

Introduction to Formal Language Theory

Course outline

Chapter 1: Basics

- Set theory
- Relations & functions
- Mathematical induction
- Graphs & trees
- Strings & languages

Chapter 2: Introduction to grammars

Chapter 3: Regular languages

- Regular grammar
- Automata
- Regular expressions

Chapter 4 :Context Free Languages

- Context free grammars
- Normal forms

Chapter 5: Push Down Automata (PDA)

- NPDA
- DPDA

Basics: outline

- Overview of languages: natural vs formal
- Review of set theory and relations
 - Set theory
 - Relations and functions
- Mathematical induction
- Graphs and trees
- Strings and languages

Overview of languages : natural Vs formal

- Language is a set of strings or sentences.
- Natural Languages
 - rules come after the language
 - evolve and develop
 - highly flexible
 - quite powerful
 - no special learning effort needed

Disadvantages

- vague
- imprecise
- ambiguous
- user and context dependent
- Ex. Amharic, English, French, ...

Overview of languages: cont'd

■ Formal Languages

- ❑ developed with strict rules
 - predefined syntax and semantics
- ❑ precise
- ❑ unambiguous
 - can be processed by machines!

Disadvantages

- ❑ unfamiliar notation
- ❑ initial learning effort
- ❑ Ex. Programming languages: Pascal, C++, ...

Overview of languages: cont'd

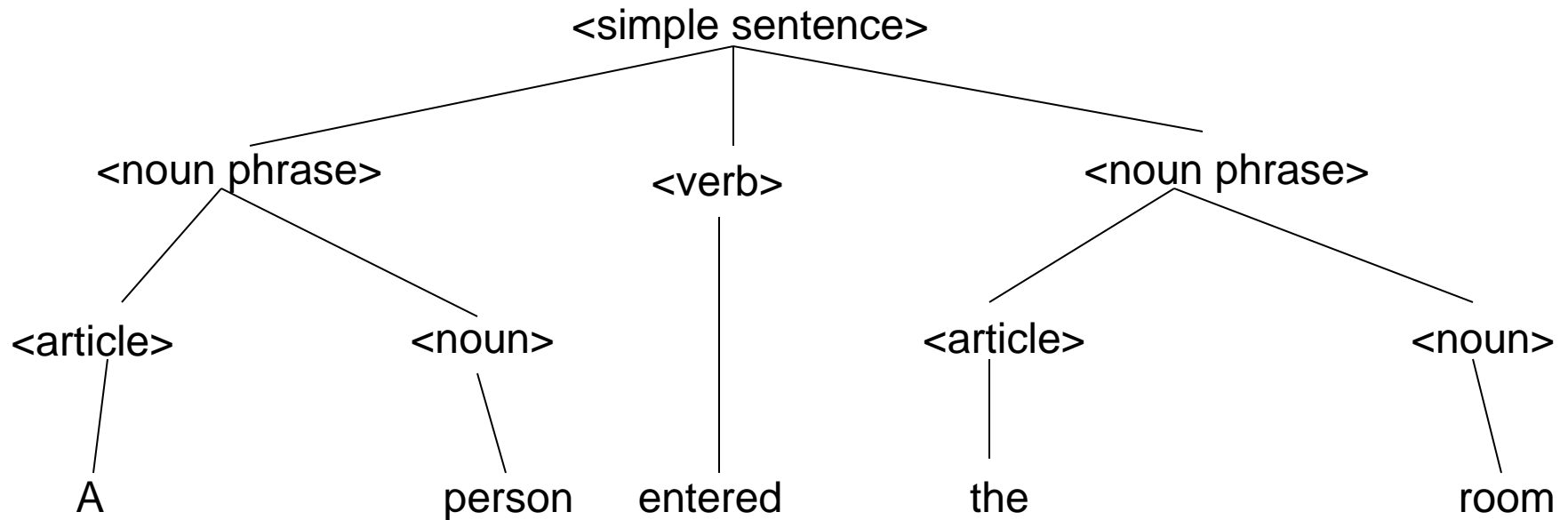
- Sentences: the basic building blocks of languages
- Sentence = Syntax + Semantics
- Grammar: the study of the structure of a sentence
- Ex:

<simple sentence> ::= <noun phrase><verb><noun phrase>

<noun phrase> ::= <article><noun>

→ A person entered the room

Overview of languages: cont'd



Derivation tree for the simple sentence: A person entered the room.

Overview of languages: cont'd

- In Pascal (as well as in many other languages), for example, an identifier is specified as follows:

$\langle \text{identifier} \rangle ::= \langle \text{letter} \rangle \mid \langle \text{letter} \rangle \{ \langle \text{letter} \rangle \mid \langle \text{digit} \rangle \}^*$

$\langle \text{letter} \rangle ::= a \mid b \mid c \dots$

$\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid \dots \mid 9$

Ex. a, x1, num, count1, ...

Review of set theory and relations

■ Sets

- A well defined collection of objects (called members or elements)
- Notation: $a \in S \rightarrow a$ is an element of the set S

■ Operation on sets

Let A and B be two sets and U the universal set

- Subset: $A \subseteq B$
- Proper subset: $A \subset B$
- Equality: $A = B$
- Union: $A \cup B$
- Intersection: $A \cap B$
- Set difference: $A \setminus B$ or $A - B$
- Complement: A' or A bar
- Cartesian product: $A \times B = \{(a,b) \mid a \in A \text{ and } b \in B\}$

Note: (a,b) is called an **ordered pair**, and is different from (b,a)

Set theory and relations: cont'd

■ Properties

Let A , B , C be sets and U the universal set

- Associative property: $A \cup (B \cup C) = (A \cup B) \cup C$
- Commutative property: $A \cup B = B \cup A$
- Demorgan's laws: $(A \cup B)' = A' \cap B'$, ...
- Involution law: $(A')' = A$

Definitions:

- Let A be a set. The **cardinality of a set A** is a measure of the "number of elements of the **set**" and denoted by $|A|$ or $\#(A)$.
- The set of all subsets of a set A is called the power set of A , denoted by 2^A .

Set theory and relations: cont'd

Definition:

Let S be a set. A collection $\{A_1, A_2, \dots, A_n\}$ of subsets of S is called a partition if $A_i \cap A_j = \emptyset$, $i \neq j$ and $S = A_1 \cup A_2 \cup \dots \cup A_n$.

Ex. $S = \{1, 2, \dots, 10\}$

Let $A_1 = \{1, 3, 5, 7, 9\}$ and $A_2 = \{2, 4, 6, 8, 10\}$, then $\{A_1, A_2\} = \{\{1, 3, 5, 7, 9\}, \{2, 4, 6, 8, 10\}\}$ is a partition of S .

Q. Find other partitions of S

■ Countability

- A finite set is countable
- If the elements of set A can be associated with $1^{\text{st}}, 2^{\text{nd}}, \dots, i^{\text{th}}, \dots$ elements of the set of Natural Numbers, then A is countable.

Note: that in this case A may not be finite.

Ex.

1. $N = \{1, 2, \dots, i^{\text{th}}, \dots\}$ is countable
2. $Z = \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\} = \{0, 1, -1, 2, -2, 3, -3, \dots\}$ is countable
3. $[0, 3]$ is uncountable (not countable)

Relations and functions

■ Relations

- Definition: A relation R is a set of ordered pairs of elements in S . (i.e is a subset of $S \times S$)

Notation: $(x, y) \in R$ or $x R y$

- Properties of relations

- Let R be a relation on a set A , then

- R is **reflexive** if for all $a \in A$, $a R a$ or $(a, a) \in R$
- R is **symmetric** if $a R b \Rightarrow b R a$
- R is **transitive** if $a R b$ and $b R c \Rightarrow a R c$, for all $a, b, c \in R$
- R is an **equivalence relation** if (a), (b) and (c) above hold.

- Let R be an equivalence relation on set A and let $a \in A$, then the equivalence class of a , denoted by $[a]$, is defined as:

$$[a] = \{b \in A \mid a R b\}$$

Relations and functions: cont'd

Examples:

Check whether the following relations are reflexive, symmetric, and transitive

1. Let R be a relation in $\{1, 2, 3, 4, 5, 6\}$ is given by $\{(1,2), (2, 3), (3, 4), (4, 4), (4, 5)\}$
2. Let R be a relation in $\{1, 2, 3, \dots, 10\}$ defined as $a R b$ if a divides b
3. Let R be defined on a set S such that $a R b$ if $a=b$
4. Let R be defined on all people in Addis Ababa by $a R b$ if a and b have the same date of birth.

Relations and functions: cont'd

■ Functions

- Definition: A function f from a set X to a set Y is a rule that associates to every element x in X a unique element in Y , which is denoted by $f(x)$.
 - The element $f(x)$ is called the image of x under f .
 - The function is denoted by $f: X \rightarrow Y$
- Functions can be defined in the following two ways:
 1. By giving the images of all elements of X
Ex. $f: \{1, 2, 3, 4\} \rightarrow \{2, 4, 6\}$ can be defined by
 $f(1) = 2, f(2) = 4, f(3) = 6, f(4) = 6$
 2. By a computational rule which computes $f(x)$ once x is given
Ex. $f: \mathbb{R} \rightarrow \mathbb{R}$ can be defined by $f(x) = x^2 + 2x + 1, x \in \mathbb{R}$ (\mathbb{R} = the set of all real numbers)

Relations and functions: cont'd

■ Let $f: A \rightarrow B$ be a function

1. f is an **into** function if $R_f \subseteq B$
2. f is an **onto** function if $R_f = B$
3. f is a **one-to-one** function
if for x_1 & $x_2 \in A$, $x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2)$
4. f is bijective (one-to-one correspondence) if it satisfies (2) and (3) above.

Ex. $f: \mathbb{Z} \rightarrow \mathbb{Z}$ is given by $f(x) = 2x$

Show that f is one-to-one but not onto.

- Definition: A set A is said to be countable iff there exists a function $f: A \rightarrow \mathbb{N}$ such that f is bijective. (\mathbb{N} =the set of natural numbers)

Mathematical induction

- Let P_n be a proposition that depends on $n \in \mathbb{Z}^+$.
Then P_n is true for all +ve n provided that:

- P_1 is true
- If P_k is true, so is P_{k+1} , for some $k \in \mathbb{Z}^+$.

Three steps:

1. Base case: verify that P_1 holds
2. Inductive hypothesis: assume that P_k holds, for some $k \in \mathbb{Z}^+$
3. Inductive step: show that P_{k+1} holds

Ex. Show that $1+2+\dots+n = n(n+1)/2$, for all $n \in \mathbb{Z}^+$.

Mathematical induction: cont'd

Solution:

Let $P_n: 1+2+\dots+n = n(n+1)/2$

Step1: for $n = 1$, P_1 holds

Step2: for some $k \in \mathbb{Z}^+$, assume P_k is true

i.e. $P_k: 1+2+\dots+k = k(k+1)/2$

Step3: WTS P_{k+1} is true

$$P_{k+1} : 1+2+\dots+k+(k+1) = (k+1)(k+2)/2$$

$$: P_k + (k+1) = (k+1)(k+2)/2$$

$$: k(k+1)/2 + (k+1) = (k+1)(k+2)/2$$

$$: [k(k+1) + 2(k+1)]/2 = (k+1)(k+2)/2$$

$$: (k+1)(k+2)/2 = (k+1)(k+2)/2$$

Therefore, P_n holds for all $n \in \mathbb{Z}^+$

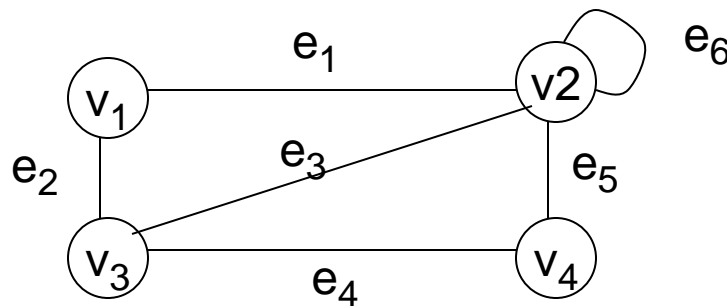
Ex. Show that $P_n = \sum_{i=1,n} (i^2) = (n+1)(n)(2n+1)/6$ for all n

Graphs and trees

■ Graphs

□ Definition: A graph (undirected graph) consists of:

- a. A non-empty set V called the set of vertices,
- b. A set E called the set of edges, and
- c. A map Φ (phi) which assigns to every edge a unique unordered pair of vertices



$$e_1 = \{v_1, v_2\}$$

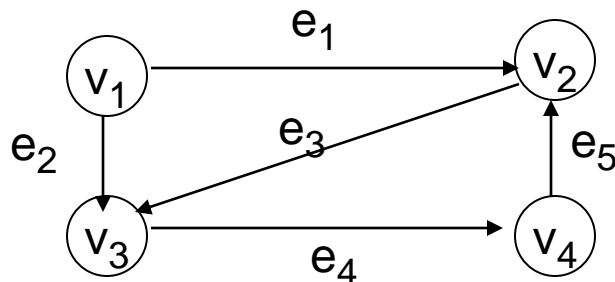
$$e_2 = \{v_1, v_3\}$$

...

$$e_6 = \{v_2, v_2\} \text{ (a self loop)}$$

Graphs and trees: cont'd

- Definition: A directed graph (digraph) consists of:
 - a. A non-empty set V called the set of vertices,
 - b. A set E called the set of edges, and
 - c. A map Φ (phi) which assigns to every edge a unique ordered pair of vertices



$e_1 = (v_1, v_2)$
 v_1 : a predecessor of v_2
 v_2 : a successor of v_1

Graphs and trees: cont'd

- Definition: The degree of a vertex v in a graph (directed or undirected) is the number of edges with v as an end vertex.

Note: that a self loop is counted twice when calculating the degree of a vertex.

Ex. In the previous graph, $\deg(v_1) = ?$ $\deg(v_2) = ?$

- Definition: A path in a graph (directed or undirected) is an alternating sequence of vertices and edges of the form $v_1 e_1 v_2 e_2 \dots e_{n-1} v_n$, beginning and ending with vertices such that e_i has v_i and v_{i+1} as its end vertices and no edge or vertex is repeated in the sequence.

The path is said to be from v_1 to v_n .

Ex. In the previous graph, $v_1 e_1 v_2 e_3 v_3 e_4 v_4$ is a path from v_1 to v_4 .

Note: that a path may be directed (if all the edges in the path have the same direction.)

Graphs and trees: cont'd

- Definition: A graph (directed or undirected) is **connected** if there is a path between every pair of vertices.

Q. Are the previous two graphs connected?

- Definition: A circuit in a graph is an alternating sequence $v_1 e_1 v_2 e_2 \dots e_{n-1} v_1$ of vertices and edges starting and ending with the same vertex such that e_i has v_i and v_{i+1} as end vertices and no edge or vertex other than v_1 is repeated.

Ex. $V_2 e_3 v_3 e_4 v_4 e_5 v_2$ is a circuit in the previous graph

Graphs and trees: cont'd

■ Trees

- Definition: A graph (directed or undirected) is called a tree if it is connected and has no circuits.

Q. Are the previous two graphs trees?

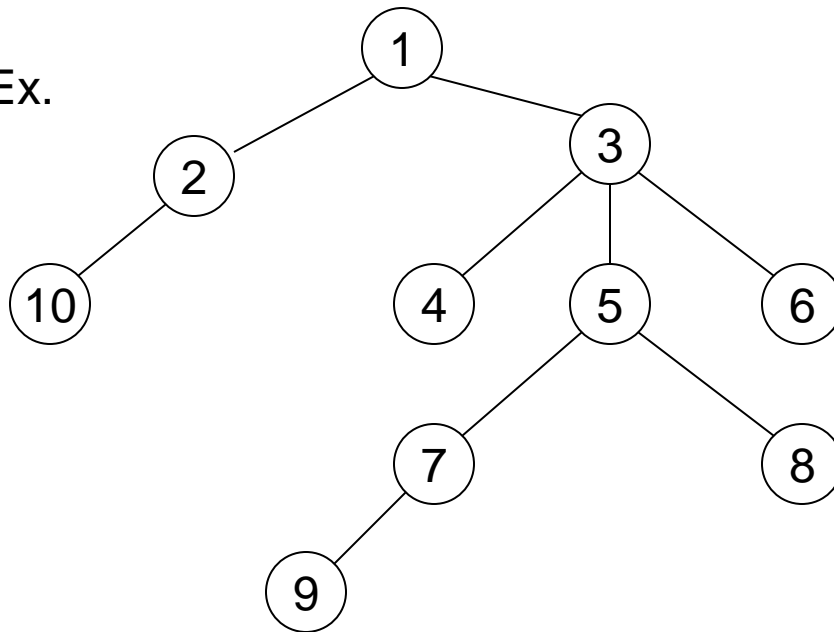
- Properties of trees:
 - In a tree there is one and only one path between every pair of vertices (nodes)
 - A tree with n vertices has $n-1$ edges
 - A **leaf** in a tree can be defined as a vertex of degree one
 - Vertices other than leaves are called internal vertices

Graphs and trees: cont'd

- Definition: An ordered directed tree is a digraph satisfying the following conditions:
 - There is one vertex called the **root** of the tree which is distinguished from all other vertices and the root has no predecessors.
 - There is a directed path from the root to every other vertex.
 - Every vertex except the root has exactly one predecessor.
(For the sake of simplicity, we refer to ordered directed trees as simply trees.)
- The number of edges in a path is called the **length** of the path.
- The height of a tree is the length of the longest path from the root.
- A vertex v in a tree is at level k if there is a path of length k from the root to the vertex v .
Q. what is the maximum possible level in a tree?
- There are several types of trees: binary, balanced binary, binary search tree, heap, general tree, ...

Graphs and trees: cont'd

Ex.



1. List the leaves.
2. List the internal nodes.
3. What is the length of the path from 1 to 9?
4. What is the height of the tree?

- Note: a path from vertex (node) n_1 to node n_k can be simply expressed as the sequence of nodes n_i , $i=1, \dots, k$ such that n_i is the parent (predecessor) of n_{i+1} ($1 \leq i \leq k$)

Strings and languages

■ Strings

- An alphabet, Σ , is a set of finite symbols.
- A string over an alphabet Σ is a sequence of symbols from Σ .
- An empty string is a string without symbols, and is denoted by λ .
- Let w be a string, then its length, denoted by $|w|$, is the number of symbols of w .

Ex. Let $\Sigma = \{0, 1\}$, the following are some strings over Σ

$w = \lambda$, $|w| = 0$; $w = 01$, $|w| = 2$; $w = 010110$, $|w| = 6$

- Given an alphabet Σ , Σ^* denotes the set of all strings (including λ) over Σ .
- $\Sigma^+ = \Sigma^* - \{\lambda\}$

Ex. $\Sigma = \{0, 1\} \Rightarrow \Sigma^* = \{\lambda, 0, 1, 01, 00, 11, 111, 0101, 0000, \dots\}$

- Σ^i is a set of strings of length i , $i = 0, 1, 2, \dots$
- Let $x \in \Sigma^*$ and $|x| = n$, then $x = a_1 a_2 \dots a_n$, $a_i \in \Sigma$

Strings and languages: cont'd

□ Operations on strings

■ Concatenation operation

□ Let $x, y \in \Sigma^*$ and $|x| = n$ and $|y| = m$. Then xy , concatenation of x and y , $= a_1a_2\dots a_nb_1b_2\dots b_m$, $a_i, b_i \in \Sigma$

□ The set Σ^* has an identity element λ with respect to the binary operation of concatenation.

Ex. $x \in \Sigma^*$, $x\lambda = \lambda x = x$

□ Σ^* has left and right cancellation

For $x, y, z \in \Sigma^*$,

$zx = zy \Rightarrow x = y$ (left cancellation)

$xz = yz \Rightarrow x = y$ (right cancellation)

□ For $x, y \in \Sigma^*$, we have $|xy| = |x| + |y|$

Strings and languages: cont'd

■ Transpose operation

- For any x in Σ^* and a in Σ , $(xa)^T = a(x)^T$

Ex. $(aaabab)^T = babaaa$

- A **palindrome** of even length can be obtained by the concatenation of a string and its transpose.
- A **prefix** of a string is a substring of leading symbols of that string.

w is a prefix of y if there exists y' in Σ^* such that $y=wy'$

Ex. $y = 123$, list all prefixes of y .

- A **suffix** of a string is a substring of trailing symbols of that string.

w is a suffix of y if there exists y' in Σ^* such that $y=y'w$

Ex. $y = 123$, list all suffixes of y .

Strings and languages: cont'd

- A **terminal** symbol is a unique indivisible object used in the generation of strings.
- A **nonterminal** symbol is a unique object but divisible, used in the generation of strings.

Ex. In English, *a*, *b*, *A*, *B*, etc are terminals and the words *boy*, *cat*, *dog*, ... are nonterminals.

In programming languages, *a*, *A*, *:*, *;*, *=*, *if*, *then*, ... are terminals

Strings and languages: cont'd

■ Languages

- Definition: A language, L , is a set (collection) of strings over a given alphabet, Σ .

- A string in L is called a sentence or word.

Ex. $\Sigma = \{0, 1\}$, $\Sigma^* = \{\lambda, 0, 1, 01, 00, 11, \dots\}$

$L_1 = \{\lambda\}$, $L_2 = \{0, 1, 01\}$ over Σ

$L_3 = \{a^n \mid n \geq 0\}$ over $\Sigma = \{a\}$

- Let L_1, L_2 be languages over Σ , then

- $L_1 L_2 = \{xy \mid x \in L_1, y \in L_2\}$

- $L\{\lambda\} = \{\lambda\}L = L$, for any language L

- $L^0 = \{\lambda\}$

- $L^1 = L$

- $L^2 = LL \equiv \{xx \mid x \in L\}$

- ...

- $L^i = L^i L^{i-1}$, for $i \geq 2$

- $L^* = \bigcup_{i=0, \infty} (L^i)$