
PARSING

Submitted by:-

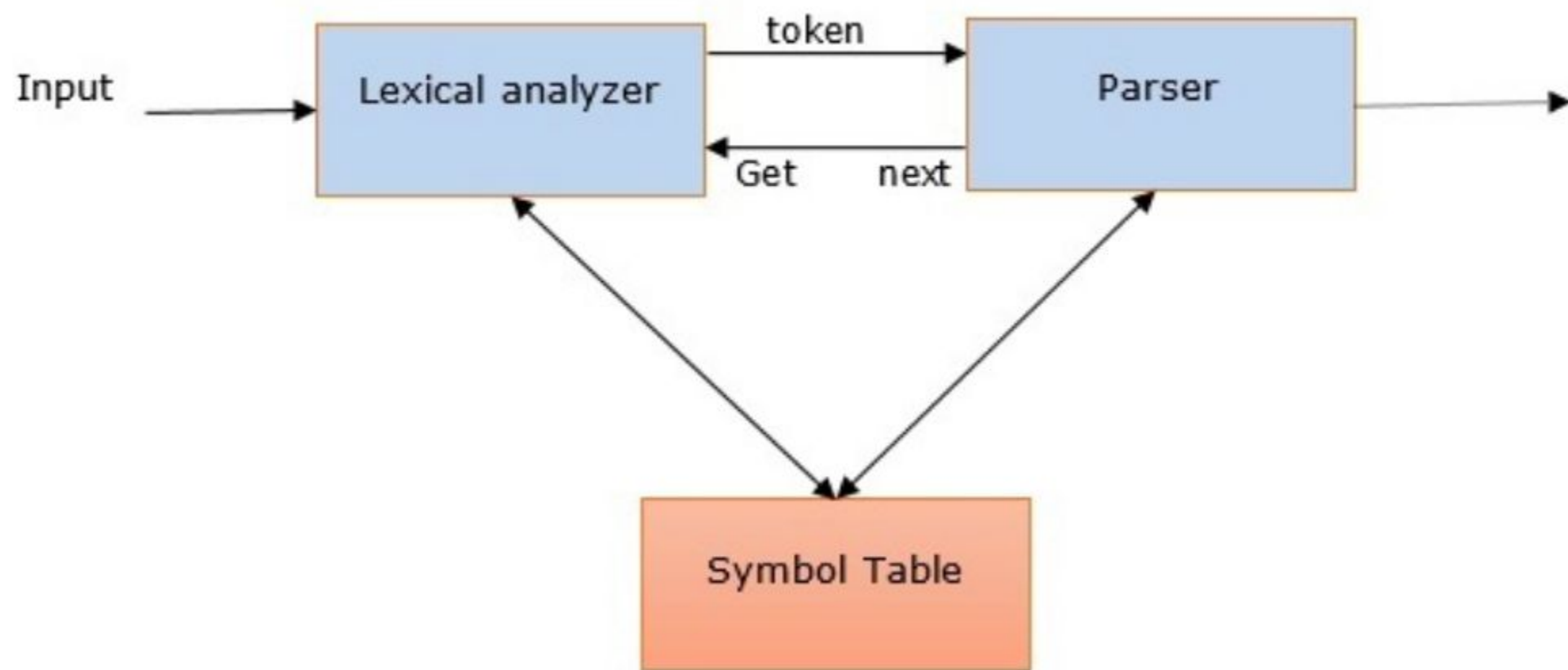
Kartika

URN-1805192 CRN-1815040 CSA2

PARSING ?

- Also called as Syntactic analysis or syntax analysis
- Syntactic analysis or parsing or syntax analysis is the third phase of NLP
- The word 'Parsing' is originated from Latin word 'pars' (which means 'part')
- Comparing the rules of formal grammar, syntax analysis checks the text for meaningfulness.
- The sentence like "Give me hot ice-cream", for example, would be rejected by parser or syntactic analyzer.
- The purpose of this phase is to draw exact meaning, or you can say dictionary meaning from the text.

It may be defined as the process of analyzing the strings of symbols in natural language conforming to the rules of formal grammar.



ROLE OF PARSER

- To report any syntax error.
- It helps to recover from commonly occurring error so that the processing of the remainder of program can be continued.
- To create parse tree
- To create symbol table.
- To produce intermediate representations (IR).

DEEP VS SHALLOW PARSING

| DEEP PARSING | SHALLOW PARSING |
|---|--|
| In deep parsing, the search strategy will give a complete syntactic structure to a sentence | It is the task of parsing a limited part of the syntactic information from the given task. |
| It is suitable for complex NLP applications | It can be used for less complex applications |
| EXAMPLES:- Dialogue Summarisation and summarisation | EXAMPLES:- Information extraction and text mining |
| It is called Full Parsing | It is called Chunking |

TYPES OF PARSING

Parsing is classified into two categories, i.e.

Top Down Parsing and Bottom-Up Parsing.

Top-Down Parsing is based on Left Most Derivation

Bottom Up Parsing is dependent on Reverse rightmost Derivation.

DERIVATION

Derivation is a sequence of production rules. It is used to get the input string through these production rules. During parsing we have to take two decisions. These are as follows:

- We have to decide the non-terminal which is to be replaced.
- We have to decide the production rule by which the non-terminal will be replaced.

Types of Derivation: Leftmost and rightmost

LEFT MOST DERIVATION

In the leftmost derivation, the input is scanned and replaced with the production rule from left to right. So in left most derivatives we read the input string from left to right.

Example:

Production rules:

1. $S = S + S$
2. $S = S - S$
3. $S = a \mid b \mid c$

INPUT

a - b + c

EXAMPLE OF LEFT MOST DERIVATION

1. $S = \mathbf{S} + S$

2. $S = \mathbf{S} - S + S$

3. $S = a - \mathbf{S} + S$

4. $S = a - b + \mathbf{S}$

5. $S = a - b + c$

RIGHT MOST DERIVATION

In the right most derivation, the input is scanned and replaced with the production rule from right to left. So in right most derivatives we read the input string from right to left.

EXAMPLE:-

1. $S = S - \mathbf{S}$
2. $S = S - S + \mathbf{S}$
3. $S = S - \mathbf{S} + c$
4. $S = \mathbf{S} - b + c$
5. $S = a - b + c$

TOP DOWN PARSING

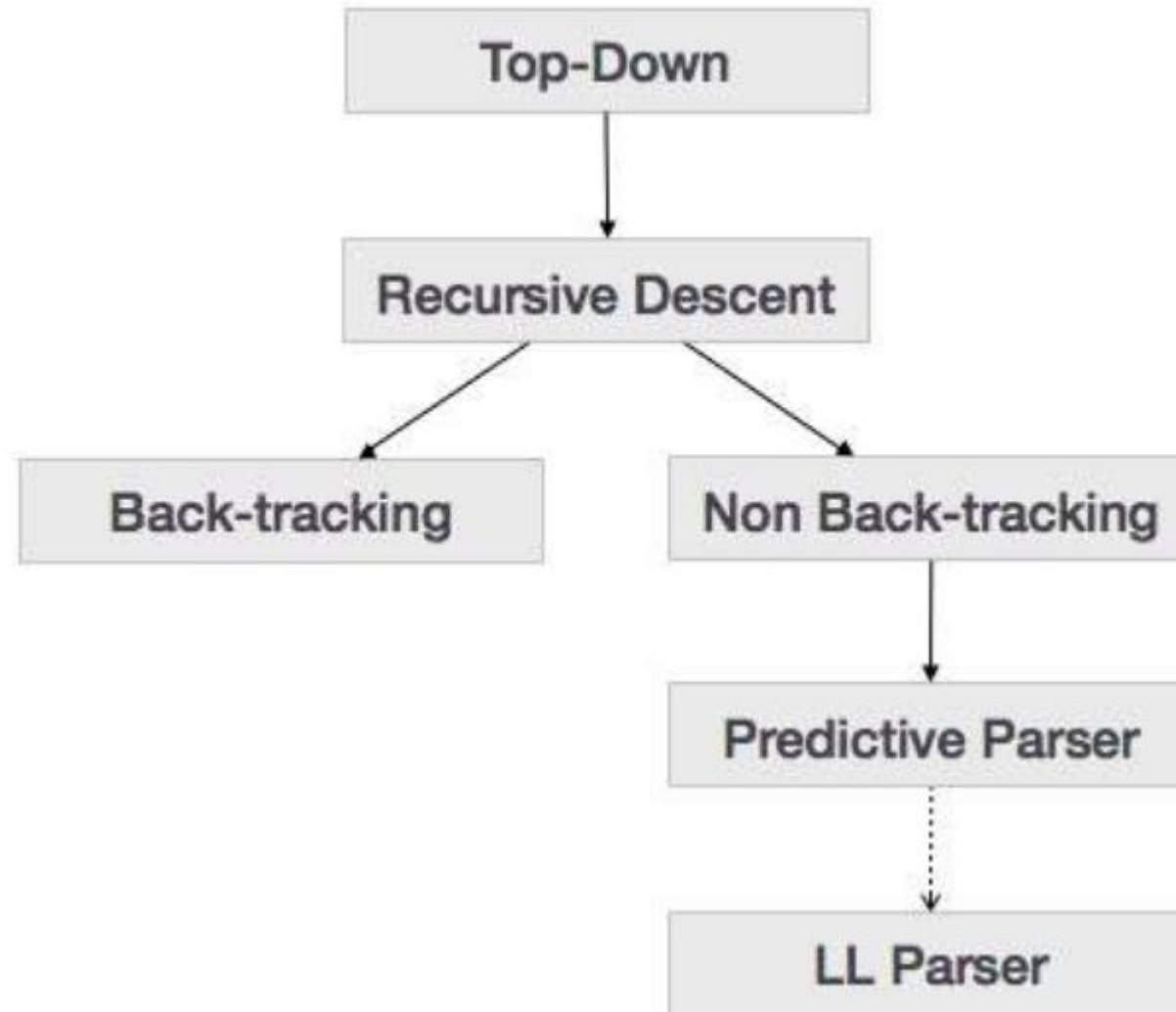
The process of constructing the parse tree which starts from the root and goes down to the leaf is Top-Down Parsing.

- Top-Down Parsers constructs from the Grammar which is free from ambiguity and left recursion.
- Top Down Parsers uses leftmost derivation to construct a parse tree.
- It allows a grammar which is free from Left Factoring.

EXAMPLE OF LEFT FACTORING-

$$S \rightarrow iEtS / iEtSeS / a$$

$$E \rightarrow b$$



WORKING OF TOP DOWN PARSER

EXAMPLE:-

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

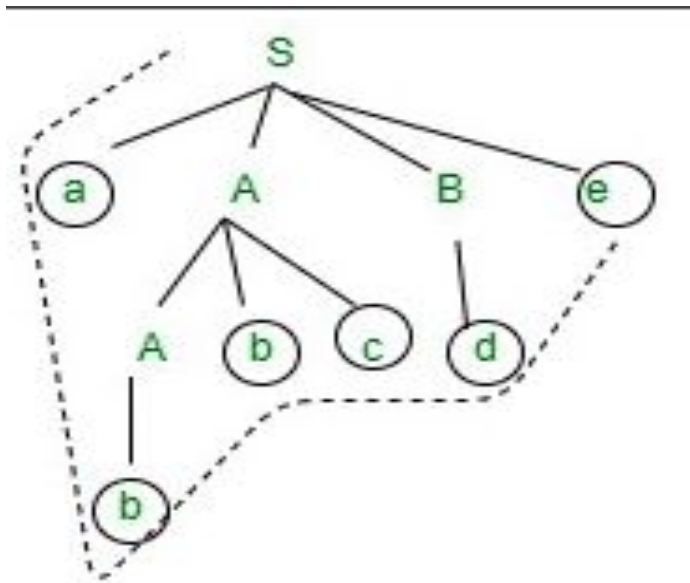
$B \rightarrow d$

Input –

abbcdc

Now, you will see that how top down approach works. Here, you will see how you can generate a input string from the grammar for top down approach.

- First, you can start with $S \rightarrow a A B e$ and then you will see input string a in the beginning and e in the end.
- Now, you need to generate $abbcde$.
- Expand $A \rightarrow Abc$ and Expand $B \rightarrow d$.
- Now, You have string like $aAbcde$ and your input string is $abbcde$.
- Expand $A \rightarrow b$.
- Final string, you will get **$abbcde$** .



BOTTOM UP PARSER OR SHIFT REDUCE PARSERS

Bottom Up Parsers / Shift Reduce Parsers

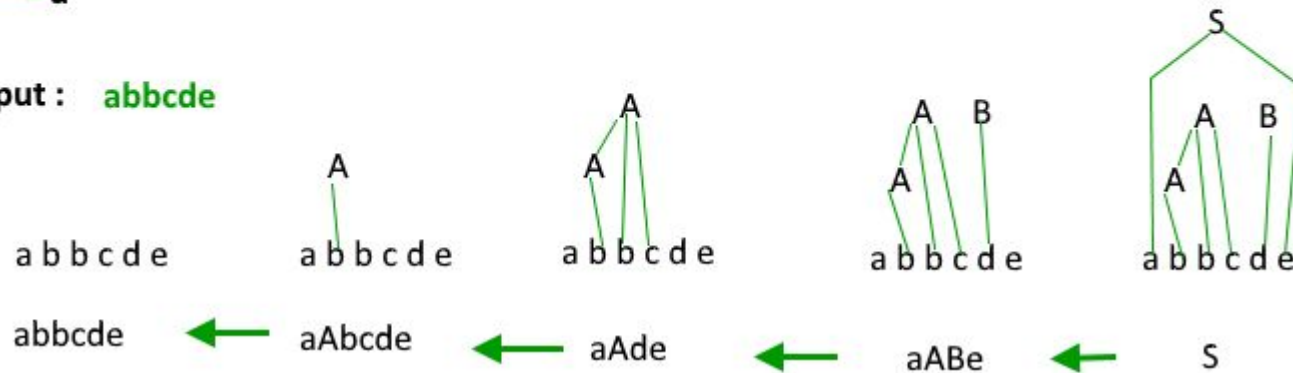
Build the parse tree from leaves to root. Bottom-up parsing can be defined as an attempt to reduce the input string w to the start symbol of grammar by tracing out the rightmost derivations of w in reverse.

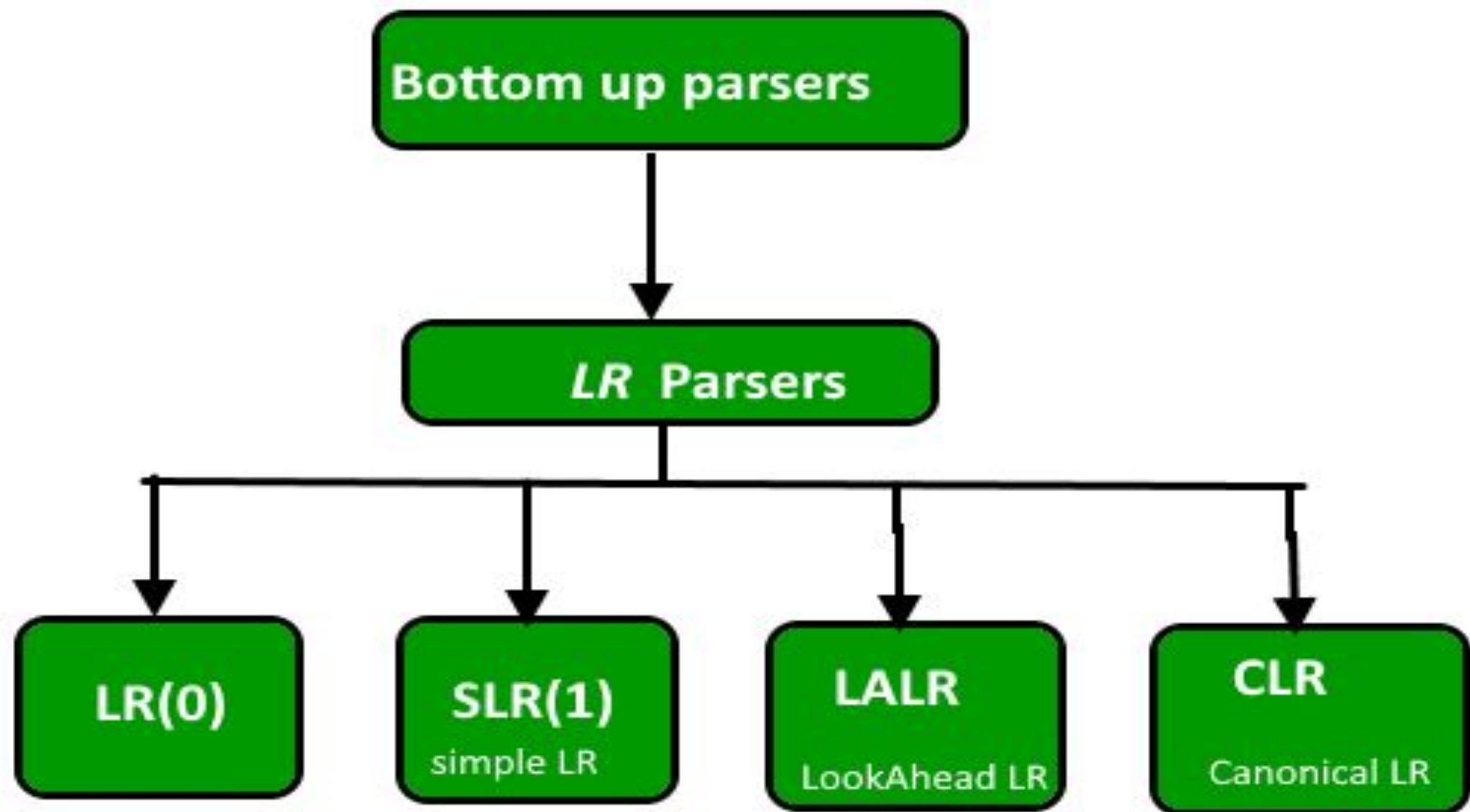
$S \rightarrow aABe$

$A \rightarrow Abc/b$

$B \rightarrow d$

Input : **abbcd e**





BOTTOM UP PARSER

There are two unique steps for bottom-up parsing. These steps are known as shift-step and reduce-step.

- **Shift step:** The shift step refers to the advancement of the input pointer to the next input symbol, which is called the shifted symbol. This symbol is pushed onto the stack. The shifted symbol is treated as a single node of the parse tree.
- **Reduce step :** When the parser finds a complete grammar rule (RHS) and replaces it to (LHS), it is known as reduce-step. This occurs when the top of the stack contains a handle. To reduce, a POP function is performed on the stack which pops off the handle and replaces it with LHS non-terminal symbol.

LR PARSER

A general shift reduce parsing is LR parsing. The L stands for scanning the input from left to right and R stands for constructing a rightmost derivation in reverse.

Benefits of LR parsing:

1. Many programming languages using some variations of an LR parser. It should be noted that C++ and Perl are exceptions to it.
2. LR Parser can be implemented very efficiently.
3. Of all the Parsers that scan their symbols from left to right, LR Parsers detect syntactic errors, as soon as possible.

S.No Top Down Parsing

Bottom Up Parsing

- | | |
|--|--|
| 1. It is a parsing strategy that first looks at the highest level of the parse tree and works down the parse tree by using the rules of grammar. | It is a parsing strategy that first looks at the lowest level of the parse tree and works up the parse tree by using the rules of grammar. |
| 2. Top-down parsing attempts to find the left most derivations for an input string. | Bottom-up parsing can be defined as an attempts to reduce the input string to start symbol of a grammar. |
| 3. In this parsing technique we start parsing from top (start symbol of parse tree) to down (the leaf node of parse tree) in top-down manner. | In this parsing technique we start parsing from bottom (leaf node of parse tree) to up (the start symbol of parse tree) in bottom-up manner. |
| 4. This parsing technique uses Left Most Derivation. | This parsing technique uses Right Most Derivation. |
| 5. It's main decision is to select what production rule to use in order to construct the string. | It's main decision is to select when to use a production rule to reduce the string to get the starting symbol. |

ISSUES IN BASIC PARSING

✓ Never explores trees that aren't potential solutions, ones with the wrong kind of root node.

✗ But explores trees that do not match the input sentence (predicts input before inspecting input).

✗ Naive top-down parsers never terminate if G contains recursive rules like $X \rightarrow X Y$ (left recursive rules).

✗ Backtracking may discard valid constituents that have to be re-discovered later (duplication of effort)

EARLEY ALGORITHM

The Earley Parsing Algorithm:

an efficient top-down parsing algorithm that avoids some of the inefficiency associated with purely naive search with the same top-down strategy

In naive search, top-down parsing is inefficient because structures are created over and over again.

Intermediate solutions are created only once and stored in a chart (dynamic programming).

Left-recursion problem is solved by examining the input.

Earley is not picky about what type of grammar it accepts, i.e., it accepts arbitrary CFGs.

function Earley-Parse(words,grammar) returns chart

 Enqueue($(\gamma \rightarrow \bullet S, [0,0])$,chart[0])

 for $i \leftarrow$ from 0 to Length(words) do

 for each state in chart[i] do

 if Incomplete?(state) and Next-Cat(state) is not POS then

 Predictor(state)

 elseif Incomplete?(state) and Next-Cat(state) is POS then

 Scanner(state)

 else

 Completer(state)

 end

 end

 return(chart)

A state consists of:

- 1) A subtree corresponding to a grammar rule $S \rightarrow NP VP$
- 2) Info about progress made towards completing this subtree $S \rightarrow NP \bullet VP$
- 3) The position of the subtree wrt input $S \rightarrow NP \bullet VP$, $[0, 3]$
- 4) Pointers to all contributing states in the case of a parser

A dotted rule is a data structure used in top-down parsing to record partial solutions towards discovering a constituent.

Earley: fundamental operations

- 1) Predict sub-structure (based on grammar)
- 2) Scan partial solutions for a match
- 3) Complete a sub-structure (i.e., build constituents)

How to represent progress towards finding an S node?

Add a dummy rule to grammar: $\gamma \rightarrow \bullet S$

This seeds the chart as the base case for recursion.

Earley's dot notation: given a production $X \rightarrow \alpha\beta$, the notation $X \rightarrow \alpha \bullet \beta$ represents a condition in which α has already been parsed and β is expected.

| | |
|-----------|--|
| Chart[0]: | \bullet astronomers saw stars with ears partial solutions ... |
| Chart[1]: | astronomers \bullet saw stars with ears partial solutions ... |
| Chart[2]: | astronomers saw \bullet stars with ears partial solutions ... |
| Chart[3]: | astronomers saw stars \bullet with ears partial solutions ... |
| Chart[4]: | astronomers saw stars with \bullet ears partial solutions ... |
| Chart[5]: | astronomers saw stars with ears \bullet |

Assumed indexing scheme:

\bullet_0 astronomers \bullet_1 saw \bullet_2 stars \bullet_3 with \bullet_4 ears \bullet_5

CYK ALGORITHM

CYK algorithm is a parsing algorithm for context free grammar.

In order to apply CYK algorithm to a grammar, it must be in Chomsky Normal Form. It uses a dynamic programming algorithm to tell whether a string is in the language of a grammar.

Algorithm :

Let w be the n length string to be parsed. And G represent the set of rules in our grammar with start state S .

1. Construct a table DP for size $n \times n$.
2. If $w = \epsilon$ (empty string) and $S \rightarrow \epsilon$ is a rule in G then we accept the string else we reject.
3. For $i = 1$ to n :
 For each variable A :
 We check if $A \rightarrow b$ is a rule and $b = w_i$ for some i :
 If so, we place A in cell (i, i) of our table.
4. For $l = 2$ to n :
 For $i = 1$ to $n-l+1$:
 $j = i+l-1$
 For $k = i$ to $j-1$:
 For each rule $A \rightarrow BC$:
 We check if (i, k) cell contains B and $(k+1, j)$ cell contains C :
 If so, we put A in cell (i, j) of our table.

5. We check if S is in $(1, n)$:
If so, we accept the string
Else, we reject.

Example –

Let our grammar G be:

$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

CYK Algorithm:- Check whether a string 'abbb' is a valid member of following CFG

* CYK applicable on CNF only (Chomsky Normal Form)
 CNF, $A \rightarrow BC$ or $A \rightarrow a$

* CYK is universal (applicable on all CNF)

Time complexity $\rightarrow O(n^3)$
 Space Complexity $\rightarrow O(n^2)$

Example

$S \rightarrow AB$

$A \rightarrow BB \mid a$

$B \rightarrow AB \mid b$

Input string \rightarrow 'abbb'
 1 2 3 4

| | (1) $\begin{smallmatrix} (4) \\ b \end{smallmatrix}$ | (2) $\begin{smallmatrix} (3) \\ b \end{smallmatrix}$ | (3) $\begin{smallmatrix} (2) \\ b \end{smallmatrix}$ | (4) $\begin{smallmatrix} (1) \\ a \end{smallmatrix}$ |
|-------|--|--|--|--|
| (1) a | (S, B) | A | (B, S) | A |
| (2) b | (S, B) | (A) | B | |
| (3) b | A | B | | |
| (4) b | B | | | |

On (12) cell or (ab)
 concatenate (11)
 and (22)
 $= aa + bb$
 $= AB$

Now S and B derives AB.

for (23) = (22)(33) = BB

for (34) = (33)(44) = BB

for (13) = (11)(23) = AA = ϕ
 \downarrow
 1 2 3
 (12)(33) = (BS) = ϕ

(BB)(SB)
 $A \neq \phi = A$

for 24 \rightarrow 2 3 4
 (22)(34)
 B $A \Rightarrow BA \Rightarrow \phi$

(23)(44)
 A B $\Rightarrow S, B$

for \rightarrow 14 \Rightarrow 1 2 3 4
 (11)(24) | (12)(34) | (13)(44)
 A SB | SB A | A B
 $\swarrow \searrow$ \downarrow \downarrow \downarrow
 AS AB | SA, BA | \downarrow
 \downarrow \downarrow \downarrow \downarrow
 ϕ S, B | ϕ ϕ | S, B

At the (1,4) cell i.e. (1 \rightarrow 4) S is there. if S is there in (1,4) cell then the string is a valid member of CFG.
 (abbb, ab, bbb) are valid.

“abbb”

$S \rightarrow AB, A \rightarrow BB|a, B \rightarrow AB|b$

Length 1 strings:

a,b,b,b

$a \Rightarrow (11)$ or (aa) , a can be derived from A

$b \Rightarrow (44)$ or (33) or (22) , b can be derived from B

Length 2 strings

$ab \Rightarrow (11)$ or (22) , $(11) \Rightarrow \mathbf{A}$ and $(22) \Rightarrow \mathbf{B}$, **concatenate** them

ie **AB** , **AB** can be derived from $\{\mathbf{S}, \mathbf{B}\}$

$bb \Rightarrow 22$ or 33 , $22 \Rightarrow \mathbf{B}$ and $33 \Rightarrow \mathbf{B}$, **concatenate** them ie **BB** can be derived from $\{\mathbf{A}\}$

“abbb”

$S \rightarrow AB, A \rightarrow BB|a, B \rightarrow AB|b$

abb (11) and (23) or 12 and 33

AA or (SB) and $B \Rightarrow SB$ and BB

“bbb” \Rightarrow 234

22 and 34 or 23 and 44

BA or $AB \Rightarrow S$ and B

“abbb” 1234

11+24 and 12+34 and 13+44

{S,B}

$S \rightarrow AB, A \rightarrow BB | a, B \rightarrow AB | b$

Input_string="abbb"

| | b (4) | b (3) | b (2) | a (1) |
|--------------|--------------|--------------|--------------|--------------|
| a (1) | {B,S} | {A} | {S,B} | {A} |
| b (2) | {B,S} | {A} | {B} | |
| b (3) | {A} | {B} | | |
| b (4) | {B} | | | |

$S \rightarrow AB|BC, A \rightarrow BA|a, B \rightarrow CC|b, C \rightarrow AB|a$

“baaba”

| | b (1) | a (2) | a(3) | b (4) | a(5) |
|-------|-------|-------|------|-------|------|
| b (1) | | | | | |
| a (2) | | | | | |
| a (3) | | | | | |
| b (4) | | | | | |
| a (5) | | | | | |

ample 2

$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

'baaba' \rightarrow input string
1 2 3 4 5

input string

$\frac{12}{ba}$

(11)(22)

(B)(AC)

BA, BC

$\downarrow \downarrow$
A S

$\frac{23}{aa}$

(22)(33)

(aa)(aa)

(AC)(AC)

AC, AA, CC

$\downarrow \downarrow \downarrow$
 $\emptyset \emptyset B$

$\frac{34}{ba}$

(33)(44)
{AC} {B}

AB BC

$\downarrow \downarrow$
S, C

$\frac{24}{ba} \Rightarrow aab$
2 3 4

(22)(34)

(AC)(SC)

(AS, AC, CC)

$\downarrow \downarrow$

B

(23)(44)

(B)(~~SA~~)

(~~BS~~, BA)

$\downarrow \downarrow$

\emptyset

$\frac{35}{aba}$

(33)(45)

(AC)(SA)

(AC, AA, CC)

$\downarrow \downarrow$

\emptyset

(34)(55)

(SC)(A, C)

SA, SC, CA

$\downarrow \downarrow$

B

$\frac{14}{1234}$

(11)(24)

(B)(B)

BB $\Rightarrow \emptyset$

(12)(34)

(SA) SC

SS, AS, AC

$\downarrow \downarrow$

\emptyset

(13)(44)

\emptyset {S, A}

= SA $\Rightarrow \emptyset$

(25) aaba
2 3 4 5

\Rightarrow (22)(35)

(AC) B \Rightarrow {C, S}

AB, BC

$\downarrow \downarrow \downarrow \downarrow$
S, C S

\rightarrow (23)(45)

(B)(SA)

(BS, BA) \Rightarrow A

{24}(45)

(B) SA \Rightarrow A

15

(11)(25) or

(12)(35) or

(13)(45) or

(14)(55)

| | 1) b | 2) a | 3) a | 4) b | 5) a |
|------|---------|---------|-------------|-------------|-----------|
| 1) b | {B} | {S, A} | \emptyset | \emptyset | {C, S, A} |
| 2) a | | {A, C} | {B} | {B} | {C, S, A} |
| 3) a | | | {A, C} | {S, C} | {B} |
| 4) b | | | | {B} | {S, A} |
| 5) a | | | | | {A, C} |

$\frac{45}{(44)(55)}$

(B)(AC)

BA, BC

$\downarrow \downarrow \downarrow \downarrow$
A S

$\frac{13}{123}$

(11)(23)

(B)(B)

BB

$\downarrow \downarrow$
 \emptyset

(12)(33)

(SA)(AC)

{SA, SC, AA,

AC}

$\downarrow \downarrow$
 \emptyset

Since there is S, start symbol
in (1,5) cell it is valid
member

We check if **baaba** is in $L(G)$:

1. We first insert single length rules into our table.

| | b | a | a | b | a |
|----------|----------|----------|----------|----------|----------|
| b | {B} | | | | |
| a | | {A,C} | | | |
| a | | | {A,C} | | |
| b | | | | {B} | |
| a | | | | | {A,C} |

Time and Space Complexity :

- **Time Complexity –**

$$O(n^3 \cdot |G|)$$

Where $|G|$ is the number of rules in the given grammar.

- **Space Complexity –**

$$O(n^2)$$

A Probabilistic Context-Free Grammar (PCFG)

| | | | | |
|----------|---------------|----|----|-----|
| S | \Rightarrow | NP | VP | 1.0 |
| VP | \Rightarrow | Vi | | 0.4 |
| VP | \Rightarrow | Vt | NP | 0.4 |
| VP | \Rightarrow | VP | PP | 0.2 |
| NP | \Rightarrow | DT | NN | 0.3 |
| NP | \Rightarrow | NP | PP | 0.7 |
| PP | \Rightarrow | P | NP | 1.0 |

| | | | |
|----|---------------|-----------|-----|
| Vi | \Rightarrow | sleeps | 1.0 |
| Vt | \Rightarrow | saw | 1.0 |
| NN | \Rightarrow | man | 0.7 |
| NN | \Rightarrow | woman | 0.2 |
| NN | \Rightarrow | telescope | 0.1 |
| DT | \Rightarrow | the | 1.0 |
| IN | \Rightarrow | with | 0.5 |
| IN | \Rightarrow | in | 0.5 |

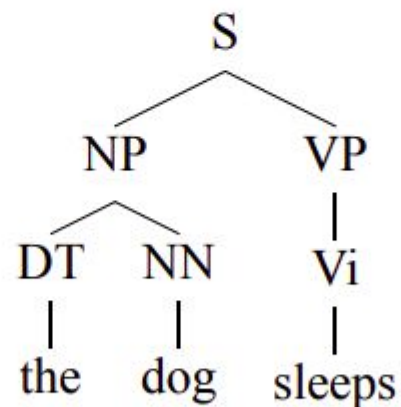
- Probability of a tree t with rules

$$\alpha_1 \rightarrow \beta_1, \alpha_2 \rightarrow \beta_2, \dots, \alpha_n \rightarrow \beta_n$$

is $p(t) = \prod_{i=1}^n q(\alpha_i \rightarrow \beta_i)$ where $q(\alpha \rightarrow \beta)$ is the probability for rule $\alpha \rightarrow \beta$.

$$\begin{aligned}
\sum_{\alpha \rightarrow \beta \in R: \alpha = \text{VP}} q(\alpha \rightarrow \beta) &= q(\text{VP} \rightarrow \text{Vi}) + q(\text{VP} \rightarrow \text{Vt NP}) + q(\text{VP} \rightarrow \text{VP PP}) \\
&= 0.3 + 0.5 + 0.2 \\
&= 1.0
\end{aligned}$$

To calculate the probability of any parse tree t , we simply multiply together the q values for the context-free rules that it contains. For example, if our parse tree t is



then we have

$$\begin{aligned}
p(t) &= q(\text{S} \rightarrow \text{NP VP}) \times q(\text{NP} \rightarrow \text{DT NN}) \times q(\text{DT} \rightarrow \text{the}) \times q(\text{NN} \rightarrow \text{dog}) \times \\
&\quad q(\text{VP} \rightarrow \text{Vi}) \times q(\text{Vi} \rightarrow \text{sleeps})
\end{aligned}$$

Properties of PCFGs

- ▶ Assigns a probability to each *left-most derivation*, or parse-tree, allowed by the underlying CFG
- ▶ Say we have a sentence s , set of derivations for that sentence is $\mathcal{T}(s)$. Then a PCFG assigns a probability $p(t)$ to each member of $\mathcal{T}(s)$. i.e., *we now have a ranking in order of probability*.
- ▶ The most likely parse tree for a sentence s is

$$\arg \max_{t \in \mathcal{T}(s)} p(t)$$

Parsing using Probabilistic Context Free Grammars

Each parse tree t_i is a sequence of context-free rules: we assume that every parse tree in our corpus has the same symbol, S , at its root. We can then define a PCFG (N, Σ, S, R, q) as follows:

- N is the set of all non-terminals seen in the trees $t_1 \dots t_m$.
- Σ is the set of all words seen in the trees $t_1 \dots t_m$.
- The start symbol S is taken to be S .
- The set of rules R is taken to be the set of all rules $\alpha \rightarrow \beta$ seen in the trees $t_1 \dots t_m$.
- The maximum-likelihood parameter estimates are

$$q_{ML}(\alpha \rightarrow \beta) = \frac{\text{Count}(\alpha \rightarrow \beta)}{\text{Count}(\alpha)}$$

where $\text{Count}(\alpha \rightarrow \beta)$ is the number of times that the rule $\alpha \rightarrow \beta$ is seen in the trees $t_1 \dots t_m$, and $\text{Count}(\alpha)$ is the number of times the non-terminal α is seen in the trees $t_1 \dots t_m$.

For example, if the rule $\text{VP} \rightarrow \text{Vt NP}$ is seen 105 times in our corpus, and the non-terminal VP is seen 1000 times, then

$$q(\text{VP} \rightarrow \text{Vt NP}) = \frac{105}{1000}$$

REFERENCES

Basics of

PCFG:-https://www.youtube.com/watch?v=wSONIMwa9rE&list=PLIQBy7xY8mbKypSJe_AjVtCuXXsdODiDi&index=3

https://www.youtube.com/watch?v=DjwH9wzCFzg&list=PLIQBy7xY8mbKypSJe_AjVtCuXXsdODiDi&index=4

CYK ALGORITHM:-

<https://www.youtube.com/watch?v=xRMn6HK84io>

<https://www.geeksforgeeks.org/cyk-algorithm-for-context-free-grammar/>

PARSING:-

https://www.tutorialspoint.com/compiler_design/compiler_design_types_of_parsing.htm

Documents: <http://www.cs.columbia.edu/~mcollins/courses/nlp2011/notes/pcfgs.pdf>