

1) Explain Top down parsing Technique with an example.

It shows how to construct an efficient known backtracking from topdown parsing called a Predictive parser.

1) first 2) follow.

first :-

To compile $\text{first}(x)$ for all grammar symbols ' x ', apply the following rules until no more terminals (α) $\in (\epsilon)$ can be added to only first .

Rules :-

* If x is a terminal then $\text{first}(x)$ is set $\{x\}$.
* If $x \rightarrow E$ is a production then add E to $\text{first}(x)$

* If x is a non-terminals a production then place a $x \rightarrow y_1, y_2, y_3, \dots, y_n$ is some i in $\text{first}(y_i)$ $\text{first}(x)$. If for

Ex :- $E \rightarrow TE'$

$E' \rightarrow +TE' | \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' | \epsilon$

$F \rightarrow (\epsilon) | id$.

1) $\text{first}(E) = \text{first}(T) = \text{first}(F) = \{c, id\}$

Non Terminal Non Terminal Non Terminal Terminal

2) $\text{first}(E') = \{+, \epsilon\}$

3) $\text{first}(T') = \{* , \epsilon\}$.

follow :- To compute $\text{follow}(A)$ apply the following rules until nothing can be added to any follow set.

Rules :-

- * \$ in $\text{follow}(S)$ where S is the start symbol and \$ is the input right and marker.
- * If there is a production in the form of $A \rightarrow d B \beta$ then everything in $\text{first}(\beta)$ except for \in is placed in $\text{follow}(B)$
- * If there is a production $A \rightarrow d B \beta$ or $A \rightarrow d B$ containing $\in (B \rightarrow \in)$ then everything in $\text{follow}(\beta)$ is in $\text{follow}(B)$

Eg: $E \rightarrow TE'$

$$\begin{aligned} E' &\rightarrow +TE'| \in \\ T &\rightarrow FT' \\ T' &\rightarrow *FT'| \in \\ F &\rightarrow (E)|\text{id.} \end{aligned}$$

$$\text{follow}(E) = \{ \$,) \}$$

$$\text{follow}(E') = \{ \$,) \}$$

$$\text{follow}(T) = \{ +, \$,) \}$$

$$\text{follow}(T') = \{ +, \$,) \}$$

$$\text{follow}(F) = \{ *, +, \$,) \}.$$

Construction of Predictive Parsing Table

I/P: Context free grammar

O/P: Predictive Parsing Table.

Method:

- 1) for each production is in the form $A \rightarrow d$ of the grammar do step 2 & 3.
- 2) for each terminal in $\text{first}(d)$ and $A \rightarrow d$ to $m[A, a]$
- 3) If " \in " is in $\text{first}(d)$ add $A \rightarrow d$ to $m[A, b]$ for each terminal b in $\text{follow}(A)$. If \in is in $\text{first}(d)$ and \$ in $\text{follow}(A)$ add $A \rightarrow d$ to $m[A, \$]$

4) make each undefined entry 'm' be error.

Predictive Parser Table | Top Down Parsing .

Non Terminals	-id	+	*	()	\$
E	$E \rightarrow TE'$		$E \rightarrow TEE'$	$E \rightarrow TE'E'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow E$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$		$T \rightarrow FTT'$	$T \rightarrow FT'$		
T'		$T' \rightarrow E$	$T' \rightarrow *FT'$		$T' \rightarrow E$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

Rules:-

- * If $x=a=\$$, the parser halts and announces successful completion of parsing.
- * If $x=a \neq \$$, then parser pops x off the stack and advances the input pointer to the next input symbol.
- * If x is a non-terminal the program consult entry $m[x,a] = x \rightarrow uvw$, the parser replaces x on top of the stack wvu .

Stack	Input	Output
\$E	id + id * id \$	-
\$E' T	id + id * id \$	$E' \rightarrow +TE'$
\$E' T' F	id + id * id \$	$T \rightarrow FT'$
\$E' T' id	id + id * id \$	$F \rightarrow id$
\$E' T'	+ id * id \$	-
\$E'	+ id * id \$	$T' \rightarrow E$
\$E' T * *	* id * id \$	$E' \rightarrow +TE'$
\$E' T	id * id \$	-
\$E' T' F	id * id \$	$T \rightarrow FT'$
\$E' T' id	id * id \$	$F \rightarrow id$

\$EIT'	*id\$
\$EIT'F*	*id\$
\$EIT'F	id\$
\$EIT'ix	id\$
\$EIT'	\$
\$EIE	\$
\$E'	\$
\$	\$

-	
T' → *FT'	
-	
F → id	
-	
T' → E	

Q) How can we convert Text to Bag of words, Text to TFIDF, Text to ngrams with suitable examples

In NLP, text data needs to be converted into numbers so that ML algorithm can understand it. One common method to do this is Bag of words (Bow) model. It turns next text like sentence, paragraph (or) document into a collection of words and counts how often each word appears but ignoring the order of words. It does not consider the order of words (or) their grammar but focus on counting how often each word appears in the text.

Key components of Bag of words :-

1) Vocabulary :- It is a list of all unique words from entire dataset each word in the vocabulary corresponds to a feature in the model.

2) Document Representation :- Each document is represented as a vector where each element shows the frequency of words from the vocabulary in that document. The frequency of each word is used as a feature for the model.

- Eg:- Sent 1 → He is a good boy
 Sent 2 → She is a good girl
 Sent 3 → Boy and girl are good.

Removing stop words,

Sent 1 → good boy

Sent 2 → good girl

Sent 3 → boy girl good

<u>Words</u>	<u>frequency</u>
--------------	------------------

good	3
boy	2
girl	2

Binary Bag of words.

$f_1 \ f_2 \ f_3 \ dp.$

good boy girl

s_1	1	1	0	1
s_2	1	0	1	1
s_3	1	1	1	1

⇒ TFIDF :-

* Spelling errors are common in user generated text (Social media, chatbot).

* Errors can be non word errors
 (These instead of their)

* Detection involves identifying words not found in reference are contextually wrong.

$$TF = \frac{\text{No. of repetition of words in sentence}}{\text{No. of words in sentence}}$$

$$IDF = \log_e \left(\frac{\text{No. of sentence}}{\text{No. of sentence containing words}} \right)$$

TF means Term frequency and

IDF means Inverse Document Frequency.

Sentence 1 :- He is a good boy

Sentence 2 :- She is a good girl

Sentence 3 :- boy girl are good.

$s_1 \rightarrow$ good boy After lowering and removing stop words:-
 $s_2 \rightarrow$ good girl removing

$s_3 \rightarrow$ boy girl are good.

$s_3 \rightarrow$ boy girl are good.

TF	S1	S2	S3
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

$$\text{sent 1} \quad f_1^{\text{good}} \quad f_2^{\text{boy}} \quad f_3^{\text{girl}}$$

$$0 \quad \frac{1}{2} \times \log_e\left(\frac{3}{2}\right) \quad 0$$

$$= 0.08$$

$$\text{Sent 2} \quad 0 \quad 0 \quad \frac{1}{2} \times \log_e\left(\frac{3}{2}\right)$$

$$= 0.08.$$

$$\text{Sent 3} \quad 0 \quad \frac{1}{3} \times \log_e\left(\frac{3}{2}\right) \quad \frac{1}{3} \times \log_e\left(\frac{3}{2}\right)$$

$$= 0.05 \quad = 0.05.$$

⇒ Go to ngrams

In NLP, ngrams are continuous sequences of n word from a given text.

A unigram considers only individual words.

A bigram considers two word sequences.

A Trigram considers three word sequences.

1) Unigram.

$$P(w) = \frac{\text{count}(w)}{\text{Total words.}}$$

2) Bigram

$$P(B|A) = \frac{\text{count}(AB)}{\text{count}(A)}$$

3) Trigram

$$P(C|A,B) = \frac{\text{count}(ABC)}{\text{count}(AB)}$$

3) Write short note on

1) first 2) follow 3) shift reduce passing
W Syntax ambiguity.

1) first :-

To compile first(x) for all grammar symbols 'x', apply the following rules until no more terminals (or) $\in (\epsilon)$ can be added to only first.

Rules:-

- If x is a terminal then $\text{first}(x)$ is set $\{x\}$
- If $x \rightarrow E$ is a production then add E to $\text{first}(x)$
- If x is a non-terminal $x \rightarrow Y_1, Y_2, Y_3, \dots, Y_n$ is a production then place a $\text{first}(x)$. If for some i in $\text{first}(Y_i)$.

2) follow.

To compute $\text{follow}(A)$ apply the following rules until nothing can be added to any followset

- ① $\$$ in $\text{follow}(S)$ where S is the start symbol and $\$$ is the input sight and marker
- ② If there is a production in the form of $A \rightarrow aB\beta$ then everything in $\text{first}(\beta)$ except for ϵ is placed in $\text{follow}(B)$
- ③ If the production $A \rightarrow aB\beta$ (a) $A \rightarrow B$ contains ϵ ($B \rightarrow E$) then everything in $\text{follow}(\beta)$ is in $\text{follow}(B)$

3) Shift Reduce Parsing.

- * A shift reduce parser is a type of bottom up parser used in NLP & compiler construction.
- * The process the input from left to right and uses a stack code and manages symbols.
- * The parsing process involves 2 main operations Shift and Reduce.
- * Shift operation :- This involves moving the next input symbol onto the stack.
- * Reduce operation :- This involves replacing a sequence of symbols on the stack with non-terminal symbol on the left hand side of rule.

4) Syntax ambiguity:-

Syntax ambiguity happens when a sentence can be parsed in more than 1 way because its structure is unclear.

for eg:- I saw the man with glasses.
This has two meanings.

- 4) what are the differences between top down parsing and bottom up parsing with examples.

Top down Parsing:-

- 1) Begins from the start symbol of the grammar and tries to derive the input string.
- 2) Derivation is from the root to leaves.
- 3) Parses tree from the top (start symbol) by expanding productions.
- 4) Used in recursive descent parsers, LL parsers
- 5) may involve backtracking (unless predictive parsing is used).

Bottom up Parsing:-

- 1) Begins from the input string and tries to reduce it to the start symbol
- 2) Derivation is from leaves to root (reverse of rightmost derivation)
- 3) Builds the parse tree from the bottom, reducing substring to non-terminals.
- 4) used in shift reduced parsers, LR parsers.
- 5) generally does not involve Backtracking

Examples for Top-down parsing:-

- ① Start with S
- ② Apply $S \rightarrow asb$: Now have asb
- ③ Apply $S \rightarrow asb$: Now have aasbb.
- ④ Apply $S \rightarrow c$: Now have aabb
- ⑤ Example for Bottom-up Parsing:-
- ⑥ look for substring matching the right side of rules
- ⑦ Replace c with nothing.

③ find asb → reduce asb to S
But we must build from bottom.

④ start reducing.

* Identify ab as a Potential Reduction

* But there's no rule $ab \rightarrow S$, so we try to reduce from aabb to S

* Pass tree construct in reverse, combining asb patterns.

5) what is meant by NLP? what are the challenges of NLP?

Natural Language Processing is a branch of AI that deals with the ability of machines to understand, analyze, manipulate and generate human language.

challenges of NLP :-

① Ambiguity :-

* words or sentences can have multiple meanings

* Types = lexical, syntactic, and semantic ambiguity.

② Sarcasm and Irony :-

* Difficult for machines to detect non literal expressions

③ Data Scarcity :-

* Lack of labelled data for many languages and domains.

④ Named Entity Recognition :-

* Identifying names, places, dates, etc., can be difficult due to variations.

⑤ code switching :-

* mixing languages in one sentence

⑥ Bias and fairness :-

* NLP models may inherit societal biases from training data.

APPLICATIONS OF NLP:-

- ① machine Translation.
eg:- Google, Translate, DeepL.
- ② chatbots and virtual assistants
eg:- Siri, Alexa, chatGPT, customer Service bots.
- ③ sentiment Analysis
Analyzing opinions in social media, reviews etc.
- ④ Speech Recognition
Converting spoken language into text.
- ⑤ Text summarization.
Creating concise summaries of longer documents.
- ⑥ Question Answering system
eg:- search engines, AI tutors.
- ⑦ Spam Detection.
Identifying and filtering out unwanted messages or emails.
- ⑧ Text classification.
Categorizing Text.