

# **Advanced Automata Theory 11**

## **Regular Languages and Learning Theory**

**Frank Stephan**

**Department of Computer Science**

**Department of Mathematics**

**National University of Singapore**

**[fstephan@comp.nus.edu.sg](mailto:fstephan@comp.nus.edu.sg)**

# Repetition: Rational Relations

**Automatic Relation:** Finite automaton reads all inputs involved at the same speed with  $\#$  supplied for exhausted inputs.

**Rational Relation:** Nondeterministic finite automaton reads all inputs individually and can read them at different speed.

The first type of automata is called **synchronous**, the second type is called **asynchronous**.

There are many relations which are rational but not automatic.

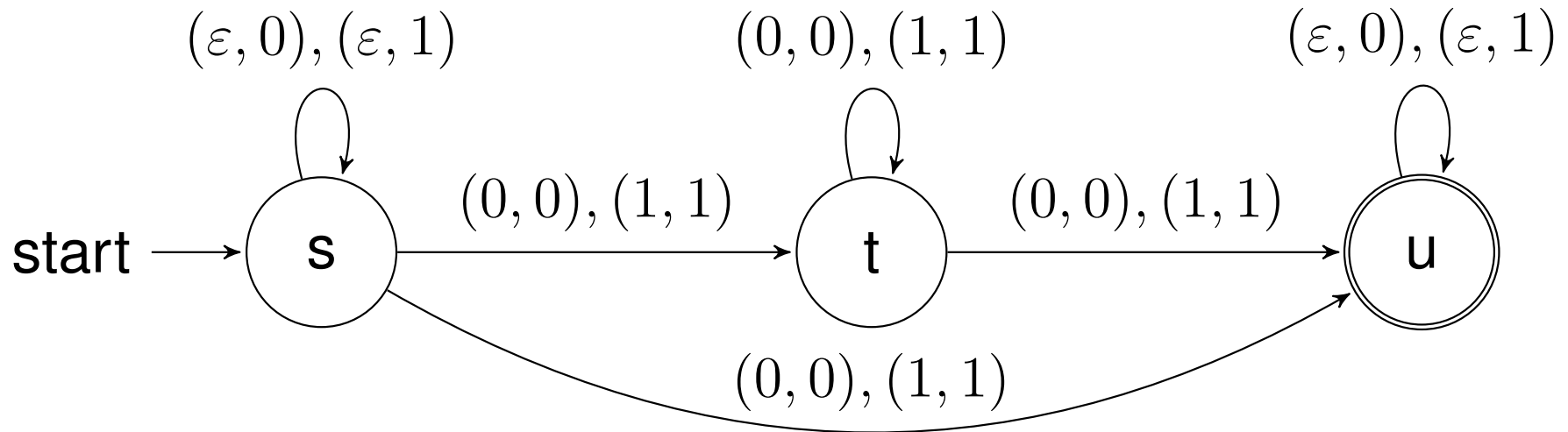
# Repetition: Formal definition

A rational relation  $R \subseteq (\Sigma^*)^n$  is given by an non-deterministic finite state machine which can process  $n$  inputs in parallel and does not need to read them in the same speed. Transitions from one state  $p$  to a state  $q$  are labelled with an  $n$ -tuple  $(w_1, w_2, \dots, w_n)$  of words  $w_1, w_2, \dots, w_n \in \Sigma^*$  and the automaton can go along this transition iff for each input  $k$  the next  $|w_k|$  symbols in the input are exactly those in the string  $w_k$  (this condition is void if  $w_k = \varepsilon$ ) and in the case that the automaton goes on this transition,  $|w_k|$  symbols are read from the  $k$ -th input word.

A word  $(x_1, x_2, \dots, x_n)$  is in  $R$  iff there is a run of the machine with transitions labelled by  $(w_{1,1}, w_{1,2}, \dots, w_{1,n})$ ,  $(w_{2,1}, w_{2,2}, \dots, w_{2,n})$ ,  $\dots$ ,  $(w_{m,1}, w_{m,2}, \dots, w_{m,n})$  ending up in an accepting state such that  $x_1 = w_{1,1}w_{2,1} \dots w_{m,1}$ ,  $x_2 = w_{2,1}w_{2,2} \dots w_{m,2}$ ,  $\dots$ ,  $x_n = w_{1,n}w_{2,n} \dots w_{m,n}$ .

# Repetition: Example 10.4

The following automaton recognises the relation of all  $(\mathbf{x}, \mathbf{y})$  where  $\mathbf{x}$  is a nonempty substring of  $\mathbf{y}$ , that is,  $\mathbf{x} \neq \varepsilon$  and  $\mathbf{y} = \mathbf{vxw}$  for some  $\mathbf{v}, \mathbf{w} \in \{0, 1\}^*$ .



In  $s$ : Parsing  $(\varepsilon, \mathbf{v})$ ;

From  $s$  to  $u$ : Parsing  $(\mathbf{x}, \mathbf{x})$ ;

In  $u$ : Parsing  $(\varepsilon, \mathbf{w})$ .

# Repetition: Transducers

**Transducers** are automata recognising rational relations where the first parameter  $x$  is called input and the second parameter  $y$  is called output.

In particular one is interested in transducers computing functions: Two possible runs accepting  $(x, y)$  and  $(x, z)$  must satisfy  $y = z$ . Given  $x$ , there is a run accepting some pair  $(x, y)$  iff  $x$  is in the domain of the function.

Two main concepts:

**Mealy machines:**

Input/Output pairs on transition-labels.

**Moore machines:**

Input on transition-labels; output in states.

# Learning Theory

Formal models of learning processes

## Query Learning

Dialogue between learner (pupil) and teacher.

Learner asks queries using a specific query language.

Teacher provides answers to these queries.

## Learning from Data

Learner reads more and more data about the concept to be learnt.

In response to the data, learner conjectures one or several hypotheses.

The last hypothesis conjectured must be correct.

# Angluin's DFA Learner

Query language – Learner can ask the following questions:

- Equivalence query (EQ): Does this DFA recognise the language  $L$  to be learnt?  
Teacher says “YES” or “NO”; when saying “NO” an  $x$  is provided on which the DFA does the opposite of what  $L$  does.
- Membership query (MQ): Is  $x$  a member of  $L$ ?  
Teacher says “YES” or “NO”.

Learner could use a list of all DFAs and ask: does the first DFA recognise  $L$ ; does the second DFA recognise  $L$ ; ...?

This needs more than  $2^n$  steps in order to deal with all languages recognised by DFA's with  $n$  states.

Angluin's algorithm needs much less queries.

# Data of Angluin's algorithm

Represent set  $S$  of states by strings  $x \in S$  on which automaton goes into corresponding states.  $x \in S$  is an accepting state iff  $x \in L$ .

Represent set  $E$  of observations so large that for each distinct  $x, y \in S$  there is an  $z \in E$  with  $L(xz) \neq L(yz)$ ; that is  $\text{row}^E(x) \neq \text{row}^E(y)$  where  $\text{row}^E(u) = \{(z, T(uz)) : z \in E\}$ .

Make a transition from  $x \in S$  on symbol  $a \in \Sigma$  to the  $y \in S$  with  $\text{row}^E(y) = \text{row}^E(xa)$ .

Maintain for this a table  $T$  of observations such that for all  $w \in S \cdot E \cup S \cdot \Sigma \cdot E$ , the entry  $T(w)$  is equal to  $L(w)$ .

$\text{DFA}(S, E, T)$  describes the DFA given by  $S, \Sigma, E, T$  as above.



# The Algorithm

Angluin's Algorithm (Slightly Modified) to learn unknown regular language  $L$  using MQs and EQs.

1. Let  $S = \{\varepsilon\}$  and  $E = \{\varepsilon\}$ ;
2. Make MQs to determine  $L(w)$  for all  $w \in S \cdot (\Sigma \cup \{\varepsilon\}) \cdot E$  and let  $T(w) = L(w)$  (where not done before);
3. Let  $\text{row}^E(u) = \{(z, T(uz)) : z \in E\}$  for all  $u \in S \cdot (\Sigma \cup \{\varepsilon\})$  and search for an  $u \in S$  and  $a \in \Sigma$  with  $\text{row}^E(ua) \neq \text{row}^E(u')$  for all  $u' \in S$ ;  
if found then let  $S = S \cup \{ua\}$  and go to 2 else go to 4;
4. Make EQ whether  $\text{DFA}(S, E, T)$  is correct;
5. If the answer is "YES" then terminate;
6. If the answer is "NO" with counter example  $w$  then let  $E = E \cup \{v : \exists u [w = uv]\}$  and go to 2.

# Verification

Let  $(\mathbf{Q}, \Sigma, \delta, \mathbf{s}, \mathbf{F})$  be the smallest dfa for  $\mathbf{L}$ .

(A) If  $\mathbf{x}, \mathbf{y} \in \mathbf{S}$  are distinct then  $\delta(\mathbf{s}, \mathbf{x}) \neq \delta(\mathbf{s}, \mathbf{y})$ .

(B) If an equivalence query is made then a new state is added eventually.

(C) Angluin's learner uses at most  $|\mathbf{Q}|$  equivalence queries.

(D) When reaching line 4,  $\mathbf{DFA}(\mathbf{S}, \mathbf{E}, \mathbf{T})$  is complete, that is, for every  $\mathbf{u} \in \mathbf{S}$  and  $\mathbf{a} \in \Sigma$  there is a  $\mathbf{u}' \in \mathbf{S}$  with  $\mathbf{row}^{\mathbf{E}}(\mathbf{u}') = \mathbf{row}^{\mathbf{E}}(\mathbf{ua})$ .

(E) If  $\mathbf{z} \in \mathbf{E}$  then  $\mathbf{DFA}(\mathbf{S}, \mathbf{E}, \mathbf{T})$  accepts  $\mathbf{z}$  iff  $\mathbf{T}(\mathbf{z}) = \mathbf{1}$ .

(F) The number of membership queries is bounded by  $|\mathbf{Q}| \cdot (\mathbf{r} + \mathbf{1}) \cdot (|\Sigma| + \mathbf{1})$  where  $\mathbf{r}$  is the sum of the lengths of all counter examples.

# Original Algorithm

Teacher has regular set  $L$  and learner makes membership and equivalence queries.

1. Initialise  $S = \{\varepsilon\}$  and  $E = \{\varepsilon\}$ .
2. For all  $w \in S \cdot E \cup S \cdot \Sigma \cdot E$  where  $T(w)$  is undefined, make MQ to find  $L(w)$  and let  $T(w) = L(w)$ .
3. If there are  $u, u' \in S$ ,  $a \in \Sigma$  and  $v \in E$  such that  $\text{row}^E(u) = \text{row}^E(u')$  and  $T(uav) \neq T(u'av)$  then let  $E = E \cup \{av\}$  and go to 2.
4. If there are  $u \in S$  and  $a \in \Sigma$  with  $\text{row}^E(ua) \neq \text{row}^E(u')$  for all  $u' \in S$  then let  $S = S \cup \{ua\}$  and go to 2.
5. Make an EQ whether  $\text{DFA}(S, E, T)$  recognises  $L$ .
6. If the answer is “YES” then terminate with  $\text{DFA}(S, E, T)$ .
7. If the answer is “NO” with counterexample  $w$  then let  $S = S \cup \{u : \exists v [uv = w]\}$  and go to 2.

# Learning from Positive Data

A text  $\mathbf{T}$  for  $\mathbf{L}$  is an infinite sequence  $\mathbf{T}(0), \mathbf{T}(1), \dots$  of strings and pause symbols  $\#$  describing  $\mathbf{L}$ :  $w \in \mathbf{L}$  iff  $w \neq \#$  and  $w = \mathbf{T}(n)$  for some  $n$ .

Learner  $\mathbf{M}$  starts with initial memory  $\mathbf{mem}_0$  and hypothesis  $\mathbf{e}_0$ . In the  $n$ -th cycle,  $\mathbf{M}$  updates  $\mathbf{mem}_n$  and input  $\mathbf{T}(n)$  to  $\mathbf{mem}_{n+1}$  and  $\mathbf{e}_{n+1}$ ; in short:  $\mathbf{M}$  maps  $(\mathbf{mem}_n, \mathbf{T}(n))$  to  $(\mathbf{mem}_{n+1}, \mathbf{e}_{n+1})$ . The values  $\mathbf{e}_1, \mathbf{e}_2, \dots$  depend on  $\mathbf{T}$ .

Learner  $\mathbf{M}$  learns  $\mathbf{L}$  iff, on every text  $\mathbf{T}$  for  $\mathbf{L}$ , almost all hypotheses  $\mathbf{e}_n$  of  $\mathbf{M}$  are the same  $\mathbf{e}$  describing  $\mathbf{L}$ .

# Example

Let  $I = \{0\}^*$  and  $L_e = \{x \in \{0, 1\}^* : |x| < |e|\}$ . For example,  $L_{00} = \{\varepsilon, 0, 1\}$  and  $L_{000} = \{\varepsilon, 0, 1, 00, 01, 10, 11\}$ .

The class  $\{L_e : e \in I\}$  can be learnt from positive data.

$\text{mem}_n$  is the longest datum seen before round  $n$  and  $\text{mem}_n = \#$  if no datum has been seen before round  $n$ .

If  $\text{mem}_n = \#$  then  $e_n = \varepsilon$  else  $e_n = 0^{k+1}$  for the length  $k$  of  $\text{mem}_n$ .

When learning  $L_{0^k}$  then after some time some string of length  $k - 1$  has shown up in the input and from then onwards,  $\text{mem}_n$  is that string and  $e_n = 0^k$ . So the learner eventually converges to the right hypothesis.

Next slide: Hypotheses made for data from sample text  $\#, \#, \varepsilon, 11, 01, 0, 1, \dots$  and the corresponding memory.

# Sample Text

$n$	first $n$ members of $T$	$mem_n$	$e_n$	$L_{e_n}$
0	—	#	$\varepsilon$	$\emptyset$
1	#	#	$\varepsilon$	$\emptyset$
2	#, #	#	$\varepsilon$	$\emptyset$
3	#, #, $\varepsilon$	$\varepsilon$	0	$\{\varepsilon\}$
4	#, #, $\varepsilon$ , 11	11	000	$\{\varepsilon, 0, 1, 00, 01, 10, 11\}$
5	#, #, $\varepsilon$ , 11, 01	11	000	$\{\varepsilon, 0, 1, 00, 01, 10, 11\}$
6	#, #, $\varepsilon$ , 11, 01, 0	11	000	$\{\varepsilon, 0, 1, 00, 01, 10, 11\}$
7	#, #, $\varepsilon$ , 11, 01, 0, 1	11	000	$\{\varepsilon, 0, 1, 00, 01, 10, 11\}$
8	#, #, $\varepsilon$ , 11, 01, 0, 1, 101	101	0000	$\{\varepsilon, 0, 1, \dots, 111\}$

If all data is from  $L_{0000}$  then the learner keeps the hypothesis  $0000$  forever.

# Learnability of Classes

**Theorem** [Gold 1967]

The class of all regular languages cannot be learnt from positive data.

**Theorem** [Angluin 1980]

An automatic family  $\{\mathbf{L}_e : e \in \mathbf{I}\}$  can be learnt from positive data iff there is for every  $e \in \mathbf{I}$  a finite subset  $\mathbf{F}_e \subseteq \mathbf{L}_e$  such that there is no  $d \in \mathbf{I}$  with  $\mathbf{F}_e \subseteq \mathbf{L}_d \subset \mathbf{L}_e$ .

**Example**

The class of  $\mathbf{L}_\varepsilon = \Sigma^*$  and  $\mathbf{L}_x = \{y \in \Sigma^* : |y| < |x|\}$  with  $x \in \Sigma^+$  does not satisfy Angluin's tell-tale condition; hence this class is not learnable.

# Learner for Tell-Tale Condition

Let  $\{L_e : e \in I\}$  be a given automatic class.

Let  $F_e$  be the tell-tale of  $L_e$ .

Now  $\text{mem}_n$  is the set of data observed before step  $n$  and updated by  $\text{mem}_{n+1} = \text{mem}_n \cup \{T(n)\}$ .

Furthermore,  $e_n$  is the length-lexicographically least  $e \in I$  with  $F_e \subseteq \text{mem}_n \subseteq L_e$ ; if such an  $e$  does not exist then  $e_n$  is a default index of  $\emptyset$ .

When learning  $L_e$ , after sufficient time the learner has seen all data of  $F_e$  and for each  $d <_{\parallel} e$ , either  $F_d \not\subseteq L_e$  or some datum in  $L_e - L_d$  has been observed.

Hence the learner will converge to  $e$ .



# Necessity of Tell-Tale Condition

Let  $\{L_e : e \in I\}$  be a given automatic class and assume that there is an  $e \in I$  such that for  $L_e$  there does not exist a finite tell-tale set.

Assume that  $M$  is a learner for the class.

Goal: Make a text  $T$  for  $L_e$  on which  $M$  outputs infinitely often a wrong conjecture.

Idea is to construct the text  $T$  inductively, starting with  $n = 0$  and alternating between 1 and 2:

1. Let  $T(n) = \min_{\Pi}(L_e - \{T(m) : m < n\})$ ,  $n = n + 1$ ;
2. Choose some  $L_d$  with  $\{T(m) : m < n\} \subseteq L_d \subset L_e$  and keep updating  $T(n) = \min_{\Pi}(L_d - \{T(m) : m < n\})$ ,  $n = n + 1$  until the conjecture  $e_n$  of  $M$  equals  $d$ .

# Automatic Learners

Automatic learner given by  $\text{mem}_0$ ,  $e_0$  and  $\text{uf} : \text{conv}(\text{mem}_n, \mathbf{x}_n) \mapsto \text{conv}(\text{mem}_{n+1}, e_{n+1})$ .  
 $\text{mem}_n$  is a string in an arbitrary alphabet;  
 $e_n$  is a hypothesis from  $\mathbf{I}$  when learning  $\{\mathbf{L}_e : e \in \mathbf{I}\}$ ;  
 $\text{uf}$  is an automatic update function of the learner.

Automatic learners can only remember few information about the past, as there is a constant  $c$  with  $|\text{mem}_{n+1}| \leq \max\{|\text{mem}_n|, |\mathbf{x}_n|\} + c$  for all  $n$ .

Explicit memory bounds (with constant  $k$ ) when learning  $\mathbf{L}_d$ :

- Word-sized:  $|\text{mem}_{n+1}| \leq \max\{|\mathbf{x}_0|, |\mathbf{x}_1|, \dots, |\mathbf{x}_n|\} + k$ .
- Hypothesis-sized:  $|\text{mem}_{n+1}| \leq |e_{n+1}| + k$ .
- Target-sized:  $|\text{mem}_{n+1}| \leq |d| + k$ .
- Constant:  $|\text{mem}_{n+1}| \leq k$ .

# Example 11.6

Assume  $L_0 = \emptyset$ ,  $L_1 = \{0\}^*$ ,  $L_2 = \{1\}^*$  and  $L_3 = \{0, 1\}^*$ .

This finite class can be learnt by tracking whether **0** and **1** have shown up in the input:

If neither **0** nor **1** are observed:  $L_0$ ;

If **0** but not **1** is observed:  $L_1$ ;

If **1** but not **0** is observed:  $L_2$ ;

If both **0** and **1** are observed:  $L_3$ .

This principle generalises to any finite class and learner needs only constant memory.

Every finite class of sets has an automatic learner.

# Example 11.7

Assume  $\Sigma = \{0, 1, 2\}$  and

$I = \{\text{conv}(v, w) : v, w \in \Sigma^* \wedge v \leq_{\text{lex}} w\} \cup \{\text{conv}(3, 3)\}$  with

$L_{\text{conv}(v,w)} = \{u \in \Sigma^* : v \leq_{\text{lex}} u \leq_{\text{lex}} w\}$  for all

$\text{conv}(v, w) \in I$ ; note that  $L_{\text{conv}(3,3)} = \emptyset$ .

Data seen so far	Hypothesis	Conjectured language
—	$\text{conv}(3,3)$	$\emptyset$
#	$\text{conv}(3,3)$	$\emptyset$
# 00	$\text{conv}(00,00)$	$\{00\}$
# 00 0000	$\text{conv}(00,0000)$	$\{00, 000, 0000\}$
# 00 0000 1	$\text{conv}(00,1)$	$\{u : 00 \leq_{\text{lex}} u \leq_{\text{lex}} 1\}$
# 00 0000 1 0	$\text{conv}(0,1)$	$\{u : 0 \leq_{\text{lex}} u \leq_{\text{lex}} 1\}$
# 00 0000 1 0 112	$\text{conv}(0,112)$	$\{u : 0 \leq_{\text{lex}} u \leq_{\text{lex}} 112\}$
# 00 0000 1 0 112 #	$\text{conv}(0,112)$	$\{u : 0 \leq_{\text{lex}} u \leq_{\text{lex}} 112\}$

# Exercises 11.8 and 11.9

## Exercise 11.8

Make an automatic learner which learns the class of all  $L_d = \{dw : w \in \Sigma^*\}$  with  $d \in \Sigma^*$ ; that is,  $I = \Sigma^*$  in this case.

## Exercise 11.9

Let  $\{L_d : d \in I\}$  be given with  $L_d \neq L_{d'}$  whenever  $d, d' \in I$  are different. Assume that an automatic learner uses this class as a hypothesis space for learning satisfying any of the above memory constraints. Let  $\{H_e : e \in J\}$  be any other automatic family containing  $\{L_d : d \in I\}$  as a subclass. Show that there is an automatic learner satisfying the same type of memory constraints conjecturing indices taken from  $J$  in place of  $I$ .

# Separation

## Theorem 11.10

Let  $I = \Sigma^*$ ,  $L_\varepsilon = \Sigma^+$  and  $L_d = \{w \in \Sigma^* : w <_{\parallel} d\}$  for  $d \in \Sigma^+$ . The class  $\{L_d : d \in I\}$  can be learnt using a word-sized memory but not using an hypothesis-sized memory.

## Word-sized Learner

Memory are length-lexicographically smallest and largest words seen so far.

If  $\varepsilon$  has not been seen, conjecture  $\Sigma^+$ . If  $\varepsilon$  and  $w$  are minimal and maximal word seen so far then let  $d$  be the successor of  $w$  and conjecture  $L_d$ .

# Theorems 11.11 and 11.12

## Theorem 11.11

There is a class which can be learnt with hypothesis-sized memory but not with target-sized memory and not with constant memory.

## Theorem 11.12

If a learner learns a class with target-sized memory then the learner's memory is also word-sized on texts for languages in the class.

In particular, if a class has an automatic learner with a target-sized memory-limitation then it also has an automatic learner with a word-sized memory-limitation.

# Exercise 11.14

Assume that  $\{L_d : d \in I\}$  is the class to be learnt and that every language in the class is finite and that for every language in the class there is exactly one index in  $I$ .

Show that if there is a learner using word-sized memory for this class, then the memory of the same learner is also target-sized. For this, show that there is a constant  $k$  such that all  $d \in I$  and  $x \in L_d$  satisfy  $|x| \leq |d| + k$  and then deduce the full result.



# Iterative Learning

A learner is iterative iff  $\text{mem}_n = e_n$ , that is, the memory is exactly the most recent hypothesis.

## Exercise 11.15

Show that there is an automatic family  $\{\mathbf{L}_d : d \in \mathbf{I}\}$  such that  $\mathbf{I}$  contains for each  $\mathbf{L}_d$  exactly one index and the  $\mathbf{L}_d$  are exactly the finite subsets of  $\{0\}^*$  with even cardinality.

Show that the class  $\{\mathbf{L}_d : d \in \mathbf{I}\}$  has an iterative learner using the given hypothesis space.

Is the same possible when the class consists of all subsets of  $\{0\}^*$  with 0 or 3 or 4 elements?

Note that an iterative learner which just conjectured an  $d \in \mathbf{I}$  must abstain from updating the hypothesis on any datum  $x \in \mathbf{L}_d$ .

# Automatic Families

Which of the following classes can be represented as automatic families? If they are automatic then provide the corresponding coding including the index set else write why they cannot be automatic families.

**Exercise 11.16:** The family of all finite subsets of  $\{0\}^* \cdot \{1\}^*$ .

**Exercise 11.17:** The family of an infinite regular set  $L$  and all subsets of up to 5 elements.

**Exercise 11.18:** The family of all sets of decimal numbers which, for some  $n > 0$ , contain exactly two digits each  $n$  times and all other digits 0 times.

# Learners

**Exercise 11.19:** Consider an automatic family  $\{L_e : e \in I\}$  which satisfies for each two distinct  $d, e \in I$  that either  $L_d \subset L_e$  or  $L_e \subset L_d$  and in which there is for each  $e$  a unique  $x_e$  such that  $x_e \in L_e$  but  $x_e \notin L_d$  for any  $d$  with  $L_d \subset L_e$ . Prove that the mapping  $e \mapsto x_e$  is automatic and construct an automatic learner for this family.

**Exercise 11.20:** Given  $\{L_d : d \in I\}$  and  $\{L_e : e \in J\}$  both satisfying the specification of Exercise 11.19, construct an automatic learner for the family of all  $K_{\text{conv}(d,e)}$  with  $d \in I, e \in J$  and  $K_{\text{conv}(d,e)} = \{0x : x \in L_d\} \cup \{1y : y \in L_e\}$ .

# Open Problems

1. Does every automatic family which has an automatic learner also have a learner with word-sized memory?
2. Does every automatic family which has a learner with hypothesis-sized memory also have a learner with word-sized memory?
3. Does every automatic family which has a learner with hypothesis-sized memory also have an iterative learner?