

①

Language Model - Grammar Based LM ; Statistical LM.

A language model (LM) is a probabilistic model that determines the likelihood of a sequence of words occurring in a given language.

It helps in generating, predicting, & analyzing text based on observed linguistic patterns. — generate human-like text.

Mathematically, given a sequence of words $w = (w_1, w_2 \dots w_n)$, a language model estimates the probability:

$$P(w) = P(w_1, w_2, \dots w_n)$$

By using different techniques (rule-based, statistical, or neural), LMs can predict the next word, complete sentences, or even generate entire documents.

Types of LM :-

Two Major types

- Grammar-Based LM
- Statistical LM.

(2)

The Importance of LMs in NLP Tasks:-

Language Models are crucial in NLP because they:

1. Enable machines to understand and generate human language.
2. Improve the performance of tasks like translation, summarization, sentiment analysis, & more.
3. Enhance user interaction in applications like chatbots & virtual assistants.

Grammar-Based Language Model.

A grammar-based language model uses predefined set of grammatical rules to construct sentences.

Components:-

1. Grammar Rules: Define syntactic structures (e.g noun phrase, verb phrase etc)
2. Parsing: The process of analyzing a sentence's structure to identify its components.

Formal Grammars:-

- * Context-free grammars,
- * Context-Sensitive grammars.

Context-free Grammars (CFGs):

These are grammars where the production rules are applied regardless of the context of a non-terminal. They are represented by the notation

$$G_1 = (N, \Sigma, P, S)$$

where,

1. Non-Terminal Symbols (N):

These are placeholders or variables that can be replaced by other symbols (both non-terminal & terminal).

2. Terminal Symbols (Σ):

These are the actual symbols or words in the language that cannot be replaced further.

3. Start Symbol (S):

This is the initial non-terminal symbol from which the production starts.

4. Production Rules (P):

These rules define how non-terminal symbols can be replaced by combination of terminal & non-terminal symbols.

The Parsing Process:-

The parsing process can be done using either top-down or bottom-up parsing methods.

Top-down parsing :

→ This method starts with start symbol S & tries to rewrite it to match the ilp sentence.

→ It applies production rules to the start symbol & continues expanding non-terminal symbols until the ilp sentence is derived.

Bottom-up parsing :-

→ This method starts with the ilp sentence & tries to construct the parse tree upwards to reach the start symbol S .

→ It identifies the parts of the ilp sentence that match the right-hand side of the production rules & combines them into non-terminal symbols until the start symbol S is reached.

Ex :- ① Top - Down.

Given the sentence "The cat sat on the mat", generate a parse tree using context-free grammar.

(5)

Step-by-step solutions:

1. Define the CFG:

$$\begin{aligned}
 S &\rightarrow NP \quad VP \\
 NP &\rightarrow Det \quad N \\
 VP &\rightarrow V \quad PP \\
 PP &\rightarrow P \quad NP \\
 Det &\rightarrow 'The' \\
 N &\rightarrow 'cat' \mid 'mat' \\
 V &\rightarrow 'sat' \\
 P &\rightarrow 'on'
 \end{aligned}$$

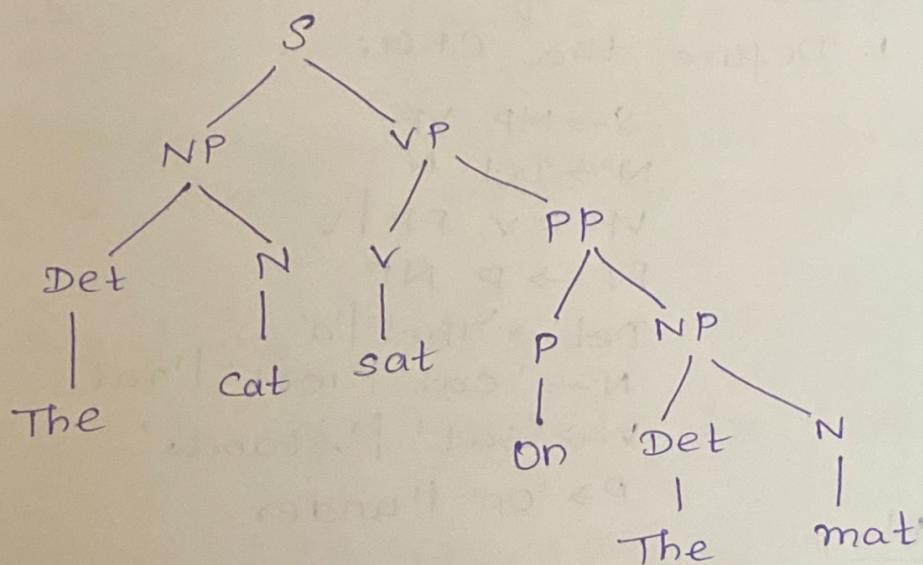
NP = Noun phrase

VP = Verb phrase

PP = preposition phrase

V = verb
 Det = Determiner
 ('The, a, they, them etc.)

2. Generate the Parse Tree:



This parse tree shows how the sentence "The cat sat on the mat" is derived from the start symbol S using the production rules of the CFG.

(6)

Ex ② :- ②

The cat sat on the ____.

How CFG will generate the correct word for the above sentence.

→ To predict the correct word to complete the sentence "The cat sat on the ___" using CFG, the model would rely on the production rules we've defined earlier.

Production rules for CFG:-

1. Define the CFG:

$$\begin{aligned} S &\rightarrow NP \quad VP \\ NP &\rightarrow Det \quad N \\ VP &\rightarrow V \quad PP \mid V \\ PP &\rightarrow P \quad NP \\ Det &\rightarrow 'The' \quad 'a' \\ N &\rightarrow 'cat' \quad 'mat' \quad 'rat' \\ V &\rightarrow 'sat' \quad 'stands' \\ P &\rightarrow 'on' \quad 'under'. \end{aligned}$$

Sohm :-

Use production rules to predict what comes next. (start with S or VP or ...)

(1) Starting with the VP (verb phrase):

$$VP \xrightarrow{V} PP$$

(2) We already have the verb "Sat":

$$V \rightarrow 'sat'$$

(3) The next part should be a Prepositional phrase (PP).

$$PP \rightarrow P \quad NP$$

(7)

(4.) We have "on" as the preposition:

$P \rightarrow 'on'$

(5) Next, we need a Noun phrase (NP)

$NP \rightarrow Det\ N$

(6.) We already have the determiner 'the':
 $Det \rightarrow 'the'$

(7) finally, we need a noun. The possible options for N are:

$N \rightarrow 'cat' | 'mat' | 'rat'$

Given the context, "The cat sat on the mat" is the most natural & grammatically correct sentence.

Therefore, the CFG would generate to complete the "mat" as the correct word to complete the sentence.

Resulting sentence:

The cat sat on the mat.

③ Example for Bottom-up Parsing :-

→ It starts with the ilp sentence and gradually works its way up to construct the parse tree until it reaches the start symbol S.

This technique is also known as shift-reduce parsing because it involves shifting ilp symbols onto a stack & reducing them based on production rules.

(8)

Let's define our CFG again:

$$\begin{aligned} S &\rightarrow NP \quad VP \\ NP &\rightarrow Det \quad N \\ VP &\rightarrow V \quad PP \mid v \\ PP &\rightarrow P \quad NP \end{aligned}$$

$$Det \rightarrow 'The' \mid 'the'$$

$$N \rightarrow 'cat' \mid 'mat'$$

$$V \rightarrow 'sat'$$

$$P \rightarrow 'on'$$

Parsing steps for "The cat sat on the mat"

1. Input Sentence:

The cat sat on the mat

2. Initialize stack & ilp buffer:

Stack: []

ilp buffer: ['The', 'cat', 'sat', 'on', 'the', 'mat']

3. shift ilp symbols onto stack:

a. shift The onto the stack:

stack: ['The']

ilp buffer: ['cat', 'sat', 'on', 'the', 'mat']

b. shift cat onto the stack:-

stack: ['The', 'cat']

ilp buffer: ['sat', 'on', 'the', 'mat']

4. Reduce by production Rule ~~'sat' → 'cat'~~:
 $Det \rightarrow 'The'$

stack: [Det, 'cat']

ilp buffer: ['sat', 'on', 'the', 'mat']

(9)

5. Reduce by production Rule $N \rightarrow 'cat'$:

stack: [Det, N]

Up buffer: ['sat', 'on', 'the', 'mat']

6. Reduce by production rule $NP \rightarrow Det\ N$:

stack: [NP]

Up buffer: ['sat', 'on', 'the', 'mat']

7. Shift sat onto the stack:

stack: [NP, 'sat']

Up buffer: ['on', 'the', 'mat']

8. Reduce by production Rule $V \rightarrow 'sat'$

stack: [NP, V]

Up buffer: ['on', 'the', 'mat']

9. Shift on onto the stack

stack: [NP, V, 'on']

Up buffer: ['the', 'mat']

10. Reduce by production rule $P \rightarrow 'on'$:

stack: [NP, V, P]

Up buffer: ['the', 'mat']

11. Shift 'the' onto the stack

stack: [NP, V, P, 'the']

Up buffer: ['mat']

12. Reduce by production rule $Det \rightarrow 'the'$:

stack: [NP, V, P, Det]

Up buffer: ['mat']

(10)

13. Shift mat onto the p_stack:

stack: [NP, V, P, Det, 'mat']

Up buffer: []

14. Reduce by production rule: $N \rightarrow \text{Det mat}$

stack: [NP, V, P, Det, N]

Up buffer: []

15. Reduce by production rule $NP \rightarrow \text{Det } N$

stack: [NP, V, P, NP]

Up buffer: []

16. Reduce by production rule $PP \rightarrow P \ NP$

stack: [NP, V, PP]

Up buffer: []

17. Reduce by production rule $\text{VP} \rightarrow V \ PP$

stack: [NP, VP]

Up buffer: []

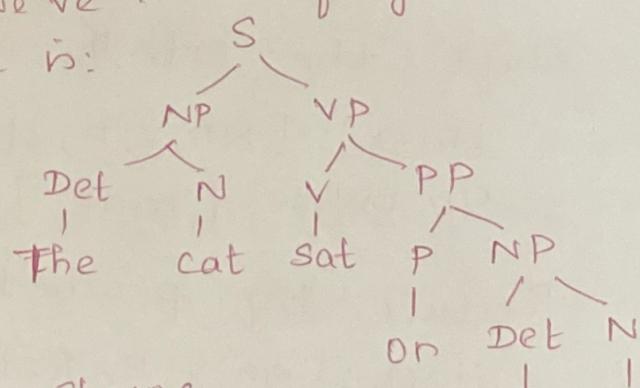
18. Reduce by production rule: $S \rightarrow NP \ VP$

stack: [S]

Up buffer: []

At this point, we've successfully parsed the sentence.

The final parse tree is:



The parse tree shows how sentence is derived from 'S' using production rules from the bottom-up approach.

Context Sensitive Grammar (CSG)

In these grammars, the production rules are dependent on the context of the non-terminal symbols.

This means that a non-terminal can be replaced by a string if it's surrounded by certain other symbols.

Ex :-

We want to generate

1. The boy runs

2. The boys run using CSG

Non-terminals: S, N_s, N_p, V_s, V_p, Det

Terminals : the, boy, boys, runs, run

Production rule:

$S \rightarrow \text{Det } N_s \text{ } V_s$

$S \rightarrow \text{Det } N_p \text{ } V_p$

$N_s \rightarrow \text{boy}$

$N_p \rightarrow \text{boys}$

$V_s \rightarrow \text{runs}$

$V_p \rightarrow \text{run}$

i) $S \rightarrow \text{Det } N_s \text{ } V_s$

This rule applies when the noun is singular (eg boy) & expect a singular verb (eg runs).

(12)

2) $S \rightarrow \text{Det } Np \text{ VP}$

This rule applies when the noun is plural (eg boys) & expect plural verb (eg run)

Ex:- 2 Generate $a^n b^n$, $n \geq 1$ using CEGI: with the help of production rule

$$S \rightarrow aCb$$

$$C \rightarrow aCb$$

$$C \rightarrow A$$

$A \rightarrow \epsilon$ (where ϵ represents empty string)

Non-terminals: S, A, C

Terminals: a, b

Sol:

Let's construct $a^3 b^3$

$$1. \quad S \rightarrow aCb$$

$$2. \quad \text{Apply } C \rightarrow aCb$$

$$aCb \rightarrow a\underline{aCb}b$$

$$3. \quad \text{Apply } C \rightarrow aCb$$

$$aacbb \rightarrow aa\underline{aCb}bb$$

$$4. \quad \text{Apply } C \rightarrow A$$

$$aaaCb \rightarrow aaaAb \rightarrow aaabbb$$

$$5. \quad \text{Apply } A \rightarrow \epsilon$$

$$aaaAb \rightarrow aaabb.$$

We generate $a^3 b^3$.
a string