# Computational Learning Theory
## Extending Patterns with the Correctness of Learning Algorithms

Akihiro Yamamoto 山本 章博

http://www.iip.ist.i.kyoto-u.ac.jp/member/akihiro/
akihiro@i.kyoto-u.ac.jp

# Contents

- What about a pair of patterns?

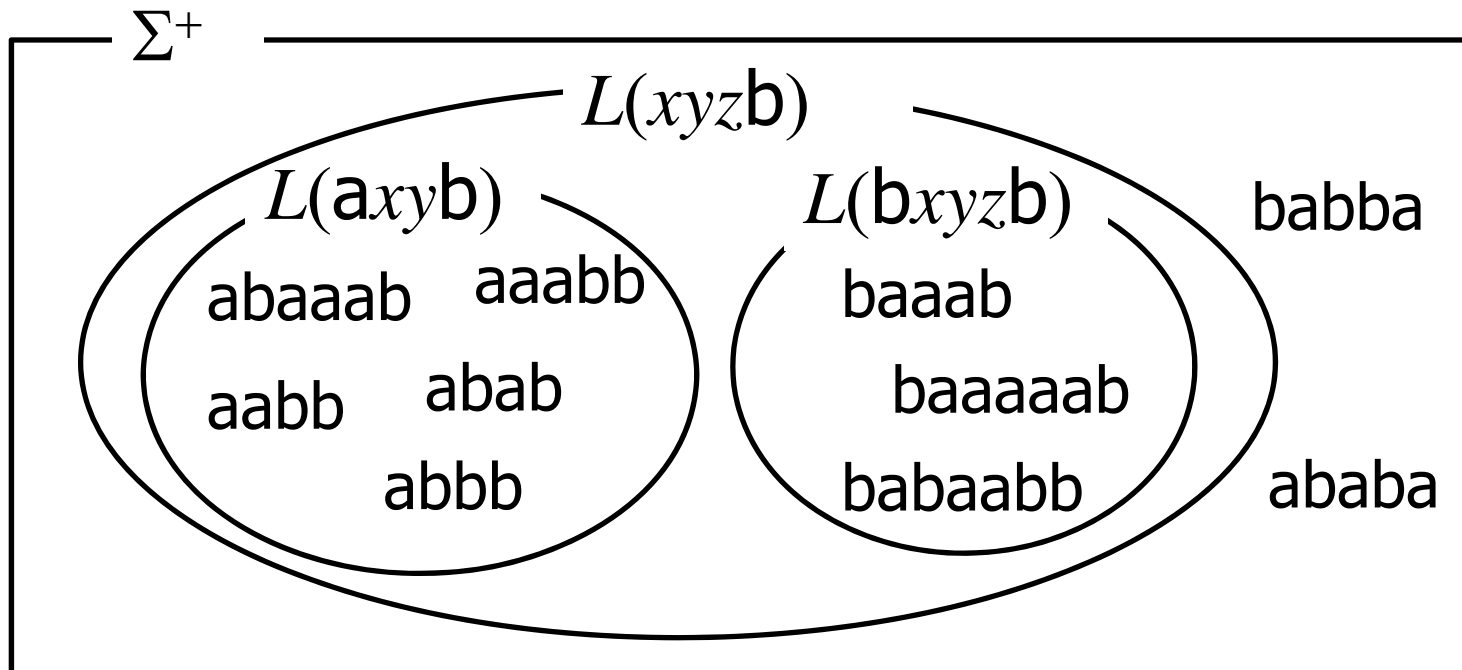- Correctness of Learning Algorithms

# What about a pair of patterns?

# Outputting a Pair of Patterns

$C = \{$babaabb, aaab, baaab, aabb, abab, baaaaab,
　　abbb, aaabb, baaab$\}$

abaaab, aabb, abab, abbb, aaabb $\in L(\mathrm{a}x\mathrm{b})$
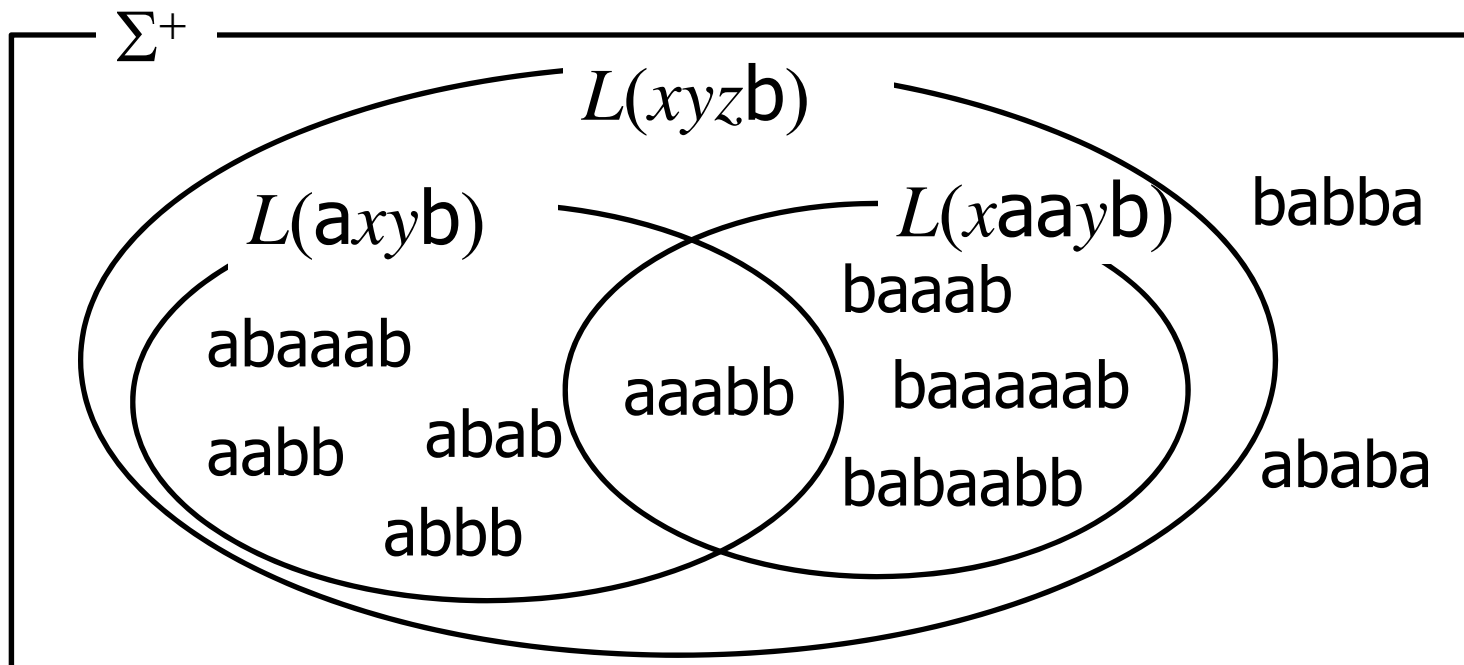
baaab, baaaaab, babaabb $\in L(\mathrm{b}xyz\mathrm{b})$
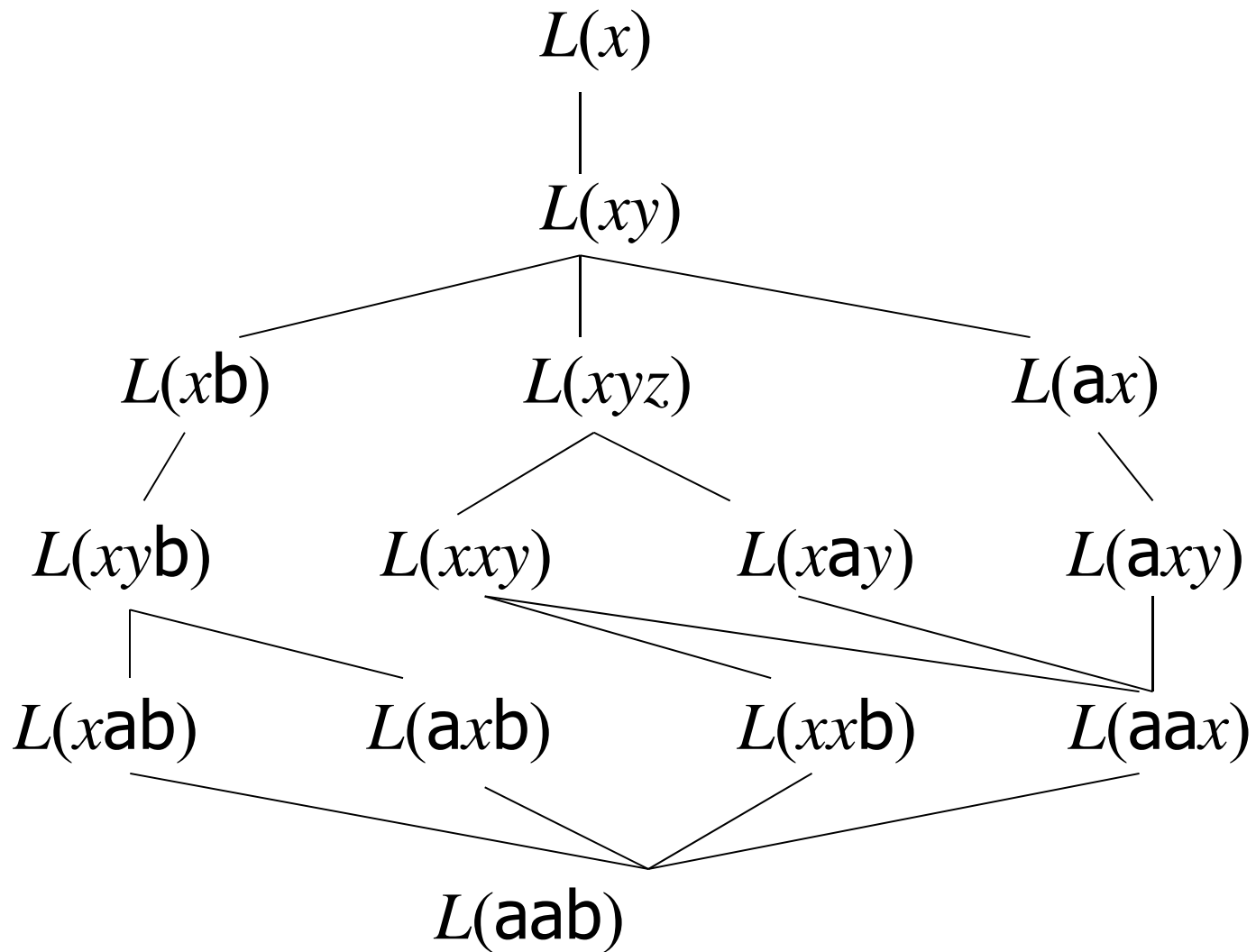
# Outputting a Pair of Patterns

$C = \{$babaabb, aaab, baaab, aabb, abab, baaaaab, abbb, aaabb, baaab$\}$

abaaab, aabb, abab, abbb, aaabb $\in L(\mathsf{a}xy\mathsf{b})$

baaab, baaaaab, aaabb, babaabb $\in L(x\mathsf{aa}y\mathsf{b})$

# Hasse Diagram

$L(x)$

$L(xy)$

$L(x\mathrm{b})$  $L(xyz)$  $L(\mathrm{a}x)$

$L(xy\mathrm{b})$  $L(xxy)$  $L(x\mathrm{a}y)$  $L(\mathrm{a}xy)$

$L(x\mathrm{ab})$  $L(\mathrm{a}x\mathrm{b})$  $L(xx\mathrm{b})$  $L(\mathrm{aa}x)$

$L(\mathrm{aab})$

# How to Construct H.D.

- A string $w$ is an instance of $\pi$ if $w = \pi\,\theta$ for some substitution $\theta$.

- A pattern $\tau$ is one step refinement of $\pi$ if $\tau = \pi\,\theta$ and $\theta$ is one of the three:

  - $\theta = \{(x, c)\}$  Replacing a variable with a symbol (character)
  - $\theta = \{(x, x_1 x_2)\}$ Spritting a variable with two new variables
  - $\theta = \{(x, z), (y, z)\}$  Unifying two variables into one

Example

$$x \rightarrow x_1 x_2 \rightarrow x_1 x_2\, x_3 \rightarrow x_1 x_1\, x_3 \rightarrow x_1 x_1 \mathsf{b}$$

# Substitution (1)

- A substitution is a set of pairs

$$\theta = \{ (x_1, \tau_1), (x_2, \tau_2), \ldots, (x_n, \tau_n) \}$$

  where $x_1, x_2, \ldots, x_n$ are distinct variables and

  $\pi_1, \pi_2, \ldots, \pi_n$ are patterns.

- Applying a substitution $\theta$ to a pattern $\pi$ is replacing every variable $x_i$ in $\pi$ with $\tau_i$ simultaneously.

  The result is denoted by $\pi\theta$.

Example

$\theta_1 = \{ (x, \text{bba}), (y, \text{ba}) \}$

$\theta_2 = \{ (x, \text{b}y\text{a}), (y, \text{a}y\text{b}) \}$

  $\text{b}x\text{a}x\text{b}\,\theta_1 = \text{bbbaabbab}, \text{b}x\text{a}x\text{b}\,\theta_2 = \text{bbyaabyab},$

  $\text{a}x\text{bb}y\text{a}\,\theta_1 = \text{abbabbbaa}, \text{a}x\text{bb}y\text{a}\,\theta_2 = \text{abyabbayba}$

# Anti-Unification

- A pattern $\pi$ is an anti-unifier of the set $C$ of positive examples if $C \subset L(\pi)$, i.e., for each positive example w, w=$\pi\theta$.

- A pattern $\pi$ is a least common anti-unifier of $C$ if $\pi$ is an anti-unifier of $C$ and no $\pi' \leq \pi$ satisfies $C \subset L(\pi')$.

Examples

The least common anti-uinifer of abaaab and aaabb is a$x$a$y$b.

The least common anti-uinifer of konnichiwa and konbanwa are kon$x$wa.

The least common anti-uinifers of konnichiwa and konbannwa are kon$x$wa and ko$x$n$y$wa.

# Notes

- The operation of "making one step refinement" can be regarded as "applying derivative"

# The learning algorithm *learn-patterns*

For $n = 1$ forever

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

    compute the list $l = \pi_1, \pi_2, \ldots, \pi_k$ of all the

    least common anti-unifications of the set of

    positives $C_n = \{s_j : \langle\, s_j, +\, \rangle$ and $j = 1, 2, \ldots, n\}$

    for each $\pi_j$ in the list $l$

        if an $e_j = \langle\, s_j, -\, \rangle$ s.t. $s_j \in L(\pi_j)$ is found

        delete $\pi_j$ from $l$

    if $l$ is not empty

        return one $\pi_i$

# In Theoretical Form

**Lemma 2** Let $\pi_1, \pi_2,\ldots,\pi_n$ be patterns. If the language $L(\pi_k)$ is minimal in $\{L(\pi_1), L(\pi_2),\ldots, L(\pi_n)\}$, then $\pi_k$ is one of the longest patterns in the list.

**Lemma 3** Let $\pi_1$ and $\pi_2$ be patterns of same length. Then $L(\pi_1) \subseteq L(\pi_2)$ if and only if $\pi_2\theta = \pi_1$.

**Note** If we do not assume $\pi_1$ and $\pi_2$ be patterns of same length, then it is not decidable whether or not $L(\pi_1) \subseteq L(\pi_2)$.

# Which pattern should be chosen?

- Let $C$ be a set of (positive) examples

1. Select all shortest examples.

2. Look for one of the minimal patterns between $x$ (a singleton variable) and the anti-unifier of the shortest examples, and return it.
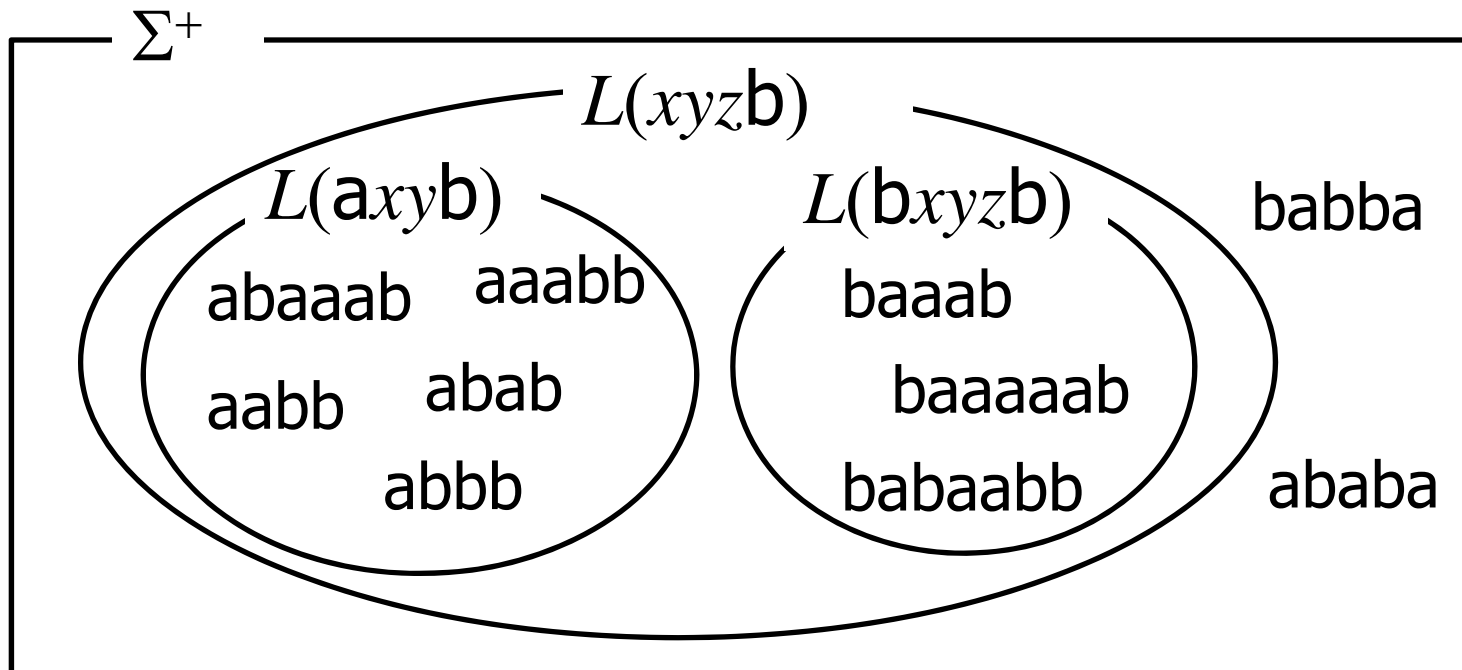
Note: If we only follow the identification-in-the-limit criterion, the second can be simplified as

2'. Return the anti-unifier of the shortest examples but this might not seem "learning".

# Outputting a Pair of Patterns

$C$ = {babaabb, aaab, baaab, aabb, abab, baaaaab,

abbb, aaabb, baaab}

abaaab, aabb, abab, abbb, aaabb $\in L(\text{a}x\text{b})$
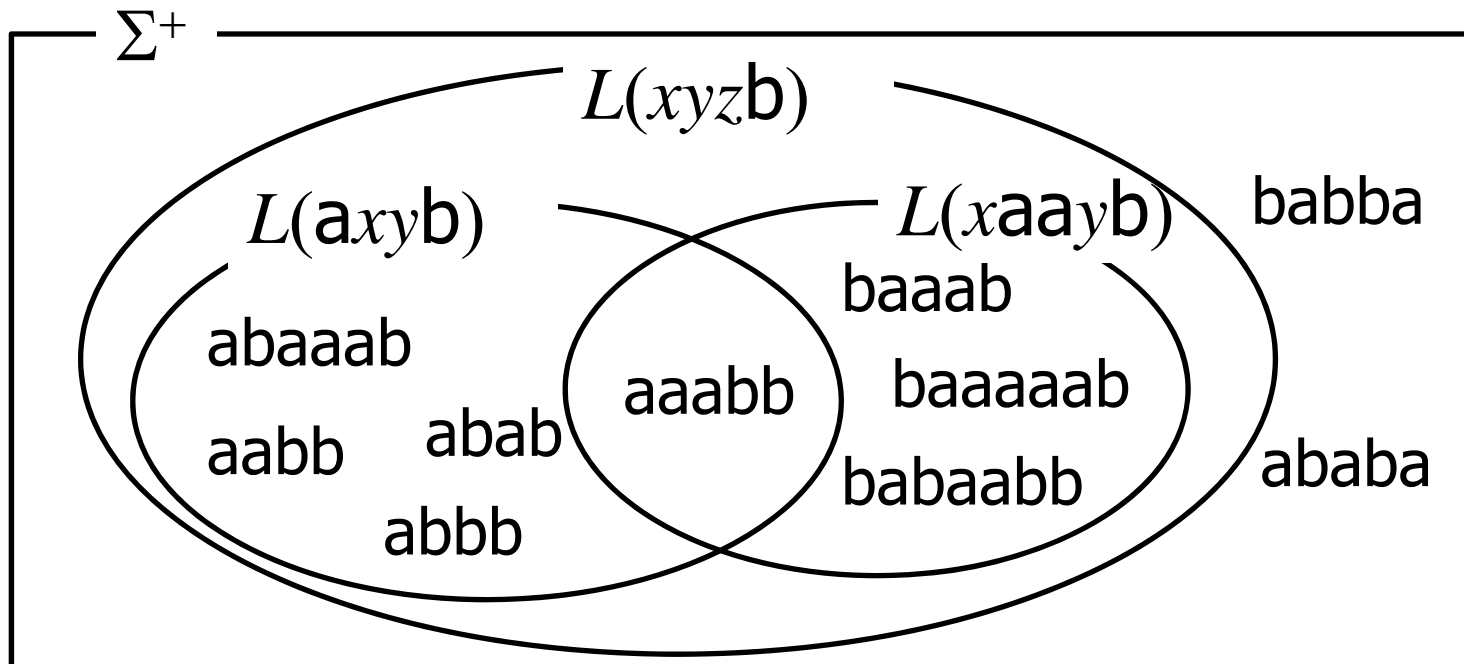
baaab, baaaaab, babaabb $\in L(\text{b}xyz\text{b})$

# Outputting a Pair of Patterns

$C = \{$babaabb, aaab, baaab, aabb, abab, baaaaab,

abbb, aaabb, baaab$\}$

abaaab, aabb, abab, abbb, aaabb $\in L(\mathsf{a}xy\mathsf{b})$

baaab, baaaaab, aaabb, babaabb $\in L(x\mathsf{aa}y\mathsf{b})$

# Downward Coverset Algorithm

Given a data set $C$

For each minimally common anti-unifier $\pi$ of $C$

    For each a pair $\pi_1$, $\pi_2$ of patterns just beneth $\pi$ in the

        Hasse Diagram

      if $L(\pi_2) \subset C - L(\pi_1)$ then

        make minimally common anti-unification $\pi_1'$

        of $C \cap L(\pi_1)$

          and

        minimally common anti-unification $\pi_2'$ of $C - L(\pi_1)$

# Correctness of Learning Algorithms
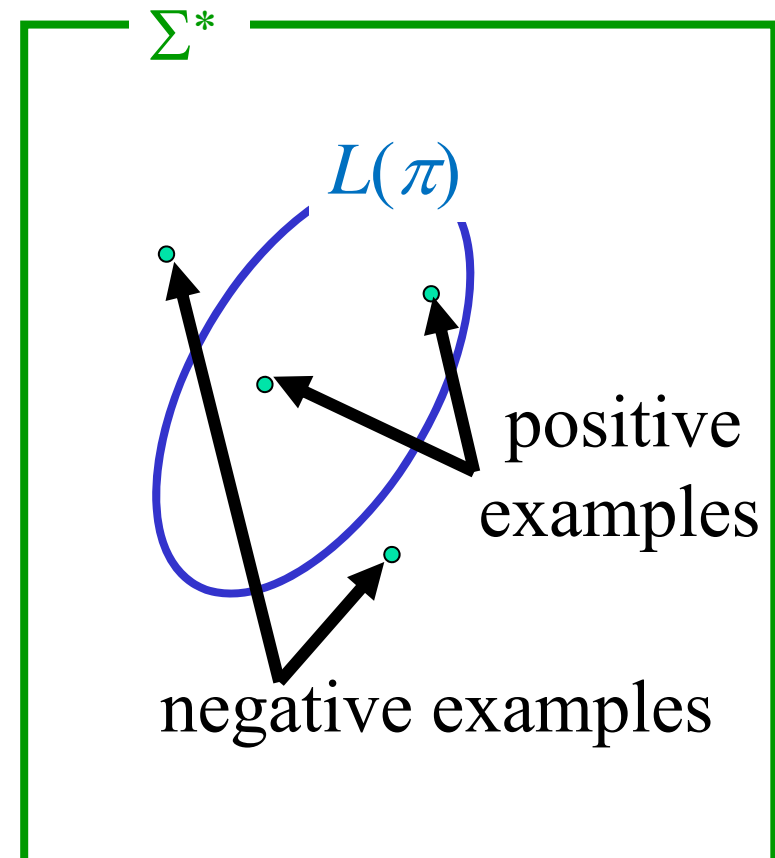
# Positive and Negative examples

$e_1, e_2, e_3, ...$

$L(\pi)$ : a language represented with a pattern $\pi$

- a positive example on $L(\pi)$ :

  $< s, +>$ for $x \in L(\pi)$

  a negative example on $L(\pi)$ :

  $< s, ->$ for $x \notin L(\pi)$

$\Sigma^*$

$L(\pi)$

positive examples

negative examples

# Abstract Classification

- A half-plane $P$ which contains $C$ (yes) and excludes $D$ (no) is to be learned

- The half-plane $P$ is represented as a pair $(w, c)$ which means the linear inequation $(w, x) + c > 0$.

  - Let $C(p) = \{x \in R^n \mid p(x)\}$ for a predicate $p$.

    Then the search space (version space) is

    $$C = \{C(\lambda x.((w, x) + c > 0)) \mid w \in R^n, c \in R^n\}.$$

    The set of parameter s are from

    $$H = \{(w, c) \mid w \in R^n, c \in R^n\}.$$

- The training examples are provided as the sets $C$ and $D$.

- A learning algorithm is provided.

# Typical evaluation method

- A learning algorithm $A$ is evaluated with test data as follows.

Step1. Let $C_*$ are set of all positive data and $D_*$ be are all negatives.

Step 2. Select subsets $C_{\text{training}} \subset C_*$ and $D_{\text{training}} \subset D_*$ for training.

Step 3. Apply $A$ to the pair $C_{\text{training}}$ and $D_{\text{training}}$ and obtain a rule $f$.

Step 4. Select subsets $C_{test}$ and $D_{test}$ make a confusion matrix.

Step 5. Calculate some measures from the confusion matrix.

# Confusion Matrix

- Every data is represented as a pair $x = <s, p>$

$p = +$ if $s \in C$ and $p = -$ if $s \in D$

|  | $C_{\text{test}}$ | $D_{\text{test}}$ |
|---|---|---|
| $\{s \in C_{\text{test}} \cup D_{\text{test}} / f(s) = 1\}$ positive | true positive | false positive |
| $\{s \in C_{\text{test}} \cup D_{\text{test}} / f(s) = 0\}$ negative | false negative | true negative |

# Measures

- Accuracy

$$\frac{|TP| + |TN|}{|TP| + |FP| + |TN| + |FN|}$$

- Error rates    1- Accuracy

- Precision (positive predictive rate)

$$\frac{|TP|}{|TP| + |FP|}$$

- Recall (coverage, true-positive rate, sensitivity)

$$\frac{|TP|}{|TP| + |FN|}$$

# When learning is correct?

- Every data is represented as a pair $x = <s, p>$

  $p = +$ if $s \in C$ and $p = -$ if $s \in D$

|  | $C_{\text{test}}$ | $D_{\text{test}}$ |
|---|---|---|
| $\{s \in C_{\text{test}} \cup D_{\text{test}} / \ f(s) = 1 \}$ positive |  | empty |
| $\{s \in C_{\text{test}} \cup D_{\text{test}} / \ f(s) = 0 \}$ negative | empty |  |

# Comparison with an Unknown Function

- Assuming an unknown discriminant function $f_*$ such that

$$C_* = \{ \, \boldsymbol{x} = <w, 1> \quad | \; f_* (w) = 1 \, \}$$

$$D_* = \{ \, \boldsymbol{x} = <w, 1> \quad | \; f_* (w) = 0 \, \}$$

we evaluate the learning algorithm $A$ by comparing its output $f$ with $f_*$.

- If every function $f$ that we treat is represented as a parameter $p$, we compare $p$ for $f$ and $p_*$ for $f_*$.

  - Every linear inequation $(\boldsymbol{w}, \boldsymbol{x}) + c > 0$ is represented as a parameter vector $(\boldsymbol{w}, c)$.

  - We evaluate $A$ with comparing $(\boldsymbol{w}, c)$ and $(\boldsymbol{w}_*, c_*)$.

# Correctness with Unknown Functions (1)

- Assuming an unknown discriminant function $f_*$, we could say that the learning algorithm $A$ is correct if the output $f$ of $A$ becomes *nearer* $f_*$ when more data are fed to $A$.

- Mathematically, consider a infinite sequence of training data sets $(C_0, D_0), (C_1, D_1), (C_2, D_2),\ldots$ such that
$$C_0 \subset C_1 \subset C_2 \subset \ldots \subset C_* \text{ and}$$
$$D_0 \subset D_1 \subset D_2 \subset \ldots \subset D_*.$$

Let $f_i$ be the output of $A$ for $C_i$ and $D_i$.

Then the algorithm $A$ is correct if $\| f_i - f_* \| \to 0$ for any of such sequences.

25

# Correctness with Unknown Functions (2)

- A similar definition of correctness could be defined: If the learning algorithm $A$ is <span style="color:red">correct</span> if $A$ outputs $f_*$ whenever an enough amount of training data are fed to $A$.

- Mathematically, consider a infinite sequence of training data sets $(C_1, D_1), (C_2, D_2), (C_3, D_3), \ldots$ such that
$$C_1 \subset C_2 \subset C_3 \subset \ldots \subset C_* \text{ and}$$
$$D_1 \subset D_2 \subset D_3 \subset \ldots \subset D_*.$$
Let $f_i$ be the output of $A$ for $C_i$ and $D_i$.
Then the algorithm $A$ is <span style="color:red">correct</span> if for <span style="color:red">each</span> of such sequences, there exists an $N$ such that $\| f_i - f_* \| = 0$ for all $n \geq N$.

# Estimation and Learning

- Estimation in statistics means to infer the value of parameters from examples.

- We assume an unknown value of $\theta$.

- The parameter $\theta$ affects the distribution of $D(\theta)$, and only finite number of data are coming from the set.

- We expect that, more data from $D(\theta)$, better conjecture $\theta^{\wedge}$ could be obtained.

- The conjecture $\theta^{\wedge}$ is (statistically) consistent if

$$\lim_{n\to\infty} E(\theta^{\wedge}) = \theta$$
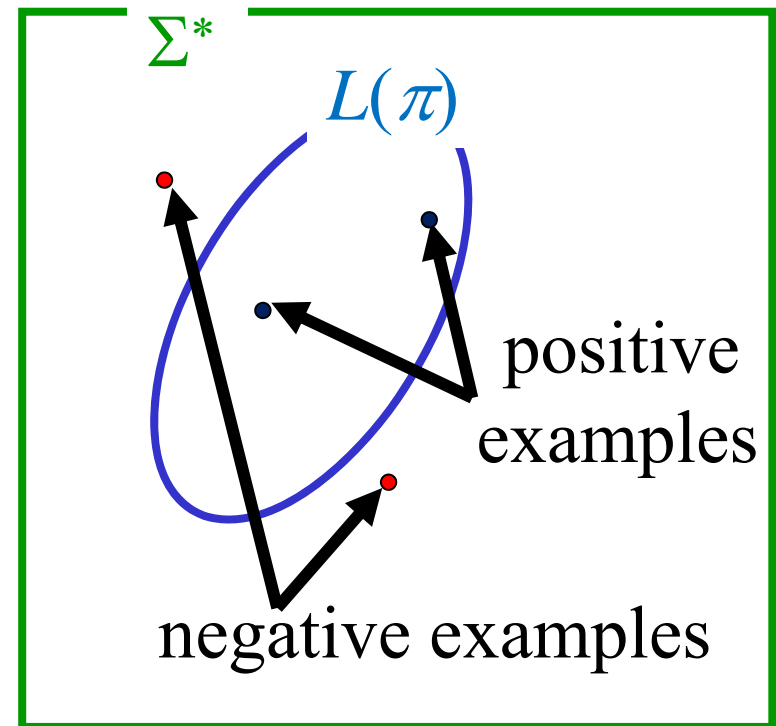
# Correctness of Learning Patterns

# Examples on $L(\pi)$
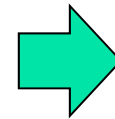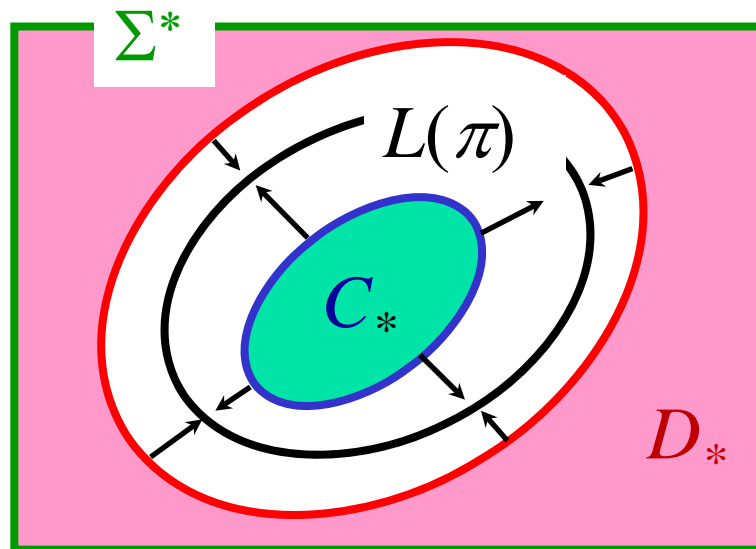
- We assume that, for an unknown pattern $\pi_*$,

$C_*$ is a finite set of positive examples on $L(\pi_*)$ and
$D_*$ is a finite set of negative examples on $L(\pi_*)$.

- a positive example on $L(\pi)$ :

$$< x, +> \text{ for } x \in L(\pi)$$

- a negative example on $L(\pi)$ :

$$< x, -> \text{ for } x \notin L(\pi)$$



$\Sigma^*$

$L(\pi)$

positive examples

negative examples

# Question

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

$$\Sigma^*$$

$$L(\pi)$$

$$C_*$$

$$D_*$$

$$\pi^\wedge \to \pi$$

# 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

$\sigma$ : <ab,+>, <aab,+>, <bbb,−>, <aaab,+>, <abba,−>, ...

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

$D_i$ = {bbb, abba}.

# Presentations

Definition A presentation of $L(\pi)$ is a infinite sequence

$$\sigma: <s_0, p_0>, <s_1, p_1>, <s_2, p_2>, \ldots$$
where $s_i \in \Sigma^*$ and $p_i = +$ or $-$ .

- $<s, +>$ is a positive example
- $<s, ->$ is a negative example

- $\sigma[n] = <s_0, p_0>, <s_1, p_1>, <s_2, p_2>, \ldots, <s_{n-1}, p_{n-1}>$

Definition A presentation $\sigma$ is complete if

any $x \in L(\pi)$ appears in $\sigma$ as a positive example $<x, +>$ at least once and

any $x \notin L(\pi)$ appears in $\sigma$ as a negative example $<x, ->$ at least once.

# Identification in the limit [Gold]

$x_1, x_2, x_3, \ldots$   ➡   $\pi_1, \pi_2, \pi_3, \ldots$

- A learning algorithm $A$ EX-identifies $L(\pi)$ in the limit from complete presentations if
  for any complete presentation $\sigma = x_1, x_2, x_3, \ldots$ of $L(\pi)$
  and the output sequence $\pi_1, \pi_2, \pi_3, \ldots$ of $A$, there exists $N$
  such that for all $n \geq N$ $\pi_n = \pi'$ and $L(\pi') = L(\pi)$

For $n = 1$ forever

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

   compute the list $l = \pi_1,\, \pi_2\,,\, \ldots,\, \pi_k$ of all the

   least common anti-unifications of the set of

   positives $C_n = \{s_j : \langle\, s_j\,,\, +\,\rangle$ and $j = 1,\, 2,\, \ldots,\, n\}$

   for each $\pi_j$ in the list $l$

      if an $e_j = \langle\, s_j\,,\, -\,\rangle$ s.t. $s_j \in L(\pi_j)$ is found

      delete $\pi_j$ from $l$

   if $l$ is not empty

      return one $\pi_i$

# Identification of patterns

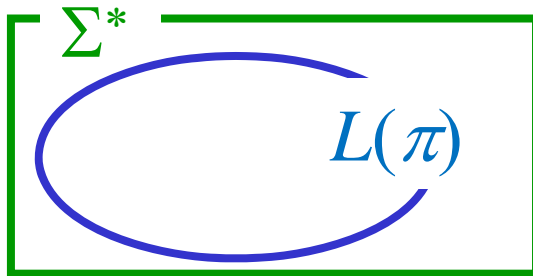Theorem The algorithm of *Learn-pattern* EX-identifies the class of all pattern languages in the

# Positive presentations

$e_1, e_2, e_3, \ldots$ → $\pi_1, \pi_2, \pi_3, \ldots$

$\Sigma^*$

$L(\pi)$

- A presentation $\sigma$ is positive if $\sigma$ consists only of positive example $< s, +>$ and any positive example occurs at least once in $\sigma$.

- A presentation $\sigma$ is complete if any $x \in L(\pi)$ appears in $\sigma$ as a positive example $< s, +>$ at least once.

# The learning algorithm *learn-patterns*

For $n = 1$ forever
    receive $e_n = \langle\, s_n\,,\, b_n\,\rangle$
    compute the list $l = \pi_1,\, \pi_2\,,\, \ldots,\, \pi_k$ of all the
    least common anti-unifications of the set of
    positives $C_n = \{s_j : \langle\, s_j\,,\, +\,\rangle$ and $j = 1, 2, \ldots, n\}$
    ~~for each $\pi_j$ in the list $l$~~
        ~~if an $e_j = \langle\, s_j\,,\, -\,\rangle$ s.t. $s_j \in L(\pi_j)$ is found~~
        ~~delete $\pi_j$ from $l$~~
    ~~if $l$ is not empty~~
        return one $\pi_i$

# Identification of patterns

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