

## UNIT 3

### Syntactic Pattern Recognition

Syntactic pattern recognition is an approach to pattern recognition that represents patterns through formal grammar rules. It is especially useful for recognizing patterns with inherent structural relationships, such as sequences or hierarchical data.

#### Overview

- Syntactic pattern recognition involves the use of grammar and language rules to describe patterns and their structural relationships.
- Patterns are treated as compositions of simpler sub patterns, which can be systematically described using grammar rules.
- Applications are common in fields such as computational linguistics, computer vision, bioinformatics, and robotics.

#### Qualifying Structure in Pattern Description and Recognition

The qualifying structure in pattern description and recognition refers to the specific rules and criteria used to define and identify patterns within data. It also refers to the organized and systematic representation of patterns based on their structural or relational properties. This approach aims to model patterns by defining their components, relationships, and rules for combining them into valid structures. It forms the basis of syntactic pattern recognition.

---

#### 1. Primitives (Basic Elements)

Primitives are the smallest, indivisible components of a pattern.

- **Definition:** Fundamental units that make up a pattern.
- **Examples:** Points, lines, curves, characters, or symbols.
- **Significance:** Serve as the building blocks for constructing complex patterns.

---

#### 2. Relationships between Primitives

Defines how primitives are related to each other spatially, temporally, or logically.

- **Types:**
  - **Spatial Relationships:** Arrangement in 2D or 3D space (e.g., proximity, alignment).
  - **Temporal Relationships:** Sequence of events or actions (e.g., order in speech recognition).
  - **Logical Relationships:** Constraints or dependencies between components.
- **Example:** In a face pattern, the eyes are positioned symmetrically around the nose.

---

#### 3. Grammatical Rules

- **Description:** Patterns are defined using a set of grammatical rules that specify how elements can be combined. These rules create a formal language for describing patterns, similar to the grammar of natural languages.
- **Example:** In programming languages, the rules for syntax (e.g., how variables, operators, and functions can be combined) are crucial for the correct interpretation of code. For instance, in Python, the grammatical rule for defining a function is:

```
def function_name(parameters):  
    # function body
```

---

#### 4. Hierarchical Structures

- **Description:** Patterns can be organized hierarchically, where complex patterns consist of simpler sub-patterns. This allows for recognizing patterns at different levels of abstraction.
  - **Example:** In image processing, a face can be considered a complex pattern that consists of simpler sub-patterns, such as eyes, nose, and mouth. Each of these sub-patterns can be recognized separately, contributing to the overall recognition of the face.
-

## 5. Attributes and Features

- **Description:** Qualifying structures include attributes or features that describe specific characteristics of the patterns. These can be quantitative (size, color) or qualitative (shape, orientation).
  - **Example:** In image recognition, features such as color histograms, edges, and corners are used to distinguish between different objects. For instance, identifying a red apple versus a green apple based on their color attributes.
- 

## 6. Contextual Rules

- **Description:** Contextual rules help qualify patterns based on their environment or context. This includes spatial relationships, temporal sequences, or the presence of other patterns.
  - **Example:** In video analysis, a contextual rule might state that a person running should be followed by a person walking, indicating a specific sequence of actions. This context is crucial for understanding behaviors in the video.
- 

## 7. Formal Descriptions

- **Description:** Patterns can be described using formal languages such as regular expressions, context-free grammars, or finite state machines. These descriptions provide a precise way to define the patterns.
- **Example:** A regular expression can describe a simple pattern for validating email addresses. For example, the regex:

`^[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-zA-Z]{2,}$` specifies the structure of valid email addresses.

---

## 8. Thresholds and Constraints

- **Description:** Qualifying structures may involve thresholds or constraints that determine when a structure is recognized as a valid pattern. This includes minimum sizes, shape ratios, or specific attribute values.
- **Example:** In anomaly detection, a threshold might be set for the number of failed login attempts. If the number exceeds a specific limit (e.g., 5 attempts), it is recognized as a potential security threat.

## Grammar-Based Approach and Applications

The Grammar-Based Approach in syntactic pattern recognition leverages formal grammar rules to identify and classify patterns within data, particularly in language processing and other domains where structured patterns are critical. Here's an overview of the approach and its applications:

### Grammar-Based Approach

1. **Concept:**
  - A grammar defines a set of rules or patterns for generating strings in a language. In syntactic pattern recognition, this means creating a formal representation of patterns that can be recognized or generated by a system.
2. **Types of Grammars:**
  - **Context-Free Grammars (CFGs):** Used widely in computational linguistics, where the syntax of a language is defined by rules that describe how symbols can be replaced or combined.
  - **Context-Sensitive Grammars (CSGs):** More complex and can capture patterns that CFGs cannot. They are used in more advanced applications where context plays a significant role.
  - **Regular Grammars:** A subset of CFGs, used for simpler pattern matching and recognition tasks.
3. **Components:**
  - **Terminals:** Basic symbols from which strings are formed.
  - **Non-terminals:** Symbols that can be replaced by combinations of terminals and other non-terminals.
  - **Production Rules:** Define how non-terminals can be replaced with other non-terminals or terminals.
  - **Start Symbol:** The initial non-terminal from which production begins.
4. **Process:**

- **Pattern Description:** The process of defining new primitives and constraints for pattern recognition.
- **Generation:** Creating strings that match the grammar rules, useful for testing or generating samples.
- **Parsing:** The process of analysing a string based on grammar rules to determine if it fits the defined patterns.
- **Recognition and Classification:** Based on the parsed data classify patterns.

## Applications

1. **Natural Language Processing (NLP):**
  - **Syntax Analysis:** Parsing sentences to understand grammatical structure and meaning.
  - **Machine Translation:** Translating text from one language to another by mapping grammatical structures between languages.
  - **Speech Recognition:** Converting spoken language into written text by recognizing syntactic patterns.
2. **Programming Languages:**
  - **Compilers:** Using grammars to parse and translate code written in programming languages into machine code or intermediate representations.
  - **Syntax Checking:** Ensuring that code adheres to grammatical rules of the programming language.
3. **Bioinformatics:**
  - **Gene Sequencing:** Recognizing patterns in genetic sequences to identify genes and other functional elements.
  - **Protein Structure Prediction:** Analyzing amino acid sequences based on known structural patterns.
4. **Information Retrieval:**
  - **Text Classification:** Categorizing documents based on their syntactic structure and content.
  - **Pattern Matching:** Identifying and extracting relevant information from large text corpora.
5. **Speech Recognition :**
  - **Pattern Recognition:** Enabling the Recognition of structure phrases and commands in Speech Processing.
6. **Image Recognition:**
  - **Pattern Recognition:** SYPR Helps to identify shapes and objects in image based on defined pattern e.g. Facial Recognition
7. **Artificial Intelligence:**
  - **Pattern Recognition:** Recognizing complex patterns in data, such as visual patterns in images or structured data in various formats.
  - **Knowledge Representation:** Using grammatical structures to represent and reason about knowledge in AI systems.
8. **Robotics Vision and Path Planning:**
  - **Pattern Recognition:** Robot uses SYPR to interpret environmental patterns navigate based on recognized landmarks.

## Advantages

- **Structured Representation:** Provides a clear and formal way to describe patterns.
- **Flexibility:** Can be adapted to various domains by modifying grammar rules.
- **Precision:** Allows for accurate recognition and generation of patterns.

## Elements of Formal Grammars, Examples of String Generation as Pattern Description:

### Elements of Formal Grammars-

1. **Terminals:**
  - **Definition:** Basic symbols from which strings are built. These are the actual characters or symbols in the strings.

- Example: In a simple grammar for arithmetic, terminals might be numbers (1, 2, 3) and operators (+, \*).
- 2. Non-terminals:**
- Definition: Symbols that can be replaced with other symbols (terminals or non-terminals) according to production rules.
  - Example: In an arithmetic grammar, non-terminals could be Expr (expression) and Term.
- 3. Production Rules:**
- Definition: Rules that define how non-terminals can be replaced with combinations of terminals and other non-terminals.
  - Example:
    - Expr can be replaced with Expr + Term or Term.
    - Term can be replaced with number or number \* Term.
- 4. Start Symbol:**
- Definition: The initial non-terminal from which the generation of strings begins.
  - Example: In our arithmetic grammar, the start symbol might be Expr.

### Example of String Generation

- Grammar:
  - Terminals: dog, barks, loudly
  - Non-terminals: Sentence, Noun Phrase, Verb Phrase
  - Production Rules:
    - Sentence  $\rightarrow$  Noun Phrase Verb Phrase
    - Noun Phrase  $\rightarrow$  dog
    - Verb Phrase  $\rightarrow$  barks | barks loudly
  - Start Symbol: Sentence

**Goal: Generate a sentence from the start symbol Sentence.**

### Steps:

1. Start with Sentence:
  - Apply the rule Sentence  $\rightarrow$  Noun Phrase Verb Phrase.
2. For Noun Phrase:
  - Apply the rule Noun Phrase  $\rightarrow$  dog.
3. For Verb Phrase:
  - Apply the rule Verb Phrase  $\rightarrow$  barks loudly.

### Combining these results:

- The Sentence becomes dog barks loudly.

### PARSING

Parsing is the process of analyzing a sequence of symbols (typically words or tokens) in order to determine their syntactic structure based on a formal grammar. The goal is to derive a syntactic tree (or parse tree) that represents how the input string adheres to grammatical rules.

### Key Concepts in Parsing:

- **Parse Tree (Syntax Tree):** A hierarchical tree-like structure that represents the syntactic structure of the input based on grammar rules.
- **Derivation:** The process of starting from the start symbol of the grammar and applying production rules to generate a string or a structure.

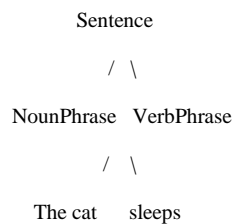
### Types of Parsing:

- **Top-Down Parsing:** This type of parser begins with the start symbol of the grammar and attempts to match the input string against possible derivations. The goal is to break down the input into smaller syntactic components.
  - **Example:** A **recursive descent parser** is a common example of top-down parsing.
- **Bottom-Up Parsing:** This method starts from the input symbols (tokens) and attempts to build up the syntactic structure by applying grammar rules in reverse. The process continues until the start symbol is reached.
  - **Example:** An **LR parser** (Left-to-right)

### How Does Parsing Work?

1. **Identify the Start:** Begin with the entire string and use rules to break it down into smaller parts.
2. **Apply Rules:** Use predefined rules (grammar) to determine how parts of the string relate to each other.
3. **Build Structure:** Construct a structured representation (like a tree) that shows how the parts fit together according to the rules.

You can visualize parsing with a simple tree diagram:



- ☐ Parsing helps break down and understand a string based on grammatical rules.
- ☐ Example Sentence: The cat sleeps.
- ☐ Grammar Rules: Define how to structure the sentence.
- ☐ Process: Divide the sentence according to the rules and build a structure.

### Parsing technique

#### 1. CYK Parsing Algorithm

The **CYK (Cocke-Younger-Kasami)** algorithm is a dynamic programming algorithm for parsing context-free grammars (CFG) in **Chomsky Normal Form (CNF)**. It is a bottom-up parsing method, which constructs a parse tree by combining smaller parts of the input step-by-step, starting from individual symbols.

#### Steps in the CYK Algorithm:

1. **Input:** The input string needs to be represented as a sequence of tokens (e.g., words in a sentence). The grammar must be in **Chomsky Normal Form (CNF)**, which means all production rules must be of the form:
  - $A \rightarrow BC$  (where A, B, C are non-terminal symbols), or
  - $A \rightarrow a$  (where A is a non-terminal and a is a terminal symbol).
2. **Initialize the Table:** Create a table (usually a 2D matrix) where each entry  $P[i][j]$  represents the set of non-terminal symbols that can generate the substring of the input from position i to position j.
3. **Fill the Table:**
  - **Base Case:** For each terminal symbol (word) in the input string, set  $P[i][i]$  to the set of non-terminals that can derive that terminal symbol. For example, if the terminal at position i is "the", and  $S \rightarrow the$  is a rule in the grammar, then  $P[i][i] = \{S\}$ .
  - **Recursive Step:** For substrings of length greater than 1, the table is filled using a dynamic programming approach:

- For each substring from  $i$  to  $j$  (with length greater than 1), check all possible ways to split the substring into two smaller parts ( $k$  from  $i$  to  $j-1$ ). For each pair of non-terminals in  $P[i][k]$  and  $P[k+1][j]$ , if there is a production rule that combines these non-terminals, then add the corresponding non-terminal to  $P[i][j]$ .
- 4. **Check for the Start Symbol:** After filling the table, check if the start symbol of the grammar (e.g.,  $S$ ) is in  $P[0][n-1]$ , where  $n$  is the length of the input string. If it is, the string can be derived from the grammar, and the input is syntactically valid.

#### Applications of the CYK Algorithm:

1. Context-Free Grammar Parsing
2. Syntax Analysis in Compilers
3. Handling Ambiguity in NLP
4. Speech Recognition
5. Machine Translation
6. Bioinformatics (RNA Secondary Structure Prediction)
7. Formal Language Recognition
8. Parsing in Formal Systems

#### Simple Example

Let's use a simple grammar and string to illustrate the CYK algorithm.

#### Grammar Rules:

1.  $S \rightarrow AB$
2.  $A \rightarrow a$
3.  $B \rightarrow b$

**String to Parse:** ab

#### Steps to Apply the CYK Algorithm:

1. **Set Up the Table:**
  - Create a table with dimensions based on the length of the string. For the string ab, we need a 2x2 table.

```

| 1 | 2 |
|---|---|
| 1 |  | 
| 2 |  | 

```

#### Fill the Table for Single Characters:

- For each character in the string, fill in which non-terminal can generate it:
  - a can be generated by A.
  - b can be generated by B.

Table after filling single characters:

#### Fill the Table for Longer Substrings:

- Now, consider substrings of length 2 (the whole string ab):
  - Check all possible splits of ab:
    - Split ab into a and b.
    - a (which can be generated by A) and b (which can be generated by B) can be combined using the rule  $S \rightarrow AB$ .

- So, S can generate ab.

Final Table:

		1		2	
		---		---	
		1		A	
		2		B	

### 1. Check the Start Symbol:

- Look at the top-right cell (which represents the whole string ab).
- If it contains the start symbol S, then the string can be generated by the grammar.

In this example, the top-right cell has S, meaning that the string ab can indeed be generated by the grammar.

## 2. ATN (Augmented Transition Network) in Parsing

An **Augmented Transition Network (ATN)** is a powerful and flexible parsing method primarily used in natural language processing. It extends the capabilities of traditional finite-state machines and context-free grammars to handle complex grammatical constructs. ATNs are widely used for syntactic analysis due to their ability to incorporate recursion, rules, and stack mechanisms, enabling more sophisticated sentence parsing.

### Key Concepts of ATNs

#### States

- **Definition:** States are points within the network that represent different stages or components of the parsing process. Each state signifies a particular syntactic unit or phase of recognition.
- **Role:** States act as checkpoints where parsing decisions are made before moving to subsequent stages.

#### 2. Transitions

- **Definition:** Transitions are links or directed arrows between states that guide the movement of the parser through the network based on input symbols, rules, and conditions.
- **Role:** Transitions dictate how the parser progresses from one state to another and serve as pathways for recognizing grammatical structures.
- **Types of Transitions:**
  - **Simple Transitions:** Move to the next state when specific input conditions are met.
  - **Recursive Transitions:** Allow re-entry into a state to recognize nested structures or recursive grammatical components

#### 3. Rules

- **Definition:** Rules act as guidelines or conditions that determine when and how transitions occur between states. They define the syntactic patterns that need to be satisfied at each stage.
- **Role:** Rules ensure that the input matches predefined grammatical structures and guide the parser in decision-making.
- **Components of Rules:**
  - Rules often include constraints or checks, such as:
    - **Part-of-speech tagging:** Ensuring a word matches the expected type (e.g., noun, verb).
    - **Lexical checks:** Verifying specific word forms.
    - **Structural constraints:** Matching multi-word phrases like adjective-noun sequences.

#### 4. Stack

- **Definition:** The stack is a memory structure used to keep track of states, transitions, and partially parsed components as the input is processed.
- **Role:** The stack enables the parser to handle:
  - **Recursive structures:** By pushing current states onto the stack before entering a new sub-state.

- **Backtracking:** Retrieving previous states when a transition fails or input does not match expected rules.
- **Mechanism:**
  - **Push:** The parser saves its current state and position on the stack before moving into a recursive transition.
  - **Pop:** The parser retrieves the saved state and continues parsing when a subcomponent is complete or if backtracking is needed.

## How Does an ATN Work?

### Start in the Initial State

Parsing begins at the starting state (e.g., S for sentence structure).

### Follow Transitions Based on Input

The parser moves through states based on input symbols and grammar rules.

*Example:* Transitioning from NP to VP on encountering a verb like "runs."

### Apply Rules

Rules ensure input matches expected patterns and guide valid transitions.

*Example:* A noun is followed by a verb to maintain sentence structure.

### Handle Stack Operations

The stack manages recursion and backtracking. Current states are pushed on entry and popped upon completion of nested structures.

*Example:* Parsing "The dog [that barked loudly] ran away."

### Complete Parsing

Parsing ends when the final state is reached, and all input is processed successfully.

*Example:* "The boy ate the apple" parses into valid components like NP and VP.

## Simple Example: Parsing a Sentence

Let's use a simple ATN to parse the sentence "The cat sleeps."

### Grammar Rules:

1. Sentence -> Noun Phrase Verb Phrase
2. Noun Phrase -> The cat
3. Verb Phrase -> sleeps

### ATN States and Transitions:

1. **States:**
  - **Start State:** Initial state where parsing begins.
  - **Noun Phrase State:** State for recognizing the noun phrase.
  - **Verb Phrase State:** State for recognizing the verb phrase.
  - **End State:** Final state indicating successful parsing.
2. **Transitions:**
  - From **Start State** to **Noun Phrase State** when encountering "The cat."
  - From **Noun Phrase State** to **Verb Phrase State** when encountering "sleeps."



- From **Verb Phrase State** to **End State** after processing the whole input.

#### Parsing the Sentence “The cat sleeps”:

1. **Start in the Initial State:**
  - Begin parsing “The cat sleeps” from the start state.
2. **Transition to Noun Phrase State:**
  - Recognize “The cat” as a noun phrase.
  - Move to the Noun Phrase state.
3. **Transition to Verb Phrase State:**
  - Recognize “sleeps” as a verb phrase.
  - Move to the Verb Phrase state.
4. **Reach the End State:**
  - All input has been processed and the end state is reached.

#### Applications of ATN Parsing:

1. Natural Language Processing (NLP)
2. Machine Translation
3. Speech Recognition
4. Information Extraction
5. Question Answering Systems
6. Grammar-Based Chatbots
7. Parsing Embedded Structures

#### ATN Diagram:

**(Start State) --The cat--> (Noun Phrase State) --sleeps--> (Verb Phrase State) --End--> (End State)**

#### Higher Dimensional Grammars:

Higher Dimensional Grammars (HDGs) extend traditional grammar formalisms to handle more complex and multidimensional structures in syntax and semantics. They are particularly useful in natural language processing, computational linguistics, and certain mathematical and AI applications where structures need to be analyzed beyond the linear or hierarchical representation provided by conventional grammars.

#### What is a Grammar?

In simple terms, a **grammar** is a set of rules that defines how sentences or strings are formed in a language. For example, in English, a basic grammar rule might be:

- **Sentence** → **Noun** + **Verb**

This means a sentence can be formed by combining a noun (like "cat") and a verb (like "runs"), producing the sentence "The cat runs."

#### What is a Higher Dimensional Grammar?

A **Higher Dimensional Grammar** generalizes this concept to work with structures beyond strings, such as:

- **2D structures** like drawings or diagrams (e.g., flowcharts).
- **3D structures** like models or shapes (e.g., 3D printed objects).
- **Graphs** representing networks (e.g., social networks, molecular structures).

#### Simple Example of a Higher Dimensional Grammar

Imagine you want to define a simple 2D drawing, like a square:

1. **Basic Components:**
  - **Point:** A dot on the paper.
  - **Line:** A straight connection between two points.
2. **Rules:**

- **Square  $\rightarrow$  4 Points + 4 Lines:**
  - Place four points on the paper.
  - Connect the points with lines to form a square.

In this example, the "grammar" isn't just about combining words but combining shapes (points and lines) according to certain rules to create a square. This is the essence of a Higher Dimensional Grammar: it allows you to define and generate complex multi-dimensional structures by applying rules.

#### Advantages of Higher Dimensional Grammars

- **Expressiveness:** They can model complex, recursive, and multidimensional structures more effectively than traditional grammars.
- **Flexibility:** HDGs can represent relationships, layouts, and networks beyond simple strings of symbols.
- **Handling Ambiguity:** By incorporating higher dimensions, ambiguities in syntax or semantics can be resolved with greater precision.
- **Real-World Applications:** Useful in fields requiring multidimensional analysis, such as natural language understanding, image processing, and knowledge representation.

#### Applications of Higher Dimensional Grammars

- **Natural Language Processing (NLP):**
  - Representing syntactic structures with dependency trees, TAGs, or hypergraphs.
  - Parsing sentences with embedded and recursive elements.
- **Image and Diagram Parsing:**
  - Understanding spatial syntax in flowcharts, tables, and diagrams.
  - Parsing two-dimensional languages (e.g., mathematical notations or visual layouts).
- **Artificial Intelligence:**
  - Parsing and modeling knowledge graphs or relational data.
  - Semantic role labeling where multidimensional relationships exist.
- **Computational Biology:**
  - Analyzing structures such as RNA or DNA chains represented as multidimensional graphs.
- **Mathematical Applications:**
  - Representing algebraic structures and geometric relationships in higher dimensions.

#### Stochastic Grammars

**Stochastic Grammars** are grammars where the production rules are not deterministic but have associated probabilities. These grammars are used in scenarios where there is uncertainty or variation in the generation process, such as in **natural language processing (NLP)**, **speech recognition**, **genetics**, and **artificial intelligence**.

#### What is a Stochastic Grammar?

1. **Grammar Basics:** A grammar consists of rules that tell you how to form valid sentences or structures. For example, a basic rule might be:
  - **Sentence  $\rightarrow$  Noun + Verb** This rule tells you that a sentence can be made up of a noun followed by a verb.
2. **Adding Probabilities:** In a stochastic grammar, each rule has a probability attached to it. For example:
  - **Sentence  $\rightarrow$  Noun + Verb [0.7]**
  - **Sentence  $\rightarrow$  Verb + Noun [0.3]** This means that there's a 70% chance the sentence will be formed as "Noun + Verb" and a 30% chance it will be "Verb + Noun".

#### Key Features of Stochastic Grammars

1. **Probabilistic Production Rules:**

- Each production rule has an associated probability, summing to 1 for a given non-terminal.
- 2. **Handling Uncertainty:**
  - Models real-world ambiguity by ranking multiple parses based on probabilities.
- 3. **Extension of CFGs:**
  - Stochastic grammars extend **Context-Free Grammars (CFGs)** by associating probabilities with rules, forming **PCFGs**.
  - Widely used in parsing and speech recognition.
- 4. **Ambiguity Resolution:**
  - Resolves ambiguity by ranking alternatives based on probabilities.

#### How Stochastic Grammars Work

1. **Defining the Grammar:**
  - Start with a formal grammar (e.g., CFG) containing **non-terminals**, **terminals**, and **production rules**.
  - Assign a probability PPP to each production rule.
2. **Parsing and Derivation:**
  - Parsing algorithms analyze input sequences.
  - The probability of a derivation is calculated by multiplying the probabilities of the applied rules:  

$$P(S \rightarrow NP VP) \times P(NP \rightarrow Det Noun) \times P(VP \rightarrow Verb) \cdot P(S \rightarrow NP VP) \times P(NP \rightarrow Det Noun) \times P(VP \rightarrow Verb)$$
3. **Ranking Parse Trees:**
  - Multiple parse trees may be generated for the input.
  - Each tree is assigned a probability based on applied rules, and the tree with the **highest probability** is selected as the best parse.
4. **Learning Probabilities:**
  - Rule probabilities are learned from training data (corpus) using statistical methods like:
    - **Maximum Likelihood Estimation (MLE).**
    - **Bayesian Inference.**
  - Frequencies of rule occurrences in the data determine their probabilities.

#### Simple Example

Imagine you're generating sentences in a very simple language:

- **Noun** → "dog" [0.5], "cat" [0.5]
- **Verb** → "runs" [0.6], "jumps" [0.4]

Here's how it works:

- **Sentence Rule:** "Sentence → Noun + Verb [1.0]"
  - This rule says that a sentence is always made up of a noun followed by a verb, with a probability of 1.0 (100%).
- **Noun Rule:** "Noun → 'dog' [0.5], 'cat' [0.5]"
  - There's a 50% chance the noun will be "dog" and a 50% chance it will be "cat".
- **Verb Rule:** "Verb → 'runs' [0.6], 'jumps' [0.4]"
  - There's a 60% chance the verb will be "runs" and a 40% chance it will be "jumps".

#### Generating a Sentence:

1. **Choose the Noun:** Randomly pick "dog" or "cat" based on the probabilities.
2. **Choose the Verb:** Randomly pick "runs" or "jumps" based on the probabilities.
3. **Form the Sentence:** Combine the chosen noun and verb to form a sentence like "The dog runs" or "The cat jumps."

#### Applications of Stochastic Grammars

1. **Natural Language Processing (NLP):**
  - Stochastic grammars are used in speech recognition and text generation to model the uncertainty in language. For instance, when a speech recognition system hears a word, it uses stochastic grammars to guess the most likely sequence of words based on probabilities.
2. **Genetic Algorithms and Evolutionary Computation:**

- In these fields, stochastic grammars can be used to generate varied solutions to optimization problems, simulating processes of natural selection and mutation with some level of randomness.
3. **Robotics and Path Planning:**
    - When robots need to navigate unpredictable environments, stochastic grammars can model possible actions and their outcomes, helping the robot decide on a course of action that has the highest probability of success.
  4. **Biological Sequence Analysis:**
    - In bioinformatics, stochastic grammars help in understanding DNA, RNA, or protein sequences by modelling the likelihood of various sequence patterns occurring, which aids in gene prediction or protein structure prediction.
  5. **Computer Music Composition:**
    - Stochastic grammars can be used to generate music by assigning probabilities to different notes or rhythms, creating pieces that are both structured and varied.

Questions:

[Chomsky Normal Form \(CNF\) in Formal Grammar](#)

## Introduction

Chomsky Normal Form (CNF) is a specific type of **Context-Free Grammar (CFG)** used in formal language theory. A grammar is in CNF if its production rules follow one of the two forms:

1.  $A \rightarrow BC$  (Where A, B, and C are non-terminal symbols, and B and C are not the start symbol)
2.  $A \rightarrow a$  (Where A is a non-terminal symbol and 'a' is a terminal symbol)

In CNF, each production must either generate a pair of non-terminals or a single terminal. This standard form is useful for simplifying parsing algorithms, especially for algorithms like **CYK** (Cocke-Younger-Kasami) for parsing context-free grammars.

## Properties of CNF

- Every production rule must be either of the form  $A \rightarrow BC$  or  $A \rightarrow a$ .
- There are no  $\epsilon$ -productions (empty string derivations), except potentially for the start symbol.
- There are no unit productions (rules of the form  $A \rightarrow B$ , where A and B are non-terminals).
- The grammar must be **non-recursive** in nature, meaning it avoids left recursion.

## Why CNF is Important

- CNF simplifies certain parsing techniques, particularly **CYK Parsing**.
- It is also useful in proving properties about languages, such as **Pumping Lemma**.

## Conversion of a CFG to CNF

To convert a given CFG to CNF, follow these general steps:

1. **Eliminate  $\epsilon$ -productions:** Remove any rules that generate the empty string ( $\epsilon$ ), ensuring that  $\epsilon$  is only generated by the start symbol.
2. **Eliminate unit productions:** Remove rules where a non-terminal directly produces another non-terminal (i.e.,  $A \rightarrow B$ ).
3. **Eliminate useless symbols:** Remove any symbols (non-terminals or terminals) that do not contribute to generating terminal strings.
4. **Convert remaining rules to CNF:**
  - If a production has more than two non-terminals (e.g.,  $A \rightarrow BCD$ ), introduce new non-terminal symbols to break the rule into binary form (e.g.,  $A \rightarrow BX$  and  $X \rightarrow CD$ ).
  - If a production has a terminal symbol and non-terminal on the right-hand side (e.g.,  $A \rightarrow aB$ ), introduce new non-terminals for terminals (e.g.,  $A \rightarrow X$  and  $X \rightarrow a$ ).

## Example of CNF Conversion

*Given CFG:*

Let's consider the following CFG:

- $S \rightarrow AB \mid a$
- $A \rightarrow BC \mid \epsilon$
- $B \rightarrow b$
- $C \rightarrow c$

#### *Step 1: Eliminate $\epsilon$ -productions*

Remove  $A \rightarrow \epsilon$  and modify other productions that use  $A$ :

- $S \rightarrow AB \mid a \mid B$
- $A \rightarrow BC$
- $B \rightarrow b$
- $C \rightarrow c$

#### *Step 2: Eliminate Unit Productions*

The production  $S \rightarrow B$  is a unit production. Replace  $B$  with its production  $b$ :

- $S \rightarrow AB \mid a \mid b$
- $A \rightarrow BC$
- $B \rightarrow b$
- $C \rightarrow c$

#### *Step 3: Convert to CNF*

- The production  $S \rightarrow AB$  is already in CNF ( $A \rightarrow BC$  form).
- The production  $S \rightarrow a$  is already in CNF ( $A \rightarrow a$  form).
- The production  $A \rightarrow BC$  is in CNF.
- The production  $B \rightarrow b$  is in CNF.
- The production  $C \rightarrow c$  is in CNF.

After these steps, the grammar is already in CNF.

#### **Final CNF Grammar:**

- $S \rightarrow AB \mid a \mid b$
- $A \rightarrow BC$
- $B \rightarrow b$
- $C \rightarrow c$

#### *b) Abstract View of the Parsing Problem*

The parsing problem in computer science and linguistics refers to the process of analyzing a sequence of symbols (often a string of characters or tokens) to determine its syntactic structure based on a formal grammar. Parsing is a key step in many applications, such as natural language processing (NLP), compilers, and interpreters. Here is an abstract view of the parsing problem:

1. **Input Sequence:**  
The parsing process begins with an input sequence, which can be a string of symbols or tokens. In NLP, this might be a sentence, and in a compiler, it could be source code.
2. **Grammar:**  
A formal grammar (e.g., Context-Free Grammar or Stochastic Grammar) defines the syntactic rules that determine how the input sequence should be parsed. These rules describe how sequences of symbols can be combined to form valid structures, such as sentences, expressions, or statements.
3. **Parsing Algorithm:**  
The parsing algorithm takes the input sequence and the grammar to generate a syntactic structure, typically a parse tree or abstract syntax tree (AST). The algorithm must determine how the input symbols match the rules of the grammar.
4. **State Transitions:**  
Parsing can be seen as a state transition process, where the parser moves from one state to another as it processes each symbol or token in the input. The state represents a particular point in the parsing process, such as recognizing a noun phrase or verb phrase in a sentence.
5. **Stack/Buffer:**  
A **stack** (or sometimes a buffer) is used to manage the state transitions during parsing. The stack keeps track of intermediate

states, especially when recursion or backtracking is involved. In some parsing approaches, the input is pushed onto the stack, and intermediate derivations are tracked until the parsing is complete.

6. **Derivation Process:**

The parser uses the grammar rules to generate derivations for the input sequence. Each rule application represents a step in the derivation process. The goal is to derive the input sequence from the starting symbol (usually the start symbol of the grammar) to reach a complete parse.

7. **Ambiguity Handling:**

In some cases, there may be multiple valid ways to parse the input sequence. The parsing algorithm must handle ambiguity, selecting the most appropriate parse tree or choosing among competing interpretations. In probabilistic parsing, the likelihood of different parses is evaluated to select the most probable one.

8. **Completion:**

Parsing is complete when the parser has consumed all input symbols, and the final structure (such as a parse tree) corresponds to the input sequence according to the grammar. If successful, the structure is usually used for further processing, such as semantic analysis, code generation, or translation.

9. **Error Handling:**

If the input cannot be parsed according to the grammar, an error must be reported. The parser may provide feedback indicating which part of the input caused the failure, enabling corrections or debugging.

## Simple Examples of String Generation as Pattern Description

Here are two simpler examples of string generation using pattern descriptions:

---

### 1. Context-Free Grammar (CFG) for Simple Sentences

A **Context-Free Grammar (CFG)** can be used to generate simple sentences consisting of a subject and a verb.

#### Example Grammar:

```
mathematica
Copy code
S → NP VP
NP → Det Noun
VP → Verb
Det → "the" | "a"
Noun → "cat" | "dog"
Verb → "runs" | "sleeps"
```

This grammar generates simple sentences like:

- **"The cat runs"**
- **"A dog sleeps"**

Here:

- **S** (Sentence) consists of a noun phrase (NP) and a verb phrase (VP).
- **NP** (Noun Phrase) is made up of a determiner (Det) and a noun (Noun).
- **VP** (Verb Phrase) consists of a verb.

---

### 2. Regular Expression for Simple Phone Numbers

A **regular expression (regex)** can be used to generate or validate simple phone numbers in the format xxx-xxx-xxxx, where x is a digit.

#### Example Regular Expression:

```
ruby
Copy code
^\d{3}-\d{3}-\d{4}$
```

This regex matches strings like:

- **"123-456-7890"**
- **"555-123-4567"**

Here:

- `\d{3}` matches three digits.
- The hyphens (-) are literal characters separating the groups of digits.
- `\d{4}` matches the final four digits of the phone number.

## Types of String Grammar

String grammar refers to a set of formal rules or production rules that describe the structure of strings (sequences of symbols). These grammars are used in various fields like linguistics, formal language theory, and programming languages to define syntactic structures. There are several types of string grammars, each with different levels of expressiveness and computational complexity. The main types of string grammars are:

---

### 1. Regular Grammar (Type 3)

#### Description:

A **Regular Grammar** is the simplest type of grammar. It consists of production rules where each rule generates a terminal string or a terminal symbol followed by a non-terminal symbol. Regular grammars can describe regular languages and are equivalent to finite automata.

#### Production Rules:

- **Right-linear:** A rule where the non-terminal appears at the end of the production. Example:  $A \rightarrow aB$
- **Left-linear:** A rule where the non-terminal appears at the beginning of the production. Example:  $A \rightarrow Ba$

#### Examples:

- **Regular Expression (Regex)** is a practical example of regular grammar, used to define simple patterns, such as strings of digits or alphabetic characters. Example: `a*` (Zero or more occurrences of 'a')

#### Applications:

- Lexical analysis in compilers
  - Text search and pattern matching
  - Describing simple languages such as identifiers and keywords in programming languages
- 

### 2. Context-Free Grammar (CFG) (Type 2)

#### Description:

A **Context-Free Grammar (CFG)** is more powerful than a regular grammar. In CFG, the left-hand side of a production rule is always a single non-terminal symbol, and the right-hand side can be a combination of terminals and non-terminals. CFGs can generate context-free languages and are widely used in programming language syntax and compiler design.

#### Production Rules:

- A single non-terminal on the left-hand side and a string of terminals and/or non-terminals on the right-hand side. Example:  $S \rightarrow aSb \mid \epsilon$

#### Examples:

- A CFG can describe balanced parentheses:

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

#### Applications:

- Parsing and syntax analysis in compilers
- Defining the syntax of programming languages
- Natural language processing (e.g., sentence structure)

### 3. Context-Sensitive Grammar (CSG) (Type 1)

A **Context-Sensitive Grammar (CSG)** is more general than a context-free grammar. In CSG, the production rules are more flexible, allowing the length of the string on the left-hand side to be greater than or equal to the string on the right-hand side. These grammars can describe context-sensitive languages, which are more complex than context-free languages.

#### Production Rules:

- A production rule where the left-hand side can have more than one non-terminal symbol, and the right-hand side can generate longer strings. Example:  $\alpha A \beta \rightarrow \alpha B \beta$  (Here, A is replaced by B in the context of symbols  $\alpha$  and  $\beta$ .)

#### Examples:

- A CSG can describe the language of the form  $\{a^n b^n c^n \mid n \geq 1\}$ , where the number of 'a's, 'b's, and 'c's is the same.

#### Applications:

- Some natural language structures that can't be expressed by CFG
- More complex syntactical structures in advanced compilers

Aspect	Grammar	Language
Definition	A formal set of rules used to define the structure of strings in a language.	A set of valid strings or sentences that can be generated or recognized by a grammar.
Purpose	To specify how words, symbols, or sentences are formed within a language.	To represent all the possible strings that can be formed based on the grammar's rules.
Components	Non-terminal symbols, terminal symbols, production rules, start symbol.	A set of strings (sequences of terminal symbols) that adhere to the grammar's rules.
Nature	A description of how language is structured.	A set of all valid strings that can be produced using the grammar.
Role	Describes the syntactic structure of the language.	Contains all valid words, phrases, or sentences generated by the grammar.
Type	A set of formal rules or a formal system.	A collection of strings that belong to the language defined by the grammar.
Example	A Context-Free Grammar (CFG) for generating simple arithmetic expressions: $E \rightarrow E + E$	$E^*E$



<b>Formality</b>	Grammars are formalized as a set of production rules (e.g., CFG, Regular Grammar, etc.).	Languages are formalized as a set of strings that are accepted by the grammar.
<b>Examples in Computation</b>	Context-Free Grammar (CFG) for a programming language, Regular Grammar for lexical analysis.	All the valid programs or sentences that can be written in a given programming language.
<b>Expressiveness</b>	Grammar's power and expressiveness depend on the type (e.g., Regular Grammar, CFG, etc.).	The language's complexity depends on the grammar used (e.g., regular languages, context-free languages).
<b>Type of Representation</b>	Typically represented as a set of rules (e.g., $S \rightarrow NP VP$ , $NP \rightarrow Det Noun$ ).	Represented as a set of strings that comply with the grammar's rules (e.g., <code>the cat sleeps</code> ).
<b>Real-World Example</b>	A CFG defining the structure of valid English sentences.	The set of valid English sentences like "The cat sleeps", "A dog runs."
<b>Abstraction Level</b>	Grammar provides an abstract description of the language's structure.	Language is the actual instantiation of all valid strings that follow a grammar's structure.
<b>Role in Computational Models</b>	In computational models (like parsers), grammars define how to analyze and generate strings.	In computational models, languages are the input/output of parsing and generation processes.
<b>Context</b>	Often used in compiler design, natural language processing, and formal language theory.	Used in areas like natural language processing, text generation, and formal language theory.
<b>Size and Scope</b>	A grammar may describe infinite possible sentences (e.g., recursive rules) with a finite set of rules.	A language may be finite or infinite depending on the grammar and the rules it uses.
<b>Dependency</b>	Grammar defines the rules that a language follows, but does not represent any concrete data.	A language is the set of strings formed by the grammar's rules, representing actual data or sentences.

## Blocks Word Description String Generation Example as Pattern Description

In the context of **String Generation**, **Pattern Description** refers to the use of grammars and rules to describe how strings (sequences of symbols) can be generated. One such method for describing strings in a formalized way is the **Block Word Description** method, which can be used to represent or generate a sequence of words or symbols based on defined patterns.

### Block Word Description

Block word description refers to using patterns and structures to describe how sequences of words can be formed. These "blocks" represent groups of words or symbols that are generated using specific production rules or patterns. This method helps in defining a formal way to generate sequences (words, phrases, or entire sentences) by stacking blocks of rules in an organized fashion.

In pattern description, blocks are typically non-terminal symbols or groups that are expanded or generated according to predefined rules.

### Example: Block Word Description in String Generation

Let's consider a simple example of generating a sentence using blocks that represent different parts of speech (POS) such as **noun phrases (NP)**, **verb phrases (VP)**, **articles (Det)**, **nouns (Noun)**, and **verbs (Verb)**.

#### Step 1: Define the Blocks (Parts of Speech)

Each block corresponds to a set of possible symbols (words). Here's how we can define each block:

- **NP (Noun Phrase):** This block can consist of a **Determiner (Det)** and a **Noun**.
  - $NP \rightarrow Det\ Noun$
- **VP (Verb Phrase):** This block consists of a **Verb** and possibly another **NP** (i.e., a direct object).
  - $VP \rightarrow Verb\ NP$
- **Det (Determiner):** This block consists of words like "the", "a".
  - $Det \rightarrow "the" \mid "a"$
- **Noun:** This block contains words like "dog", "cat".
  - $Noun \rightarrow "dog" \mid "cat"$
- **Verb:** This block consists of action words like "chases", "catches".
  - $Verb \rightarrow "chases" \mid "catches"$

### *Step 2: Combine Blocks to Form Sentences*

Now, by using the rules to generate strings, we can combine these blocks to form sentences.

- **Sentence Generation:** We start with a Sentence (S) block and expand it using the rules for **NP** and **VP**.
  - $S \rightarrow NP\ VP$
- **Expand NP:** An **NP** consists of a **Det** and a **Noun**.
  - $NP \rightarrow Det\ Noun$
  - For example:  $NP \rightarrow "the" "dog"$
- **Expand VP:** A **VP** consists of a **Verb** and another **NP**.
  - $VP \rightarrow Verb\ NP$
  - For example:  $VP \rightarrow "chases" "the\ cat"$

Combining all these, we get the final sentence:

- $S \rightarrow NP\ VP \rightarrow "the\ dog" "chases" "the\ cat"$

### *Step 3: Variations of Sentence Generation*

By altering the components chosen in each block, we can generate other valid sentences based on the same pattern.

- "a cat chases the dog"
- "the dog catches the cat"
- "a cat catches a dog"