFG: parsing, HTML, XML

Grammar for nested parens

$$S \rightarrow (S) \mid SS \mid \epsilon$$

This is ambiguous



Two different parse trees for $()()()$.

$$S \rightarrow (S)S \mid \epsilon$$

Unambiguous.

Can we show that everything generated by this grammar is a collection of properly nested parens?

Can we show that every properly nested string of parens is generated by this grammar?

Def A sequence of parens is properly nested if ~~the~~ (1) the total number of left (&) are equal & (2) in any prefix the number of (\geq number of).

Prop 1 Any string generated by G₂ is properly nested.

Proof By induction on the length of the derivation.

$$S \rightarrow (S)S \xrightarrow{*} (\alpha)\beta.$$

$$\begin{array}{ccc} \downarrow & & \downarrow \\ \alpha & & \beta \end{array}$$

By IH : ~~first prop~~ α, β are both properly nested.

So $(\alpha)\beta$ must have equal number of (&).

The prefixes of $(\alpha)\beta$ are :

- (i) ((ii) (α') where $\alpha' \leq \alpha$ (iii) (α) (iv) $(\alpha)\beta'$ where $\beta' \leq \beta$.

In each case the property (2) holds.

Prop 2 ~~See~~ If α is properly nested then $S \xrightarrow{*} \alpha$.

Proof If α is properly nested then the first symbol must be (so $\alpha = (\alpha'$

& The leading (must be matched by a) somewhere.

Two possibilities $\alpha = (\alpha')(\beta)$ or $\alpha = (r)\delta$

~~Now it must be~~ Now the matching) is the first point where $N_L - N_R = 0$ so r must be properly matched & (r) is properly matched so must δ be then

$$S \rightarrow (S)S \quad S \rightarrow r \quad \& \quad S \rightarrow \delta \text{ must be derivable}$$

$$\begin{array}{ccc} \downarrow & & \downarrow \\ r & & \delta \end{array} \text{ gives a derivation of } (r)\delta.$$

6

~~CFL~~ Closure properties:

- (1) L_1, L_2 are CFLs then $L_1 \cup L_2$ is a CFL
 $S \rightarrow S_1 | S_2$
- (2) $L_1 \cdot L_2$ is a CFL $S \rightarrow S_1 S_2$
- (3) L^* is a CFL Add new start symbol S' & new rules
 $S' \rightarrow SS' | \epsilon$.

CFLs are NOT closed under intersection:

~~S~~ $a^*b^*c^*$ $L_1 = \{a^n b^n c^m \mid n, m \geq 0\}$

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

$L_1 \cap L_2$ is not context-free.

~~CFLs~~ CFLs are not closed under complement. If they were they would have to be closed under \cap .

$L = \{a^n b^n c^n \mid n \geq 0\}$ is not a CFL but \bar{L}

\bar{L} is a CFL.

If L is a CFL & R is a regular language
then $L \cap R$ is a CFL.