



# Computational Learning Theory

## Regular Expression vs. Monomilas

---

Akihiro Yamamoto 山本 章博

<http://www.iip.ist.i.kyoto-u.ac.jp/member/akihiro/>  
[akihiro@i.kyoto-u.ac.jp](mailto:akihiro@i.kyoto-u.ac.jp)



# Contents

---

- What about a regular expressions?
- Learning in the Limit
- General Theory of Learning from Positive Data



# What About Regular Expressions?

---



# Regular Expressions (1)

---

- Regular expression was invented by S. Kleene, a mathematician, to represent sets in mathematics.
- Some interfaces of operating systems employ regular expressions in order to sets of files, etc.
  - Such an interface is sometimes called a “shell” in UNIX-orgined operating systems, e.g. Ubuntu.
  - Regular expressions can be used in the command window (prompt) of MS Window systems.

```
$ ls *.c
```

```
$ ls [abc]*.c
```



## Regular Expressions (2)

---

- Some commands in UNIX-originated operating systems also employ regular expressions in order to represent patterns of strings.
  - Examples of such editors are `ed`, `sed`, `vi`, `more`, ...
- Some programming languages based on manipulating characters and strings also employ regular expressions.
  - Examples of such languages are `awk`, `perl`, `python`,....

```
import re
regex = r'ab+'
text = "abbabbabaaabb"
pattern = re.compile(regex)
matchObj = pattern.match(text)
```



## Regular Expressions (3)

---

- From the history of their usage, so many variations and modifications are invented and introduced into particular commands or languages.
- The simplest regular expressions are constructed of characters **a**, **b**, **c** and operations **|**, **\***, where a string  $w$  of characters represents the set  $\{w\}$  and
  - $(R / S)$  represents the union  $R \cup S$
  - $(RS)$  represents the set of catenations  
$$\{wv \mid w \in R \text{ and } v \in S \}$$
  - $(R^*)$  represents the set of Kleene Closure of  $R$



# Kleene Closure

---

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

$$L^n = L L^{n-1} = \{ uv \mid u \in L, v \in L^{n-1} \} \quad (n \geq 1)$$

$$L^* = \{\varepsilon\} \cup L \cup L^2 \cup L^3 \cup \dots = \bigcup_{n=1}^{\infty} L^n$$

- Sometimes the set  $L^* - \{\varepsilon\}$  is denoted by  $L^+$ .

## Example

$$L = \{aa, ab\}$$

$$L^2 = \{aaaa, aaab, abaa, abab\}$$

$$L^3 = \{aaaaaa, aaaaab, aaabaa, aaabab, abaaaa, \dots\}$$

...

$$L^* = \{\varepsilon, aa, ab, aaaa, aaab, abaa, abab, aaaaaa, \dots\}$$



# Examples

---

Regular Expression  $R$

Set of strings  $L(R)$

$aababb$

$\{aababb\}$

$(aab)|(abb)$

$\{aab, abb\}$

$a(ab)^* b$

$\{w \mid w = aub \text{ and } u \in \{ab\}^*\}$   
 $= \{ab, aabb, aababb, \dots\}$

$a(a \mid b)^* b$

$\{w \mid w = aub \text{ and } u \in \{a, b\}^*\}$   
 $= \{ab, aab, abb, aaab, aabb, \dots\}$

$a(aa \mid bb)^* b$

$\{w \mid w = aub \text{ and } u \in \{\{aa\} \cup \{bb\}\}^*\}$   
 $= \{aaab, aabb, aaaaab, aaaabb, aabaab, aababb, aaaaaaab, \dots\}$





# Defining languages with patterns

---

- A language defined with a pattern  $\pi$  is  $\{\sigma \mid \sigma = \pi\theta \text{ for some non-empty grounding substitution } \theta\}$   
The language is denoted by  $L(\pi)$ .

## Example

$$L(axb) = \{a**ab**, a**bb**, a**aab**, a**abb**, a**abab**, a**abbb**, \dots\}$$

$$L(ayb) = \{a**ab**, a**bb**, a**aab**, a**abb**, a**abab**, a**abbb**, \dots\}$$

$$L(bxaxb) = \{b**aaab**, b**babbb**, \\ b**aaaaab**, b**abaabb**, b**baabab** b**bbabbb**, \\ b**aaaaaaab**, \dots\}$$

$$L(bxayb) = \{b**aaab**, b**aabbb**, b**aaaaab**, b**aaabb**, b**aabab**, \dots \\ b**baab**, b**babbb**, b**baaab**, b**baabb**, b**babab**, \dots \\ b**aaaaab**, b**aaaabb**, b**aaaaaab**, b**aaaabb**, \dots \\ b**baaab**, b**baabb**, b**baaaaab**, b**baaabb**, \dots\}$$



# RE vs. Monomials in Learning

---

- While both regular expressions and monomials represents data set of strings, they are different when we treat them in machine learning.
- Assume the case that an *unknown (hidden) representation  $R$*  is learned from training examples *in the limit*.
  - If we adopt a **regular expression** to represent  $R$ , we **cannot** learn  $R$  only from **only positive** examples, i.e. unsupervised learning.
  - If we adopt a **monomial** to represent  $R$ , we **can** learn  $R$  only from **only positive** examples.



# Learning in the Limit

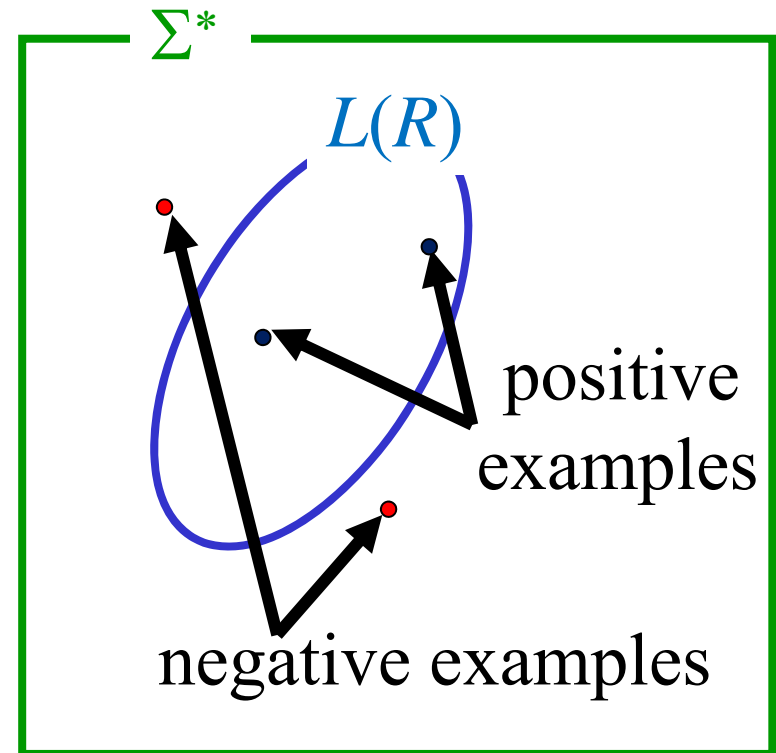
---

# Examples on $L(R)$

- We assume that, for an **unknown** rule  $R_*$ ,  
 $C_*$  is a finite set of positive examples **on  $L(R_*)$**  and  
 $D_*$  is a finite set of negative examples **on  $L(R_*)$** .

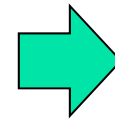
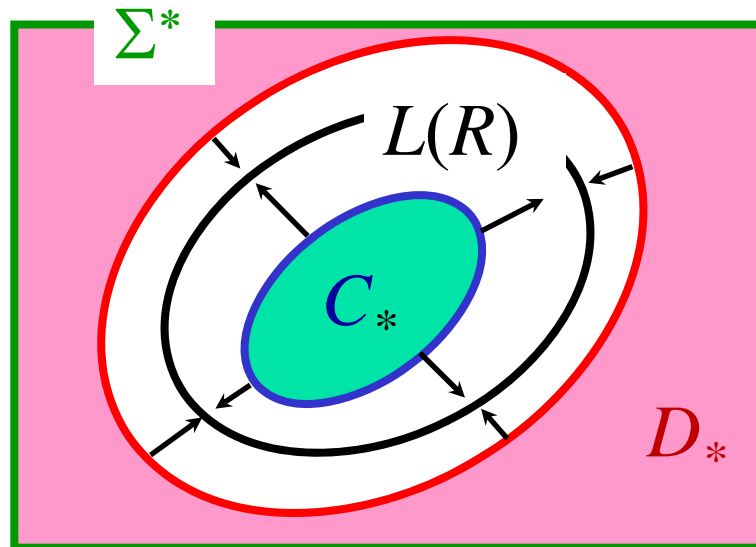
$L(R)$  : the set represented by  
the representation  $R$

- a positive example **on  $L(R)$**  :  
 $\langle x, + \rangle$  for  $x \in L(M)$
- a negative example **on  $L(R)$**  :  
 $\langle x, - \rangle$  for  $x \in \overline{L(M)}$



# Question

- If we give more and more (negative and positive) examples on  $L(R_*)$  to an learning algorithm, does it eventually conjecture the unknown  $R_*$ ?
- We have to give mathematical definitions of
  - giving **more and more** examples, and
    - or giving examples **many enough**
  - conjecturing  $M$  **eventually**.



$$R^{\wedge} \rightarrow R$$



# Assumption

---

- Without loss of generality, we may assume that learning algorithm takes examples in  $C_*$  and  $D_*$  **one by one**.
- In the situation that both  $C_i$  and  $D_i$  grow, we assume that an infinite **sequence**  $\sigma$  of strings marked with either + or -, and some **truncation** of  $\sigma$  corresponds to  $C_i$  and  $D_i$ .

## Example

$\sigma: \langle ab, + \rangle, \langle aab, + \rangle, \langle bbb, - \rangle, \langle aaab, + \rangle, \langle abba, - \rangle, | \dots$

$$C_i = \{ab, aab, aaab\},$$

$$D_i = \{bbb, abba\}.$$



# Presentations

---

**Definition** A **presentation** of  $L(R)$  is a **infinite** sequence

$$\sigma: \langle s_0, p_0 \rangle, \langle s_1, p_1 \rangle, \langle s_2, p_2 \rangle, \dots$$

where  $s_i \in \Sigma^*$  and  $p_i = +$  or  $-$ .

- $\langle s, + \rangle$  is a positive example
- $\langle s, - \rangle$  is a negative example
- $\sigma[n] = \langle s_0, p_0 \rangle, \langle s_1, p_1 \rangle, \langle s_2, p_2 \rangle, \dots, \langle s_{n-1}, p_{n-1} \rangle$

**Definition** A presentation  $\sigma$  is **complete** if

any  $x \in L(R)$  appears in  $\sigma$  as a positive example  $\langle x, + \rangle$   
at least once and

any  $x \in \overline{L(R)}$  appears in  $\sigma$  as a negative example  $\langle x, - \rangle$   
at least once.

# Identification in the limit [Gold]



$x_1, x_2, x_3, \dots$



$R_1, R_2, R_3, \dots$

- A learning algorithm  $A$  **EX-identifies**  $L(R)$  **in the limit from complete presentations** if for any complete presentation  $\sigma = x_1, x_2, x_3, \dots$  of  $L(R)$  and the output sequence  $R_1, R_2, R_3, \dots$  of  $A$ , there exists  $N$  such that for all  $n \geq N$   $R_n = R'$  and  $L(R') = L(R)$
- A learning algorithm  $A$  **BC-identifies**  $L(R)$  **in the limit from complete presentations** if for any complete presentation  $\sigma = x_1, x_2, x_3, \dots$  of  $L(R)$  and the output sequence  $R_1, R_2, R_3, \dots$  of  $A$ , there exists  $N$  such that for all  $n \geq N$   ~~$R_n = R'$~~  and  $L(R_n) = L(R)$





# A Well-known Result on RE

---

**Theorem** For every set  $L(R)$  represented by a regular expression  $R$ , there exists a unique minimal expression  $R'$  such that  $L(R)=L(R')$ .



## Embedding the Modified Generate-and-Test Algorithm into the Framework

---

Assume a procedure of enumerating all RE so that the enumeration  $R_0, R_1, R_2, \dots, R_i, \dots$  satisfies

$$|R_0| \leq |R_1| \leq |R_2| \leq \dots \leq |R_i| \leq \dots$$

Input  $\sigma = x_1, x_2, \dots$ : presentation (an infinite sequence)

Initialize  $k = 0$  /\*  $R_0$  is the simplest RE \*/

**for**  $N = 1, 2, \dots$

$\sigma[N] = x_1, x_2, \dots, x_N$

**forever**

let  $k' = k$

**for**  $n = 1, 2, \dots, N$ ,

**if** ( $x_n \in C$  and  $x_n \notin L(R_{k'})$ ) or ( $x_n \in D$  and  $x_n \in L(R_{k'})$ )

replace  $k$  with  $k + 1$

**if**  $k' = k$

terminate and output  $R_k$



# On the Generate-and-Test Algorithm

---

**Theorem** For any regular expression  $R_*$ ,

the modified generate-and-test algorithm EX-identifies  $L(R_*)$  in the limit from complete presentations.

**Proof** Let  $\sigma$  be an any complete presentation on  $L(R_*)$ .

Let  $R_N$  be the output of the algorithm for the input  $\sigma[N]$ .

If  $L(R_*) \neq L(R_N)$ , then there must be a string  $x \in \Sigma^*$

$(x \in L(R_*) \text{ and } x \notin L(R_N)) \text{ or } (x \notin L(R_*) \text{ and } x \in L(R_N)).$

Since  $\sigma$  is complete,  $x$  must be appears in the sequence with the sign  $+$  if  $x \in L(R)$  or otherwise with  $-$ .

This means that  $R_N$  must be replaced with another expression, at latest, when  $x$  appears in  $\sigma$ .

Once the algorithm outputs  $R_N$  s.t.  $L(R_*) = L(R_N)$ , it never changes the output afterwards.



## Revised version of *learn-patterns*

- ~~Fix an effective enumeration of patterns on  $\Sigma \cup X$ :~~

~~$\pi_1, \pi_2, \dots,$~~

$k = 1, \pi = \pi_1$

for  $n = 1$  forever

receive  $e_n = \langle s_n, b_n \rangle$

while (  $0 \leq \exists j \leq n$

$(e_j = \langle s_j, + \rangle$  and  $s_j \notin L(\pi)$ ) ~~and~~

~~$(e_j = \langle s_j, - \rangle$  and  $s_j \in L(\pi)$ )~~

$\pi = \pi'$  for an appropriate  $\pi'$ ;  $k ++$

output  $\pi$

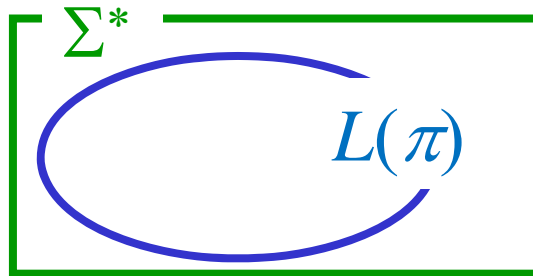
# Positive Presentations



$e_1, e_2, e_3, \dots$



$\pi_1, \pi_2, \pi_3, \dots$



- A **presentation** of  $L(\pi)$  is a **infinite** sequence consisting of positive and negative example.
- A presentation  $\sigma$  is **positive** if  $\sigma$  consists only of positive example  $\langle s, + \rangle$  and any positive example occurs at least once in  $\sigma$ .

# Identification in the limit [Gold]



$s_1, s_2, s_3, \dots$



$\pi_1, \pi_2, \pi_3, \dots$

- A learning algorithm  $A$  **EX-identifies**  $L(\pi)$  **in the limit from positive presentations** if for any positive presentation  $\sigma = s_1, s_2, s_3, \dots$  of  $L(g)$  and the output sequence  $\pi_1, \pi_2, \pi_3, \dots$  of  $A$ , there exists  $N$  such that for all  $n > N$   $\pi_n = \pi'$  and  $L(\pi') = L(\pi)$
- A learning algorithm  $A$  **BC-identifies**  $L(\pi)$  **in the limit from positive presentations** if for any positive presentation  $\sigma = s_1, s_2, s_3, \dots$  of  $L(g)$  and the output sequence  $\pi_1, \pi_2, \pi_3, \dots$  of  $A$ , there exists  $N$  such that for all  $n > N$   ~~$\pi_n = \pi'$~~  and  $L(\pi_n) = L(\pi)$



# Identification in the limit [Gold]

---

- A learning algorithm  $A$  **EX-identifies** a class  $C$  of languages **in the limit from positive presentations** if  $A$  EX-identifies every language in  $C$  in the limit from positive presentations.
- A learning algorithm  $A$  **BC-identifies** a class  $C$  of languages **in the limit from positive presentations** if  $A$  BC-identifies every language in  $C$  in the limit from positive presentations.

# Anti-Unification of Strings

- For a set  $C$  of strings of same length

$$s_1 = c_{11} c_{12} \dots c_{1i} \dots c_{1k}$$

$$s_2 = c_{21} c_{22} \dots c_{2i} \dots c_{2k}$$

...

$$s_n = c_{n1} c_{n2} \dots c_{nj} \dots c_{nk}$$

the anti-unification of  $C$  is a pattern

$$\pi = \gamma(c_{11}c_{21}\dots c_{n1})\gamma(c_{12}c_{22}\dots c_{n2}) \dots \gamma(c_{1k}c_{2k}\dots c_{nk})$$

where

$$\gamma(c_1c_2\dots c_n) = \begin{cases} c & \text{if } c_1 = c_2 = \dots = c_n = c \\ x_{i(c_1c_2\dots c_n)} & \text{otherwise.} \end{cases}$$

and  $i(c_1c_2\dots c_n)$  is the “index” of  $c_1c_2\dots c_n$ .





# Identification of patterns

---

**Theorem** The revised algorithm of *Learn-pattern* with computing an anti-unification EX-identifies the class of all pattern languages in the limit from positive presentations.



# A Negative Result

---

**Theorem** [Gold] There is no learning algorithm which identifies any regular expression from positive data.

# A Negative Result (2)

- We construct a positive presentation  $\sigma$  of  $L((ab)^*)$  in the following manner.
- Let  $e_1$  be a string in  $L$ . Since the regular expression  $e_1$  is also in  $\mathbf{C}$  and  $A$  must identify  $\{e_1\}$ . So the first  $N_1$  examples of  $\sigma$  are all  $e_1$ , until “ $A$  identifies the regular expression  $e_1$ .”

$$\exists N_1 \forall n > N_1 h_n = e_1$$



$e_1, e_1, \dots, e_2, \dots$   
 $\underbrace{\hspace{10em}}_{N_1+1}$

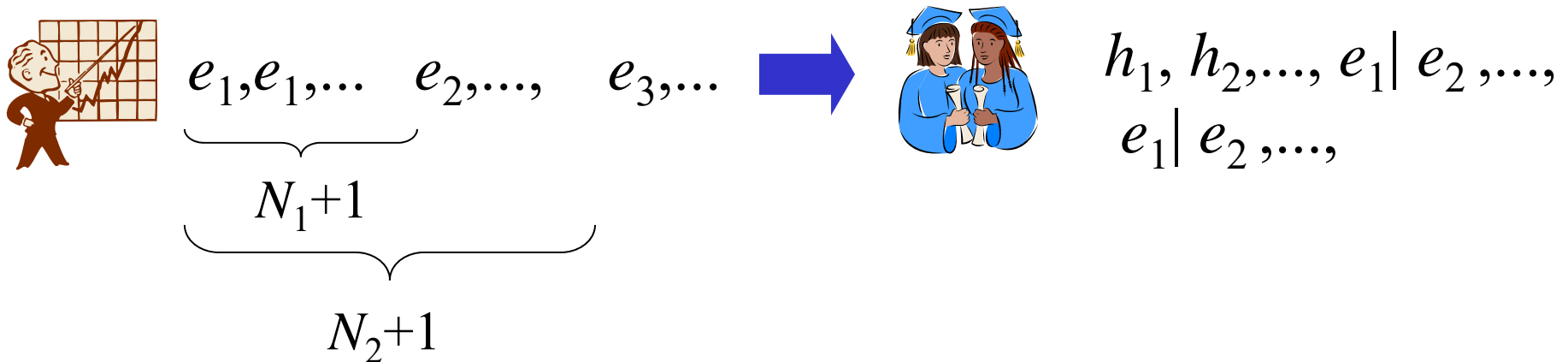


$h_1, h_2, h_3, \dots, e_1, e_1, \dots$

# A Negative Result (3)

- Let the  $(N_1+1)$ -th example be  $e_2$  which is different from  $e_1$ .
- Since  $\mathbf{C}$  contains  $e_1|e_2$ , the learning algorithm  $A$  identifies  $e_1|e_2$  in the limit.

$$\exists N_1 \forall n > N_2 > N_1 R_n = e_1|e_2$$





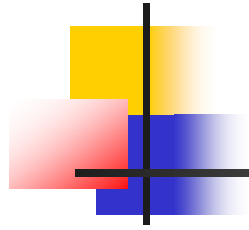
## A Negative Result (4)

---

- Let the  $(N_2+1)$ -th example be  $e_3$  which is different from both of  $e_1$  or  $e_2$ .
- Since  $C$  contains  $e_1 | e_2 | e_3$ ,  $A$  identifies  $e_1 | e_2 | e_3$  in the limit.

$$\exists N_3 \forall n > N_3 > N_2 > N_1 h_n = e_1 | e_2 | e_3$$

- The language  $L = \{e_1, e_2, e_3, e_4, \dots\}$  is a infinite and  $A$  cannot identify  $L$ .



# General Theory of Learning from Positive Data



# GCD and Learning

---

A class of languages in  $\mathbf{N}$  :

$$L(\mathbf{N}) = \{L(m) \mid m \in \mathbf{N}\}$$

$$L(m) = \{0\underbrace{1\dots 1}_n 0 \mid n \bmod m = 0\}$$

$$L(m) = \{n \in \mathbf{N} \mid n \bmod m = 0\}$$

A class of languages in  $\mathbf{Z}$  :

$$L(\mathbf{N}) = \{L(m) \mid m \in \mathbf{N}\}$$

$$L(m) = \{\underbrace{1\dots 1}_n \mid n \bmod m = 0\} \cup \{0\underbrace{1\dots 1}_n \mid n \bmod m = 0\}$$

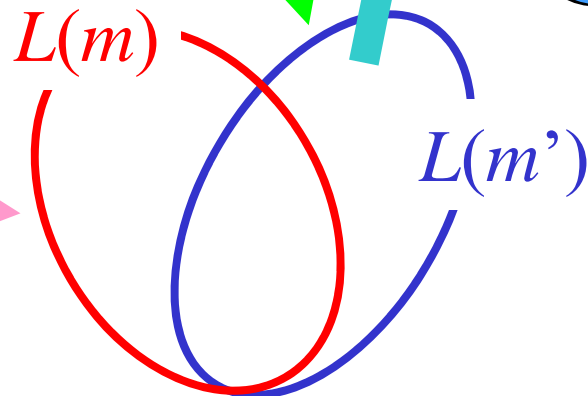
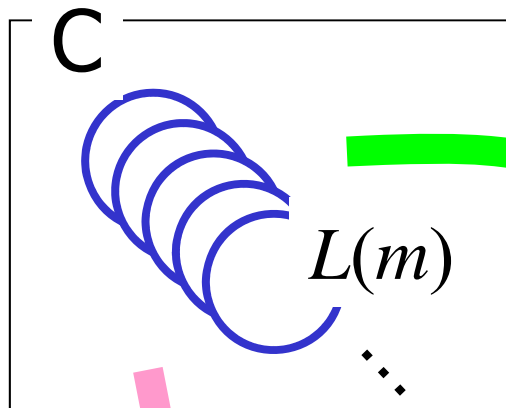
$$L(m) = \{n \in \mathbf{Z} \mid |n| \bmod m = 0\}$$

# GCD and Learning



- $L(\mathbf{N}) = \{L(m) \mid m \in \mathbf{N}\}$   
 $L(m) = \{01\dots10 \mid n \bmod m = 0\}$

Positive presentation  
72, 48, 60, ..., 12, ...



Compute the GCD  
of  $s_1, s_2, \dots, s_k$   
with Euclidean  
Algorithm

Conjecture  
72, 24, 12, ..., 12, ...





# Proving that $L(\mathbf{N})$ is identifiable

---

- For every  $n \in \mathbf{N}$ , the characteristic set of  $L(m)$  in  $L(\mathbf{N})$  is  $\{m\}$ , that is,  $\{m\} \subseteq L(m')$  implies  $L(m) \subseteq L(m')$ .

- To see this, assume that  $\{m\} \subseteq L(m')$ .

This is equivalent to  $m \in L(m')$  and from the definition of  $L(m')$ ,  $m = k' m'$  for some  $k' \in \mathbf{N}(\mathbf{Z})$ .

- $L(m) = \{n \in \mathbf{N} \mid n \bmod m = 0\}$  (  $\{n \in \mathbf{Z} \mid |n| \bmod m = 0\}$  ).

Let  $n$  be any element in  $L(m)$ . Then, from the definition, there exists  $k \in \mathbf{N}(\mathbf{Z})$  such that  $n = k m$ . For the  $k'$  and  $k$ , it holds that  $n = k k' m'$ . This means  $n \in L(m')$ , and therefore  $L(m) \subseteq L(m')$ .

# Analysis of Patterns (1)

Example  $\pi = axxbbyaa$

$L(axxbbyaa)$

$=\{aaabbaaa, aaabbbaa, abbbbaaa, abbbbaa, aaaaabbaaa, aaaaabbaa, aababbaaa, aababbbbaa, \dots, abaabaabbbbababaa, \dots\}$

- Using examples as long as  $\pi$  :

$aaabbaaa, aaabbbaa, abbbbaaa, abbbbaa$

$\theta_1 = \{(x,a), (y,a)\}$   $\theta_2 = \{(x,a), (y,b)\}$   $\theta_3 = \{(x,b), (y,a)\}$   $\theta_4 = \{(x,a), (y,b)\}$

We can know that the 2nd, 3rd, and 6th positions must be variables.

The variable at the 6th position is different from those at the 2nd and 3rd.



## Analysis of Patterns (2)

---

- Any language  $L(\pi')$  containing the four strings must be a superset of  $L(\pi)$ .

aaabbaaa, aaabbbbaa, abbbbaaa, abbbbaaa

$\theta_1 = \{(x,a), (y,a)\}$   $\theta_2 = \{(x,a), (y,b)\}$   $\theta_3 = \{(x,b), (y,a)\}$   $\theta_4 = \{(x,a), (y,b)\}$

- If  $\pi'$  and  $\pi$  are of same length,  $\pi'$  has more variables than  $\pi$ .
- If  $\pi'$  is shorter than  $\pi$ ,  $\pi'$  has at least one variable with which some substring of  $\pi$  longer than 2 must be replaced.



# Characteristic Set of $L(\pi)$

---

- Let  $\pi$  be a pattern which contains variables  $x_1, x_2, \dots, x_n$ .

Consider the following substitutions:

$$\theta_a = \{(x_1, \mathbf{a}), (x_2, \mathbf{a}), \dots, (x_n, \mathbf{a})\},$$

$$\theta_b = \{(x_1, \mathbf{b}), (x_2, \mathbf{b}), \dots, (x_n, \mathbf{b})\},$$

$$\sigma_1 = \{(x_1, \mathbf{a}), (x_2, \mathbf{b}), \dots, (x_n, \mathbf{b})\},$$

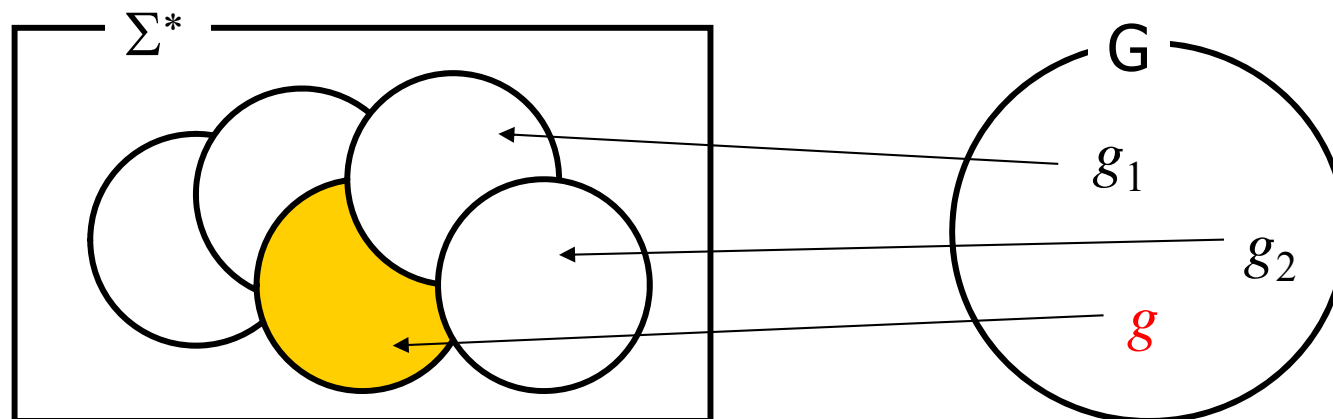
...

$$\sigma_n = \{(x_1, \mathbf{b}), (x_2, \mathbf{b}), \dots, (x_n, \mathbf{a})\}$$

- The set  $\{p\theta_a, p\theta_b, p\sigma_1, \dots, p\sigma_n\}$  is a characteristic set of  $L(\pi)$ .

# A General Framework of Learning

- A class of formal languages  $L(\mathbf{G})$  indexed with  $\mathbf{G}$
  - $\mathbf{G}$ : A set of expressions such that each expression in  $\mathbf{G}$  represents one language in  $L(\mathbf{G})$ , and every language in  $L(\mathbf{G})$  is represented by at least one expression in  $\mathbf{G}$ .
    - We assume that There is an algorithm which determines whether or not  $w \in L(g)$  for every string  $w \in \Sigma^*$  and  $g$ .
- Examples of  $\mathbf{G}$  : a set of finite state automata, a set of CFGs, a set of patterns,...





## C2: The Characteristic Set Property

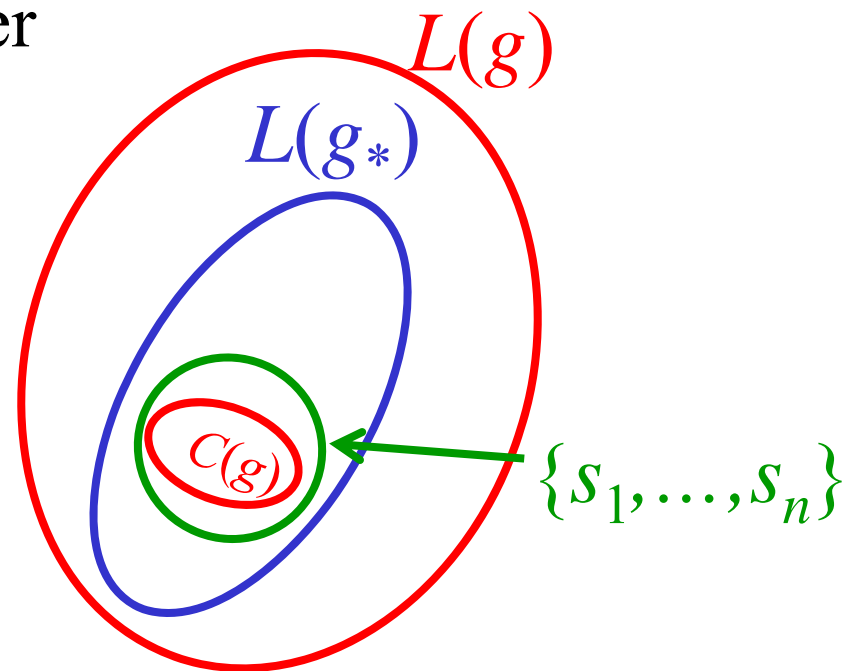
- A subset  $C(g)$  of a language of  $L(g)$  is a **characteristic set** of  $L(g)$  in  $L(\mathbf{G})$  if
  - (1)  $C(g)$  is a finite set and
  - (2) for every  $L(g') \in L(\mathbf{G})$   $C(g) \subseteq L(g')$  implies
$$L(g) \subseteq L(g')$$

**Theorem [Kobayashi]** A class  $L(\mathbf{G})$  of languages is identifiable in the limit from positive presentation **if** every language  $L(g)$  in  $L(\mathbf{G})$  has a characteristic set  $C(g)$  in  $L(\mathbf{G})$ .

# Which grammar should be chosen?

- Choose  $g$  such that  $C(g) \subseteq \{s_1, \dots, s_n\}$ 
  - The examples are from  $L(g_*)$ , that is,  $\{s_1, \dots, s_n\} \subseteq L(g_*)$ .  
and therefore  $C(g) \subseteq L(g_*)$ . From the definition of characteristic sets, this implies  $L(g) \subseteq L(g_*)$ .

So over generalization never happens.





# EC1: The Finite Tell-tale Property

- A subset  $T(g)$  of a language of  $L(g)$  is a **finite tell-tale** of  $L(g)$  in  $L(G)$  if
  - (1)  $T(g)$  is a finite set and
  - (2)  $T(g) \subseteq L(g') \subsetneq L(g)$  for **no**  $L(g') \in L(G)$  other than  $L(g)$

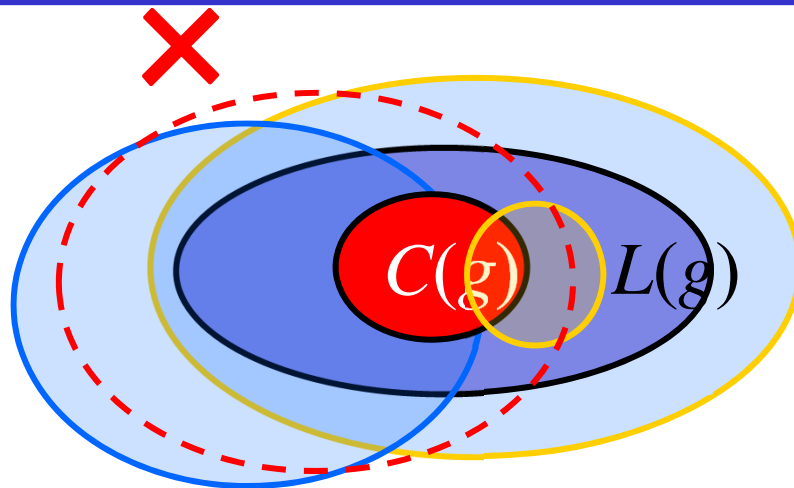
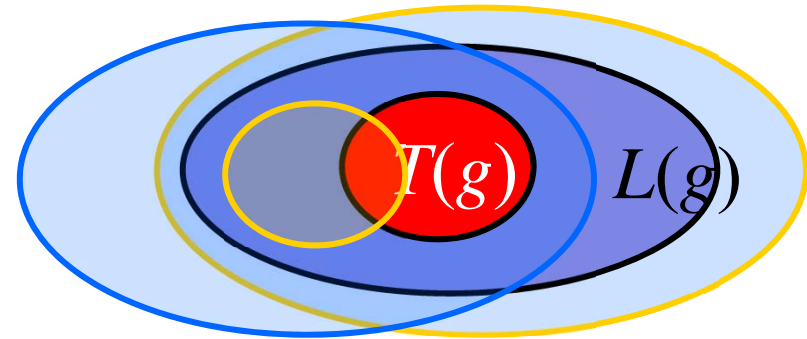
**Theorem [Angluin]** A class  $L(G)$  of languages is identifiable in the limit from positive presentation **if and only if** every language  $L(g)$  in  $L(G)$  has a finite tell-tail  $T(g)$  in  $L(G)$  and there is a procedure which generates elements of  $T(g)$  when the grammar  $g$  is given as an input.



# Tell-tales and Characteristic Sets

Finite Tell-tale  $T(g)$  of  $L(g)$ :

- $T(g) \subseteq L(g)$  ( $T$  is a finite set)
- For no  $L(g') \in \mathbf{L}(\mathbf{G})$  other than  $L(g)$ ,  $T(g) \subseteq L(g') \subseteq L(g)$



Characteristic set  $C(g)$  of  $L(g)$ :

- $T(g) \subseteq L(g)$  ( $T$  is a finite set)
- For every  $L(g') \in \mathbf{L}(\mathbf{G})$   
 $C(g) \subseteq L(g')$  implies  $L(g) \subseteq L(g')$



## Analysis of Patterns (3)

---

**Lemma 1** For every string  $s$ , there are only finite number of pattern languages containing  $s$ .

**Proof.** If  $s \in L(\pi)$ , then  $|s| \geq |\pi|$ .

**Example** The languages containing  $s = \mathbf{aab}$  are

$L(\mathbf{aab})$ ,

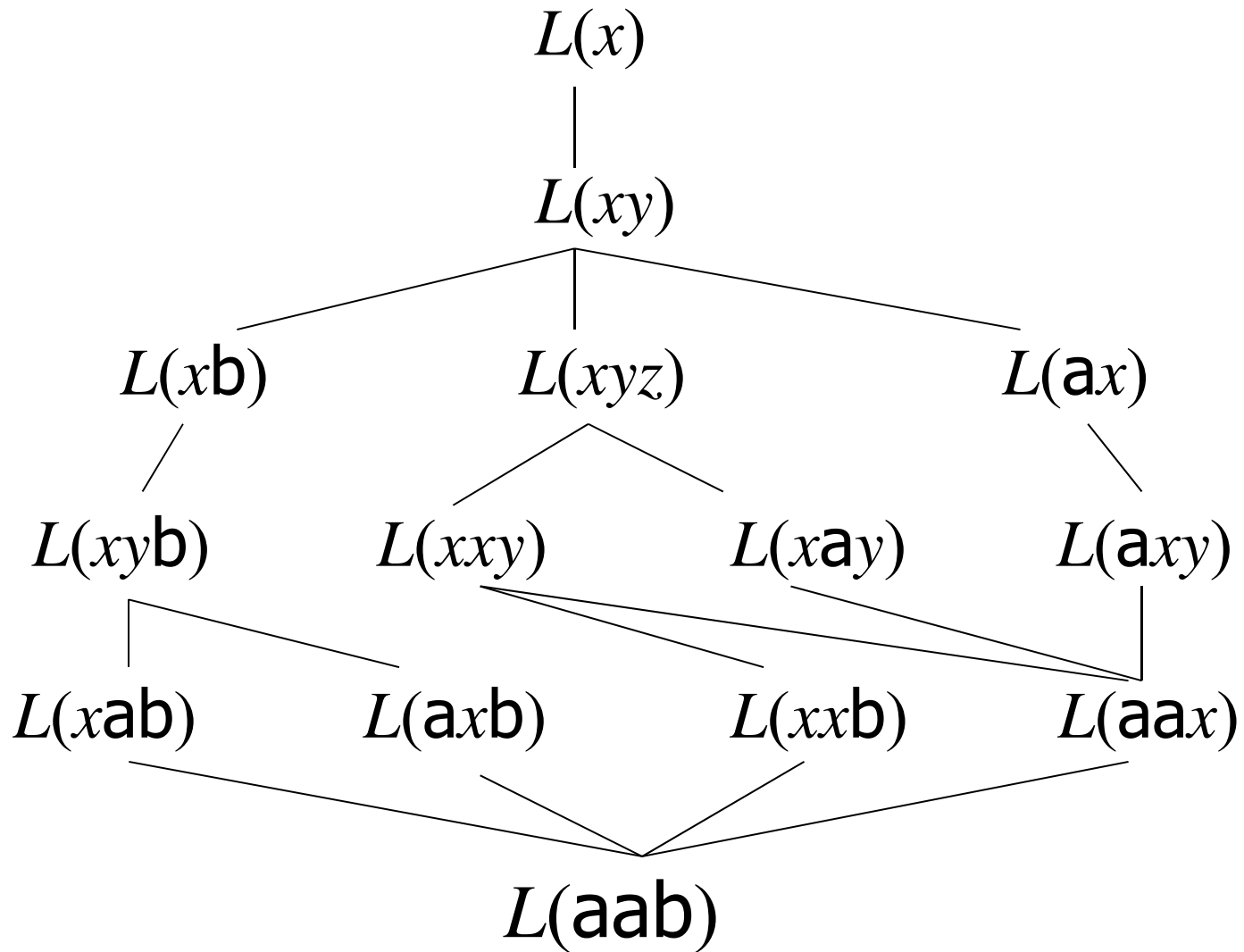
$L(\mathbf{xab})$ ,  $L(\mathbf{axb})$ ,  $L(\mathbf{aax})$ ,  $L(\mathbf{xxb})$ ,  $L(\mathbf{xb})$ ,  $L(\mathbf{ax})$ ,  $L(\mathbf{x})$ ,

$L(\mathbf{xyb})$ ,  $L(\mathbf{xay})$ ,  $L(\mathbf{axy})$ ,  $L(\mathbf{xxy})$ ,  $L(\mathbf{xy})$ ,

$L(\mathbf{xyz})$ ,



# Hasse Diagram





## C4: Finite thickness

---

- A class  $L(G)$  of languages has **the finite thickness** if for all  $w \in \Sigma^*$  there are only a finite number of languages in  $L(G)$  which contain  $w$ .

**Theorem [Angluin]** A class  $L(G)$  of languages is identifiable in the limit from positive presentation **if and only if**  $L(G)$  of languages has the finite thickness.



# $L(\mathbf{N})$ has the Finite Thickness

---

- From the finite thickness condition:

$L(\mathbf{N}) = \{L(m) \mid m \in \mathbf{N}\}$  has the finite thickness property.

- From the fact

$$\text{GCD}(e_1, e_2, \dots, e_k) \geq \text{GCD}(e_1, e_2, \dots, e_k, e_{k+1})$$

and the following property:

Let  $a_1, a_2, \dots, a_n, \dots$  be a infinite sequence of natural numbers satisfying that

$$a_n \geq a_{n+1} \text{ for all } n \geq 1.$$

Then there is  $N \geq 1$  such that  $a_n = a_{n+1}$  for all  $n \geq N$ .



## C3: Finite Elasticity

- A class  $L(\mathbf{G})$  of languages has the **infinite elasticity** if there is an infinite sequence of strings  $w_0, w_1, w_2, \dots$ , and an infinite sequence languages in  $L(\mathbf{G})$   $L(g_0), L(g_1), L(g_2)$  such that
$$\{w_0, w_1, \dots, w_{n-1}\} \subseteq L(g_n) \text{ and } w_n \notin L(g_n) \text{ for every } n \geq 1.$$
A class  $L(\mathbf{G})$  of languages has the **finite elasticity** if it does not have the infinite elasticity.

**Th. [Wright]** A class  $L(\mathbf{G})$  of languages is identifiable in the limit from positive presentation **if**  $L(\mathbf{G})$  has the finite elasticity.

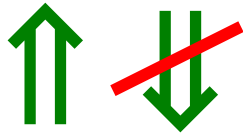


# Relation among the conditions

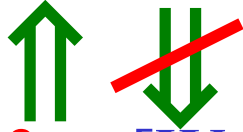
---

$U$  : a class of languages

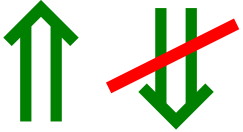
- EC1 (necessary and sufficient) [Angluin]



- C2: [Kobayashi]



- C3: [Wright]



- C4: [Angluin]



# Announcement

---

- The lectures on 28<sup>th</sup> November follow the Schedule for Monday.
- The next lecture of this course is on 5<sup>th</sup> December.