

1. Grammar Formalisms

- L CFG
- L Grammar rules
- L Treebanks
- L Normal forms for Grammar
 - CNF
 - GNF

2. Dependency - Based Syntax.

- L Dependency Grammar
- L Feature structures
- L Unification of feature structure

3. Parsing Techniques

- syntactic parsing
- Ambiguity
- Dynamic programming parsing
- Shallow parsing.

4. Probabilistic Parsing.

1. PCFG
2. P CYK Alg.
3. Probabilistic Lexicalized CFGs

Context Free Grammar

CFG

NLP is the capability of computer software to understand the natural language.

↓

* There are variety of languages in the world.

* Each language has its own structure, like SVO, SOV

English

S V O

SOVX.

I eat mango.

↓ called

Grammar

↓ Has

certain set of rules

Other Lang :-

S O V

O S V

Natural Language

as input (sentences)

→ what is allowable
→ what is not allowed

NLP slw

Set of Rules
"Grammar"

Parsing

Taking out
meaning from
input structure
sentence

LHS

* Single
symbol

* Non
terminal

RHS.

* Single/
multiple
symbol

* Non
terminal/
terminal

CFG

Ex :-

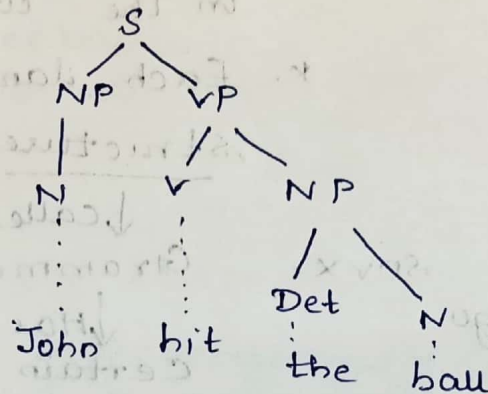
John hit the ball

$S \rightarrow NP \quad VP$

$VP \rightarrow V \quad NP$

$NP \rightarrow N$

$NP \rightarrow Det$
 L.H.S. R.H.S.



CFG :-

Context free grammar is a formal grammar which is used to generate all possible strings in a given formal language.

$G = (V, T, P, S)$ G = Grammar

V = Non-terminal set

T = Terminal set

P = production rule

S = Start symbol.

Rule :-

$\alpha \rightarrow \beta$

$\alpha \rightarrow$ Single Variable

$\beta \in (V + T)^*$

$T = \{ \text{this, that, the, book, flight, we, read, John, ball} \dots \}$ All the words in vocabulary & terminal

$V = \{ S, NP, N, VP, V, Det, Noun, Aux \dots \}$

P = {

$S \rightarrow NP \quad VP$

$S \rightarrow Aux \quad NP \quad VP$

$S \rightarrow VP$

$NP \rightarrow Det \quad Noun$

$VP \rightarrow Verb$

$VP \rightarrow Verb \quad NP$

}

Det \rightarrow the | a | this | that

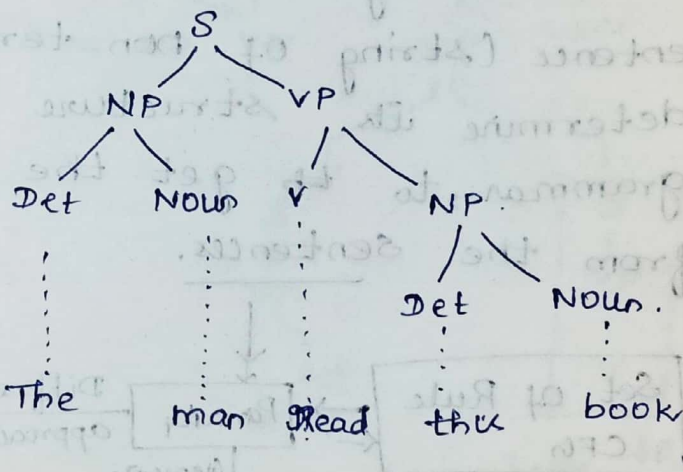
Noun \rightarrow book | flight | John | ball

Verb \rightarrow book | include | Read

Aux \rightarrow does | is

Ex:-

The man read this book

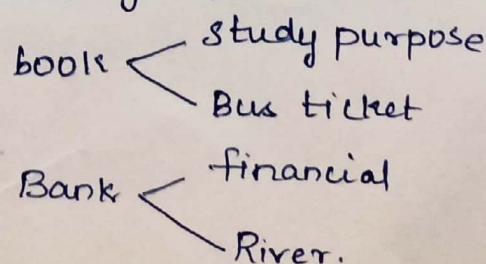


Properties of CFG:-

\Rightarrow Set of possible derivation

\Rightarrow string SET^*

\Rightarrow Each string is language generated by CFG may have more than one derivation (Ambiguity)



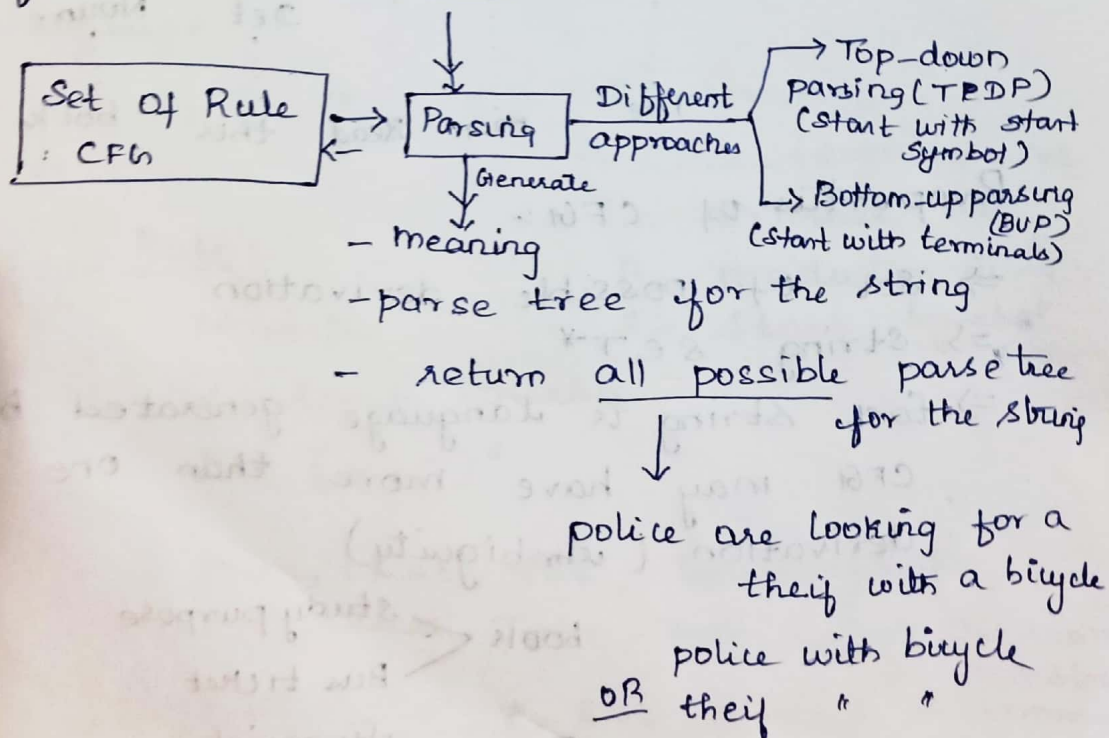
Summary :-

⇒ CFG is a list of rules that define the set of all well formed sentences in a language.

⇒ Each rule has a left hand side, which identified a syntactic category, & a right hand side define its alternative components parts.

CFG Parsing :-

Parsing is a method of analyzing a sentence (string of non-terminals) to determine its structure according to the grammar, to get the meaning out from the sentences.



Ex:-

5

String OR sentence given:-

book that flight.

Rules Given.

$S \rightarrow NP \quad VP$

$S \rightarrow VP$

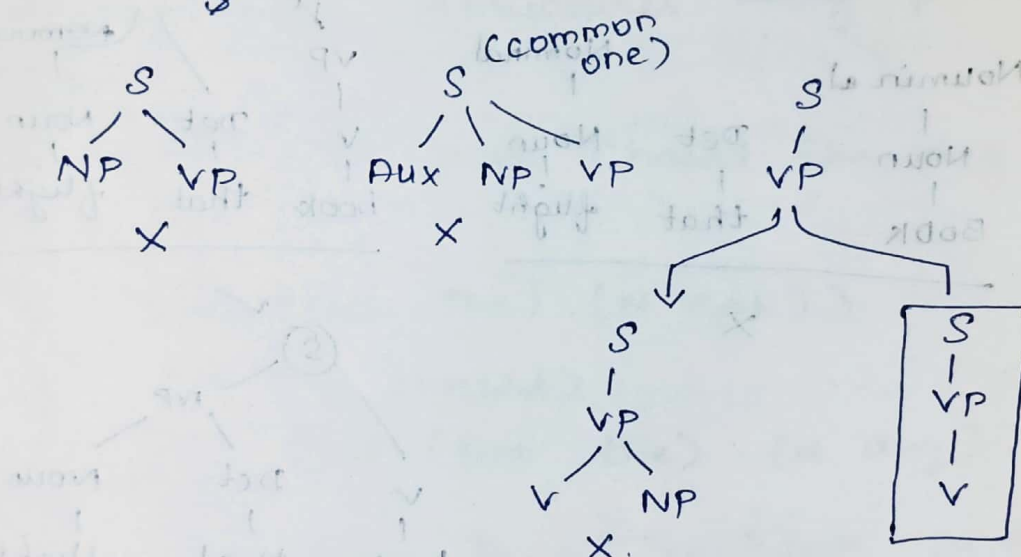
$NP \xrightarrow{(det)} ART \quad N$

$NP \xrightarrow{(det)} ART \quad ADI \quad N$

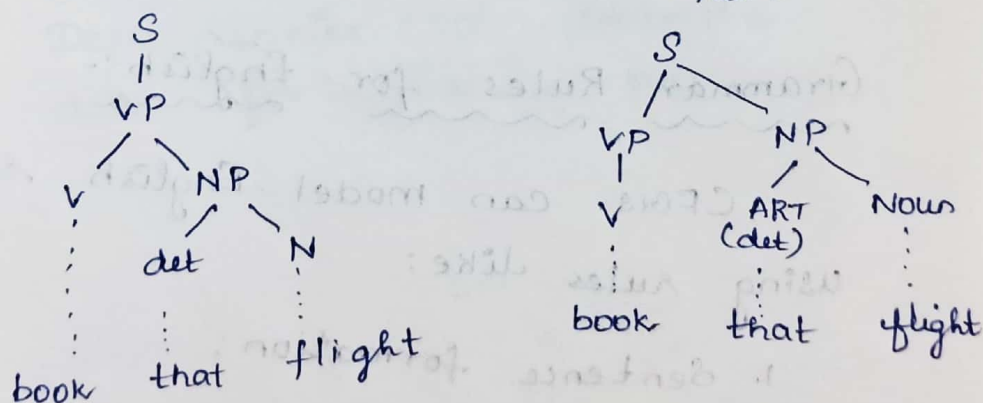
$VP \rightarrow VP$

$VP \rightarrow X' \quad NP$

Construct parse tree:- TDP



take this



6

Bottom-up parsing :- (BUP)

1. Start with the input text
2. Derive the text from Rule

Noun Det Noun Verb Det Noun
| | | | | |
Book that flight Book that flight

3. Each of these can be derived from non-terminals

Noun

Noun

Book

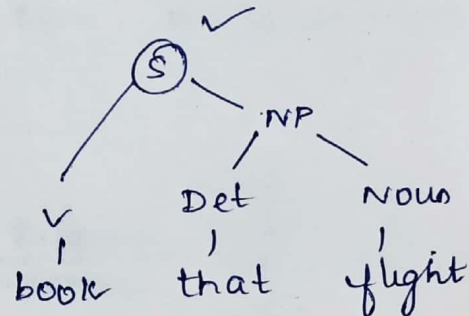
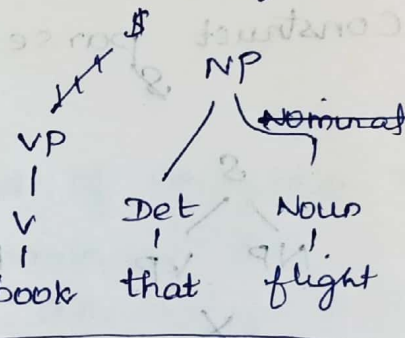
Det

that

Noun

Noun

flight



Grammar Rules for English :-

CFGs can model English syntax using rules like;

1. Sentence formation:

$S \rightarrow NP \ VP$

2. Noun phrase:

$NP \rightarrow \text{Det } N \mid N.$

3. Verb phrase:

$VP \rightarrow V \ NP \mid V$

For complex structures, additional rules handles adjective phrases (AP), prepositional phrases (PP) & Subordination.

Treebanks :-

A treebank is a corpus where sentences are manually annotated with syntactic structures using parse trees.

Ex:- "The cat Chased the dog"

(S

(NP (Det The) (N cat))

(VP (V chased)

(NP (Det the) (N dog))))

Treebanks like Penn Treebank and Universal Dependencies (UD) provide annotated data for training parsers).

Normal Forms for Grammar. (8)

Grammars can be transformed into normal forms for efficient parsing.

① Chomsky Normal Form (CNF):

Each rule is of the form

$$A \rightarrow BC \text{ or } A \rightarrow a,$$

where $A, B, C \Rightarrow$ non-terminals

$a \Rightarrow$ terminal.

② Greibach Normal Form (GNF):

Rules have the form $A \rightarrow a\alpha$, where a is terminal & α is a sequence of non-terminals.

Ex CNF conversion

from

$$S \rightarrow NP \quad VP$$

$$VP \rightarrow V \quad NP$$

$$NP \rightarrow \text{Det} \quad N$$

To CNF:

$$S \rightarrow X \quad VP$$

$$X \rightarrow NP$$

$$VP \rightarrow V \quad Y$$

$$Y \rightarrow NP$$

$$NP \rightarrow \text{Det} \quad N.$$

Chomsky Normal Form CNF.

9

In Chomsky Normal Form (CNF) we have a restriction on the length of RHS; which is, elements in RHS should either be two variables or a Terminal.

A CFG is in CNF if the productions are in the following forms:

$$\begin{aligned} A &\rightarrow a \\ A &\rightarrow BC \end{aligned}$$

where A, B, C are non-terminal and a is a terminal.

Steps to convert a given CFG to CNF.

Step 1:- If the start symbol S occurs on some right side, create a new start symbol s' and a new production $s' \rightarrow S$

Step 2: Remove Null productions.
(Using the Null production) Rem

Step 3: Remove Unit productions

Step 4: Replace each production

$$\begin{aligned} A &\rightarrow B_1 \dots B_n \text{ where } n > 2, \\ \text{with } A &\rightarrow B_1 C \text{ where } C \rightarrow B_2 \dots B_n \end{aligned}$$

Repeat this step for all productions having 2 or more symbols on the right side.

Step 5: If the right side of any production is in the form $A \rightarrow aB$ where 'a' is a terminal and A & B are non-terminals, then the production is replaced by $A \rightarrow XB$ & $X \rightarrow a$.

Repeat this step for every production which is of the form $A \rightarrow aB$.

Ex :-

Convert the following CFG to CNF.

P: $S \rightarrow ASA \mid aB$, $A \rightarrow B \mid S \mid \epsilon$, $B \rightarrow b \mid \epsilon$

Soln :-

1) Since S appears in RHS, we add a new state s' and $s' \rightarrow S$ is added to the production.

P: $s' \rightarrow S$, $S \rightarrow ASA \mid aB$,
 $A \rightarrow B \mid S \mid \epsilon$,
 $B \rightarrow b \mid \epsilon$

2) Remove the null productions:
 $B \rightarrow \epsilon$ and $A \rightarrow \epsilon$

After Removing $B \rightarrow \epsilon$:

P: $s' \rightarrow S$
 $S \rightarrow ASA \mid aBa$, $A \rightarrow B \mid S \mid \epsilon$, $B \rightarrow b$

After Removing $A \rightarrow \epsilon$:

P: $s' \rightarrow S$, $S \rightarrow ASA \mid aB \mid a \mid AS \mid SA \mid S$,
 $A \rightarrow B \mid S$, $B \rightarrow b$

3) Remove the Unit Production: (11)

$S \rightarrow S, S' \rightarrow S, A \rightarrow B \text{ \& } A \rightarrow S:$

After Removing $S \rightarrow S:$

$P: S' \rightarrow S, S \rightarrow ASA|AB|a|AS|SA,$
 $A \rightarrow B|S, B \rightarrow b.$

After Removing $S' \rightarrow S:$

$P: S' \rightarrow (ASA|AB|a|AS|SA),$
 $S \rightarrow ASA|AB|a|AS|SA,$
 $A \rightarrow B|S$
 $B \rightarrow b$

After Removing $A \rightarrow B:$

$P: S' \rightarrow ASA|AB|a|AS|SA,$
 $AS \rightarrow ASA|AB|a|AS|SA,$
 $A \rightarrow b|S$
 $B \rightarrow b$

After Removing $A \rightarrow S:$

$P: S' \rightarrow ASA|AB|a|AS|SA,$
 $S \rightarrow ASA|AB|a|AS|SA$
 $A \rightarrow b|ASA|AB|a|AS|SA$
 $B \rightarrow b.$

~~After Rem~~

4) Now find out the productions that has more than Two variables in RHS

$S' \rightarrow ASA, S \rightarrow ASA \text{ \& } A \rightarrow ASA$

After removing these, we get:

$P: S' \rightarrow AX|AB|a|AS|SA$
 $S \rightarrow AX|AB|a|AS|SA$
 $A \rightarrow b|AX|AB|a|AS|SA$
 $B \rightarrow b, X \rightarrow SA$

5) Now change the productions (12)

$S' \rightarrow aB$, $S \rightarrow aB$ and $A \rightarrow aB$

Finally we get:

P: $S \rightarrow Ax | YB | a | AS | SA$

$S \rightarrow Ax | YB | a | AS | SA$

$A \rightarrow b | Ax | YB | a | AS | SA$

$B \rightarrow b$

$X \rightarrow SA$

$Y \rightarrow a$

Ex:- 2.

Given CFG, convert to CNF.

$S \rightarrow NP \mid VP$

$VP \rightarrow V \mid NP \mid V$

$NP \rightarrow Det \mid N \mid N$

$Det \rightarrow \text{"the"} \mid \text{"a"}$

$N \rightarrow \text{"cat"} \mid \text{"dog"}$

$V \rightarrow \text{"chased"} \mid \text{"barked"}$

convert to CNF.

Soln:-

1. Break long rules into binary rules

$S \rightarrow NP \mid VP$ (Already in CNF)

$VP \rightarrow V \mid NP$ (" " ")

$NP \rightarrow Det \mid N$ (" " ")

$VP \rightarrow V$ (Needs modification) \rightarrow

Introduce new non-terminal X :

$VP \rightarrow V \mid X$

$X \rightarrow \epsilon$

2. Ensure all right-hand sides contain at most two symbols.

Final CNF Grammar:

$$\begin{aligned} S &\rightarrow NP \ VP \\ VP &\rightarrow V \ NP \mid V \ X \\ X &\rightarrow \epsilon \\ NP &\rightarrow Det \ N \\ Det &\rightarrow "the" \mid "a" \\ N &\rightarrow "cat" \mid "dog" \\ V &\rightarrow "chased" \mid "barked". \end{aligned}$$

CNF - Use :-

→ CNF is a restricted form of CFG that simplifies parsing alg & efficiently processes sentences in NLP.

- * → Mandatory for CYK parsing
- * → Improves parsing efficiency.
- * → Reduces ambiguity in grammars.
- * → Works well with probabilistic model (PCFGs)
- * → Standardizes grammar for NLP applications.