

THE NATIONAL UNIVERSITY
of SINGAPORE



School of Computing
Computing 1, 13 Computing Drive, Singapore 117417

TRB1/10

On Automatic Families

*Sanjay Jain, Yuh Shin Ong, Shi Pu and
Frank Stephan*

January 2010

Technical Report

Foreword

This technical report contains a research paper, development or tutorial article, which has been submitted for publication in a journal or for consideration by the commissioning organization. The report represents the ideas of its author, and should not be taken as the official views of the School or the University. Any discussion of the content of the report should be sent to the author, at the address shown on the cover.

OOI Beng Chin
Dean of School

On Automatic Families

Sanjay Jain, Yuh Shin Ong, Shi Pu and Frank Stephan

*S. Jain: School of Computing, National University of Singapore
Block COM1, 13 Computing Drive, Singapore 117417
Email: sanjay@comp.nus.edu.sg*

*Y. S. Ong: School of Computing, National University of Singapore
Block COM1, 13 Computing Drive, Singapore 117417
Email: yuhshin@comp.nus.edu.sg*

*S. Pu: School of Computing, National University of Singapore
Block COM1, 13 Computing Drive, Singapore 117417
Email: pushi@nus.edu.sg*

*F. Stephan: Department of Mathematics, National University of Singapore
Block S17, 10 Lower Kent Ridge Road, Singapore 119076
Email: fstephan@comp.nus.edu.sg*

This paper summarises previous work on automatic families. It then investigates a natural measure which exists inside every automatic family: the size of a regular language in this family is just the length of its index. This measure satisfies various properties similar to those of Kolmogorov complexity. Furthermore, there is a quite natural extension of it to all regular languages such that, for each automatic family, the new general measure differs from that of the least index only by up to a constant. However, this generalisation fails to have several properties similar to Kolmogorov complexity. Furthermore, a characterisation is given regarding when a class of languages is a subclass of an automatic family.

Sanjay Jain and Frank Stephan are supported in part by NUS grant numbers R252-000-308-112 and R146-000-114-112.

1. Introduction

Automatic families are uniformly regular families of sets which are mainly used in inductive inference (learning theory) as a way to represent hypothe-

sis spaces with decidable first-order theory. For such spaces, one can decide whether they are explanatorily learnable in the limit and have effective procedures to generate such learners. Furthermore, they have also been used to study various other learnability notions. Automatic families are closely related to automatic graphs and automatic structures. All these notions are based on the notion of a finite automaton. The present work formalises the notion of a size of a language with respect to an automatic class. It also relates this size to various measures of the size of a language derived from its minimal deterministic finite automaton. Furthermore, an overview of related properties and applications of automatic families is given.

The basic notion of the field is that of a finite automaton. Informally, a finite automaton could be considered as an algorithm which reads a word from the left to the right and has a constant-sized memory; when it reaches the right end of the word, the algorithm says either “accept” or “reject”. This informal description permits to write algorithms representing finite automata quite well; on the other hand, when one wants to prove properties of finite automata, it is better to write them as a finite set of states together with a state-transition relation, a starting state and a set of final states. The automaton then reads the symbols of the word to be checked from the left to the right and each time changes its state based on a transition-relation which is a set of triples of the form (old state, current symbol, new state). Ideally, for every old state and current symbol, there is a unique triple in the table; an automaton with this property is called a “deterministic finite automaton” (dfa) and every finite automaton can be brought into this form. Now, for studying automatic structures, the key generalisation is to permit automata to track various inputs simultaneously, provided that the reading of these inputs is synchronised. Informally, the algorithm consists of one big loop which first reads, from each of those inputs which are not already exhausted, one symbol and then updates its internal memory accordingly; the memory is restricted to a constant amount of bits. Here, the algorithm knows when an input is exhausted and when not, so this piece of information can be taken into account. When every input is exhausted, the algorithm says “accept” or “reject”. The fact that all inputs are read at the same speed is important, otherwise the model would become too powerful and certain undecidability problems of the theory of such models would arise. Mathematically, one can also map this back to the case of a single input which is formed as the convolution of the inputs. The convolution $conv(x, y)$ of two inputs x and y would consist of symbols which are pairs of the corresponding symbols from x and y ; in the case that one of

these input words is shorter than the other, it is brought up to the same length by appending the special symbol \diamond sufficiently often; the symbol \diamond is not contained in any of the input alphabets used for x and y . One can define *conv* on more than two arguments similarly. The next algorithm is an example which decides the lexicographic order of two strings x and y ; the algorithm uses an underlying ordering $<$ on Σ :

- (1) If the input x is exhausted then accept and terminate.
- (2) If the input y is exhausted then reject and terminate.
- (3) Read the current symbol a of the input x and the current symbol b of the input y .
- (4) If $a < b$ then accept and terminate.
- (5) If $b < a$ then reject and terminate.
- (6) Go to (1).

Note that this algorithm also accepts if both words are equal; so it accepts if $x \leq_{lex} y$ and rejects if $y <_{lex} x$. Furthermore, the algorithm produces an early decision, when possible. So when comparing 25880 with 2593 the algorithm would already make the decision after having read 258 from 25880 and 259 from 2593. Such an early decision is permitted for easier formulations of the algorithms, although formally the algorithm has to scan the full input.

Besides automatic relations like the lexicographic ordering, one can also define automatic functions where the automaton recognises a function f as follows: it reads convolution of x and y and accepts iff x is in the domain of f and $y = f(x)$. An automatic family is a specific type of structure. Here one has the domain D and an index domain I which are both regular sets; furthermore there is an automatic relation on $I \times D$ defining a family $\{L_i : i \in I\}$ such that $x \in L_i$ iff the underlying relation contains $conv(i, x)$. Such a family defines a class of subsets of D and two families have the same range iff they contain the same languages. In many cases it is also handy to have that the indexing is a one-one indexing, that is, $L_i \neq L_j$ for different $i, j \in I$. However, this property is not mandatory within the present work. An important result from the early days of automatic structures [8, 9, 15] is that whenever a relation or function is first-order definable using other automatic relations or functions then it is itself automatic. Furthermore, the first-order theory of automatic structures is decidable. These two results give this field some importance in model checking and similar applications and the two results are also frequently used when constructing learners in inductive inference. The next examples of automatic families illustrate the

concept.

Example 1.1: Let $\Sigma = \{0, 1, 2\}$ and $D = \Sigma^*$.

The family of all $L_i = \{x \in D : x \text{ extends } i\}$ with $i \in \Sigma^*$ is an automatic family; an example member of it is $L_{001} = \{001, 0010, 0011, 0012, 00100, 00101, \dots, 00122, 001000, 001001, \dots\}$. The automaton recognising the family accepts $\text{conv}(i, x)$ iff every symbol of i appears at the same position in x .

The family $L_{\text{conv}(i,j)} = \{x : i <_{lex} x <_{lex} j\}$ of all $i, j \in \Sigma^*$ with $i <_{lex} j$ is also automatic; note that in some special cases like $j = i0$ the language $L_{\text{conv}(i,j)}$ is empty. The automaton recognising the family is mainly checking whether x is between the two boundaries i and j which are coded up in the index $\text{conv}(i, j)$.

The family of all L_{i3a3b} with $i3a3b \in I = \Sigma^* \cdot \{3\} \cdot \Sigma \cdot \{3\} \cdot \Sigma$ and $L_{i3a3b} = \{ixaybz : x, y, z \in \Sigma^*\}$ is automatic. For example, $L_{00123231}$ is given by the regular expression $0012 \cdot (0+1+2)^* \cdot 2 \cdot (0+1+2)^* \cdot 1 \cdot (0+1+2)^*$. The automaton recognising the automatic family has to memorise the content of a and b when reaching the corresponding position in the input $\text{conv}(i3a3b, u)$ as a and b can occur in u much later than in $i3a3b$. This is possible as a finite automaton can store constant amount of memory.

Example 1.2: Using a certain regular domain D , one can also code the natural numbers with addition, the relation $<$ and a predicate F recognising the Fibonacci numbers [25]. If now $\phi(x, i, j, k)$ is a first-order formula with four inputs defined using $+$, $<$, F and natural numbers then the family of all $L_{\text{conv}(i,j,k)} = \{x \in D : \phi(x, i, j, k) \text{ is true}\}$ is an automatic family where the index-domain is the set of all convolutions of three elements of D . An example for such a formula ϕ is $\phi(x, i, j, k) \Leftrightarrow \exists y \exists z [i < x + y < j \wedge x + x = z + z + z + 2 \wedge F(x + y + y + k)]$.

Note that the above example is more in the traditional style of automatic structures where all aspects of coding can be freely chosen in order to meet the specification. Automatic families usually are a bit more fixed; here one wants to find an automatic indexing of a given class of regular languages; that is, while the indexing is considered to be “chosen”, the languages inside the class and the domain D are more considered as “given”. The results in the next section will establish various facts on the indexings which show that the indexings are not completely free, but some aspects of them are determined by the class which has to be represented by the automatic family. See [4, 5, 8, 9, 15, 22, 23, 25] for more information on automatic structures.

2. The size of languages inside a family

Having the concept of an automatic family, one can use its indexing in order to introduce a measure for the size of the languages inside the given family.

Definition 2.1: Given an automatic family \mathcal{L} and a language $R \in \mathcal{L}$, let $d_{\mathcal{L}}(R) = \min\{|i| : i \in I \wedge L_i = R\}$ be the size of R .

The next result shows that the size depends only up to an additive constant on the chosen automatic family; so enlarging the underlying family or just changing its indexing has not much impact on the size of a languages inside the family.

Proposition 2.2: Let $\mathcal{L} = \{L_i : i \in I\}$ and $\mathcal{H} = \{H_j : j \in J\}$ be two automatic families. Then there is a constant c such that $d_{\mathcal{L}}(R)$ and $d_{\mathcal{H}}(R)$ differ at most by c for all $R \in \mathcal{L} \cap \mathcal{H}$.

Proof: The basic idea is to look at the set $O = \{\text{conv}(i, j) : i \in I \text{ and } j \in J \text{ and } L_i = H_j \text{ and } i, j \text{ are the length-lexicographically first indices of their respective sets}\}$. Note that the length-lexicographic order $<_l$ is automatic. Now O is first-order definable:

$$\begin{aligned} \text{conv}(i, j) \in O \Leftrightarrow & \quad \forall x \in D [x \in L_i \Leftrightarrow x \in H_j] \\ & \wedge \forall i' [L_{i'} = L_i \Rightarrow i \leq_l i'] \\ & \wedge \forall j' [H_{j'} = H_j \Rightarrow j \leq_l j']. \end{aligned}$$

By a result of Khoussainov and Nerode [15] the set O is regular. Note that any two members $\text{conv}(i, j), \text{conv}(i', j') \in O$ satisfy $(i = i' \vee j = j') \Rightarrow (i = i' \wedge j = j')$. Thus, due to the pumping lemma, there is a constant c such that, for every pair $\text{conv}(i, j) \in O$, the length of i and j differ by at most the constant c . \square

This result is a bit parallel to the corresponding result in the field of Kolmogorov complexity [18] that the Kolmogorov complexity of an object depends only up to a constant on the underlying universal machine. The main difference is that here the measures $d_{\mathcal{L}}$ are only defined on a subfamily \mathcal{L} of the regular languages and not on all of them; this invokes some problems and in the following it is investigated to which degree one can overcome these problems. Before doing this in the next sections, first a striking parallel to Kolmogorov complexity is pointed out: Boolean operations and images of sets under functions essentially have the complexity of the input sets.

Theorem 2.3: *For every automatic family \mathcal{L} (with nonempty indices) and every automatic function f there is a further automatic family \mathcal{H} such that for all $L_i, L_j \in \mathcal{L}$ the following statements hold:*

- $L_i \cup L_j \in \mathcal{H}$ and $d_{\mathcal{H}}(L_i \cup L_j) \leq \max\{d_{\mathcal{L}}(L_i), d_{\mathcal{L}}(L_j)\}$,
- $L_i \cap L_j \in \mathcal{H}$ and $d_{\mathcal{H}}(L_i \cap L_j) \leq \max\{d_{\mathcal{L}}(L_i), d_{\mathcal{L}}(L_j)\}$,
- $D - L_i \in \mathcal{H}$ and $d_{\mathcal{H}}(D - L_i) \leq d_{\mathcal{L}}(L_i)$,
- $\{f(x) : x \in L_i\} \in \mathcal{H}$ and $d_{\mathcal{H}}(\{f(x) : x \in L_i\}) \leq d_{\mathcal{L}}(L_i)$.

Proof: Let \mathcal{H} contain the following sets for $i, j \in I$:

- $H_{conv(0,i,j)} = L_i \cup L_j$,
- $H_{conv(1,i,j)} = L_i \cap L_j$,
- $H_{conv(2,i,i)} = D - L_i$ and
- $H_{conv(3,i,i)} = \{f(x) : x \in L_i\}$.

Then one defines that $J = \{conv(0, i, j), conv(1, i, j), conv(2, i, i), conv(3, i, i) : i, j \in I\}$ and $\mathcal{H} = \{H_k : k \in J\}$. In order to see that \mathcal{H} is an automatic family, note that one can construct finite automata A_0, A_1, A_2, A_3 such that A_0 accepts $conv(conv(0, i, j), x)$ iff $i, j \in I$ and $x \in L_i \cup L_j$; furthermore, A_1 accepts $conv(conv(1, i, j), x)$ iff $i, j \in I$ and $x \in L_i \cap L_j$; A_2 accepts $conv(conv(2, i, j), x)$ iff $i, j \in I$, $i = j$ and $x \in D - L_i$; A_3 accepts $conv(conv(3, i, j), x)$ iff $i, j \in I$, $i = j$ and there is an $y \in L_i$ such that $f(y) = x$. Note that the automata A_0, A_1, A_2, A_3 exist as the relations are first-order defined from other automatic relations. Now there is an automaton which accepts $conv(k, x)$ whenever one of the automata A_0, A_1, A_2, A_3 accepts $conv(k, x)$. Hence the family \mathcal{H} is automatic.

Using that no string in I is empty, one sees that the convolutions $conv(0, i, j)$ and $conv(1, i, j)$ have at most the length $\max\{|i|, |j|\}$ and that $conv(2, i, i), conv(3, i, i)$ have length i . Thus the automatic family \mathcal{H} covers the desired sets and satisfies the conditions on the sizes of the indices. \square

Note that in Theorem 2.3 one can just add 1 to each bound on the right side in order to get rid of the “nonempty indices” in the statement of the theorem. One cannot get rid of the change from \mathcal{L} to \mathcal{H} as one cannot assume that \mathcal{L} is closed under Boolean operations. Theorem 3.7 below formulates the results using a common measure for the size of the sets and of their Boolean combinations.

3. Universal Complexity Measures

In the following, let A_R be the smallest deterministic finite automaton accepting R ; A_R has to be complete, that is, for every state p and every symbol $a \in \Sigma$ there is exactly one state q such that A goes from p to q on input a . Let $d_{dfa}(R)$ denote the number of states of A_R and $d_{run}(R)$ denote the maximum n such that, for some input word x , A_R on x goes through n different states. Note that $d_{run}(R) \leq d_{dfa}(R)$. The next result establishes an inequality for the converse direction. This inequality witnesses that, for each n , there are only finitely many R with $d_{run}(R) = n$. Thus d_{run} satisfies some minimum requirement for measuring the size of a language adequately.

Proposition 3.1: *Let $R \subseteq \Sigma^*$ be a regular language. If Σ has at least two members then*

$$d_{dfa}(R) \leq \frac{|\Sigma|^{d_{run}(R)} - 1}{|\Sigma| - 1}$$

else $d_{dfa}(R) = d_{run}(R)$.

Proof: Let A_R be the minimal deterministic finite automaton recognising R . One can look at the finite tree T of all runs of A_R in which no state is visited twice. The height of this tree T is at most $d_{run}(R) - 1$. Furthermore, every state in A occurs in this tree T as it can be reached by a repetition-free run. In the case that Σ has exactly one element, T has $d_{run}(R)$ members and $d_{dfa}(R) = d_{run}(R)$. In the case that Σ has at least two members, the formula

$$d_{dfa}(R) \leq |T| \leq |\Sigma|^0 + |\Sigma|^1 + \dots + |\Sigma|^{d_{run}(R)-1} = \frac{|\Sigma|^{d_{run}(R)} - 1}{|\Sigma| - 1}$$

provides an upper bound on the number of members of T and thus on the value $d_{dfa}(R)$. \square

Remark 3.2: The exponential bound in the case of the alphabet Σ having at least two symbols looks large, but the gap cannot be made much smaller. For example, if $L_n = \{xx : x \in \Sigma^n\}$ then $d_{run}(L_n) = 2n + 1$ while $d_{dfa}(L_n) \geq |\Sigma|^n$, as for every $x \in \Sigma^n$ the derivative $L_n[x] = \{x\}$ is encoding x .

The following connections hold between the size based on automatic families and these two measures.

Theorem 3.3: *For every automatic family \mathcal{L} there is a constant c such that $d_{dfa}(R) \leq d_{\mathcal{L}}(R) \cdot c + 1$ for all $R \in \mathcal{L}$.*

Proof: Let A be the automaton recognising the automatic class, that is, A accepts $\text{conv}(i, x)$ whenever $i \in I$ and $x \in L_i$. Let c be the number of states of A . Let n be the length of the index i of some L_i . Now one constructs an automaton $\text{Comb}(A, i)$ which recognises the language L_i . The states of $\text{Comb}(A, i)$ are the disjoint union of sets Q_0, Q_1, \dots, Q_n , where Q_1, Q_2, \dots, Q_n each consist of a copy of the states of A and Q_0 consists of a copy of the starting state of A ; the unique state in Q_0 is the starting state of $\text{Comb}(A, i)$. For a state p in $\text{Comb}(A, i)$ let p' denote the state in A from which it is copied. A state $p \in Q_m$ is accepting iff there is a word $x \in L_i$ such that A after processing $|x|$ symbols of $\text{conv}(i, x)$ is in the state p' , $m < n \Rightarrow |x| = m$ and $m = n \Rightarrow |x| \geq m$. There are two types of transitions: For $m < n$ there is a transition from $p \in Q_m$ to $q \in Q_{m+1}$ on input $a \in \Sigma$ iff A goes from p' to q' on input $\text{conv}(b, a)$ where b is the symbol of i at position m . Furthermore, for $p, q \in Q_n$ and $a \in \Sigma$ there is a transition on $\text{Comb}(A, i)$ from p to q on input a iff there is a transition from p' to q' in A on input $\text{conv}(\diamond, a)$.

One can summarise the construction as follows: the automaton $\text{Comb}(A, i)$ simulates on input x the behaviour of A on the input $\text{conv}(i, x)$. The verification that this simulation is indeed doing what it should is left to the reader. Note that $\text{Comb}(A, i)$ has $c \cdot n + 1$ states and thus satisfies the required bound. \square

Remark 3.4: This multiplicative constant c is indeed needed: If \mathcal{L} is the class of all finite languages consisting of up to c strings, then each automaton accepting the language $L_n = \{0^d 1^n 0^d : d < c\}$ needs at least $c \cdot n$ states in order to memorise the number of 0s and then to count the number of 1s before comparing the number of 0s after the block of 1s with the memorised value. Also, one can easily verify that — up to an additive constant — $d_{\mathcal{L}}(L_n) = n$.

The next result is the main contribution of this paper. It shows that d_{run} is a measure which meets the expectation in at least one point: given any automatic class \mathcal{L} , the measures $d_{\mathcal{L}}$ and d_{run} coincide on \mathcal{L} up to a constant.

Theorem 3.5: *For every automatic family \mathcal{L} there is a constant c such that, for all $R \in \mathcal{L}$, the values $d_{\text{run}}(R)$ and $d_{\mathcal{L}}(R)$ differ from each other by at most c .*

Proof: Let \mathcal{L} be the given automatic family and let c be the number of states of the automaton A recognising the family. Let $R \in \mathcal{L}$.

Now it is shown that $d_{run}(R) \leq d_{\mathcal{L}}(R) + c$. For every index $i \in I$, one constructs the automaton $Comb(A, i)$ as described in Theorem 3.3. In any run, this automaton passes through at most $|i| + c$ states.

Now it is shown that $d_{\mathcal{L}}(R)$ is bounded by $d_{run}(R)$ plus a constant independent of R . For $i \in I$, let $L_i[x] = \{y : xy \in L_i\}$ be the left derivative of the language L_i . Now, if $|x| \geq |i|$ then one can produce the following automaton B accepting the left derivative $L_i[x]$:

- The set of states of B equals the set of states of A ;
- The starting state of B is the state of A after having read the first $|x|$ symbols of $conv(i, x)$;
- The state transition of B from p to q on a symbol a occurs iff A goes on the one-symbol word $conv(\diamond, a)$ from p to q ;
- The accepting states of B and A are the same.

Hence the language $L_i[x]$ is recognised by an automaton containing only c states. Let B_q be the automaton B with the starting state q . Now let n be the least positive natural number such that, for every $x \in \{0, 1\}^n$, $L_i[x]$ is recognised by some automaton B_q . Note that the choice of q only depends on the state in which A is after having read first $|x|$ symbols of $conv(i, x)$. Now let t code the table which maps each of the states p of A reached after processing the first n symbols of $conv(i, x)$ to the state q identifying the automaton B_q which recognises $L_i[x]$. There are only c^c such possible tables; hence one can use a proper alphabet to encode the possible entries of t .

In short: j is the shortest index of the form $conv(k, t)$ such that k is the shortest possible nonempty prefix of i for which there are c finite automata B_1, B_2, \dots, B_c , each of c states, such that, for all $x \in \{0, 1\}^{|k|}$, B_b recognises $L_i[x]$ if A is in the state b after processing $conv(k, x)$. The j is a function of i and $H_j = L_i$. Let J be the set of such j formed using indices $i \in I$. One can verify that the family $\{H_j : j \in J\}$ is automatic and J is regular.

Now let $i \in I$ and $j \in J$ be given with $L_i = H_j$. It follows, using the pumping lemma, that every word in L_i of length $d_{run}(L_i)$ is of the form uvw such that $L_i[uv^{|i|}w] = L_i[uvw]$. Hence $L_i[uvw]$ is recognised by one of the automata B_q and so $|j| \leq d_{run}(L_i)$. Thus it holds, for all $R \in \mathcal{L}$, that $d_{run}(R) \geq d_{\mathcal{H}}(R)$. Therefore $d_{run}(R)$ differs from $d_{\mathcal{L}}(R)$ only by a constant independent of R . \square

Corollary 3.6: *For every automatic family \mathcal{L} there is a constant c such that $d_{\mathcal{L}}(R) \leq d_{dfa}(R) + c$ for all $R \in \mathcal{L}$.*

The counterpart to Theorem 2.3 is the following.

Theorem 3.7: *For every automatic family \mathcal{L} and every automatic function f there is a constant c such that, for all $L, H \in \mathcal{L}$, it holds that*

- $d_{run}(L \cup H) \leq \max\{d_{run}(L), d_{run}(H)\} + c$,
- $d_{run}(L \cap H) \leq \max\{d_{run}(L), d_{run}(H)\} + c$,
- $d_{run}(D - L) \leq d_{run}(L) + c$ and
- $d_{run}(\{f(x) : x \in L\}) \leq d_{run}(L) + c$.

However the first, second and fourth condition of this list do not hold without the restriction to an automatic family.

Proof: The main part of this result follows from the Proposition 2.2 and Theorems 2.3 and 3.5. So the rest of the proof is to show that these connections do not hold in general, except for the third condition which is just obtained by interchanging acceptance and rejection inside D .

Let $n \in \{1, 2, 3, \dots\}$. For the first two conditions one considers the languages $(\Sigma^n)^*$ and $(\Sigma^{n+1})^*$. While $d_{run}((\Sigma^n)^*) = n$ and $d_{run}((\Sigma^{n+1})^*) = n + 1$, it holds that $d_{run}((\Sigma^n)^* \cup (\Sigma^{n+1})^*) = n(n + 1)$ and $d_{run}((\Sigma^n)^* \cap (\Sigma^{n+1})^*) = n(n + 1)$. For the fourth condition, assume $\Sigma = \{0, 1, 2\}$ and consider the automatic function f which interchanges 1 and 2 at every second occurrence of one of these digits. So $f(001001001001) = 001002001002$ and $f(12121212221212121) = 1111111122222222$. Let $L = (0^n 1)^*$. Then $\{f(x) : x \in L\} = (0^n 10^n 2)^* + (0^n 10^n 2)^* \cdot 0^n 1$. While $d_{run}(L) = n + 2$ it holds that $d_{run}(\{f(x) : x \in L\}) = 2n + 3$. \square

Remark 3.8: Besides d_{run} , one could also look at the following measure: $d_{rf}(R)$ is the largest number n such that there is an input $x \in \Sigma^*$ on which the minimal automaton A_R for R goes through exactly n states without repeating any of them.

Note that $d_{rf}(R) \leq d_{run}(R) \leq d_{dfa}(R)$ for every regular language R and the proof of Proposition 3.1 gives directly that

$$d_{dfa}(R) \leq \frac{|\Sigma|^{d_{rf}(R)} - 1}{|\Sigma| - 1}$$

for the case that Σ has at least two elements. Furthermore, Theorem 3.5 holds also with d_{rf} in place of d_{run} .

Although there are many parallel results for these two measures, d_{run} and d_{rf} are not identical. Consider the language $R = \{0^n, 1^n\}^*$ with $n \geq 2$. The minimal automaton R_A has the states reached by 0^m with $m < n$, 1^m

with $m < n$ plus one rejecting state which is never left; as the initial state is double counted in this list, there are in total $2n$ states and $d_{dfa}(R) = 2n$. On the word $0^n 1^n 01$ all states are visited and therefore $d_{run}(R) = 2n$. However $d_{rf}(R) = n + 1$ and this maximum number is taken on the input $0^{n-1}1$.

4. Characterising Automatic Families

The central question of this section is: When is a class \mathcal{L} of regular languages a subclass of some automatic family. The answer is that every language in the class must be representable by an automaton of a specific form.

Theorem 4.1: *A class \mathcal{L} is a subclass of an automatic family iff there is a constant c such that every $R \in \mathcal{L}$ is accepted by a deterministic finite automaton whose states can be partitioned into sets $Q_0, Q_1, Q_2, \dots, Q_n$ satisfying the following conditions: each set Q_m has up to c states; Q_0 consists exactly of the starting state; if there is a transition from a state p to a state q , then there are r, r' with $p \in Q_r, q \in Q_{r'}$ and $r' = \min\{r + 1, n\}$.*

Proof: Given an automatic class \mathcal{L} with an automaton A recognising the class, one can construct, for each $i \in I$, the automaton $Comb(A, i)$ to recognise L_i as in Theorem 3.3; it is easy to see that the automaton is of the above form.

For the other way round, consider the class \mathcal{H} of all languages which are accepted by an automaton of the above form with a given fixed constant c . The indices j of \mathcal{H} consist of symbols which code up Q_0, Q_1, \dots, Q_n . For each Q_m it is coded which of the states of Q_m (numbered as $1, 2, \dots, c$) are accepting and what transitions are there from Q_m to $Q_{\min\{m+1, n\}}$ based on various inputs from Σ . Without loss of generality state 1 from Q_0 is the starting state and that does not need to be coded. Now $J = \Gamma^*$ where

$$\Gamma = \{\text{reject}, \text{accept}\}^c \times \{1, 2, \dots, c\}^{\{1, 2, \dots, c\} \times \Sigma},$$

that is, where each symbol in Γ codes the acceptance of the c states in Q_m plus the transition table to $Q_{\min\{m+1, n\}}$. Without loss of generality, the empty string just codes the empty language. This convention permits to avoid a domain check for the index j .

As a finite automaton is the same as an algorithm working from the left to the right through the word with constant memory, the algorithm is now given more explicit than it would be in the case of an automaton. It runs in stages and it memorises information which can be stored in constantly

many bits; note that the number of these bits depends on the value of c but not of the value of n . For easier readability, this information is memorised in variables which are initialised in the first two steps. The algorithm reads in steps (2) and (4) the symbols unless the end of the corresponding inputs (j and x) is reached in which case the last value is not overwritten. If j is empty, then the algorithm rejects all x .

- (1) Variables: status $\in \{\text{reject}, \text{accept}\}$; state $\in \{1, 2, \dots, c\}$; $a \in \Sigma$ (current input symbol to be processed); b (table of current Q_m).
Let state = 1 and b be a code such that 1 is a rejecting state and all transitions are from 1 to 1.
- (2) If there is some symbol of the code of j to be read
then read b
else let b unchanged.
- (3) If the current value of state according to b is an accepting state
then let status = accept
else let status = reject.
- (4) If there is some symbol of the input word to be read
then read a
else accept/reject according to status and terminate.
- (5) Decode from b the new value of state in dependence of a and of the current value of state.
- (6) Go to (2).

The verification is left to the reader, as it is straightforward but lengthy. Note that \mathcal{L} is automatic iff the set $\{j \in J : H_j \in \mathcal{L}\}$ is regular. \square

The characterisation from Theorem 4.1 could be put into a more general form.

Theorem 4.2: *A class \mathcal{L} is a subclass of an automatic family iff there is a constant c such that $R \in \mathcal{L}$ iff there is an n such that:*

- For every m there are at most c different derivatives $R[x]$ with $x \in \Sigma^m$;
- There are at most c different derivatives $R[x]$ with $x \in \Sigma^* \wedge |x| \geq n$.

Here the derivative $R[x]$ is the set $\{y \in \Sigma^* : xy \in R\}$.

In learning theory, an important and well-studied family is that of the pattern languages [2, 17, 24]. Here a pattern is a string consisting of constants and variables, where the language generated by the pattern is the set of all those words which are obtained by replacing variables by strings. If the

strings replacing the variables are permitted to be empty one speaks of “erasing pattern languages” and if these strings are not permitted to be empty one speaks of “non-erasing pattern languages”. The learnability of pattern languages has been extensively studied and while there is an explanatory learner for the class of non-erasing pattern languages [2], such a learner does not exist in the case of the erasing pattern languages [21]. In particular, the class of “regular pattern languages” is quite important and one might ask what the automatic counterpart of it is. Here Shinohara [24] defined that a pattern is called a *regular pattern* iff it contains every variable at most once. So if $\Sigma = \{0, 1, 2\}$ and the pattern is $i = 01121x121y112z$ then the language generated by this pattern is given by the regular expression $01121 \cdot (0+1+2)^* \cdot 121 \cdot (0+1+2)^* \cdot 112 \cdot (0+1+2)^*$; in the non-erasing case, one would have to replace “ $(0+1+2)^*$ ” by “ $(0+1+2) \cdot (0+1+2)^*$ ” as every variable consists of at least one letter.

Theorem 4.3: *A class \mathcal{L} of erasing regular pattern languages is contained in an automatic family iff there is a constant c such that, in every pattern of a language in the family, there are at most c constants after the occurrence of the first variable.*

Proof: Given the constant c , it is first shown that there is an automatic family containing all the erasing regular pattern languages with up to c constants after the first occurrence of a variable. Note that there are only finitely many regular patterns which start with a variable and contain up to c constants; the reason is that a double variable xy has the same effect as a single variable in the case that variables can take the empty word and are not repeated. Now let Γ be the set of all these patterns and assume that Γ is coded such that it is disjoint to Σ . One chooses as an indexing the set $\Sigma^*\Gamma$. If L_g is the language generated by $g \in \Gamma$ then extend this definition to $L_{ig} = i \cdot L_g$ for $i \in \Sigma^*$. The family of all L_{ig} is automatic as the finite automaton recognising $\{conv(ig, x) : x \in L_{ig}\}$ accepts $conv(ig, x)$ iff $x = iy$ for some y and $y \in L_g$; note that the latter can be checked as there are only finitely many languages of the form L_g and so one can combine the corresponding finite automata to get one automaton for checking whether $y \in L_g$.

For the converse direction, assume that an automatic family $\{L_i : i \in I\}$ of erasing regular pattern languages is given. As shown in the proof of Theorem 3.5, there is a constant c such that, for every $i \in I$, there are at most c different derivatives $L_i[x]$ with $|x| \geq |i|$. Now assume that $i \in I$ codes a pattern language of the form $ux_0a_1x_1 \dots a_nx_n$ where $u \in \Sigma^*$,

$a_1, \dots, a_n \in \Sigma$, x_0 is a variable and x_1, \dots, x_n are each either a variable or the empty string. Furthermore, let $b \in \Sigma - \{a_1\}$. Now choose k larger than the length of i . Then $ub^k a_1 \dots a_n$ is a string such that there are n different derivatives $L_i[ub^k a_1 \dots a_n]$ with shortest word $a_{m+1} \dots a_n$, where $m \in \{1, 2, \dots, n\}$. Furthermore, as $|ub^k| \geq |i|$ the inequality $n \leq c$ holds and therefore the pattern language has at most c constants after the occurrence of the variable x_0 . \square

As the counterpart for non-erasing regular pattern languages, one just postulates that there are at most c symbols (counting both, symbols from Σ and variables) after the occurrence of the first variable. Note that every regular pattern generates a regular language [24]. The following remark [20] shows that the converse direction fails in certain cases.

Remark 4.4: The non-erasing pattern $xyxz$ generates the regular language $\cup_{a \in \Sigma} a \Sigma^* a \Sigma^*$ which, for alphabets of size 2 or more, cannot be generated by a regular pattern.

If Σ has at least four symbols, then every erasing pattern language which is regular is generated by a pattern which does not have repetitions of the variables. However for alphabets of size 2 and 3 there are also erasing pattern languages which are regular sets but need repetitions of variables to be generated. For $\Sigma = \{0, 1\}$ the pattern $tuv1uwwy1xyz$ generates the language $(0+1)^* \cdot 1 \cdot (0+1)^* \cdot 1 \cdot (0+1)^* - 10 \cdot (00)^* \cdot 1$, which is not generated by any regular pattern.

5. Applications of Automatic Families in Learning Theory

Automatic families are a special case of indexed families [1, 16] which are a widely studied subject in inductive inference. Here an indexed family $\{L_i : i \in I\}$ is uniformly recursive, that is, the domain D , the set of admissible indices I and the indexing $i, x \mapsto L_i(x)$ are all recognisable by Turing machines. So automatic families are just the restriction obtained when replacing Turing machines by finite automata with possibly several inputs. Gold [7] formalised the notion of learning and created the following notion of explanatory learning: A family $\{L_i : i \in I\}$ is explanatorily learnable iff there is a recursive learner which — in parallel — reads more and more data about the language R to be learnt from a text T and outputs a sequence e_0, e_1, e_2, \dots of indices which syntactically converges to one index $i \in I$ with $R = L_i$. Here a text is an infinite sequence of words containing every member of the set to be learnt but not any nonmember; the words in the

text can be in arbitrary and adversary order and repetitions are permitted. The interested reader might find more background information on learning theory in the standard textbooks of inductive inference [13, 19]. Angluin [1] gave a criterion on the learnability of an indexed family which — in the case of automatic families — can be simplified to the following one [11].

Theorem 5.1: *An automatic family $\mathcal{L} = \{L_i : i \in I\}$ is explanatorily learnable iff for every $i \in I$ there is a constant c such that, for all $j \in I$, the equality $L_i = L_j$ holds whenever $\{x \in L_i : |x| \leq d_{\mathcal{L}}(L_i) + c\} \subseteq L_j \subseteq L_i$.*

In other words, one can build a learner which — when learning R — always conjectures the least i such that all data in L_i shorter than $|i| + c$ have already been observed but no data outside L_i .

Based on this observation, Jain, Luo and Stephan [11] formulated the notion of an automatic learner which is less powerful than a recursive learner. An automatic learner has a long term memory which stores all relevant information about the data observed; this long term memory is a string like any input word. The learner is then given by an update function

$$F: (\text{old long term memory, current data item}) \mapsto (\text{new long term memory, new hypothesis})$$

Starting from an initial value for its long term memory, the learner reads in each round a current datum and updates its memory and the hypothesis according to F . The update function F has to be automatic. Jain, Luo and Stephan [11] showed that not every learnable automatic class can also be learnt by an automatic learner. Indeed even quite simple classes like the one given by $L_i = \Sigma^* - \{i\}$ (where $D = I = \Sigma^*$) is learnable only if the alphabet consists of one symbol only; in the case of an alphabet with at least two symbols this class is no longer learnable by an automatic learner, as the automatic learner cannot memorise enough information about the data observed. Special cases considered are where the long term memory cannot be longer than the longest datum observed so far, where the long term memory is the last hypothesis conjectured (iterative learning) and where the long term memory consists of up to c data items observed previously (bounded example memory); the ability for an automatic learner to learn depends heavily on the nature of such restrictions imposed. Ong [20] investigated the learnability of automatic families of pattern languages and related classes.

Jain, Martin and Stephan [12] considered the setting of robust learning [6, 14, 26] and asked when every translation of an automatic family \mathcal{L} is

learnable, where a translation is given by a first-order definable operator which preserves inclusions between among all languages as well as non-inclusions among languages from the class. An example is $\Phi(L) = \{x : \exists y \in L [y \neq x]\}$ and the underlying family \mathcal{L} contains \emptyset , every singleton $\{x\}$ and the full set D . It is easy to see that \mathcal{L} is learnable. However the given translation of \mathcal{L} is not learnable, as the learner would have to converge to $\Phi(D)$ after having seen finitely many data items and then the input language could be $\Phi(\{y\}) = D - \{y\}$ for some y which has not yet been observed in the input. General characterisations are given for the question of the following type: “For which classes are all translations learnable under a given criterion?” Besides the standard criteria from inductive inference, the paper also looks at query learning. Here a learner asks a teacher in a given query language about the class, for example the learner produces an i and asks whether the language L_i is a superset of the language R to be learnt. The learner asks finitely many queries and has then to conjecture the correct index of the input language. Angluin [3] started to investigate the learnability of regular languages from queries. She showed that the class of all regular languages can be learnt using membership and equivalence queries in polynomial time where, when learning R , the time bound depends on $d_{dfa}(R)$ and the largest counter example observed. The following result of Jain, Martin and Stephan [12] links explanatory learning, query learning and robustness.

Theorem 5.2: *The following conditions are equivalent for an automatic family $\mathcal{L} = \{L_i : i \in I\}$.*

- (a) *For every $i \in I$ there is a bound b such that, for all $j \in I$ with $L_j \subset L_i$, there is a $k \in I$ with $k \leq_U b$, $L_j \subseteq L_k$ and $L_i \not\subseteq L_k$.*
- (b) *Every translation of \mathcal{L} is explanatorily learnable.*
- (c) *Every translation of \mathcal{L} can be learnt using superset queries and membership queries.*
- (d) *\mathcal{L} can be learnt using superset queries.*

Note that in condition (d) it does not matter whether one learns only \mathcal{L} using superset queries or every translation of \mathcal{L} using superset queries. The reason for this is that translations do not change the inclusion structure of the automatic family. This result shows that there are also connections to robust notions of query learning.

Jain, Luo, Semukhin and Stephan [10] investigated the question on what can be said on the learnability of uncountable families. For this they use ω -automatic families, where the indices of the sets are ω -words and where a

nondeterministic Büchi-automaton checks whether a finite word x belongs to the language defined by an ω -word. Also here the Büchi automaton is fed with the convolution of the input word and the ω -word serving as the index of the language. As there are uncountably many indices, it is no longer possible for the learner to come up with the correct index after finite time; therefore the model is adjusted to a verification game. So the learner reads in parallel a text consisting of all the words in the language and an ω -index; the output is a sequence of Büchi automata which converges to one fixed automaton. Then this automaton has to accept the ω -word iff it is an index for the language observed. Also in this setting, learnability is equivalent to Angluin's tell-tale condition [1]. However it is necessary for this result that the learner has the right to choose the indexing; otherwise the criterion is more restrictive. Also further other criteria are transferred to this model and it is shown that one can abstain from adjusting the indexing if one requires vacillatory learning, where the learner in the limit oscillates between finitely many Büchi automata, where each of them accepts the observed ω -word iff it is an index for the language to be learnt.

References

1. Dana Angluin. Inductive inference of formal languages from positive data. *Information and Control* 45:117–135, 1980.
2. Dana Angluin. Finding patterns common to a set of strings. *Journal of Computer and System Sciences*, 21:46–62, 1980.
3. Dana Angluin. Learning regular sets from queries and counterexamples. *Information and Computation*, 75:87–106, 1987.
4. Achim Blumensath. *Automatic structures*. Diploma thesis, RWTH Aachen, 1999.
5. Achim Blumensath and Erich Grädel. Automatic structures. *15th Annual IEEE Symposium on Logic in Computer Science (LICS)*, pages 51–62, IEEE Computer Society, 2000.
6. Mark Fulk. Robust separations in inductive inference. *Proceedings of the 31st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 405–410, St. Louis, Missouri, 1990.
7. E. Mark Gold. Language identification in the limit. *Information and Control* 10:447–474, 1967.
8. Bernard R. Hodgson. *Théories décidables par automate fini*. Ph.D. thesis, University of Montréal, 1976.
9. Bernard R. Hodgson. Décidabilité par automate fini. *Annales des sciences mathématiques du Québec*, 7(1):39–57, 1983.
10. Sanjay Jain, Qinglong Luo, Pavel Semukhin and Frank Stephan. Uncountable automatic classes and learning. *Algorithmic Learning Theory, Twentieth International Conference, ALT 2009, Porto, Portugal, October 3-5, 2009*.

- Proceedings. Springer LNAI 5809:293–307, 2009. Technical Report TRB1/09, School of Computing, National University of Singapore, 2009.
11. Sanjay Jain, Qinglong Luo and Frank Stephan. Learnability of automatic classes. Technical Report TRA1/09, School of Computing, National University of Singapore, 2009.
 12. Sanjay Jain, Eric Martin and Frank Stephan. *Identification in the limit and automatic translators*. Manuscript, 2010.
 13. Sanjay Jain, Daniel N. Osherson, James S. Royer and Arun Sharma. *Systems That Learn*. MIT Press, 2nd Edition, 1999.
 14. Sanjay Jain and Frank Stephan. *A tour of robust learning. Computability and Models. Perspectives East and West*. Edited by S. Barry Cooper and Sergei S. Goncharov. Kluwer Academic / Plenum Publishers, University Series in Mathematics, pages 215–247, 2003.
 15. Bakhadyr Khoussainov and Anil Nerode. Automatic presentations of structures. *Logical and Computational Complexity*, (International Workshop LCC 1994). Springer LNCS 960:367–392, 1995.
 16. Steffen Lange and Thomas Zeugmann. Language learning in dependence on the space of hypotheses. *Proceedings of the Sixth Annual Conference on Computational Learning Theory (COLT)*, pages 127–136. ACM Press, 1993.
 17. Steffen Lange and Rolf Wiehagen. Polynomial time inference of arbitrary pattern languages. *New Generation Computing*, 8:361–370, 1991.
 18. Ming Li and Paul Vitányi. *An Introduction to Kolmogorov Complexity and Its Applications*. Third Edition, Springer, 2008.
 19. Daniel Osherson, Michael Stob and Scott Weinstein. *Systems That Learn, An Introduction to Learning Theory for Cognitive and Computer Scientists*. Bradford — The MIT Press, Cambridge, Massachusetts, 1986.
 20. Yuh Shin Ong. *Automatic Structures and Learning*. Final Year Project, School of Computing, National University of Singapore, 2010.
 21. Daniel Reidenbach. A non-learnable class of E-pattern languages. *Theoretical Computer Science*, 350:91–102, 2006.
 22. Sasha Rubin. *Automatic Structures*. Ph.D. Thesis, Univ. of Auckland, 2004.
 23. Sasha Rubin. Automata presenting structures: a survey of the finite string case. *The Bulletin of Symbolic Logic*, 14:169–209, 2008.
 24. Takeshi Shinohara. Polynomial time inference of extended regular pattern languages. *RIMS Symposia on Software Science and Engineering, Kyoto, Japan, Proceedings*. Springer LNCS 147:115–127, 1982.
 25. Wai Yean Tan. *Automatic Structures*. Honours Year Thesis, Department of Mathematics, National University of Singapore, 2008.
 26. Thomas Zeugmann. On Bärzdiņš’ Conjecture. *Proceedings of the International Workshop on Analogical and Inductive Inference (AII’86)*, Springer LNCS 265:220–227, 1986.