

Advanced Automata Theory 12

Open Problems in Automata Theory

Frank Stephan

Department of Computer Science

Department of Mathematics

National University of Singapore

fstephan@comp.nus.edu.sg

Repetition: 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.

Repetition: 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.

Repetition: Learning from 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 mem_0 and hypothesis e_0 . In the n -th cycle, \mathbf{M} updates mem_n and input $\mathbf{T}(n)$ to mem_{n+1} and e_{n+1} ; in short: \mathbf{M} maps $(\text{mem}_n, \mathbf{T}(n))$ to $(\text{mem}_{n+1}, e_{n+1})$. The values e_1, 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 e_n of \mathbf{M} are the same e describing \mathbf{L} .

Repetition: Learnability

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.

Repetition: Automatic Learners

Automatic learner given by mem_0 , e_0 and $\text{uf} : \text{conv}(\text{mem}_n, \mathbf{x}_n) \mapsto \text{conv}(\text{mem}_{n+1}, e_{n+1})$.
 mem_n is a string in an arbitrary alphabet;
 e_n is a hypothesis from \mathbf{I} when learning $\{\mathbf{L}_e : e \in \mathbf{I}\}$;
 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: $|\text{mem}_{n+1}| \leq k$.

Repetition: 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 #	$\text{conv}(0,112)$	$\{u : 0 \leq_{\text{lex}} u \leq_{\text{lex}} 112\}$

Repetition: Limitations

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

Overview

Open Problems from Work in Singapore

These problems are not that famous but also not that difficult. There is some chance to get them out with some good ideas.

Famous Open Problems

Some of these problems have been open for many years or decades. They are difficult to get out. However, once, the undergraduate student Róbert Szelepcsényi from Slovakia (at that time ČSSR) solved one problem of this type by showing that context-sensitive classes are closed under complement; he and Immerman from the USA solved this problem independently.

Open Problems from Lecture 11

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?

Many-One and One-One Reducibilities

One says a set A is many-one reducible to B via f iff $\forall x [A(x) = B(f(x))]$. If f is one-one then A is one-one reducible to B via f .

The set A is many-one / one-one reducible to B iff there is a function f witnessing this of a corresponding type (like being automatic or rational).

Incomparable sets: For alphabet $\{0, 1\}$, the sets $\{0\}^*$ and $\{0\}^* \cdot \{1\} \cdot \{0, 1\}^*$ are incomparable with respect to automatic one-one reducibility.

Automatic Reducibility

Reducibilities without base set

Let $\mathbf{A} \leq_{\text{au}} \mathbf{B}$ denote that there is an automatic function \mathbf{f} such that $\forall \mathbf{x}, \mathbf{y} \in \mathbf{A} [\mathbf{x} \neq \mathbf{y} \Rightarrow \mathbf{f}(\mathbf{x}) \neq \mathbf{f}(\mathbf{y}) \wedge \mathbf{f}(\mathbf{x}) \in \mathbf{B}]$.

Similarly one writes $\mathbf{A} \leq_{\text{tr}} \mathbf{B}$ for the corresponding definition where \mathbf{f} is any function computed by a finite transducer.

Example

$\{0\}^* \cup \{1\}^* \leq_{\text{tr}} \{2\}^*$ and $\{0\}^* \cup \{1\}^* \not\leq_{\text{au}} \{2\}^*$.

A regular set \mathbf{A} has size $\Theta(n^k)$ iff there is $c > 0$ with

$$\forall n [n^k/c - c \leq |\{\mathbf{x} \in \mathbf{A} : |\mathbf{x}| \leq n\}| \leq n^k \cdot c + c]$$

and if \mathbf{A} is infinite and does not have size $\Theta(n^k)$ for any k then \mathbf{A} is exponential-sized.

Results of Tan Wai Yean

Theorem [Tan 2010]

Let A, B be regular sets.

1. The sets A, B are comparable for tr -reducibility:
 $A \leq_{\text{tr}} B$ or $B \leq_{\text{tr}} A$. Furthermore, $A \leq_{\text{tr}} B$ if one of the following conditions holds:
 - A, B are both finite and $|A| \leq |B|$;
 - A is finite and B infinite;
 - A has size $\Theta(n^k)$ and B has size $\Theta(n^h)$ with $k \leq h$;
 - B is exponential-sized.
2. If A is polynomial-sized or finite then $A \leq_{\text{au}} B$ or $B \leq_{\text{au}} A$. If A is of size $\Theta(n^k)$, B is of size $\Theta(n^h)$ and $k < h$ then $A \leq_{\text{au}} B$ and $B \not\leq_{\text{au}} A$.

Exercise and Open Problem

Exercise 12.4

Make an automatic one-one function which maps the domain $A = 0^*(1^* \cup 2^*)$ to a subset of $B = (0000)^*(1111)^*(2222)^*$, that is, show that $A \leq_{\text{au}} B$.

Open Problem [Tan 2010]

Are there regular sets A, B such that $A \not\leq_{\text{au}} B$ and $B \not\leq_{\text{au}} A$?

Remark

There are context-free sets A, B with $A \not\leq_{\text{au}} B$ and $B \not\leq_{\text{au}} A$: $A = \{0\}^*$ and $B = \{x2y : x, y \in \{0, 1\}^* \wedge |x| = |y|\}$.

The XX-Problem

For the following, consider the fixed one-one reduction $x \mapsto xx$ and $x \mapsto x(x^{mi})$.

Open Problem [Zhang 2013]

1. Given a regular language A , is there a regular language B such that, for all x , $A(x) = B(xx)$?
2. Given a context-free language A , is there a context-free language B such that, for all x , $A(x) = B(xx)$?

Zhang [2013] solved the reverse directions: If B is regular then the set $A = \{x : xx \in B\}$ is also regular. However, the set $B = \{0^n 1^n 2^m 0^m 1^k 2^k : n, m, k \in \mathbb{N}\}$ is context-free while the corresponding set $A = \{x : xx \in B\} = \{0^n 1^n 2^n : n \in \mathbb{N}\}$ is properly context-sensitive.

The XM-Problem

Theorem [Fung 2014]

Assume that A is a regular set and let $B = \{u : \text{there are an odd number of pairs } (y, z) \text{ with } u = yz, y \in A \text{ and } z^{mi} \in A\}$. Then B is regular and for all x , $x \in A \Leftrightarrow x \cdot x^{mi} \in B$.

Open Problem

Given a context-free set A , is there a context-free set B such that for all x , $x \in A \Leftrightarrow x \cdot x^{mi} \in B$?

Ordered Groups

An ordered group $(G, +, <, 0)$ satisfies the group and order axioms and for all x, y, z that $x < y$ implies $x + z < y + z$ and $z + x < z + y$.

Theorem [Jain, Khoussainov, Stephan, Teng and Zou 2014]
If an ordered group has a regular domain and fully automatic addition then the group is commutative.

Question [Khoussainov] Is there a fully automatic copy of $(\mathbb{Z}, +, 0)$ such that the set $\{x \in \mathbb{Z} : x > 0\}$ is not regular?

Theorem [Jain, Khoussainov, Stephan, Teng and Zou 2014]
There is a fully automatic copy of $(\{x \cdot 6^y : x, y \in \mathbb{Z}\}, +, 0)$ in which neither the order nor the set of positive numbers is automatic.

Synchronising words

Definition

Given a dfa with states Q and transition function δ , a word w is called **synchronising** iff $\delta(p, w) = \delta(q, w)$ for all $p, q \in Q$.

Example

The dfa which accepts all decimal numbers which are multiples of **3** does not have a synchronising word.

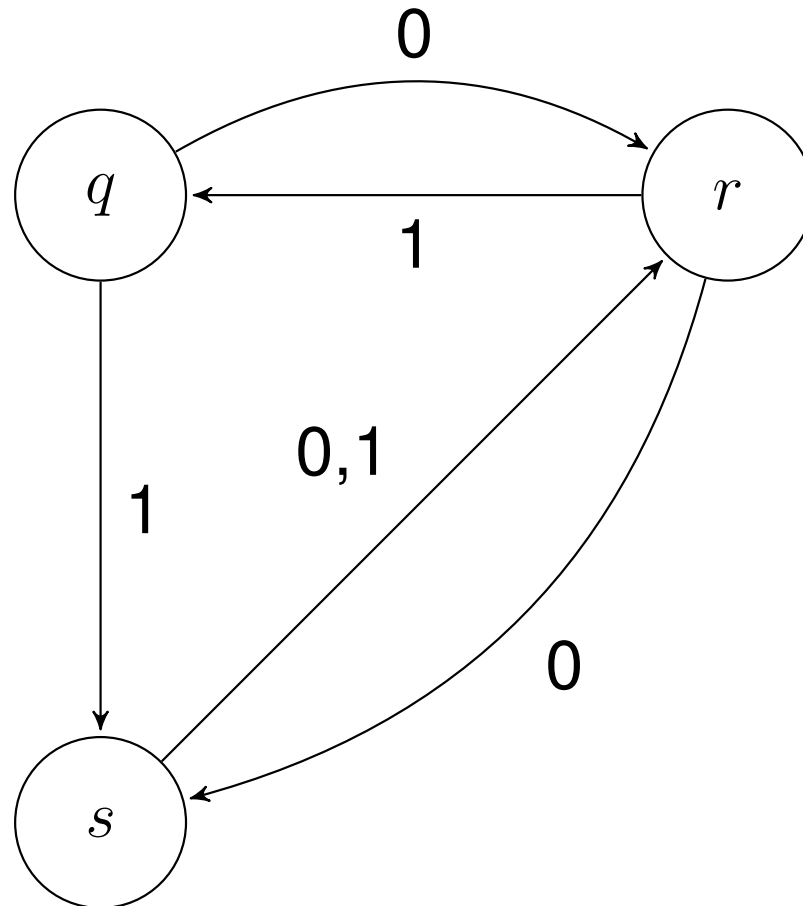
Example

The dfa of two states which accepts a decimal number iff it ends with **0, 2, 4, 6, 8** and rejects a decimal number iff it ends with **1, 3, 5, 7, 9** has a synchronising word.

Theorem [Černý 1964]

For each n there is a complete dfa with n states which has a synchronising word of length $(n - 1)^2$ and no shorter ones.

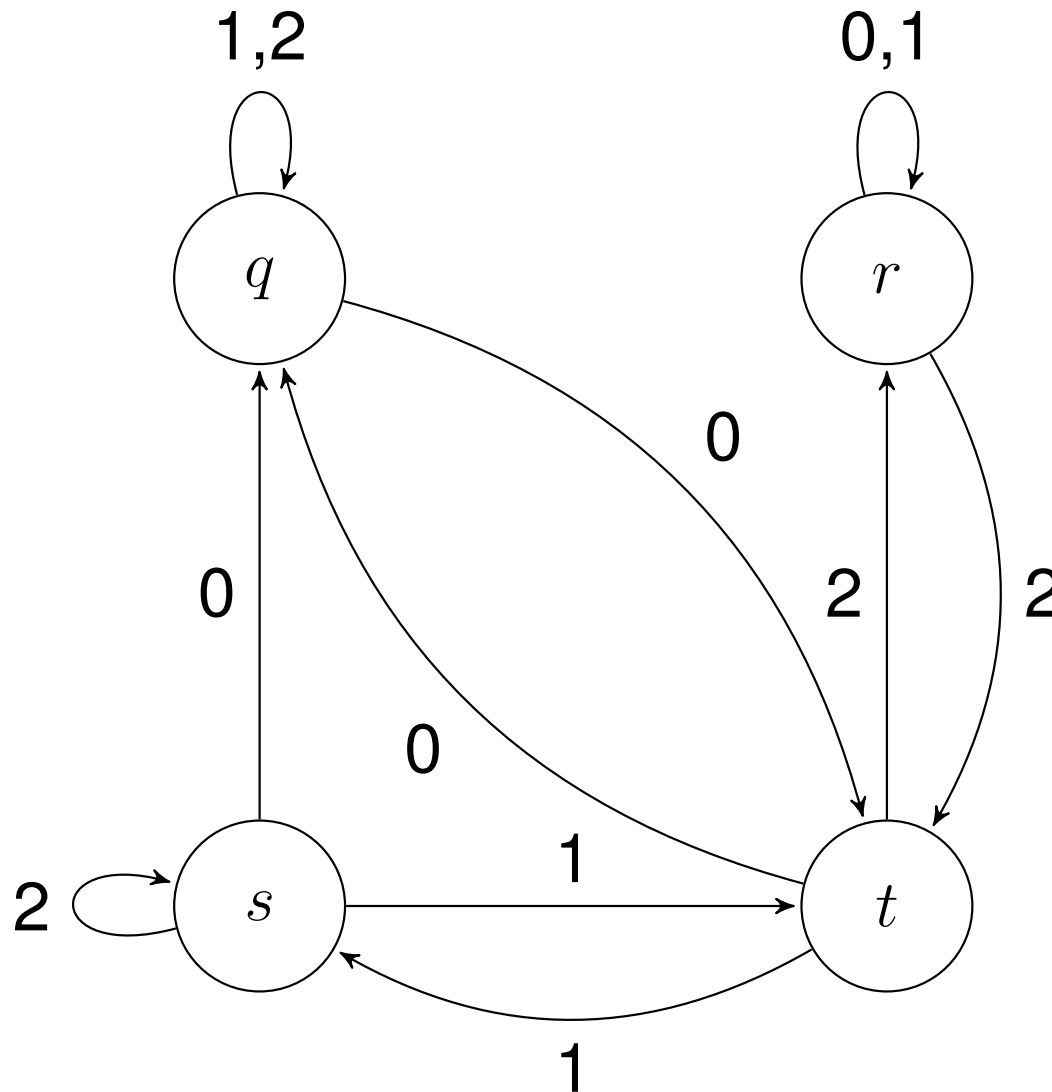
Example for 3 States, Length 4



The word **0110** is a synchronising word which sends all states to **r**.

Example for 4 States, Length 9

Alphabet $\{0, 1, 2\}$ and synchronising word **012020120**.



Upper Bounds

Černý's lower bound is quadratic, all known upper bounds are cubic or worse.

Theorem [Frankl 1982; Klyachko, Rystov and Spivac 1987; Pin 1983]. Assume a complete dfa has n states and has a synchronising word. Then it has a synchronising word not longer than $(n^3 - n)/6$.

Example

If $n = 3$ and the dfa has a synchronising word then there is a synchronising word of length up to 4.

Černý's Conjecture

If a dfa of n states has a synchronising word then it has one not longer than $(n - 1)^2$.

Exercise 12.13

Prove Černý's conjecture for $n = 4$.

Star Height

Consider regular expressions formed by union, concatenation and Kleene star and count how often stars are nested:

- Let S_0 contain all finite languages, note that S_0 is closed under union and concatenation;
- For each n , let S_{n+1} contain all languages which can be formed by taking unions and concatenations of languages of the form L or L^* with $L \in S_n$.

The star-height of a regular language L is the minimal n with $L \in S_n$.

Theorem

For each n there is a regular language of star height n .

Examples

Here some examples for the lowest levels.

- The language $L_0 = \{0, 11, 222, 3333\}$ is finite and has star-height **0**;
- The language $L_1 = \{00, 11\}^*$ has star-height **1**;
- The language $L_2 = (\{00, 11\}^* \cdot \{22, 33\} \cdot \{00, 11\}^* \cdot \{22, 33\})^*$ has star-height **2**.

Generalised Star-Height

Let T_0 contain the all sets formed from Σ^* and all finite languages by iterated use of intersection, union, concatenation and set-difference. Let T_{n+1} be all sets formed from languages L and L^* with $L \in T_n$ by iterated use of intersection, union, concatenation and set-difference. The languages $L \in T_n$ which are not in any T_m with $m < n$ have generalised star-height n .

Open Problem

For which levels $n > 1$ do there exist regular languages of generalised star-height n ?

Examples

- The language $\{0, 1\}^*$ has generalised star-height **0**, as

$$\{0, 1\}^* = \Sigma^* - \bigcup_{a \in \Sigma - \{0, 1\}} \Sigma^* a \Sigma^*;$$

- L_2 from above has generalised star-height **1**, as

$$L_2 = \{00, 11, 22, 33\}^* \cap (\{0, 1\}^* \cdot \{22, 33\} \cdot \{0, 1\}^* \cdot \{22, 33\})^*$$

and so L_2 is the intersection of two languages of generalised star-height **1**;

- $L_3 = \{w : w \text{ does not have a substring of the form } v\}$ for a fixed v is of generalised star-height **0** as

$$L_3 = \Sigma^* - \Sigma^* \cdot v \cdot \Sigma^*;$$

- $L_4 = \{w : w \text{ has an even number of } 0\}$ is of generalised star-height **1**.

Quiz

Determine the generalised star-height of the following languages over the alphabet $\{0, 1, 2\}$ – it is zero or one:

1. $\{00, 11, 22\}^* \cdot \{000, 111, 222\}^*$;
2. $\{0, 1\}^* \cdot 2 \cdot \{0, 1\}^* \cdot 2 \cdot \{0, 1\}^*$;
3. $(\{0, 1\}^* \cdot 2 \cdot \{0, 1\}^* \cdot 2 \cdot \{0, 1\}^*)^*$;
4. $(\{0, 1\}^* \cdot 2 \cdot \{0, 1\}^*)^*$;
5. $(\{0, 1\}^+ \cdot 22)^*$;
6. $(\{0, 1\}^* \cdot 22)^*$;
7. $((00)^+ \cdot 11)^+ \cdot 22^+$.

Recognising Languages

Deciding Membership

Let L be a language and let $n = |w|$ be the length of an input word w . For various types of L there are methods to decide whether $w \in L$.

If L is regular then a dfa can check in linear time whether a word $w \in L$.

If L is context-free then one can check in time $O(n^3)$ whether $w \in L$; this algorithm of Cocke, Kasami and Younger has been improved to $O(n^{2.38})$ by Coppersmith and Winograd.

If L is context-sensitive then one can check in time $O(c^n)$ for some c whether $w \in L$. Furthermore, there are non-deterministic algorithms running in space $O(n)$ and deterministic algorithms running in space $O(n^2)$.

Questions

Open Problem

What is the best time complexity to decide the membership of a context-free language?

Open Problem

Can the membership in a given context-sensitive language be decided in deterministic linear space?

This is related to the overall problem whether a non-deterministic algorithm using space $s(n)$ can be made into a deterministic algorithm using space $O(s(n))$; the best bound known is $O(s^2(n))$.

Open Problem

Given an automatic group G with a finite generator set F , let $L = \{w \in F^* : w \text{ and } \varepsilon \text{ represent the same member of } G\}$. Is L solvable in logarithmic space?

Isomorphism Problems

In certain cases one does not ask how difficult a problem is but only whether it has at all an algorithmic solution. One can, for example, for automatic well-orderings $(\mathbf{A}, <)$ and $(\mathbf{B}, <)$ decide with an algorithm whether they are isomorphic; similarly for automatic linear orders. One cannot do this for automatic equivalence relations.

Open Problem

- Assume that $(\mathbf{A}, \text{Succ}_{\mathbf{A}}, \mathbf{P}_{\mathbf{A}})$ and $(\mathbf{B}, \text{Succ}_{\mathbf{B}}, \mathbf{P}_{\mathbf{B}})$ are automatic structures such that $(\mathbf{A}, \text{Succ}_{\mathbf{A}})$ and $(\mathbf{B}, \text{Succ}_{\mathbf{B}})$ are isomorphic to the natural numbers with successor and that $\mathbf{P}_{\mathbf{A}}$ and $\mathbf{P}_{\mathbf{B}}$ are regular predicates (subsets) on \mathbf{A} and \mathbf{B} . Can one decide whether $(\mathbf{A}, \text{Succ}_{\mathbf{A}}, \mathbf{P}_{\mathbf{A}})$ and $(\mathbf{B}, \text{Succ}_{\mathbf{B}}, \mathbf{P}_{\mathbf{B}})$ are isomorphic?
- Can one decide whether the commutative fully automatic groups $(\mathbf{A}, +)$ and $(\mathbf{B}, +)$ are isomorphic?

Complexity of Deciding Games

One studies the complexity of deciding games (or other problems) in dependence of parameters. For parity games, the parameters taken here are the number n of nodes and the number m of values.

- McNaughton 1993: $O((kn)^{m+1})$ for some k .
- Browne, Clarke, Jah, Long and Marrero 1997: $O(n^2 \cdot (2n/m)^{(m+3)/2})$.
- Jurdinski, Patterson and Zwick 2006/2008: $n^{k \cdot \sqrt{n}}$ for some k .
- Schewe 2007/2016: $n^2 \cdot (k \cdot n \cdot m^{-2})^{m/3}$ for some k .
- Calude, Jain, Khoussainov, Li and Stephan 2016: $O((m/\log(n))^4 \cdot n^{3.45+\log(m/\log(n)+3)})$.

Open Problem: Can parity games be decided in polynomial time?

Winning Statistics 1

Modify parity game to parity game with winning statistics. Winning statistics are vectors $\mathbf{b}_0 \mathbf{b}_1 \dots, \mathbf{b}_k$ with $k = \lceil \log(n) + 2 \rceil$ with the following meaning: $\mathbf{b}_i > \mathbf{0}$ stands for the observation of 2^i nodes in the play so far such that between any of these nodes the highest value was of the player's parity and \mathbf{b}_i is the largest value observed with starting the end of the sequence. If $\mathbf{b}_i, \mathbf{b}_j > \mathbf{0}$ and $j < i$ then the sequence for \mathbf{b}_j can only start after the 2^i nodes of the previous sequence and also the node with value \mathbf{b}_i have been observed.

Winning statistics indicate whether a player will eventually win in the case that the winner plays a memoryless winning strategy.

Winning Statistics 2

Initialisation: All b_i of a winning statistics are 0 .

Update rule: For each new node with value b , choose the largest i which can do one of the following:

- b and b_0, b_1, \dots, b_{i-1} have the players parity and b_i does not;
- $0 < b_i < b$.

If found, let $b_i = b$ and all $b_j = 0$ for all $j < i$.

Theorem

If Anke plays a memoryless winning strategy her winning statistics will eventually have $b_k > 0$ what will indicate that she has a winning strategy while Boris' winning statistics will never indicate a win. Similarly for Boris and his winning statistics.

Reduction to Survival Games

Game graphs of survival game consist of all nodes $(\mathbf{v}, \mathbf{p}, \mathbf{w})$ with node \mathbf{v} of parity game, player \mathbf{p} to move and current winning statistics \mathbf{w} of Boris. Move from $(\mathbf{v}, \mathbf{p}, \mathbf{w})$ to $(\mathbf{v}', \mathbf{p}', \mathbf{w}')$ is possible iff there is edge from \mathbf{v} to \mathbf{v}' in parity game, $\mathbf{p} \neq \mathbf{p}'$ and on move with value of \mathbf{v}' , winning statistics \mathbf{w} of Boris are updated to \mathbf{w}' and \mathbf{w} is not already won for Boris. Anke loses if game gets stuck.

There is an algorithm to determine the winner of a survival game which needs time linear in number of edges of the survival game, so linear in time $\mathbf{n}^2 \cdot \mathbf{n}^{\log(\mathbf{m})+4}$ where $\mathbf{n}^{\log(\mathbf{m})+4}$ is an upper bound on the number of winning statistics for Boris. This upper bound is based on the fact that Winning statistics consist of $\lceil \log(\mathbf{n}) + 3 \rceil$ numbers with $\lceil \log(\mathbf{m} + 1) \rceil$ bits.

More careful counting gives slightly better bound.

Exercises 12.25-12.28

Let $A = \{0\}^* \cdot \{1\}^*$, $B = \{00\}^* \cdot \{11\}^* \cdot \{22\}^*$;
 $C = \{00, 11, 22\}^*$; $D = \{0, 1\}^* - \{1\}^*$.

Exercise 12.25

How are the above sets A, B, C, D ordered by \leq_{au} ? Provide the reductions where they exist.

Exercise 12.26

How are the above sets A, B, C, D ordered by \leq_{tr} ? Provide the reductions where they exist.

Exercise 12.27

Find for A, B, C regular sets A', B', C' such that for all x ,
 $x \in A \Leftrightarrow xx \in A'$ and $x \in B \Leftrightarrow xx \in B'$ and
 $x \in C \Leftrightarrow xx \in C'$.

Exercise 12.28

Provide a regular set E such that E' exists as in 12.27, but E' is neither E nor $E \cdot E$.

Exercises 12.29-12.31

Determine for the following languages **F**, **G**, **H** the minimal complete dfa and either determine its smallest synchronising word or show that there is no synchronising word.

Exercise 12.29

Do the above for $\mathbf{F} = (\{0\}^* \cdot \{1\} \cdot \{0\}^* \cdot \{1\})^*$.

Exercise 12.30

Do the above for $\mathbf{G} = (\{0\}^+ \cdot \{1\} \cdot \{0\}^+ \cdot \{1\})^*$.

Exercise 12.31

Do the above for $\mathbf{H} = (\{0\}^+ \cdot \{1\} \cdot \{0\}^+ \cdot \{1\})^* \cup (\{0, 1\}^* \cdot \{11\} \cdot \{0, 1\}^* \cdot \{00\}) \cup (\{1\} \cdot \{0, 1\}^* \cdot \{00\})$.

Final Examination

- Final examination counts 60 marks; midterm counts 30 marks and homework counts 10 marks.
- The final examination covers all 12 lectures.
- There are 10 questions with six marks each equally distributed over the material.
- Some material in the lecture notes was not covered or only sketched during lecture hours; this is additional add-on for those who are interested.
- When revising the material, look also at the exercises in order to test whether you understood the material and also check out the self-tests.
- There will be no questions covering topics which are not in the lecture notes.
- The final examination is a closed book assessment.