# 01204211 Discrete Mathematics
## Lecture 7a: Languages and regular expressions[1]
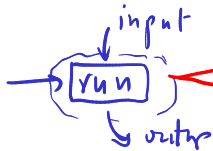
Jittat Fakcharoenphol

August 22, 2023

# What is computation?

- เขียนโปรแกรม python → run
  - input
  - output

model ของการคำนวณ



เปลี่ยน model ในนี้

- Functional prog. languages
  Haskell, Lisp,
  JavaScript

→ classical physics

→ quantum physics → "quantum computer"

เปลี่ยน model ในนี้

Models of computations → Machines

Machines
- finite automata ⊆
- Turing machines

λ-calculus
(functional prg)

grammars
वाक्य [pattern]
*
- regular expression
- context-free grammar

ภาษา
Languages = specifications

อักฤน จัน ทั้งัน

natural language

NLP.

↳ _formal_ language.

set ของ __strings__ ที่เป็นของชุดๆ

ex. { w | w แทน เลขจำนวน }

{ p | p เป็น จำนวนเฉพาะ เป็นกลุ่ม ด เลข—ๆ }

{ p | p เป็น palindrome }

{ s | s เป็น สตริง ที่ประกอบด้วย '(' ')' ที่ถูกต้องปูงๆๆ }

# Formal definition: strings

Intuitively, a string is a *finite* sequence of symbols. However, to be able to formally prove properties of strings we need a precise definition.

Let a finite set $\Sigma$ be the **alphabet**. (E.g., for bit strings, $\Sigma = \{0, 1\}$; for digits, $\Sigma = \{0, 1, \ldots, 9\}$; for English string $\Sigma = \{a, b, \ldots, z\}$.)
The following is a recursive definition of strings.

good        g · o · o · d

## Recursive definition of strings

A **string** $w$ over alphabet $\Sigma$ is either

► the empty string $\varepsilon$, or

► $a \cdot x$ where $a \in \Sigma$ and $x$ is a string.

The set of all strings over alphabet $\Sigma$ is denoted by $\Sigma^*$.

# Review: more recursive definitions

## Lengths

For a string $w$, let $|w|$ be the length of $w$ defined as

$$|w| = \begin{cases} 0 & \text{when } w = \varepsilon \\ 1 + |x| & \text{when } w = a \cdot x \end{cases}$$

*(handwritten annotation: base case)*

## Concatenation

For strings $w$ and $z$, the concatenation $w \cdot z$ is defiend recursively as

$$w \cdot z = \begin{cases} z & \text{when } w = \varepsilon \\ a \cdot (x \cdot z) & \text{when } w = a \cdot x \end{cases}$$

# Review: proving facts about strings

## Lemma 1

*For strings $w$ and $x$, $|w \cdot x| = |w| + |x|$.*

## Proof.

จ: พิสูจน์ โดยทฤษฎีด้วย.

Induction Hypothesis (IH): " สำหรับ ทุก string $s$ ที่ $|s| < |w|$
$$|s \cdot x| = |s| + |x|. \text{"}$$

__Case 1__: $w = \varepsilon$ $\qquad |w \cdot x| = |x| \qquad$ นิยามของ•
$$= 0 + |x|$$
$$= |w| + |x| \qquad \text{นิยามของ length} \quad \checkmark$$

__Case 2__: $w = a \cdot y$
โดย $a \in \Sigma$. ··· $|w \cdot x| = |(a \cdot y) \cdot x| = |a \cdot (y \cdot x)|$
$$= 1 + |y \cdot x| \qquad \text{นิยามของ length}$$
$$= \underline{1 + |y| + |x|} \qquad \text{จาก IH}$$

# Formal languages

A **formal language** is a set of strings over some finite alphabet $\Sigma$.
Examples:

$\Sigma = \{0, 1\}$ , $\Sigma^* = \{$ strings ทั้ง.

0

$\{\omega \in \Sigma^* \mid \omega$ ขึ้นต้นด้วย $1 \}$ ex. $\{1, 10, 1011, \ldots$

$\{\omega \in \Sigma^* \mid \omega$ เป็น จำนวน เลขคี่ $\}$ ex $\{10, 11, 101, 111, \ldots$
repr.

$\Sigma = \{a, b\}$

$\{\omega \in \Sigma^* \mid \omega$ เป็น palindrome $\}$ ex. $\{a$ aba aaa , abba
bab , bbabb …

$\emptyset$ , $\{\varepsilon\}$

# Careful...

These are different languages: $\emptyset, \{\varepsilon\}$
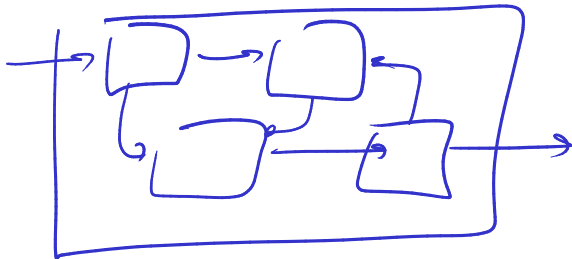And $\varepsilon$ is not a language.

# How to describe languages?

① ระบุกฎ กติกา ธรรมชาติ + mathematical notation

$$\{ w \mid \text{———} \}$$

② วิธีที่ ระบุแล้ว

↳ ทำให้ คำนวณได้ง่าย ....

⓪ $\{0, 01, 11, 1001\}$    ไม่ต้องมีบทกติกา.   ⋮ ⋮

# Composition

# Combining languages

If $A$ and $B$ are languages over alphabet $\Sigma$.

▶ Basic set operations: $A \cup B$, $A \cap B$, $\bar{A} = \Sigma^* \setminus A$.

set minus

▶ Concatenation: $A \cdot B$.

$$A \cdot B = \{ x \cdot y \mid x \in A, y \in B \}$$

$A = \{ HELLO, GOODBYE \}$

$B = \{ PITA, THAKSIN, SETTHA, CHU \}$

$C = \{ 1, 2, \ldots \}$

$A \cdot C = \ldots$

▶ Kleene closure or Kleene star: $A^*$.

$$
\begin{array}{l}
w \in A^* \text{ iff} \\
① \quad w = \varepsilon \\
② \quad w = x \cdot y \quad \text{เมื่อ } x \in A \text{ และ } y \in A^*
\end{array}
$$

ex $\{0,1\}^*$
$= \{ \varepsilon, 0, 1, 00, 10, 01, 11, \ldots \}$

Also $A^+ = A \cdot A^*$

# Examples

# Regular languages

## Definition: regular languages

A language $L$ is **regular** if and only if it satisfies one of the following conditions:

- ▶ $L$ is empty;    $\emptyset$
- ▶ $L$ contains one string (can be the empty string $\varepsilon$);    {here}
- ▶ $L$ is a union of two regular languages;
- ▶ $L$ is the concatenation of two regular languages; or
- ▶ $L$ is the Kleene closure of a regular language.

Examples  $\Sigma = \{0, 1\}$  ~~{0111...}~~

$\phi$

$\{\varepsilon\}$ , $\{0\}$ , $\{010\}$ , $\{10\}$ , $\{0001\}$ , $\{1001\}$ , $\{011111\}$

$\phi \cup \{0\}$ , $\underline{\{010\} \cup \{10\}}$ , $\{0001\} \cup (\{010\} \cup \{10\})$

$\{0\} \bullet \{010\}$ , $\phi \bullet \{10\}$ , $\{\varepsilon\} \bullet \{10\}$ , $(\{010\} \cup \{10\}) \bullet (\{1001\} \cup \{10\})$

$\underset{\|}{}$ $\{00010\}$     $\underset{\|}{}$ ~~$\phi$~~     $\underset{\|}{}$ $\{10\}$

$$\left[ \left[ (\{01\} \cup \{10\}) \bullet \{1\} \right]^* \bullet \{0\} \bullet (\{01\} \cup \{10\}) \right]^*$$

# Regular expressions

Let $\Sigma = \{0, 1\}$. Consider

$$(((\{01\} \cup (\{1\} \cdot (\{0\} \cup \{10\}))) \cup (\{00\} \cdot (\{1\})^*)) \cdot ((\{0\} \cdot \{0\}) \cdot \{1\})$$



regular
expression
tree

$(\{0\} \cdot \{0\}) \cdot \{1\}$

# Regular expressions

Regular language

$$(( \{01\} \cup (\{1\} \cdot (\{0\} \cup \{10\}))) \cup (\{00\} \cdot (\{1\})^*)) \cdot ((\{0\} \cdot \{0\}) \cdot \{1\})$$

is represented as

$$(01 + 1(0 + 10) + 00(1)^*)001$$

Regular expressions

- ▶ omit braces around one-string sets
- ▶ use $+$ instead of $\cup$
- ▶ omit $\cdot$
- ▶ follow the precedence: Kleene star operator $*$, $\cdot$ (implicitly), and $+$.

*Remark: $+$ and $\cdot$ are associative, i.e., $(A + B) + C = A + (B + C)$ and $(A \cdot B) \cdot C = A \cdot (B \cdot C)$.*

0010

เลขที่

$$0 + 1(0+1)^* \, \underline{0}$$

$$(0+1)^*$$

$$00^* \qquad 000^* \quad 00 \, ; \, 000 \qquad\qquad (00)^*$$

0, 00, 0w

All strings over $\{0, 1\}$ except $\underline{010}$.

$$\Sigma^{1*} - \{010\} \qquad \underline{1010}$$

$$\varepsilon + \underline{1}(0+1)^* + \underline{0}\left(\varepsilon + \underline{0}(0+1)^*\right)$$
$$+ \underline{0.1}\left(\varepsilon + 1(0+1)^*\right)$$
$$+ 010\ \underline{(0+1)}(0+1)^*$$

- $\varepsilon$
- يبدأ $\underline{1}$
- تبدأ $\underline{00}$ ....
- تبدأ $\underline{01}$
- او $010 + ....$

$$(01 + 1)^* \dashrightarrow \varepsilon,\ 01,\ 1$$
$$011\quad 0101\quad 011111$$

$$(ab + c)^*$$

$$\varepsilon, ab, abab, ababc, abcab, ccc, abccabccab, ....$$

$$(10 + 0)^* \quad 0, \varepsilon, 00010, 001010010, ....$$

# Subexpressions

$$0\,1\,1\,0\,(1+0)\,(1+0\,1)^*$$

$\{0110\}$     $\{01\}\cdot\{10\}$     $\{0\}\cdot\{1\}\cdot\{1\}\cdot\{0\}$

$$\underbrace{\left(0110(1+0)\right)}_{\text{sub expression}}\cdot\underbrace{(1+01)}_{\text{subexpression}}{}^*$$

$$0110\cdot(1+0)$$

# Regex is everywhere

# Proofs about regular expressions - structural induction

R

(I.H.) [____] ถือว่า เป็น พวกๆ subexpr ของ R

Case 1 :  R = ∅   ✓

Case 2 :  R = s   ✓

Case 3 :  R = S + T      เมื่อ  S & T เป็น reg ex.

case 4 :  R = S · T

case 5 :  R = S*

## Lemma 2

*Every regular expression that does not use the symbol $\emptyset$ represents a non-empty language.*

**Proof.**

$\Sigma = \{0,1\}$

$\emptyset$

$0 + 1$

$0^*$

$\varepsilon$

*Every regular expression that does not use the symbol $\emptyset$ represents a non-empty language.*

$R$

$P(R)$

$Q(R)$

$\forall R \left[ P(R) \Rightarrow Q(R) \right]$

**Proof.**

→ Let $R$ be a regular expression that does not use the symbol $\emptyset$. We prove by (structural) induction that $R$ represents a non-empty language.

*Every regular expression that does not use the symbol $\emptyset$ represents a non-empty language.*

**Proof.**
Let $R$ be a regular expression that does not use the symbol $\emptyset$. We prove by (structural) induction that $R$ represents a non-empty language.

**Induction hypothesis:** Every subexpression of $R$ that does not use the symbol $\emptyset$ represents a non-empty language.

## Lemma 2

*Every regular expression that does not use the symbol ∅ represents a non-empty language.*

**Proof.**
Let $R$ be a regular expression that does not use the symbol ∅. We prove by (structural) induction that $R$ represents a non-empty language.

**Induction hypothesis:** Every subexpression of $R$ that does not use the symbol ∅ represents a non-empty language.

*Case 1:* $R = \emptyset$.    ✗ contradiction    ข้อแย้งกับ เป็นไปบริๆเพราะ ว่า $R$ ไม่มี $\emptyset$

### Lemma 2

*Every regular expression that does not use the symbol $\emptyset$ represents a non-empty language.*

**Proof.**

Let $R$ be a regular expression that does not use the symbol $\emptyset$. We prove by (structural) induction that $R$ represents a non-empty language.

**Induction hypothesis:** Every subexpression of $R$ that does not use the symbol $\emptyset$ represents a non-empty language.

*Case 1:* $R = \emptyset$.

*Case 2:* $R$ is a single string.

R บ่ง non-empty language ของ พวกเรา

## Lemma 2

*Every regular expression that does not use the symbol $\emptyset$ represents a non-empty language.*

**Proof.**

Let $R$ be a regular expression that does not use the symbol $\emptyset$. We prove by (structural) induction that $R$ represents a non-empty language.

**Induction hypothesis:** Every subexpression of $R$ that does not use the symbol $\emptyset$ represents a non-empty language.

*Case 1:* $R = \emptyset$. ✗

*Case 2:* $R$ is a single string.

$$L(R) \neq \emptyset$$

**Proof.** (cont.2/4)
*Case 3:* $R = S + T$ for some regular expressions $S$ and $T$.

เนื่องจากว่า $R$ ไม่มี symbol $\emptyset$ , $S$ และ $T$ ดี ไม่มี $\emptyset$ ด้วย

เพราะ $S$ & $T$ เป็น subexpression ของ $R$.

ดังนั้น $S$ represent non-empty language จาก I.H.

และ $T$ represent non-empty lang จาก I.H.

เมื่อเอา $R$ เป็น

ดังนั้น language ที่เป็น $R$ ดี represent non-empty lang ด้วย.

union กับ non-empty
lang

$$L(R) \neq \emptyset$$

$$L(S) \neq \emptyset$$
$$L(T) \neq \emptyset$$

**Proof.** (cont.3/4)
*Case 4:* $R = S \cdot T$ for some regular expressions $S$ and $T$.

$\Rightarrow$ S & T เป็น non-empty languages

$\underline{L_S}$ & $\underline{L_T}$

- เป็น language ที่แทนด้วย $S$ ว่า A.
- เป็น language ที่แทนด้วย $T$ ว่า B.

ต้องแสดง $R$ ที่แทน language $A \cdot B$ ——— $L(S) \cdot L(T)$

สมมติว่า A ไม่ empty ดี $x \in A$
B ไม่ empty ดี $y \in B$

ต้องแสดง $x \cdot y \in A \cdot B$ (ภาษาเป็น $\emptyset$)

นั้นคือ language ที่แทนด้วย $R$ ดังไม่ empty

$\begin{cases} \swarrow x \text{ ใน } L(S) \\ y \text{ ใน } L(T) \end{cases}$

$x \cdot y \in L(S) \cdot L(T)$

$L(R) \neq \phi$

**Proof.** (cont.4/4)

*Case 5:* $R = S^*$ for some regular expression $S$.

แต่ว่า จาก $\varepsilon \in$ language ที่ออกแต่ว่า $\underline{S^*}$

$\quad R$ ออกแต่ language ที่ไม่ empty.

$$\varepsilon \in L(S^*)$$

ดังน

$$L(R) \neq \emptyset$$

**Proof.** (cont.4/4)

*Case 5:* $R = S^*$ for some regular expression $S$.

In every case, the language $L(R)$ is non-empty.

เวลาที่ เราเจอ reg ex
$R$ ต่างๆ
ให้ $L(R)$ เป็น
language ที่
represent โดย $R$

ex: $1$     $L(1) = \{1\}$

$1(0+1)$     $L(1(0+1))$

$= \{10, 11\}$

## Lemma 3

*Every non-empty regular language is represented by a regular expression that does not use the symbol $\emptyset$.*

$$\forall L [ \underbrace{Q(L) \Rightarrow P(L)}]$$

$$L \neq \phi$$

$$\hookrightarrow \exists R [L(R) = L \quad \text{b.o.:}$$
$$R \text{ 7ud symbol } \phi ]$$

$$\boxed{\neg Q(L) \vee P(L)}$$

### Lemma 3

*Every non-empty regular language is represented by a regular expression that does not use the symbol $\emptyset$.*

Let $R$ be a regular expression.

## Lemma 3

*Every non-empty regular language is represented by a regular expression that does not use the symbol $\emptyset$.* αυση, regular language $L$ ...

Let $R$ be a regular expression. We prove that if $L(R) \neq \emptyset$, then there exists a regular expression $R'$ such that $L(R) = L(R')$ and $R'$ does not contain $\emptyset$.

## Lemma 3

*Every non-empty regular language is represented by a regular expression that does not use the symbol $\emptyset$.*

Let $R$ be a regular expression. We prove that if $L(R) \neq \emptyset$, then there exists a regular expression $R'$ such that $L(R) = L(R')$ and $R'$ does not contain $\emptyset$.

We prove by induction. What should the induction hypothesis be?

สำหรับ ทุกๆ subexpression $S$ ของ $R$

ถ้า $L(S) \neq \emptyset$ จะมี reg ex $S'$ ที่ $S'$ ไม่มี symbol $\emptyset$

โดยที่ $L(S') = L(S)$

**I.H.:** For every subexpression $S$ of $R$, if $L(S) \neq \emptyset$, there exists an $\emptyset$-free regular expression $S'$ such that $L(S) = L(S')$.

**I.H.:** For every subexpression $S$ of $R$, if $L(S) \neq \emptyset$, there exists an $\emptyset$-free regular expression $S'$ such that $L(S) = L(S')$.

## What are the cases that we have to consider?

Case 1: $R = \phi$

Case 2: $R \neq \emptyset$ single sbn.

Case 3: $R = S + T$     สมมุติ regex $S$ & $T$

Case 4: $R = S \cdot T$     สมมุติ ของ regex $S \cdot T$

Case 5: $R = S^*$     สมมุติ ของ regex $S'$

▷ אם $L(R) \neq \phi$, אז יש $R'$ אז' $L(R) = L(R')$ כזאת: $R'$ ללא symbol $\phi$.

$$(\phi + 1) \cdot 0^*$$

$$1 \cdot 0^*$$

**I.H.:** For every subexpression $S$ of $R$, if $L(S) \neq \emptyset$, there exists an $\emptyset$-free regular expression $S'$ such that $L(S) = L(S')$.

---

**What are the cases that we have to consider?**

▶ $R = \emptyset$ , $L(R) = \phi$ ✗

▶ $R$ is a single string. ✓

▶ $R = \overline{S + T}$ for some regular expressions $S$ and $T$.

▶ $R = \check{S} \cdot \check{T}$ for some regular expressions $S$ and $T$.

⋯▶ $R = S^*$ for some regular expression $S$.

אם $L(S) = \phi$ , $L(T) \neq \phi$

ב'case 4 cases.

$\Rightarrow$
$L(S) \neq \phi$, $L(T) \neq \phi$
יש $S'$ בנ: $T'$ אז'ם כל $\phi$
ואז $R' = S' + T'$

case 1: $L(S) = \phi \rightarrow$ אז $\boxed{R' = \varepsilon}$

case 2: $\boxed{L(S) \neq \phi}$ יש $S'$ אז $L(S') = L(S)$ בנ:
$S'$ ללא $\phi$     $R' = (S')^*$

HELLO$^{®}$ = OLLEH

(E-ex1-6) For string $w$, the reversal $w^R$ is defined recursively as follows:

$\{\varepsilon, 0, 01, 100\}^{®}$

$$w^R = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ x^R \bullet a & \text{if } w = ax \text{ for some symbol } a \text{ and some string } x \end{cases}$$

$= \{\varepsilon, 0, 10, 001\}$

For a language $L$, the reversal of $L$ is defined as

$$L^R = \{w^R \mid w \in L\}.$$

You may assume the following facts.

- $L^* \bullet L^* = L^*$ for every language $L$.
- $(w^R)^R = w$ for every string $w$.
- $(x \bullet y)^R = y^R \bullet x^R$ for all strings $x$ and $y$.

อ้าง L เป็น regular language

$L^R$ เป็น regular — .

$(L^*)^R = (L^R)^*$

A

 สมมุว language $L$ ถูก represent ด้วย

vejular expression $R$ . จะพิสูจน์ว่า $L^R$ เป็น regular language

ฉันจะสร้าง reg. ex. $S$ ที่ $L(S) = L^R$

Prove that $(L^R)^* \subseteq (L^*)^R$.