



Computational Learning Theory

Correctness of Algorithms for Learning Finite State Automata

Akihiro Yamamoto 山本 章博

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



Note

- There might be several automata consistent with given C and D .
- In such automata there is a minimum one in the sense that the number of states in it is minimum.
- Unfortunately it is proved that the problem of finding a **minimum** automaton consistent with given C and D is NP-hard.
 - The activity of a learning algorithm should not be evaluated (justified) only on the viewpoint of optimization.
 - Even though it were not ensured that the algorithm returns the best solution, the algorithm could work as “learning”.



A Simple Generate-and-Test Algorithm

Assume a procedure of enumerating all FA so that the enumeration $M_0, M_1, M_2, \dots, M_i, \dots$ satisfies

$$P(M_0) \leq P(M_1) \leq P(M_2) \leq \dots \leq P(M_i) \leq \dots$$

Let the input data x_1, x_2, \dots, x_N

Initialize $M = M_0$ as an automaton consisting of one state

let $k = 0$

forever

let $k' = k$

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

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

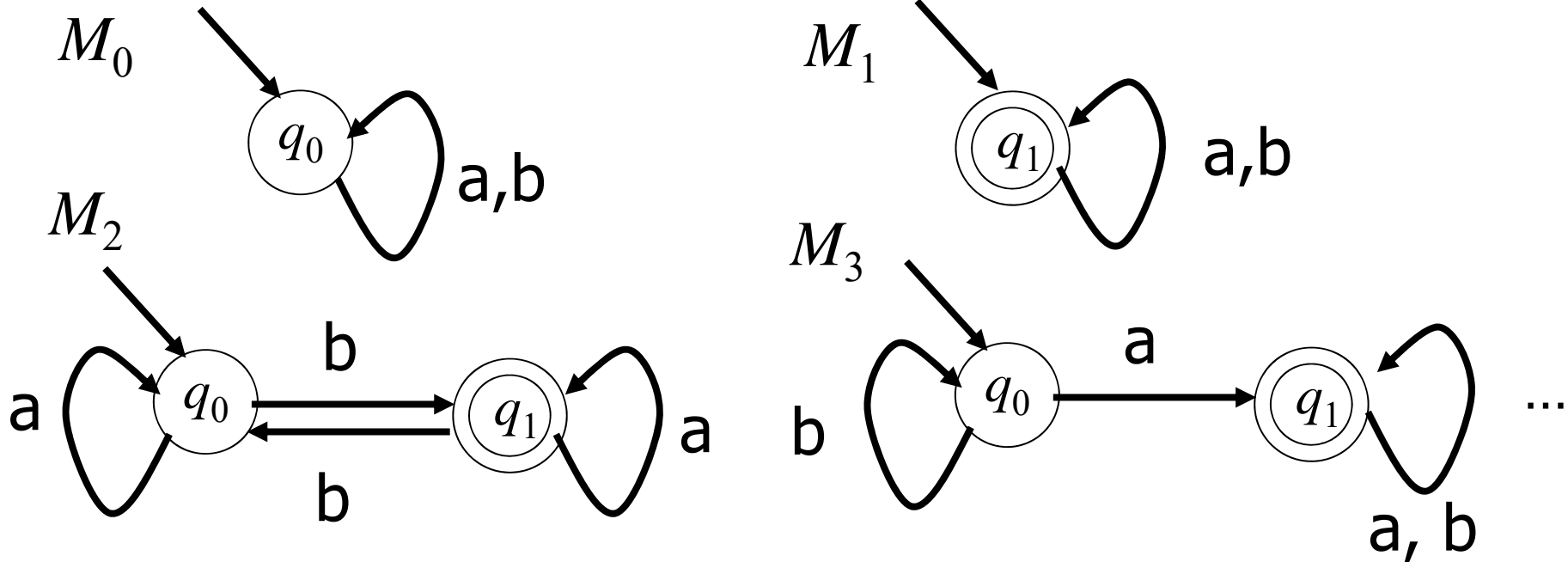
replace k with $k + 1$

if $k' = k$

terminate and output M_k

Example

- Enumeration of FAs



Input:

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

Output:

$M_1 \quad M_1 \quad M_3 \quad M_3 \quad \dots$

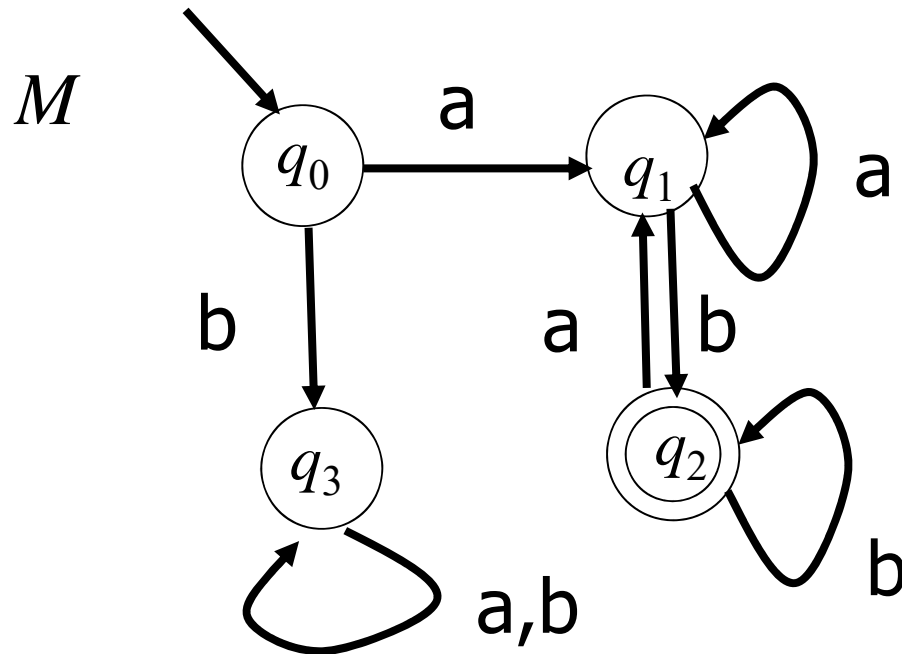
Example(cont.)

Input:

$\sigma: \dots, \langle \text{bbb}, - \rangle, \langle \text{aaab}, + \rangle, \langle \text{abba}, - \rangle, \dots \dots \dots$

Output:

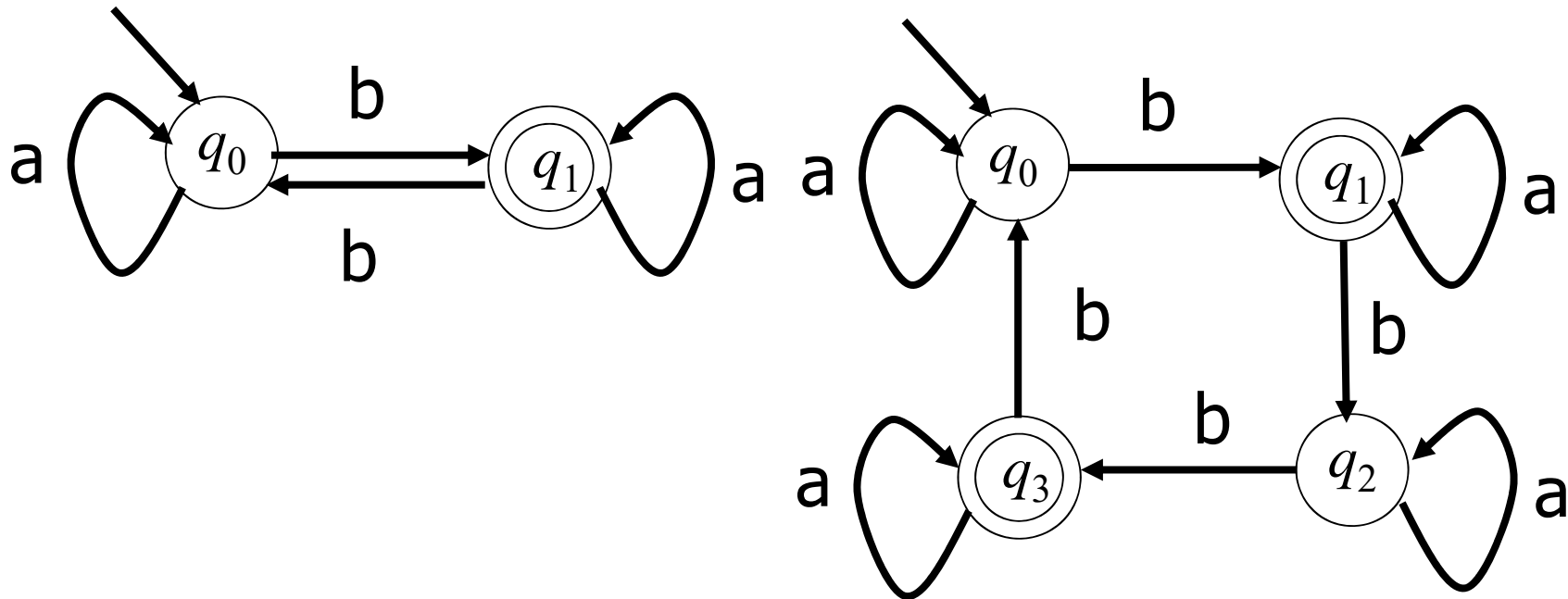
$\dots M_3 M_3 \dots \dots M' M' M'$



$$L(M') = L(M)$$

Note

- For a finite automaton M , there might be exists another automaton such that $L(M) = L(M')$.





Evaluation of Learning Algorithm

Support Vector Machine

- **Input:** a set of numerical data

$$\{(\mathbf{x}_1, c_1), (\mathbf{x}_2, c_2), \dots, (\mathbf{x}_m, c_m)\} \quad \mathbf{x}_i \in \mathbf{R}^n$$

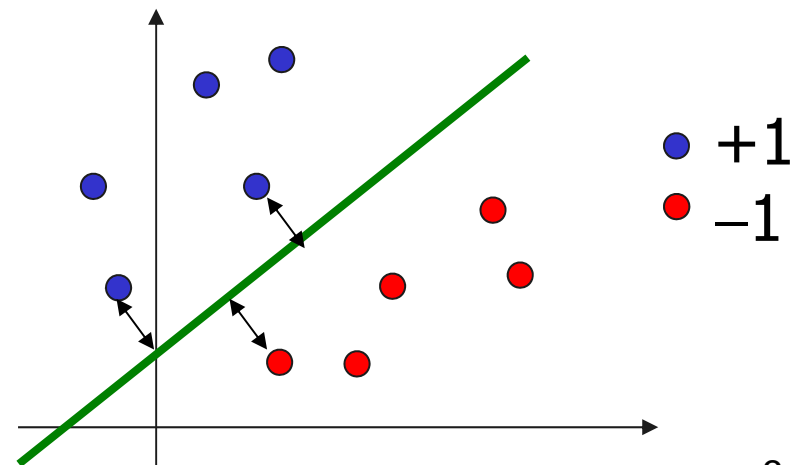
where each $c_i \in \{+1, -1\}$ is a class signal for \mathbf{x}_i

Output: find a linear function (hyper-plane)

$$f(\mathbf{x}) = \sum w_i \mathbf{x}_i \cdot \mathbf{x} + b$$

which $\text{sign}(f(\mathbf{x}_i)) = y_i$ for all i and

maximize the margin $\min_{1 \leq i \leq m} d(f, \mathbf{x}_i)$





General Learning

- Recently a machine learning method is recognized as one to find

$$\operatorname{argmin}_{f \in H} (\sum_{\mathbf{x} \in D} \operatorname{Loss}(f, \mathbf{x}) + \lambda P(f))$$

where

$\operatorname{Loss}(f, \mathbf{x})$ is a loss function and

$P(f)$ is a penalty function.

- This definition is declarative.
- This course we introduce some of the instances of $\operatorname{Loss}(f, \mathbf{x})$ and $P(f)$.



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 (\mathbf{w}, c) which means the linear inequation $(\mathbf{w}, \mathbf{x}) + c > 0$.

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

Then the search space (version space) is

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

The set of parameter s are from

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

- 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.

Step 1. Let C_* be set of all positive data and D_* be 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_{training} and D_{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 $\mathbf{x} = \langle w, s \rangle$
 $s = +$ if $w \in C$ and $s = -$ if $w \in D$

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



Measures

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

- Error rates $1 - \text{Accuracy}$

- Precision (positive predictive rate)

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

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

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



Comparison with an Unknown Function

- Assuming an unknown discriminant function f_* such that

$$C_* = \{ \mathbf{x} = \langle \mathbf{w}, 1 \rangle \mid f_*(\mathbf{w}) = 1 \}$$

$$D_* = \{ \mathbf{x} = \langle \mathbf{w}, 1 \rangle \mid f_*(\mathbf{w}) = 0 \}$$

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 $(\mathbf{w}, \mathbf{x}) + c > 0$ is represented as a parameter vector (\mathbf{w}, c) .
 - We evaluate A with comparing (\mathbf{w}, c) and (\mathbf{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 $(C_0, D_0), (C_1, D_1), (C_2, D_2), \dots$ such that

$$C_0 \subset C_1 \subset C_2 \subset \dots \subset C_* \text{ and}$$

$$D_0 \subset D_1 \subset D_2 \subset \dots \subset D_*.$$

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

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



Correctness with Unknown Functions (2)

- A similar definition of correctness could be defined:
If the learning algorithm A is **correct** if
 A outputs f_* whenever an enough amount of training data are fed to A .
- Mathematically, consider a infinite sequence of training data sets $(C_1, D_1), (C_2, D_2), (C_3, D_3), \dots$ such that
$$C_1 \subset C_2 \subset C_3 \subset \dots \subset C_*$$
and
$$D_1 \subset D_2 \subset D_3 \subset \dots \subset D_*$$
.
Let f_i be the output of A for C_i and D_i .
Then the algorithm A is **correct** if for **each** 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 θ .
- The parameter θ 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 θ^\wedge could be obtained.
- The conjecture θ^\wedge is (**statistically**) **consistent** if

$$\lim_{n \rightarrow \infty} E(\theta^\wedge) = \theta$$



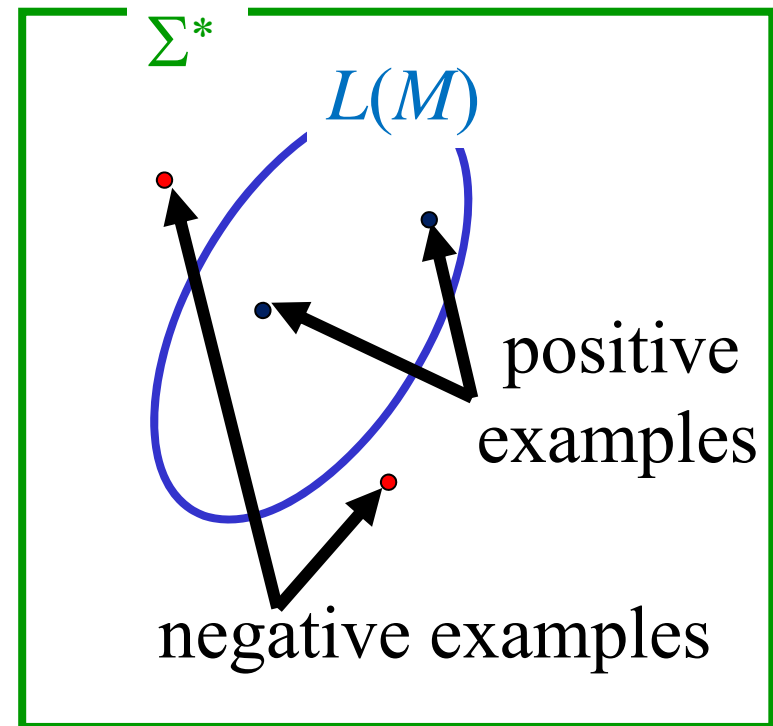
Correctness of Learning Automata

Examples on $L(M)$

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

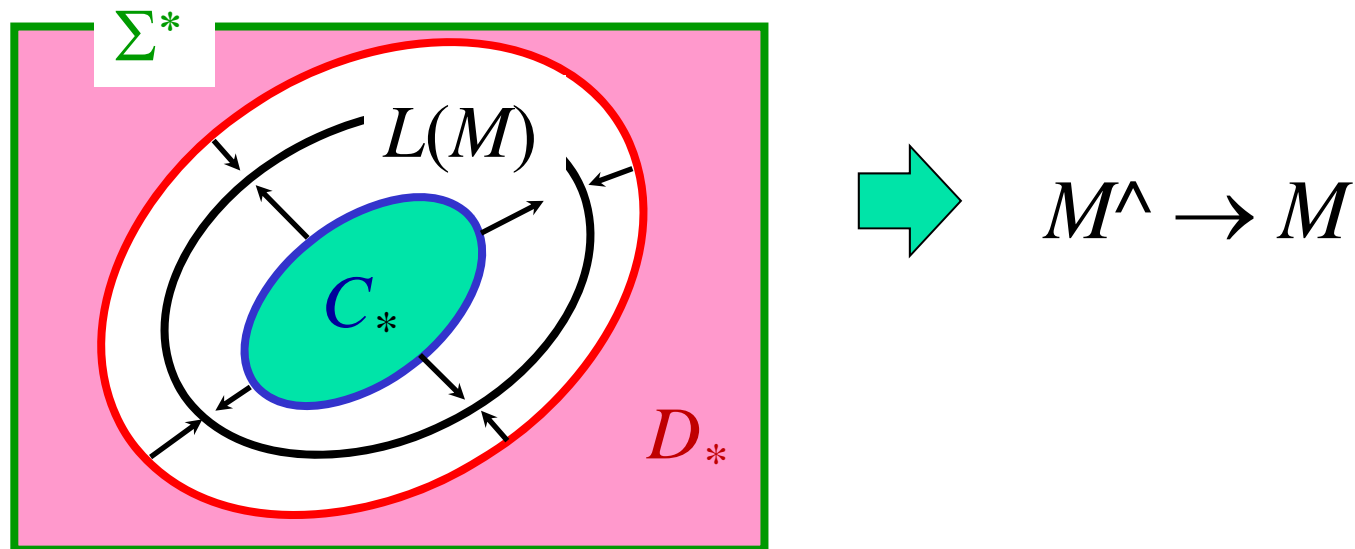
$L(M)$: a language accepted by a finite state automaton M

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



Question

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





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** σ of strings marked with either + or -, and some **truncation** of σ 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 $L(M)$ 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 σ is **complete** if

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

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

Identification in the limit [Gold]



x_1, x_2, x_3, \dots



M_1, M_2, M_3, \dots

- A learning algorithm A **EX-identifies** $L(M)$ **in the limit from complete presentations** if for any complete presentation $\sigma = x_1, x_2, x_3, \dots$ of $L(M)$ and the output sequence M_1, M_2, M_3, \dots of A , there exists N such that for all $n \geq N$ $M_n = M'$ and $L(M') = L(M)$
- A learning algorithm A **BC-identifies** $L(M)$ **in the limit from complete presentations** if for any complete presentation $\sigma = x_1, x_2, x_3, \dots$ of $L(M)$ and the output sequence M_1, M_2, M_3, \dots of A , there exists N such that for all $n \geq N$ ~~$M_n = M'$~~ and $L(M_n) = L(M)$



A Well-known Result on FA

Theorem For every language $L(M)$ accepted by a finite state automaton M , there exists a unique minimal automaton M' such that $L(M)=L(M')$, where “minimal” means that the number of states in M' is minimal in such automata.



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

Assume a procedure of enumerating all FA so that the enumeration $M_0, M_1, M_2, \dots, M_i, \dots$ satisfies

$$P(M_0) \leq P(M_1) \leq P(M_2) \leq \dots \leq P(M_i) \leq \dots$$

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

Initialize $k = 0$ /* M_0 as an automaton consisting of one state*/

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 $x_n \notin L(M_{k'})$) or ($x_n \in D$ and $x_n \in L(M_{k'})$)

replace k with $k + 1$

if $k' = k$

terminate and output M_k



On the Generate-and-Test Algorithm

Theorem For any finite state automaton M_* on Σ ,

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

Proof Let σ be an any complete presentation on $L(M_*)$.

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

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

$(x \in L(M_*) \text{ and } x \notin L(M_N)) \text{ or } (x \notin L(M_*) \text{ and } x \in L(M_N)).$

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

This means that M_N must be replaced with another automaton, at latest, when x appears in σ .

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

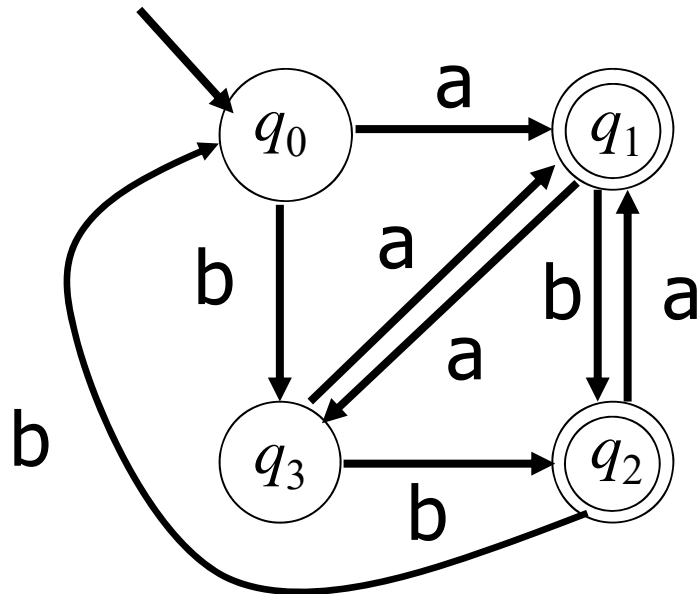


Data Sets Enough to Output Hidden Automata

Minimal Test Sets

- A set $S \subset \Sigma^*$ is a **minimal test set** for a FA M if for each state q of M , there exists exactly one string x such that $\delta(q_0, x) = q_i$.

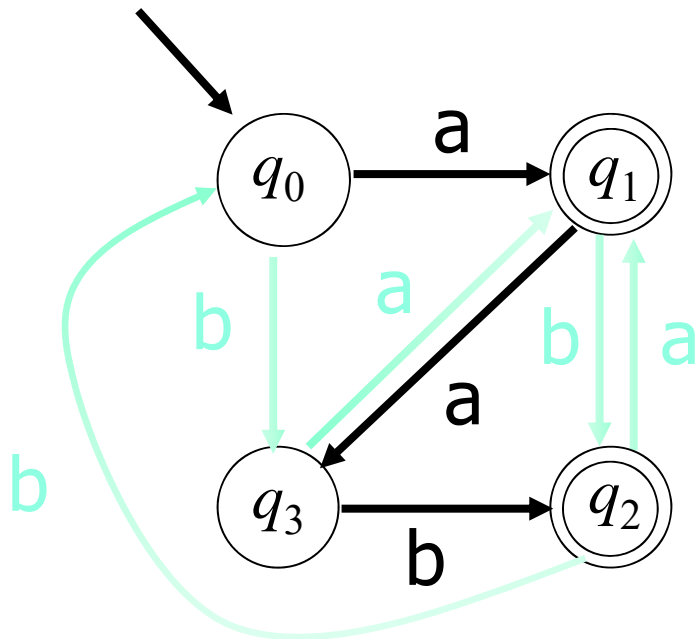
Example Examples of test sets of M are $S_1 = \{\varepsilon, a, aa, aab\}$ and $S_2 = \{\varepsilon, a, ab, b\}$.



Minimal Test Sets

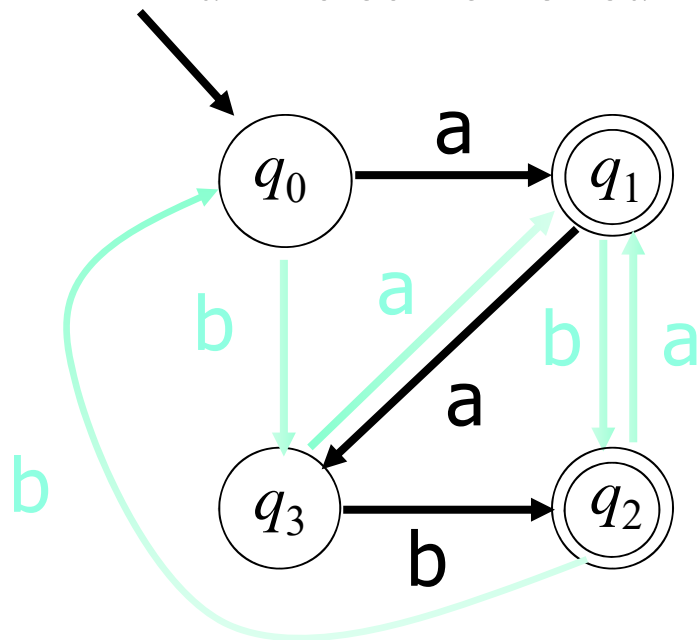
- Intuitively, a test set gives a “skelton” of the finite state automaton.
 - But the set is not sufficient to identify the FA.

Example Examples of test sets of M are $S_1 = \{\varepsilon, a, aa, aab\}$ and $S_2 = \{\varepsilon, a, ab, b\}$.



Prefix closed Test Sets

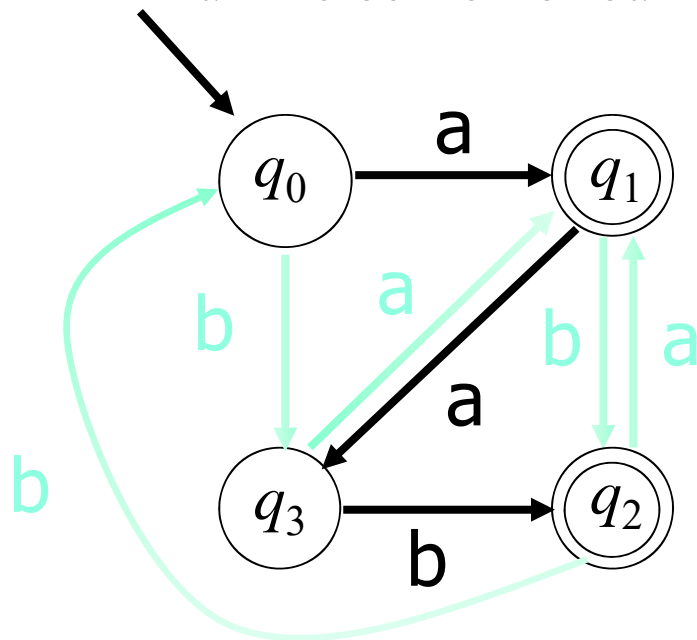
- A set of strings S is **prefix closed (suffix closed)** if and only if every prefix (resp. suffix) of every member of S is also a member of S .
- Intuitively, a prefix closed minimal test set gives a “skelton” of the finite state automaton.
 - But the set is not sufficient to identify the FA.



Example Both $S_1 = \{\varepsilon, a, aa, aab\}$ and $S_2 = \{\varepsilon, a, ab, b\}$ are prefix closed.

Prefix closed Test Sets

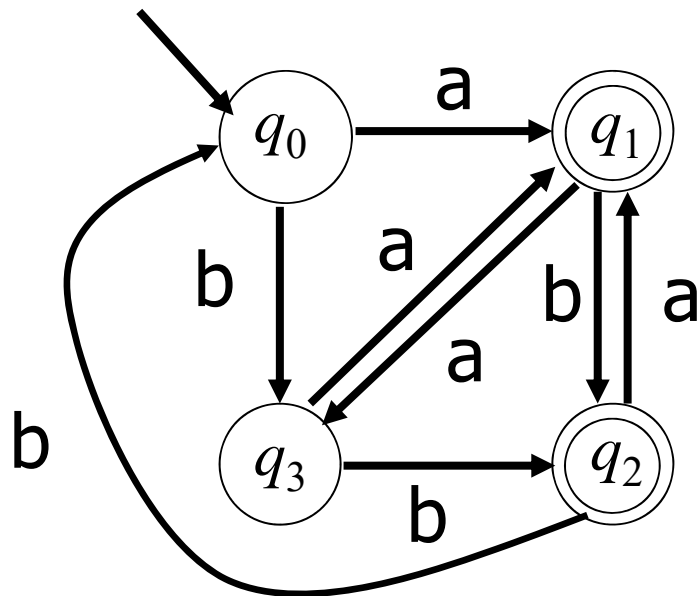
- A set of strings S is **prefix closed (suffix closed)** if and only if every prefix (resp. suffix) of every member of S is also a member of S .
- Intuitively, a prefix closed minimal test set gives a “skelton” of the finite state automaton.
 - But the set is not sufficient to identify the FA.



Example Both $S_1 = \{\varepsilon, a, aa, aab\}$ and $S_2 = \{\varepsilon, a, ab, b\}$ are prefix closed.

Fixing an Order

- We fix one ordering for listing elements of a set.
- **Example** Following the lexicographic ordering, elements of $S_1 = \{\varepsilon, a, aa, aab\}$ is listed as ε, a, aa, aab

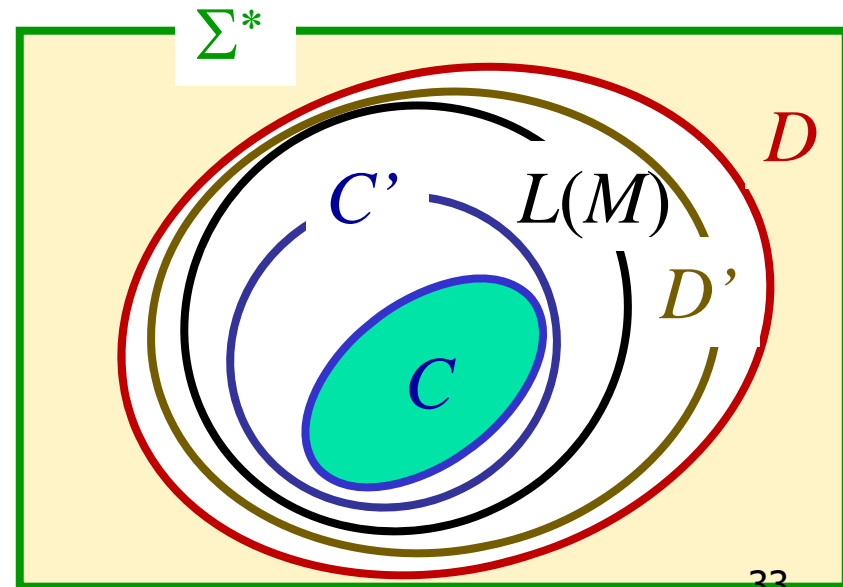


Characteristics Examples

- Assume an algorithm A which learns FA.
- Assume that we treat only minimal FA.
- A pair (C, D) of sets of examples is **characteristic** for a FA M if for **any** pair (C', D') of examples such that

$$C \subset C' \subset L(M) \text{ and } D \subset D' \subset \overline{L(M)}$$

the algorithm A returns M .



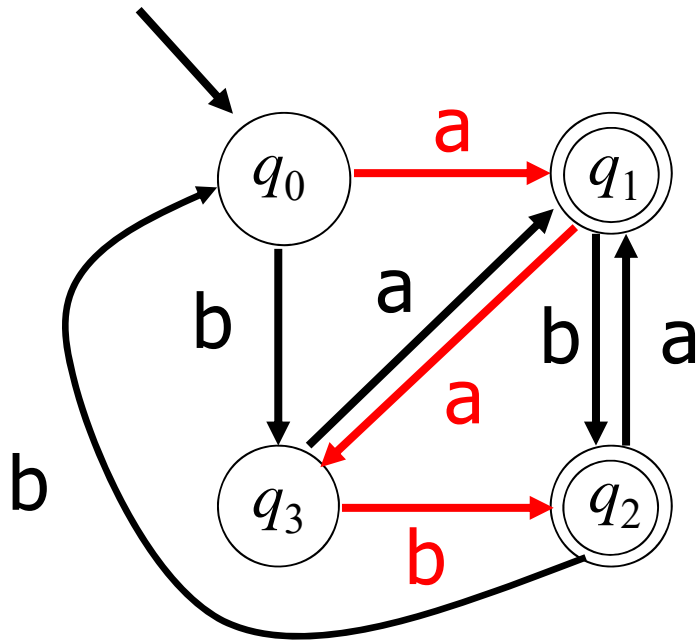
Observation table

- An observation table (S, E, T) :
 - S : a prefix closed set $S \subset \Sigma^*$
 - E : a suffix closed set $E \subset \Sigma^*$
 - $T : (S \cup S \Sigma)E \rightarrow \{0, 1\}$
 - $S \Sigma = \{ sa \mid s \in S \text{ and } a \in \Sigma \}$
 - The element of the position (s, w) shows whether or not the automaton M accepts sw .

		E	
		ε	b
S	ε	0	0
	a	1	1
	aa	0	1
	aab	1	0
$S \Sigma$	b	0	1
	ab	1	0
	aaa	1	1
	$aaba$	1	1
	$aabb$	0	0

Observation table

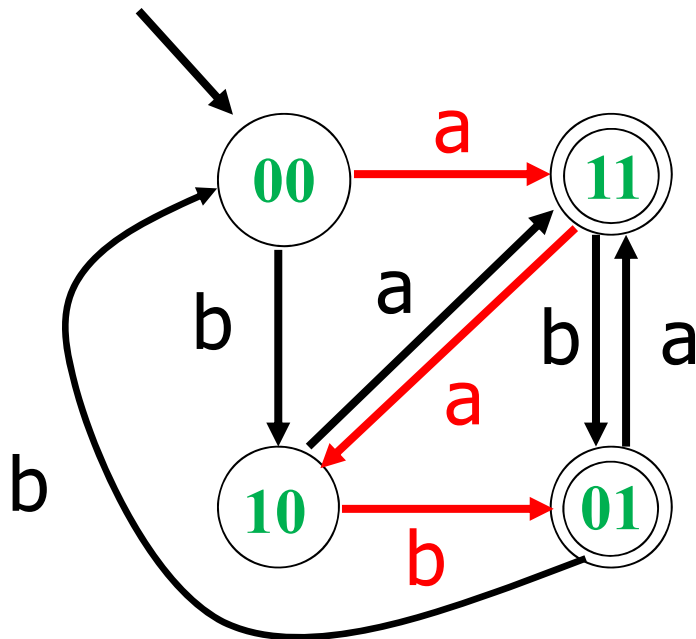
- An observation table (S, E, T) :
 - S : a prefix closed set $S \subset \Sigma^*$
 - E : a set $E \subset \Sigma^*$
 - $T : (S \cup S\Sigma)E \rightarrow \{0, 1\}$



		E	
		ε	b
S	ε	0	0
	a	1	1
	aa	0	1
	aab	1	0
$S\Sigma$	b	0	1
	ab	1	0
	aaa	1	1
	$aaba$	1	1
	$aabb$	0	0

Observation table

- An observation table (S, E, T) :
 - S : a prefix closed set $S \subset \Sigma^*$
 - E : a set $E \subset \Sigma^*$
 - $T : (S \cup S\Sigma)E \rightarrow \{0, 1\}$



		E	
		ϵ	b
S	ϵ	0	0
	a	1	1
	aa	0	1
	aab	1	0
$S\Sigma$	b	0	1
	ab	1	0
	aaa	1	1
	$aaba$	1	1
	$aabb$	0	0



How to construct the table

Input : a minimal FA A

Output : The characteristic set of polynomial size

$S :=$ the minimal test set of A , $E := \{ \varepsilon \}$, $S' := S\Sigma - S$,

Generate (S, E, T) ;

while there exists $w, v \in S$ s.t. $\text{row}(w) = \text{row}(v)$ but

$T(wc, e) \neq T(vc, e)$ for some $c \in \Sigma$ and $e \in E$

$E := E \cup \{ae\}$;

Generate (S, E, T) ;

end while

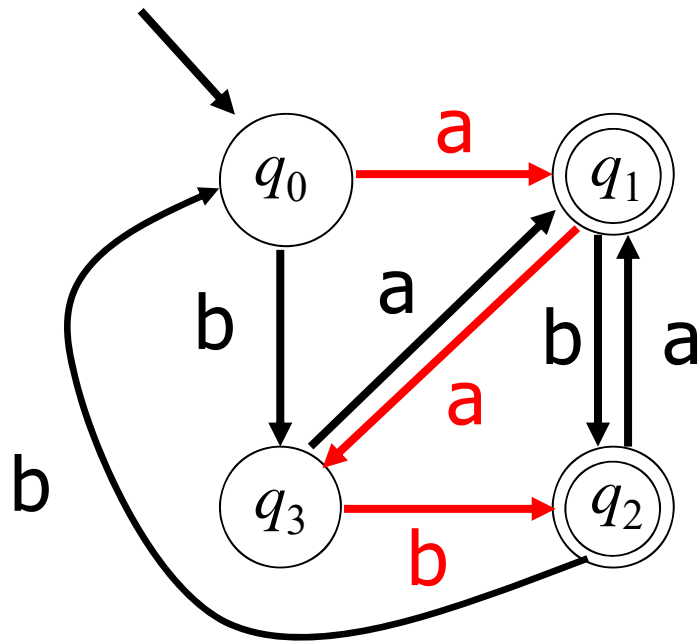
$C = \{ we \mid w \in S \cup S\Sigma, e \in E, \text{ and } T(wc, e) = 1 \}$

$D = \{ we \mid w \in S \cup S\Sigma, e \in E, \text{ and } T(wc, e) = 0 \}$

return (C, D) ;

Example

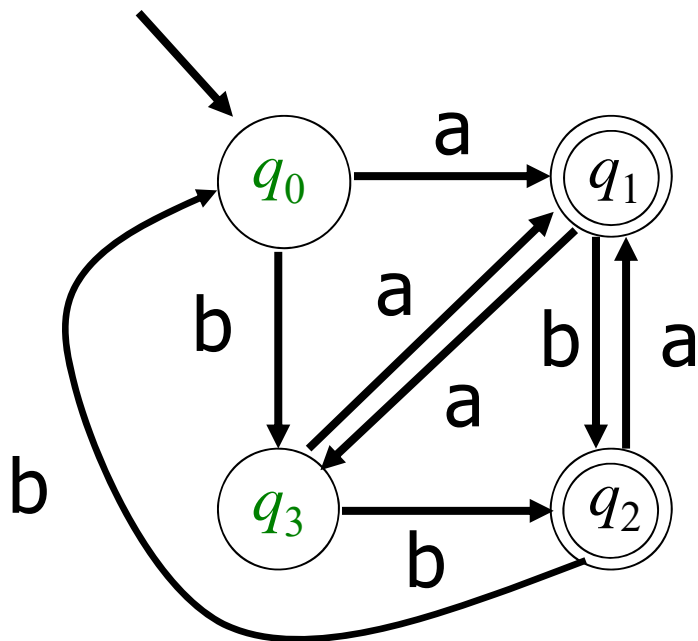
- $S = \{\varepsilon, a, aa, aab\}$
- $S\Sigma = \{a, aa, aaa, aaba, b, ab, aab, aabb\}$
- $E = \{\varepsilon\}$.



	E	
	ε	
S	ε	0
	a	1
	aa	0
	aab	1
$S\Sigma$	b	0
	ab	1
	aaa	1
	aaba	1
	aabb	0

Example

Because $T(\varepsilon, \varepsilon) = T(\mathbf{aa}, \varepsilon)$, check whether or not $T(\mathbf{a}, \varepsilon) = T(\mathbf{aaa}, \varepsilon)$, and whether or not $T(\mathbf{b}, \varepsilon) = T(\mathbf{aab}, \varepsilon)$.

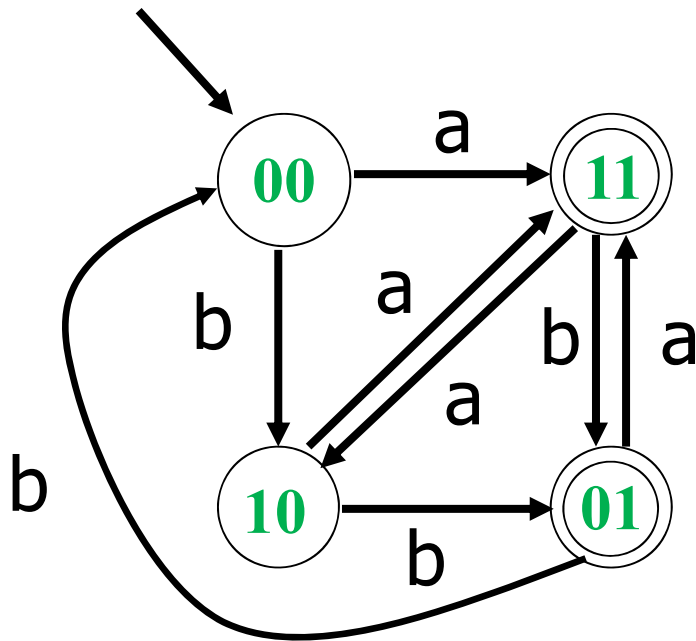


	ε
ε	0
a	1
aa	0
aab	1
b	0
ab	1
aaa	1
aaba	1
aabb	0

Labels on the left: S (bracketed over rows 2-5), $S \Sigma$ (bracketed over rows 6-9). Label on top: E (bracketed over column 2).

Example

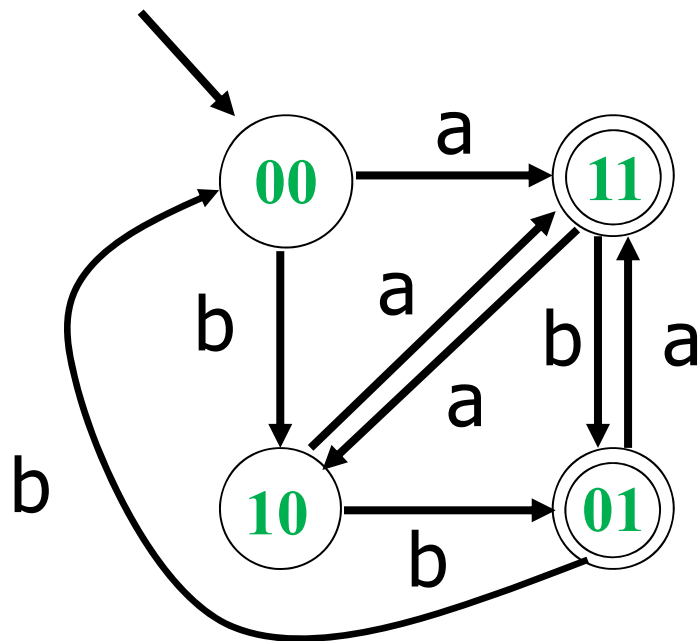
- $E := E \cup \{b\}$
- Fill all of the new elements of the extended table.



		E	
		ε	b
S	ε	0	0
	a	1	1
	aa	0	1
	aab	1	0
$S \Sigma$	b	0	1
	ab	1	0
	aaa	1	1
	aaba	1	1
	aabb	0	0

Example

- There is no w and v in the S part s.t. $\text{row}(w) = \text{row}(v)$, end the loop.
- $C = \{a, ab, bb, aaa, aab, aaab, aaba, aabab\}$
- $D = \{\varepsilon, b, aa, abb, aabb, aabbb\}$



	E	
	ε	b
ε	0	0
a	1	1
aa	0	1
aab	1	0
b	0	1
ab	1	0
aaa	1	1
$aaba$	1	1
$aabb$	0	0

S (rows 1-4)
 $S \Sigma$ (rows 5-8)



Consistent Table

- An observation table (S, E, T) is **consistent** if and only if for every pair $w, v \in S$ such that $\text{row}(w) = \text{row}(v)$, $\text{row}(wc) = \text{row}(vc)$ for any $c \in \Sigma$.
 - Intuitively, in a consistent table, every row in the S part can be regarded as one state of an automaton.

Proposition A consistent table T represents an automaton M such that, for $w \in S \cup S\Sigma$ and $e \in E$, M accepts we if and only if $T(w, e) = 1$.



Characteristic Examples

Theorem Suppose T be the table obtained above method from M . Then the pair (C, D) where

$$C = \{we \mid w \in S \cup S \Sigma \text{ and } e \in E \text{ and } T(w, e) = 1\}$$

$$D = \{we \mid w \in S \cup S \Sigma \text{ and } e \in E \text{ and } T(w, e) = 0\}$$

is characteristic w.r.t. the generate-and-test algorithm and M .



The Myhill-Nerode Theorem

Theorem The following three statements are equivalent:

(1) The language L is accepted by some finite automaton.

(2) L is the union of some equivalence classes of a **right invariant** equivalence relation of finite index.

(3) Let equivalence relation R_L be defined by: $x R_L y$ if and only if for all $z \in \Sigma^*$ xz is in L iff yz is in L . Then R_L is finite index.

- An equivalence relation R is **right invariant** iff $x R y$ implies $xz R yz$ for all $z \in \Sigma^*$.
- The index of equivalence relation R is the number of equivalence classes.

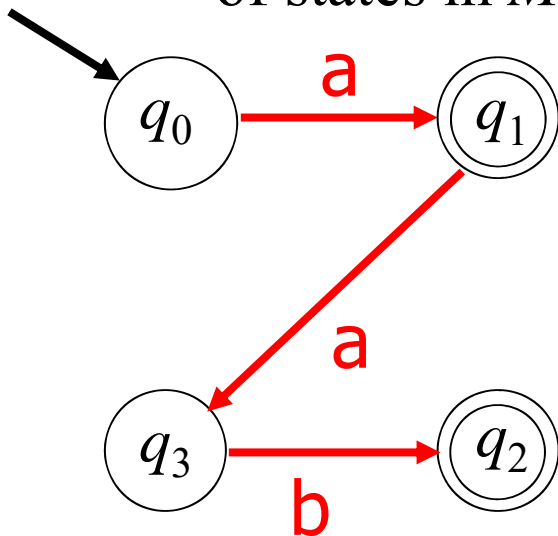
Why Characteristic Set?

- Let M' be an automaton the number of whose states is minimal and less than that of M . Then it holds that

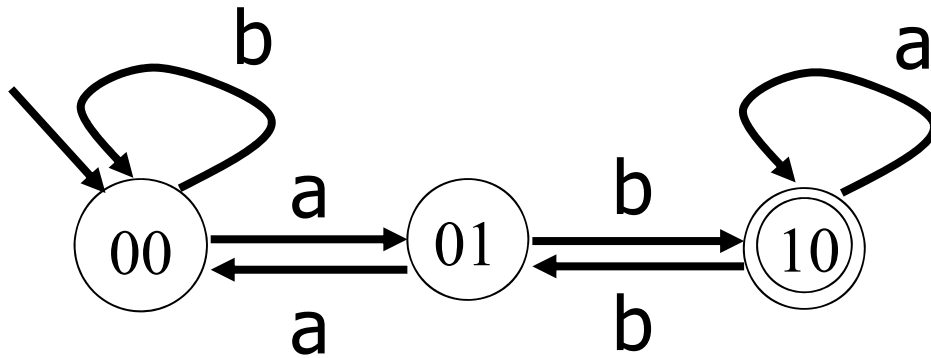
$$C \cap \overline{L(M')} \neq \emptyset \text{ or } D \cap L(M') \neq \emptyset.$$

- Proof: Since S is a test set minimal and prefix closed, we can construct a prefix tree T such that there is a one-to-one and on-to mapping between the set of nodes in T and the set of states in M . If $C \subseteq L(M')$ and $D \subseteq L(M')$, then some

node in M' must corresponds to more than two nodes in T . However from Myhill-Nerode's theorem, the set E avoids such correspondence.

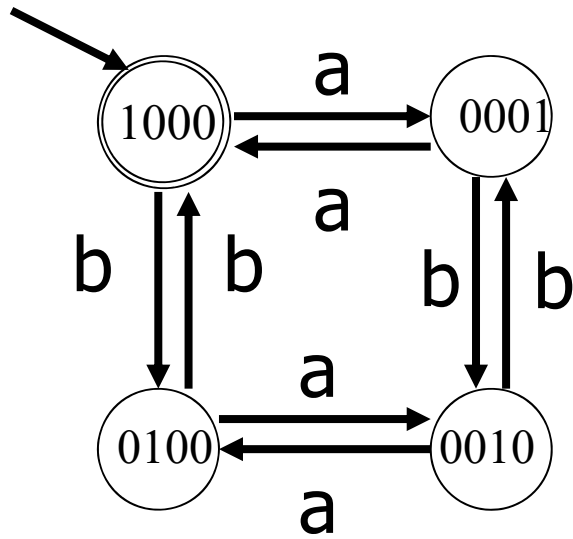


Example 2



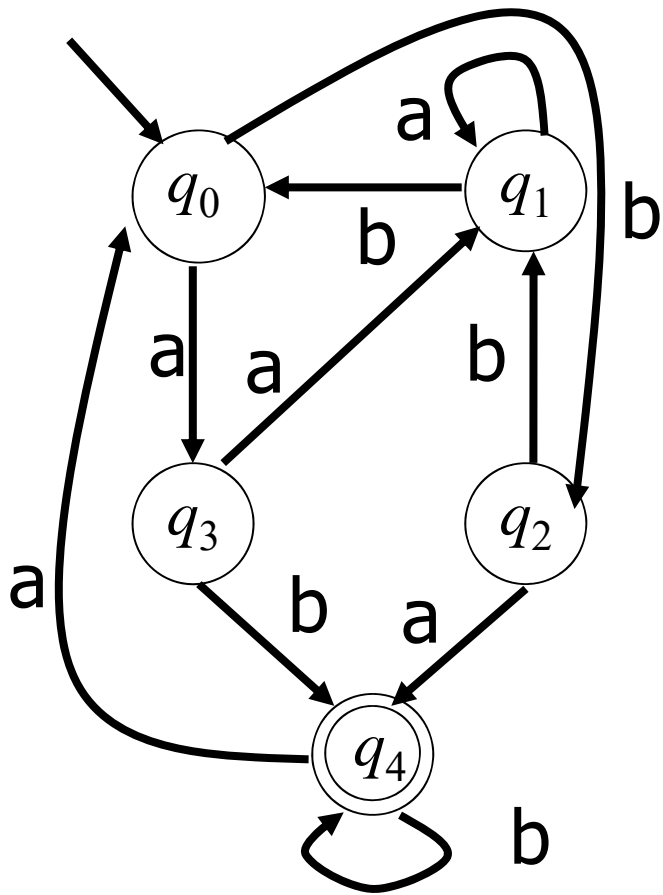
	ϵ	b
ϵ	0	0
a	0	1
ab	1	0
b	0	0
aa	0	0
aba	1	0
abb	0	1

Example 3



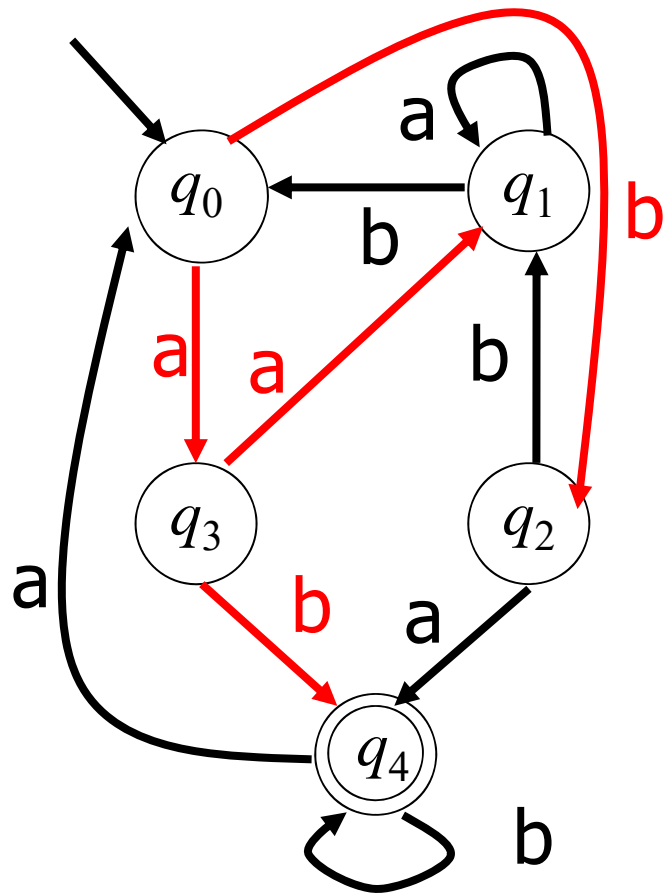
		E			
		ϵ	b	ab	bab
S	ϵ	1	0	0	0
	a	0	0	0	1
	b	0	1	0	0
$S \Sigma$	ab	0	0	1	0
	aa	1	0	0	0
	ba	0	0	1	0
	bb	1	0	0	0
	aba	0	1	0	0
	abb	0	0	0	1

Example 4



		E			
		ϵ	b	a	ab
S	ϵ	0	0	0	1
	a	0	1	0	0
	aa	0	0	0	0
	ab	1	1	0	0
	b	0	0	1	1
$S \Sigma$	ba	1	1	0	0
	bb	0	0	0	0
	aaa	0	0	0	0
	aab	0	0	0	1
	aba	0	0	0	1
	abb	1	1	0	0

Example 4



		E			
		ϵ	b	a	ab
S	ϵ	0	0	0	1
	a	0	1	0	0
	aa	0	0	0	0
	ab	1	1	0	0
	b	0	0	1	1
$S \Sigma$	ba	1	1	0	0
	bb	0	0	0	0
	aaa	0	0	0	0
	aab	0	0	0	1
	aba	0	0	0	1
	abb	1	1	0	0