## Generation of DFA Minimization Problems

Gregor Hans Christian Sönnichsen

February 18, 2020

#### Abstract

The theory of deterministic finite automatons (DFAs) is a classical topic of computer science-related courses. A typical task for students is to minimize a DFA by deleting irrelevant elements. However generation of those DFAs that shall be minimized is often done manually by the exercise instructor. This work presents approaches to automize the generation of DFA minimization tasks.

# $\begin{array}{c} \textbf{Generierung von DFA Minimierungsproblemen} \\ \textbf{-} \\ \textbf{Zusammenfassung} \end{array}$

 $Zusammen fassung. \dots \\$ 

# Table of Contents

1	Intr	oducti	on	2				
	1.1	Prelim	inaries	2				
		1.1.1	Deterministic Finite Automatons	2				
		1.1.2	Minimal DFAs	3				
		1.1.3	Practical Isomorphy of DFAs	4				
		1.1.4	Equivalent and distinguishable state pairs	4				
		1.1.5	The minimization algorithm	5				
		1.1.6	m-ComEquivPairs	6				
		1.1.7	Another equivalence automaton	6				
	1.2	_	rements analysis	6				
		1.2.1	Example of a DFA minimization task for students	7				
		1.2.2	Discussion of difficulty adjustment possibilities	7				
	1.3	Appro	ach and general algorithm	9				
2	Gen	erating	g minimal DFAs	10				
	2.1	Using	trial and error	10				
		2.1.1	Ensuring $A_{test}$ is minimal and $D(A_{test})$ is correct	11				
		2.1.2	Ensuring $A_{test}$ is planar	11				
		2.1.3	Ensuring $A_{test}$ is new	12				
		2.1.4	Option 1: Generating $A_{test}$ via Randomness					
		2.1.5	Option 2: Generating $A_{test}$ via Enumeration					
	2.2	Altern	ative approach: Building $m(i)$ bottom up	15				
	2.3	3 Related research on DFA generation						
	2.4	2.4 Empirical and combinatorial results						
3	Extending minimal DFAs 19							
	3.1	Creatin	ng equivalent state pairs	19				
		3.1.1	Adding outgoing transitions					
		3.1.2	Adding ingoing transitions	20				
		3.1.3	The algorithm					
		3.1.4	Creating equivalent state pairs does not change $D$	21				
	3.2	Adding	g unreachable states	24				
4	Con	Conclusion						
$\mathbf{A}$	Bibliography							
В	Erklärung							

### Chapter 1

### Introduction

Automata theory is recommended as part of a standard computer science curriculum [12, pp. 5-6]. As other such theories it provides the chance to gain a precise cognitive model yielding new perspectives on problems and givens. This may thus lead to increased problem solving skills and more accurate thinking.

A typical task in automata theory is the minimization of a given deterministic finite automaton (DFA). The classic textbook "Introduction to automata theory, languages, and computation" by Hopcroft et. al. [15] presents a practicable minimization algorithm. In this work we will confine ourselves to look at DFA minimizations using that algorithm.

In an introduction course to theoretical computer science minimization tasks are thus likely to occur in supplementary exercises or exams. As of the creation of such tasks, one may assume, that it is done mostly manually. Automation would yield here the following advantages:

- freeing time for other things, e.g. research, helping students face-to-face, designing the whole exercise sheet
- generation of tasks which lie in a well-defined range
- increased predictability and consistency of the generated task properties, which can be adjusted accurately through various parameters
- saves human operators from the generating task which involves monotonous work

**Gregor:** Delete or find extern from wikipedia Engagement on this topic promises moreover increased clarification which kind of minimization tasks can be generated, and where difficulties of those tasks lie.

This work aims to provide theoretical foundations for a DFA minimization task generator. What requirements a user has towards such a program will be discussed in a short requirements analysis. Based on this work a DFA minimization generator has been implemented. It can be found at https://github.com/bt701607/Generation-of-DFA-Minimization-Problems.

#### 1.1 Preliminaries

We start with defining preliminary theoretical foundations.

#### 1.1.1 Deterministic Finite Automatons

A 5-tuple  $A = (Q, \Sigma, \delta, s, F)$  with Q being a finite set of states,  $\Sigma$  a finite set of alphabet symbols,  $\delta \colon Q \times \Sigma \to Q$  a transition function,  $s \in Q$  a start state and  $F \subseteq Q$  final states

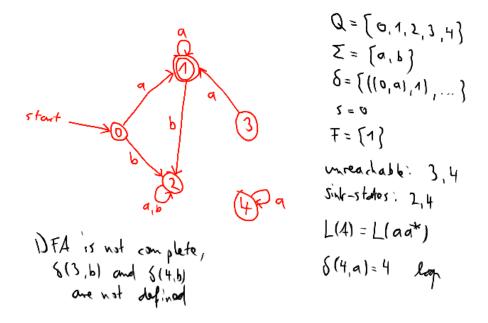


Figure 1.1: An example DFA and its properties.

is called deterministic finite automaton (DFA) [15, p. 46]. From now on  $\mathfrak A$  shall denote the set of all DFAs.

We say  $\delta(q, \sigma) = p$  is a transition from q to p using symbol  $\sigma$ . We define the extended transition function  $\delta^* : Q \times \Sigma^* \to Q$  of a DFA  $A = (Q, \Sigma, \delta, s, F)$  as:

- $\delta^*(q,\varepsilon) = q$
- $\delta^*(q, w\sigma) = \delta(\delta^*(q, w), \sigma)$  for all  $q \in Q, w \in \Sigma^*, \sigma \in \Sigma$

Then, the language of that DFA is defined as  $L(A) = \{ w \mid \delta^*(w) \in F \}$  [15, pp. 49-50. 52].

Given a state  $q \in Q$ . With  $d^-(q)$  we denote the set of all *ingoing* transitions  $\delta(q', \sigma) = q$  of q. With  $d^+(q)$  we denote the set of all *outgoing* transitions  $\delta(q, \sigma) = q'$  of q [9, pp. 2-3]. If a transition is of the form  $\delta(q, \sigma) = q$ , then we say that q has a *loop*.

**Definition 1.** We say a state q is (un-)reachable in a DFA A, iff there is (no) a word  $w \in \Sigma^*$  such that  $\delta^*(s, w) = q$ .

If all states of a DFA A are reachable, then we say A is accessible [9, p. 2].

A DFA is called *complete* iff for all states, every symbol of the alphabet is used on an outgoing transition:  $\forall q \in Q \colon \forall \sigma \in \Sigma \colon \exists p \in Q \colon \delta(q, \sigma) = p$ . Note, that every incomplete DFA can be converted to a complete one by adding a so called *dead state* [15, p. 67]. The resulting automaton has the same language.

#### 1.1.2 Minimal DFAs

This section closely follows [24, pp. 42-45]. We call a DFA A minimal, if there exists no other automaton with the same language using less states. With  $\mathfrak{A}_{min}$  we shall denote the set of all minimal DFAs.

The Nerode-relation  $\equiv_L \subseteq \Sigma^* \times \Sigma^*$  of a language L with alphabet  $\Sigma$  is defined as follows:

$$x \equiv_L y \Leftrightarrow_{def} \forall z \in \Sigma^* : (xz \in L \Leftrightarrow yz \in L)$$

The Nerode-relation of a DFA A is the the Nerode-relation of its language:  $\equiv_{L(A)}$ . If the context makes it clear, than we will shorten the notation of a equivalence class  $[x]_{\equiv_L}$  with [x].

The equivalence class automaton  $A_L = (Q_L, \Sigma_L, \delta_L, s_L, F_L)$  to a regular language L with alphabet  $\Sigma$  is defined as follows:

- $Q_L = \{ [x] \mid x \in \Sigma^* \}$
- $\Sigma_L = \Sigma$
- $\delta_L([x], \sigma) = [x\sigma], \ \forall x \in \Sigma^*, \ \forall \sigma \in \Sigma$
- $s = [\varepsilon]$
- $\bullet \ F = \{ \ [x] \mid x \in L \ \}$

**Theorem 1.** Given a language L, then the equivalence class automaton  $A_L$  is minimal.

#### 1.1.3 Practical Isomorphy of DFAs

Given two DFAs  $A_1 = (Q_1, \Sigma_1, \delta_1, s_1, F_1)$  and  $A_2 = (Q_2, \Sigma_2, \delta_2, s_2, F_2)$ . We say  $A_1$  and  $A_2$  are practical isomorph, iff:

- $|Q_1| = |Q_2|, |\Sigma_1| = |\Sigma_2|$  and
- there exists a bijection  $\phi \colon \Sigma_1 \to \Sigma_2$  such that:
- there exists a bijection  $\pi: Q_1 \to Q_2$  such that:

$$\pi(s_1) = s_2$$

 $\forall q \in Q_1 : (q \in F_1 \iff \pi(q) \in F_2)$ 

$$\forall q \in Q_1 \colon \forall \sigma \in \Sigma_1 \colon \pi(\delta_1(q,\sigma)) = \delta_2(\pi(q),\phi(\sigma)))$$

Note that practical isomorphy between two DFAs  $A_1$ ,  $A_2$  does not imply  $L(A_1) = L(A_2)$ . This would be given, if  $\Sigma_1 = \Sigma_2$  were the case (see [24, p. 45]). However the language of such DFAs is equivalent except for an exchange of alphabet symbols:

$$\{ \phi(\sigma_0) \dots \phi(\sigma_n) \mid \sigma_0 \dots \sigma_n \in L(A_1) \} = L(A_2)$$

In [24, p. 45] we can find the following statement:

**Theorem 2.** Every minimal DFA is unique (has a unique language) except for renaming of states.

We can thus formulate the following corollary:

**Corollary 1.** Every minimal DFA is unique (has a unique language) except for practical isomorphy.

Gregor: Write down practical isomorphy test, argue its correctness.

#### 1.1.4 Equivalent and distinguishable state pairs

**Definition 2** (Equivalent and Distinguishable State Pairs). [15, p. 154] A state pair  $q_1, q_2 \in Q$  of a DFA  $A = (Q, \Sigma, \delta, s, F)$  is called *equivalent*, iff  $\sim_A (q_1, q_2)$  is true, whereas

$$q_1 \sim_A q_2 \Leftrightarrow_{def} \forall z \in \Sigma^* : (\delta^*(q_1, z) \in F \Leftrightarrow \delta^*(q_2, z) \in F)$$

If  $q_0 \not\sim_A q_1$ , then  $q_0$  and  $q_1$  are called a *distinguishable* state pair. The relation  $\sim_A$  is an equivalence relation

#### 1.1.5 The minimization algorithm

This minimization algorithm MINIMIZEDFA works in four major steps, removing essentially states in such a way, that no unreachable states and no equivalent state pairs are left

1. Compute all unreachable states via breadth-first search.

```
1: function ComUnreachables(A)
 2:
          U \leftarrow Q \setminus \{s\}
                                                                                         ▶ undiscovered states
          O \leftarrow \{s\}
                                                                                               ▷ observed states
 3:
          D \leftarrow \{\}
                                                                                             4:
          while |O| > 0 do
               N \leftarrow \{ p \mid \exists q \in O \ \sigma \in \Sigma \colon \ \delta(q, \sigma) = p \land p \notin O \cup D \} 
               U \leftarrow U \setminus N
 7:
               D \leftarrow D \cup O
 8:
               O \leftarrow N
 9:
10:
          return U
```

2. Remove all unreachable states and their transitions.

```
1: function REMUNREACHABLES(A, U)
2: \delta' \leftarrow \delta \setminus \{ ((q, \sigma), p) \mid q \in U \lor p \in U \}
3: return (Q \setminus U, \Sigma, \delta', s, F \setminus U)
```

3. Compute all equivalent state pairs ( $\sim_A$ ). Inspired by Schöning [24, p. 46] and Martens [26, p. 17].

```
1: function ComEquivPairs(A)

2: M \leftarrow \{(p,q), (q,p) \mid p \in F, q \notin F\}

3: do

4: M' \leftarrow \{(p,q) \mid (p,q) \notin M \land \exists \sigma \in \Sigma \colon (\delta(p,\sigma), \delta(q,\sigma)) \in M\}

5: M \leftarrow M \cup M'

6: while M' \neq \emptyset

7: return Q^2 \setminus M
```

Note that Comequiveral requires its input automaton to be complete. **Gregor:** Why?

4. Merge all equivalent state pairs, which are exactly those in  $\sim_A$ . Inspired by Högberg [17, p. 10].

```
1: function REMEQUIVPAIRS(A, \sim_A)
                                                                                    \triangleright [\cdot]_{\sim_A} shall be abbreviated [\cdot]
           Q_E \leftarrow \emptyset
 2:
           \delta_E \leftarrow \emptyset
 3:
           F_E \leftarrow \emptyset
 4:
           for q in Q do
 5:
                Add [q] to Q_E
 6:
                for \sigma in \Sigma do
 7:
 8:
                      \delta_E([q], \sigma) = [\delta(q, \sigma)]
                if q \in F then
 9:
                      Add [q] to F_E
10:
           return (Q_E, \Sigma, \delta_E, [s], F_E)
11:
```

Note that RemequivPairs creates complete automatons.

```
1: function MINIMIZEDFA(A_{task})
2: A_{re} \leftarrow \text{REMUNREACHABLES}(A_{task}, \text{COMUNREACHABLES}(A_{task}))
3: A_{sol} \leftarrow \text{REMEQUIVPAIRS}(A_{re}, \text{COMEQUIVPAIRS}(A_{re}))
4: return A_{sol}
```

This DFA minimization algorithm has been found by Hopcroft [14] and was established in teaching by Hopcroft et al. [15, pp. 154-164].

**Theorem 3.** [15, pp. 162-164] MINIMIZEDFA computes a minimal DFA to its input DFA.

#### 1.1.6 *m*-ComEquivPairs.

When looking at COMEQUIVPAIRS, one notes, that it computes distinct subsets of  $Q \times Q$  on the way. Indeed, one could write the algorithm in such a way, that these subsets are explicitly computed in form of a function  $m: \mathbb{N} \to \mathcal{P}(Q \times Q)$ :

```
1: function m-ComEquivPairs(A)

2: i \leftarrow 0

3: m(0) \leftarrow \{(p,q),(q,p) \mid p \in F, q \notin F\}

4: do

5: i \leftarrow i+1

6: m(i) \leftarrow \{(p,q),(q,p) \mid (p,q) \notin \bigcup m(\cdot) \land \exists \sigma \in \Sigma \colon (\delta(p,\sigma),\delta(q,\sigma)) \in m(i-1)\}

7: while m(i) \neq \emptyset

8: return \bigcup m(\cdot)
```

Using this redefinition, we can easier refer to the state pairs marked in a certain iteration. We will use both variants in exchange.

We will denote the number of iterations done by ComEquivPairs on an DFA A as D(A).

#### 1.1.7 Another equivalence automaton

Consider step 3 and 4 of MINIMIZEDFA. A view on these algorithms reveal, that they are essentially collapsing each  $\sim$ -equivalence class of a DFA A into one state.

**Definition 3** ( $\sim_A$ -equivalence automaton). We will call the automaton  $E_A = (Q_E, \Sigma_E, \delta_E, s_E, F_E)$  created by ComEquivPairs and RemEquivPairs(A) the  $\sim_A$ -equivalence automaton of A. It has the following properties:

```
 \bullet \ Q_E = \{ \ [q] \mid q \in Q \ \} 
 \bullet \ \Sigma_E = \Sigma 
 \bullet \ \delta_E([q], \sigma) = [\delta(q, \sigma)], \forall q \in Q, \sigma \in \Sigma 
 \bullet \ s_E = [s]_{\sim} 
 \bullet \ F_E = \{ \ [q] \mid q \in Q \ \}
```

Note, that we know by theorem 3 that  $E_A$  is guaranteed to be a minimal automaton to A, if A has no unreachable states.

#### 1.2 Requirements analysis

Now that we have introduced all necessary basic definitions, we shall do a short analysis of an example DFA minimization task and its sample solution, as it could have been given to students in an introductory course to automata theory.

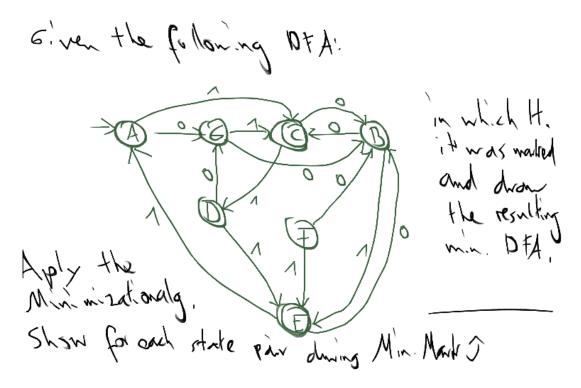


Figure 1.2: An example DFA minimization task.

#### 1.2.1 Example of a DFA minimization task for students

Gregor: search for typical task in standard text books

Figures 1.2 and 1.3 show such a task and solution. The students are confronted with a task DFA  $A_{task}$ . Firstly, unreachable states have to be eliminated, we then gain the reachable DFA  $A_{re}$ . Secondly equivalent state pairs of  $A_{re}$  are merged such that the minimal solution DFA  $A_{sol}$  is found. The table T displayed in figure 1.3 is nothing else but a visualization of the function m, whereas  $T(q_0, q_1) = i \Leftrightarrow (q_0, q_1) \in m(i)$ .

We do some rather formal statements and requirements. Firstly, we can state that

- $A_{re} = \text{RemUnreachables}(A_{task}, \text{ComUnreachables}(A_{task}))$  and
- $A_{sol} = \text{RemEquivPairs}(A_{re}, \text{ComEquivPairs}(A_{re}))$

Therefore  $A_{sol}$  is minimal regarding  $A_{re}$  and  $A_{task}$ . Secondly the languages of  $A_{task}$ ,  $A_{re}$  and  $A_{sol}$  are be equal. We know that Comequiveral requires  $A_{re}$  to be complete and that Remequiveral creates complete DFAs, so  $A_{sol}$  is complete too. Furthermore we know that every state of  $A_{re}$  is reachable since it is the output of Remunreachables.

#### 1.2.2 Discussion of difficulty adjustment possibilities

Concerning the execution of MINIMIZEDFA we find that its difficulty can be classified through various classification numbers.

ComEquivPairs-depth  $(D(A_{task}))$ . Consider the computation of the sets m(i) in ComEquivPairs. Determining m(0) is quite straightforward, because it consists simply of tests whether two states are in  $F \times Q \setminus F$  (see 0, line 3). Determining m(1) is less easy: The rule for determining all m(i), i > 0 is different to that for m(0) and more complicated (see 0, line 6). Determining m(2) requires the same rule. It shows nonetheless a students understanding of the terminating behavior of ComEquivPairs: It does not

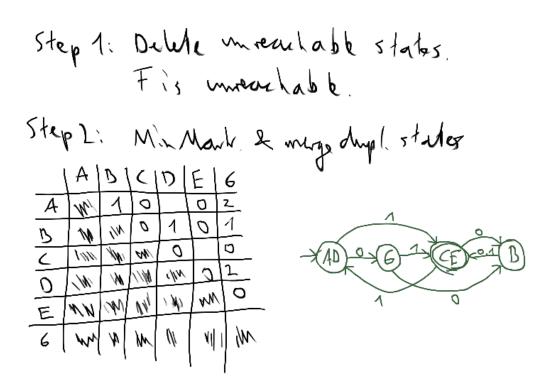


Figure 1.3: Solution to the DFA minimization task in fig. 1.2.

stop after computing m(1), but only when no more distinguishable state pairs were found. Concerning the sets m(i), i > 2 however no additional understanding can be shown.

It would therefore be sensible if  $D(A_{task})$  could be adjusted for example by parameters  $m_{min}$ ,  $m_{max}$  which give lower and upper bounds for that value.

Number of states  $(Q_{sol}, Q_{eq}, Q_{unr})$ . To control the number of states in  $A_{task}, A_{re}$  and  $A_{sol}$ , we will introduce three parameters:  $Q_{sol}, Q_{eq}, Q_{unr} \in \mathbb{N}$ . These parameters get their meaning by the following equations:

$$\begin{aligned} |Q_{sol}| &= \mathcal{Q}_{sol} \\ |Q_{re}| &= \mathcal{Q}_{sol} + \mathcal{Q}_{eq} \\ |Q_{task}| &= \mathcal{Q}_{sol} + \mathcal{Q}_{eq} + \mathcal{Q}_{unr} \end{aligned}$$

It is sensible to have  $Q_{unr} > 1$ ,  $Q_{eq} > 1$ , such that REMUNREACHABLES and REMEQUIV-PAIRS will not be skipped. To not make the task trivial,  $Q_{sol} > 2$  is sensible. An exercise instructor will find it useful, to control exactly how big  $Q_{unr}$ ,  $Q_{eq}$  and  $Q_{sol}$  are: The higher  $Q_{unr}$ ,  $Q_{eq}$ , the more states have to be eliminated and merged. The higher  $Q_{sol} + Q_{eq}$ , the more state pairs have to be checked during ComEquivPairs.

**Gregor:** Concept for min max.

**Alphabet size.** The more symbols the alphabet of  $A_{task}$ ,  $A_{re}$  and  $A_{sol}$  has (note how MINIMIZEDFA does not change the alphabet), the more transitions have to be followed when checking whether  $(\delta(q,\sigma),\delta(p,\sigma)) \in m(i-1)$  is true for each state pair p,q.

Number of final states. Since most DFAs in teaching have about 1 to 3 final states, so being able to set a number of final states, that is familiar to students, might be a reason. We devote the parameter f to this number.

Uniqueness of  $A_{sol}$ . For example for an exam it would be sensible to be able to generate a task, so  $A_{sol}$  and  $A_{task}$ , that has not been generated before. We will use the criterion of practical isomorphy to distinguish between used and unused DFAs. This is a sensible criterion, since state and symbol labels are not relevant to understanding the minimization algorithm and can easily be adjusted afterwards.

If  $A_{sol}$  is already not practical isomorph to all previously generated DFAs, then  $A_{task}$  will automatically be different too, since  $A_{sol}$ ,  $A_{task}$  describe a unique language then (see cor. 1).

Completeness of  $A_{task}$ . Even though Comune Comune and Remune Achables and Remune Chables do not require their input DFA  $A_{task}$  to be complete, it is sensible to build it that way. The implications of the completeness-property are - in comparison to the other concepts involved here - rather subtle. This is especially due to its purely representational nature, a DFA has the same language and D-value, whether it is represented in its complete form or not. Nonetheless we shall introduce a parameter c, that determines if there exist unreachable states, that make  $A_{task}$  incomplete. Thus an exercise lecturer could showcase this matter on a DFA and generate according exercises.

**Planar drawing of**  $A_{task}$ . A graph G is planar if it can be represented by a drawing in the plane such that its edges do not cross. Such a drawing is then called planar drawing of G. A visual aid for students would be given, if the task DFA were planar and presented as a planar drawing. In this work libraries and parameters  $p_1, p_2 \in \{0, 1\}$  (toggling planarity of  $A_{sol}, A_{task}$ ) will be used to allow the option of planarity, but neither ensuring planarity nor planar drawing will be investigated further theoretically.

Maximum degree of any state in  $A_{task}$ . The degree deg(q) of a state  $q \in Q$  in a DFA A is defined as  $deg(q) = |d^-(q)| + |d^+(q)|$ , so the total number of transitions in which q participates. By capping the maximum degree for all states, the graphical representation of the DFA would be more clear. In this work the inclusion of a maximum degree parameter is omitted.

### 1.3 Approach and general algorithm

In this work we will first build the solution DFA (step 1), and - based on that - the task DFA by creating equivalent states and adding unreachable states (step 2). Both steps will fulfill all criteria chosen above and are covered in depth in chapter 2 respectively chapter 3.

We will see that D and L of both DFAs will be set when building  $A_{sol}$ . We know that creating equivalent states and adding unreachable does not change  $L(A_{task})$  in comparison to  $A_{sol}$ , else MinimizeDFA would not work (a minimal DFA has in particular the same language as the original DFA). However we must ensure, that adding those states does not change D. Since unreachable states are eliminated before Comequiveralers is applied, we need only to prove, that creating equivalent states does not change the D-value. We will do this during the discussion of step 2, more specifically in section 3.1.4.

At the beginning of chapter 2 and 3, we will provide formal problem definitions for both steps, that specify precisely all requirements. Here we shall content ourselves with the definition of the main algorithm:

- 1: function GenerateDFAMINIMIZATIONPROBLEM( $Q_{sol}, a, f, m_{min}, m_{max}, p_1, p_2, Q_{eq}, Q_{unr}, c$ )
- 2:  $A_{sol} \leftarrow \text{GENERATENEWMINIMALDFA}(Q_{sol}, a, f, m_{min}, m_{max}, p_1)$
- 3:  $A_{task} \leftarrow \text{ExtendMinimalDFA}(A_{sol}, p_2, \mathcal{Q}_{eq}, \mathcal{Q}_{unr}, c)$
- 4: return  $A_{sol}, A_{task}$

### Chapter 2

# Generating minimal DFAs

We seek algorithms for generation of minimal DFAs that fulfilling the conditions defined in the requirements analysis section 1.2.2. We formally subsume these conditions via the GenerateNewMinimalDFA-problem:

**Definition 4** (GenerateNewMinimalDFA).

Given:

```
\mathcal{Q}_{sol} \in \mathbb{N} number of states a \in \mathbb{N} alphabet size f \in \mathbb{N} number of final states m_{min}, m_{max} \in \mathbb{N} lower and upper bound for D-value p \in \{0,1\} planarity-bit
```

<u>Task:</u> Compute, if it exists, a solution DFA  $A_{sol}$  with

- $|Q_{sol}| = Q_{sol}$ ,  $|\Sigma_{sol}| = a$ ,  $|F_{sol}| = f$
- $m_{min} \le D(A_{sol}) \le m_{max}$
- $A_{sol}$  being planar iff p=1
- $\bullet$   $A_{sol}$  being new

We consider different approaches to solve this problem, of which those using trial-and-error will be discussed most broadly.

Note that the presented algorithms will not be able to compute all of  $\mathfrak{A}_{min}$  since we are going to exclude minimal DFAs that are practical isomorph to already found ones.

### 2.1 Using trial and error

We will develop an algorithm that makes partly use of the trial-and-error paradigm to find matching DFAs. The approach here is as follows:

Firstly a test DFA  $A_{test}$  is generated by use of either randomness or enumeration. Alphabet size and number of (final) states will already be correct. On this DFA then tests will be executed, to check if it is minimal, planar (if wished) and new. If this is the case,  $A_{test}$  will be returned, if not, new test DFAs are generated until all tests pass.

By constructing test DFAs with already correct alphabet size and number of (final) states we are able to subdivide the search space of DFAs in advance into much smaller pieces which are in particular finite.

**Gregor:** How much smaller? Why now finite?

```
1: function BuildnewMinimalDFA-1 (Q_{sol}, a, f, m_{min}, m_{max} \in \mathbb{N}, p \in \{0, 1\})
        while True do
 2:
            generate DFA A_{test} with |Q|, |\Sigma|, |F| matching Q_{sol}, a, f
 3:
            if A_{test} not minimal or not m_{min} \leq D(A_{test}) \leq m_{max} then
 4:
                continue
 5:
            if p = 1 and A_{test} is not planar then
 6:
                continue
 7:
            if A_{test} is not new then
 8:
                continue
 9:
10:
            return A_{test}
```

We will complete this algorithm by resolving how the tests in lines 4, 6 and 8 work and by showing two methods for generation of automatons with given restrictions of |Q|,  $|\Sigma|$  and |F|.

#### 2.1.1 Ensuring $A_{test}$ is minimal and $D(A_{test})$ is correct

In order to test, whether  $A_{test}$  is minimal, we could simply use the minimization algorithm and compare the resulting DFA and  $A_{test}$  using an isomorphy test. However it is sufficient to ensure, that no equivalent or unreachable states exist. **Gregor:** minimality planarity complete under isomorphy

To get  $D(A_{test})$ , we have to run ComEquivPairs entirely anyway. Hence we can combine the test for equivalent states with computing the DFAs D-value:

```
1: function HasEquivalentStates(A)
          depth \leftarrow 0
 2:
          M \leftarrow \{(p,q), (q,p) \mid p \in F, q \notin F\}
 3:
          do
 4:
               depth \leftarrow depth + 1
 5:
               M' \leftarrow \{(p,q) \mid (p,q) \not\in M \land \exists \sigma \in \Sigma \colon (\delta(p,\sigma),\delta(q,\sigma)) \in M\}
 6:
                M \leftarrow M \cup M'
 7:
          while M' \neq \emptyset
 8:
          hasDupl \leftarrow |\{(p,q) \mid p \neq q \land (p,q) \notin M\}| > 0
 9:
          return hasDupl, depth
10:
```

Since ComEquivPairs basically computes all distinguishable state pairs  $\not\sim_A$ , we test in line 9, whether there is a pair of distinguishable states not in  $\not\sim_A$ .

Regarding the unreachable states, we can just use Comune and test whether the computed set is empty:

```
1: function HasUnreachableStates(A)
2: return |ComUnreachables(A)| > 0
```

#### 2.1.2 Ensuring $A_{test}$ is planar

There exist several algorithms for planarity testing of graphs. In this work, the library  $pygraph^1$  has been used, which implements the Hopcroft-Tarjan planarity algorithm. More information on this can be found for example in this [18] introduction from William Kocay. The original paper describing the algorithm is [13].

<sup>1</sup>https://github.com/jciskey/pygraph

#### 2.1.3 Ensuring $A_{test}$ is new

In our requirements we stated, that we wanted the generated solution DFA to be new, meaning not practically isomorph to any previously generated solution DFA. This implies the need of a database, that allows saving and loading DFAs. We name this database DB1. Assuming the database is relational, the following scheme is proposed:

```
|Q_A| |\Sigma_A| |F_A| D(A) isPlanar(A) encode(A)
```

With this scheme we can fetch once all DFAs matching the search parameters. Thus we need not fetch all previously found DFAs every time, but only those that are relevant. Afterwards we must only check whether any practical isomorphy test on the current test DFA and one of the fetched DFAs is positive. If any test DFA passes all tests and is going to be returned, then we have to save that DFA in the database.

A more concrete specification of this proceeding is shown below, embedded in the main algorithm:

```
1: function BUILDNEWMINIMALDFA-2 (Q_{sol}, a, f, m_{min}, m_{max}, p)

2: l \leftarrow all DFAs in DB1 matching Q_{sol}, a, f, m_{min}, m_{max}, p

3: while True do

4: ...

5: if A_{test} is practical isomorph to any DFA in l then

6: continue

7: save A_{test} and its respective properties in DB1

8: return A_{test}
```

#### 2.1.4 Option 1: Generating $A_{test}$ via Randomness

We now approach the task of generating a random DFA whereas alphabet and number of (final) states are set.

Corollary ?? tells us, that the states names are irrelevant for the minimality of a DFA, therefore we will give our generated DFAs simply the states  $e_1, \ldots, e_{\mathcal{Q}_{sol}}$ . For alphabet symbols this is not given. But since we **Gregor:** TODO minimality and planarity complete under isomorphy

We can state, that our start state is  $e_1 \in Q$ , since we can apply an isomorphism to every DFA, such that its start state is renamed to  $e_1$ .

The remaining elements that need to be defined are  $\delta$  and F. The set of final states is supposed to have a size of f and be a subset of Q. Therefore we can simply choose randomly f distinct states from Q.

The transition function has to make the DFA complete, so we have to choose an "end" state for every combination in  $Q \times \Sigma$ . There is no restriction as to what this end state shall be, so given  $e \in Q$  and  $\sigma \in \Sigma$  we can randomly choose an end state from Q.

With defining how to compute  $\delta$  we have covered all elements of a DFA.

```
1: function BUILDNEWMINIMALDFA-3A (Q_{sol}, a, f, m_{min}, m_{max}, p)

2: l \leftarrow all DFAs in DB1 matching Q_{sol}, a, f, m_{min}, m_{max}, p

3: Q \leftarrow \{e_1, \dots, e_{Q_{sol}}\}

4: \Sigma \leftarrow \{\sigma_1, \dots, \sigma_a\}

5: while True do

6: \delta \leftarrow \emptyset
```

```
for e in Q do
 7:
                   for \sigma in \Sigma do
 8:
                       e' \leftarrow \text{random chosen state from } Q
 9:
                       \delta \leftarrow \delta \cup \{((e, \sigma), e')\}
10:
11:
              s \leftarrow e_1
              F \leftarrow \text{random sample of } f \text{ states from } Q
12:
              A_{test} \leftarrow (Q, \Sigma, \delta, s, F)
13:
              if A_{test} not minimal or not m_{min} \leq D(A_{test}) \leq m_{max} then
14:
15:
16:
              if p = 1 and A_{test} is not planar then
17:
                   continue
18:
              if A_{test} is isomorph to any DFA in l then
                   continue
19:
20:
              save A_{test} and its respective properties in DB1
              return A_{test}
21:
```

#### 2.1.5 Option 2: Generating $A_{test}$ via Enumeration

The second method of test DFA generation is based on the idea, that instead of randomly generating F and  $\delta$ , we could just enumerate through all possible final state sets and transition functions.

Both enumerations are finite, given  $Q_{sol}$  and a. Having a requirement of f final states, then  $Q_{sol}$  choose f is the number of possible F-configurations. On the other hand there are  $Q_{sol}^{Q_{sol}*a}$  possible  $\delta$ -configurations: We have to choose one of  $Q_{sol}$  possible end states for every combination in  $Q \times \Sigma$  - so  $Q_{sol}*a$  times.

Again we will call our states and symbols w.l.o.g.  $e_1, \ldots, e_{\mathcal{Q}_{sol}}$  resp.  $\sigma_1, \ldots, \sigma_a$ . We will represent the state of an enumeration with two fields  $F_F$  and  $F_\delta$ . The first field shall have  $\mathcal{Q}_{sol}$  Bits, whereas Bit  $F_F[i] \in [0,1]$  represents the information, whether  $e_i$  is a final state or not. The second field shall have  $\mathcal{Q}_{sol} * a$  entries containing state indices, such that entry  $F_\delta[i*a+j] = k, k \in [1, \mathcal{Q}_{sol}]$  says, that  $\delta(e_i, \sigma_j) = e_k$ . These semantics are illustrated in figure 2.1.

Given an enumeration state  $b_f, b_t$  and  $Q_{sol}, a, f$  we will then compute the next DFA based on this state as follows. We will treat both fields as numbers,  $F_f$  as binary and  $F_\delta$  as  $Q_{sol}$ -ary. To get to the next DFA, we will first increment  $F_\delta$  by 1. If  $F_\delta = 1 \dots 1$ , then we increment  $F_F$  until it contains f ones (again) and set  $F_\delta$  to  $0 \dots 0$ . This behavior is summarized in the following algorithm: **Gregor:** Clarify what happens at 11111...

```
1: function IncrementEnumProgress (F_F, F_\delta, \mathcal{Q}_{sol}, a, f)
2:
        add 1 to (F_{\delta})_{\mathcal{Q}_{sol}}
        if F_{\delta} = 0 \dots 0 then
3:
              while \#_1(F_F) \neq f do
                                                                           \triangleright if the number of 1s in F_F is not f
4:
                  add 1 to (F_F)_2
                  if F_F = 0 \dots 0 then
6:
7:
                       \mathbf{return} \perp
                  F_{\delta} = 0 \dots 0
8:
        return F_F, F_\delta
```

The example in figure 2.2 illustrates such increments.

Based on the incremented bit-fields the new DFA can be build according to the semantics defined above:

Eiven: 
$$q = 4$$
,  $a = 2$ ,  $f = 3$ 

example  $b_f$ :

1101

50

 $b_f [0] = 1$ 
 $b_f [1] = 1$ 
 $b_f [2] = 0$ 
 $b_f [3] = 1$ 
 $f = \{q_3, q_2, q_3\}$ 

Start

 $f = \{q_3, q_2, q_3\}$ 

So e.g.

 $f = \{q_3, q_2, q_3\}$ 

Figure 2.1: Example for two possible configurations of the fields  $F_F$  and  $F_\delta$  given  $\mathcal{Q}_{sol}$ , a and f. Below the corresponding DFA is drawn.

Figure 2.2: The upper half shows how a  $F_{\delta}$ -increment results in a change in the resulting DFAs transition function:  $\delta(e_3, \sigma_1) = e_1$  becomes  $\delta(e_3, \sigma_1) = e_2$ . The lower half shows what happens, if  $F_{\delta}$  has reached its end.

```
1: function DFAFROMENUMPROGRESS (F_F, F_\delta, \mathcal{Q}_{sol}, a, f)
            Q \leftarrow \{e_1, \dots, e_{\mathcal{Q}_{sol}}\}
            \Sigma \leftarrow \{\sigma_1, \ldots, \sigma_a\}
 3:
 4:
            for i in [1, \ldots, \mathcal{Q}_{sol}] do
 5:
                  for j in [1,\ldots,a] do
 6:
                        \delta(e_i, \sigma_j) = e_{F_{\delta}[i*a+j]}
 7:
 8:
            F \leftarrow \{ e_i \mid i \in [1, \dots, \mathcal{Q}_{sol}] \land F_F[i] = 1 \}
 9:
            return (Q, \Sigma, \delta, s, F)
10:
```

The initial field values are each time 0...0. Note how construction and use of these fields results in DFAs with correct alphabet size and number of (final) states. An enumeration can finish either because a matching DFA has been found or all DFAs have been enumerated.

Once the enumeration within a call of BuildnewMinimalDFA has been finished, it is reasonable to save the enumeration progress (meaning the current content of  $F_F$ ,  $F_\delta$ ), such that during the next call enumeration can be resumed from that point on. The alternative would mean, that the enumeration is run in its entirety until that point again, whereas all so far found DFAs would be found to be not new. Thus we introduce a second database DB2 with the following table:

$$|Q_A|$$
  $|\Sigma_A|$   $F_F$   $F_\delta$ 

We reduce the enumeration room for each calculation.

```
1: function BUILDNEWMINIMALDFA-3B (Q_{sol}, a, f, m_{min}, m_{max}, p)
         l \leftarrow \text{all DFAs in DB1 matching } \mathcal{Q}_{sol}, a, f, m_{min}, m_{max}, p
         F_F, F_\delta \leftarrow \text{load enumeration progress for } \mathcal{Q}_{sol}, a, f, p \text{ from DB2}
 3:
         while True do
 4:
             if F_F, F_\delta is finished then
 5:
                  save F_F, F_\delta
 6:
                  return \perp
 7:
             A_{test} \leftarrow \text{next DFA based on } F_F, F_{\delta}
 8:
             if A_{test} not minimal or not m_{min} \leq D(A_{test}) \leq m_{max} then
 9:
                  continue
10:
             if p = 1 and A_{test} is not planar then
11:
12:
13:
             if A_{test} is isomorph to any DFA in l then
                  continue
14:
             save F_F, F_\delta in DB2
15:
16:
             save A_{test} and its respective properties in DB1
             return A_{test}
17:
```

### 2.2 Alternative approach: Building m(i) bottom up

Build m from m-ComEquivPairs iteratively. (Why would this basically result in running ComEquivPairs all the time?)

#### 2.3 Related research on DFA generation

Nicaud provides an overview of results on random generation and combinatorial properties of DFAs in [21]. We will outline relevant related research.

Nicaud's summary indicates, that research has focused on randomized generation of accessible, but not minimal DFAs so far. In the following we will sketch some approaches that have come up.

Using the recursive method. Champarnaud and Paranthoën [9] continue ideas started by Nicaud in his thesis [20]. Let  $\mathfrak{F}_{n,m}$  be the set of extended m-ary trees of order n. These trees are characterized by a partitioning  $V = N \uplus L$  with |N| = n and  $v \in N \Rightarrow d^+(v) = m$  and  $v \in L \Rightarrow d^+(v) = 0$ . We define the following set of tuples using s = n(m-1):

$$\mathfrak{R}_{\mathfrak{m},\mathfrak{n}} = \{ (k_1, \dots, k_s) \in \mathbb{N}^s \mid \forall i \in [2, s] \colon k_i \ge \left\lceil \frac{i}{m-1} \right\rceil \text{ and } k_i \ge k_{i-1} \}$$

In [9, p. 6] it is shown that there exists a bijection  $\varphi$  between  $\mathfrak{F}_{n,m}$  and  $\mathfrak{R}_{m,n}$  which maps to  $k_i$ ,  $i \in [1, s]$  of a tuple the number of leaves visited before the *i*th leaf in a tree. The connection to accessible DFAs is established by proving that "transition structures<sup>2</sup>" with |Q| = n,  $|\Sigma| = m$  reduced to the set of the smallest paths from the s to each other state are in bijection with extended m-ary trees of order n (see [9, p. 8]).

As a consequence they are able to construct a random generation of accessible complete DFAs using the "recursive method" from [22] which generates n-tuples [9, p. 10]. Nicaud states in his survey that the algorithm's runtime is  $\mathcal{O}(n^2)$  but notes, that generation of DFAs with more than "a few thousand states" is practically hard to do [21, pp. 10-11].

Almeida et. al. [1, 2, 23] present and implement methods using a string-encoding of DFAs for exact enumeration and random generation of DFAs. Nicaud [21, p. 11] states in a remark, that this approach uses the same recursive method and differs only in the DFA encoding.

Using Boltzmann sampler. Bassino, David and Nicaud present and implement a more efficient random generator of accessible complete DFAs in [4, 6]. Their idea is based on so called Boltzmann samplers. This framework of samplers is characterized in particular by the fact that the size of its generated objects are not fixed but in an interval around a given input size - this stands in opposition to most random generators in literature [11, p. 2].

In [6] the authors use a Boltzmann sampler to generate set partitions that are shown to be in bijection with so called box diagrams [6, p. 8] which are in turn in bijection to accessible complete DFAs [6, p. 4]. They thus acquire an average runtime complexity of  $\mathcal{O}(n^{3/2})$  for a single random generation.

Using a rejection algorithm. Carayol and Nicaud [8] give a simple algorithm with the same runtime complexity. They use a result stating that the size of accessible DFAs is concentrated around some computable value. In the end random possibly inaccessible DFAs of a specific size are generated, of which afterwards all unreachable states are deleted. This is thus essentially a rejection algorithm with clever generation of test DFAs. They furthermore show that allowing approximate sampling with the number of states being in  $[n - \varepsilon \sqrt{n}, n + \varepsilon \sqrt{n}]$  results in linear expected runtime.

<sup>&</sup>lt;sup>2</sup>Those are essentially DFAs without final state sets.

Others and comparison to algorithm presented in this work. In his survey Nicaud mentions a paper by Bassino and Sportiello [3] that yields random generation of accessible DFAs in expected linear time. This work will not be discussed further here.

In this work we use a rejection algorithm that generates test DFAs either by randomization or by enumeration. Both methods implement a naive approach. The generated test DFAs are not necessary minimal and in particular not necessary accessible as in [8]. The enumeration method uses encodings of DFAs similar to those used by Almeida et. al. [23].

#### 2.4 Empirical and combinatorial results

Concerning combinatorial properties of DFAs, several authors (e.g. [6, 10, 16]) consider a work from Vyssotsky [25] in the Bell laboratories to be the first on this subject. A contribution by Korshunov [19] is often cited in this regard, for he firstly "determines an asymptotic estimate of the number of accessible complete and deterministic n-state automata over a finite alphabet" [5].

Implementations (e.g. [1, 4]) of various random and enumeration generation methods have given rise to several empirical observations concerning the number of minimal DFAs, their fraction among all DFAs and so forth.

Domaratzki, Kisman, and Shallit [10] give some asymptotic estimates and explicit computations for the number of each several types of languages and automata that are distinct. The here relevant results have been subsumed and extended in [2, p. 8] by means of exact enumeration and are confirmed in [4].

$ \Sigma (k)$	Q (n)	$ \mathcal{A}_{min,n,k} $	$ \mathcal{A}_{n,k} $	Minimal %
k=2	2	24	64	0.38
	3	1028	5832	0.18
	4	56014	1048576	0.05
	5	3705306	312500000	0.01
	6	286717796	139314069504	0.0
	7	25493886852	86812553324672	0.0
k = 3	2	112	256	0.44
	3	41928	157464	0.27
	4	26617614	268435456	0.1
	5	25184560134	976562500000	0.03
k=4	2	480	1024	0.47
	3	1352732	4251528	0.32
	4	7756763336	68719476736	0.11
k=5	2	1984	4096	0.48
	3	36818904	114791256	0.32

Figure 2.3: Table depicting the amount of minimal complete DFAs among all complete DFAs for various sizes of  $Q, \Sigma$ . The numbers of minimal DFAs (bold numbers) are taken from [2, p. 8].

In table 2.3 we use these results to determine the ratios of minimal complete DFAs among all complete DFAs for given |Q| and  $|\Sigma|$ . The number of all DFAs is computed as follows:

$$|\mathcal{A}_{n,k}| = \underbrace{n^{n*k}}_{\text{\#possible }\delta\text{'s}} * \underbrace{2^n}_{\text{\#possible sets }F}$$

Thus we gain an insight into how probable the generation of a distinct minimal test DFA is without applying further constraints. For our proposed default parameters  $n \in [4-5]$ 

and  $k \in [2-3]$  the probabilities of successful generation range from 1% to 5%. Practical tests have shown that this leads to sufficient short run times for our implementation.

Further interesting results in this area include the determination of the fraction of minimal automata among accessible complete DFAs [5] and asymptotic estimates for the number of states that a random minimized DFA has [7].

### Chapter 3

# Extending minimal DFAs

We firstly define a formal problem for extending a minimal DFA  $A_{sol}$  to a task DFA  $A_{task}$  based on our requirements analysis (see 1.2.2):

**Definition 5** (ExtendMinimalDFA).

Given:

```
A_{sol} = (Q, \Sigma, \delta, s, F) \in \mathfrak{A}_{min} solution DFA  \mathcal{Q}_{eq} \in \mathbb{N}  number of states creating equivalent state pairs  \mathcal{Q}_{unr} \in \mathbb{N}  number of unreachable states  p \in \{0,1\}  planarity-bit  c \in \{0,1\}  completeness-bit
```

Task: Compute, if it exists, a task DFA A<sub>task</sub> with

- $Q_{task} = Q_{sol} \cup \{r_1, \dots, r_{Q_{eq}}, u_1, \dots, u_{Q_{unr}}\}$
- $r_1, \ldots, r_{Q_{eq}}$  each creating an equivalent state pair
- $u_1, \ldots, u_{\mathcal{Q}_{unr}}$  unreachable
- $\Sigma_{task} = \Sigma_{sol}, \, s_{task} = s_{sol}, \, F_{task} \subseteq F_{sol}$
- $A_{task}$  being planar iff p=1
- $A_{task}$  being complete iff c = 1
- $A_{sol}$  being isomorph to MinimizeDFA $(A_{task})$

In order to fulfill these requirements we will deduce for both kinds of states how they may be added by examining their desired properties. We will show for the action of adding equivalent states, that this does not change a DFAs D-value.

#### 3.1 Creating equivalent state pairs

Step 3 and 4 of the minimization algorithm are concerned with detection and elimination of equivalent state pairs. We now want to add states  $r_1, \ldots, r_{Q_{eq}}$  to a DFA  $A_{sol}$ , gaining  $A_{re}$  with  $Q_{re} = Q_{sol} \cup \{r_1, \ldots, r_{Q_{eq}}\}$ , such that each of these states is equivalent to a state in  $A_{re}$ . Note that, for reasons of clarity, we are going to abbreviate from now on  $A_{re} = A$ ,  $Q_{re} = Q$ ,  $\sim_{A_{re}} = \sim_A$  etc.

Consider the properties  $r_1, \ldots, r_{Q_{eq}}$  must have. They are equivalent to states  $o_1, \ldots, o_{Q_{eq}}$  of A.

$$\exists r_1, \dots, r_{\mathcal{Q}_{eq}} \in Q \colon \exists o_1, \dots, o_{\mathcal{Q}_{eq}} \in Q \colon \forall i \in [1, \mathcal{Q}_{eq}] \colon r_i \sim_A o_i$$

But we know also and in particular, that each of them is equivalent a state e of  $A_{sol}$ .

$$\exists r_1, \ldots, r_{\mathcal{Q}_{eq}} \in Q \colon \forall i \in [1, \mathcal{Q}_{eq}] \colon \exists e \in Q_{sol} \colon r_i \sim_A e$$

In our algorithm, we will choose the state e for each state we add.

#### 3.1.1 Adding outgoing transitions

Regarding the outgoing transitions of any  $r_i$  equivalent to a state e, we are directly restricted by the relationship  $\forall \sigma \in \Sigma \colon [\delta(r_i, \sigma)]_{\sim_A} = [\delta(e, \sigma)]_{\sim_A}$ . Thus, when adding some  $r_i$ , we have to choose for each symbol  $\sigma \in \Sigma$  at exactly one transition (completeness requirement for A) from the following set:

$$O_{e,\sigma} = \{ ((r_i,\sigma),q) \mid q \in [\delta(e,\sigma)]_{\sim_A} \}$$

Since the solution DFA is complete and since every here added state gets a transition for every alphabet symbol, we know that every  $O_{e,\sigma} \neq \emptyset$ .

**Gregor:** Why does this not affect the eq. class of any other state?

#### 3.1.2 Adding ingoing transitions

First of all, we know, that  $r_i$  is reachable, since every state of A must be reachable, so we need to give  $r_i$  at least one ingoing transition. Doing this, we have to ensure, that any state q, that gets such an outgoing transition to  $r_i$  remains in its  $\sim$ -equivalence class.

Thus a fitting state q has to have a transition to some state in  $[r_i]_{\sim_A} = [e]_{\sim_A}$  already. So, given a state q with  $\delta(q, \sigma) = p$  and  $p \in [e]_{\sim_A}$ , we can set  $\delta(q, \sigma) = r_i$  and thus "steal" q its ingoing transition.

We see here, that q must have at least 2 ingoing transitions, else it would become unreachable. Thus we summarize:

$$I_e = \{ ((q, \sigma), p) \mid \delta(q, \sigma) = p \land p \in [e] \land d^-(p) \ge 2 \}$$

Choose at least one  $((q, \sigma), p) \in I_e$ , remove  $((q, \sigma), p)$  from  $\delta$  and add  $((q, \sigma), r_i)$ .

These finding lead us to a general requirement regarding the choice of a state e for an  $r_i$ : The equivalence class of any e has to contain at least one state with at least 2 ingoing transitions (see fig. 3.1). We establish the following notion to pin down this restriction:

$$duplicatable(q) \Leftrightarrow_{def} (\exists p \in [q]_{\sim_A} : |d^-(p)| \ge 2)$$

The number of duplicatable states in any accessible DFA A is 0 for  $|\Sigma| \leq 1$  (due to the restriction  $|d^-(p)| \geq 2$ ) and greater than 0 for  $|\Sigma| > 1$  due to the pigeonhole principle: An accessible complete DFA has  $|Q||\Sigma|$  transitions which have to be spread across |Q| states.

#### 3.1.3 The algorithm

- 1: function CreateEquivalentStatePairs $(A, \mathcal{Q}_{eq})$
- 2:  $Q \leftarrow Q_{sol}$
- 3:  $\delta \leftarrow \delta_{sol}$
- 4:  $F \leftarrow F_{sol}$
- 5:  $K \leftarrow \{ \{q\} \mid q \in Q \}$   $\triangleright$  tracks the equivalence classes of A

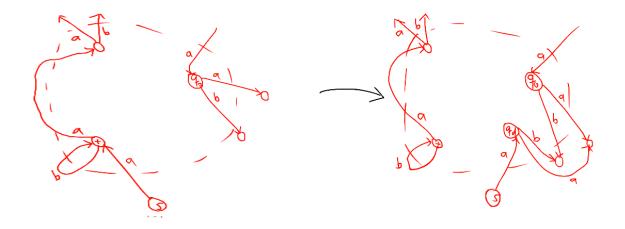


Figure 3.1: If an equivalence class (here denoted by the states in the dashed area) contains a state with 2 or more ingoing transitions (in this case p), then a state equivalent to any of the classes states may be added. Here r is equivalent to o and is "stealing" the ingoing transition  $\delta(q, a)$  from p.

```
k(q) = C such that q \in C and C \in K
                                                                              \triangleright returns the equivalence class to q
 6:
          in(q) = |d^{-}(q)| for all q \in Q
                                                                                 \triangleright tracks the number of ingoing t.
 7:
          for i in [1, \mathcal{Q}_{eq}] do
 8:
               for q in Q do
                                                                                         \triangleright find a duplicatable state e
 9:
                    if in(q) \geq 2 then
10:
11:
                         e \leftarrow \text{random chosen state from } k(q)
                         break
12:
               r_i \leftarrow \text{unused state label}
                                                                                    \triangleright create to e equivalent state r_i
13:
               Add r_i to Q
14:
               Add r_i to k(e)
15:
                                                                                                                \triangleright add d^+(r_i)
16:
               for \sigma in \Sigma do
                    \delta(r_i, \sigma) = \text{random chosen state from } k(\delta(e, \sigma))
17:
               P \leftarrow \{ ((s, \sigma), t) \in \delta \mid t \in k(e), in(t) \geq 2 \}
                                                                                                                \triangleright add d^-(r_i)
18:
               C \leftarrow \text{random nonempty subset of } P
19:
20:
               for ((s,\sigma),t) in C do
                    in(t) \leftarrow in(t) - 1
21:
                    in(r_i) \leftarrow 1
22:
                    \delta(s,\sigma) = r_i
23:
          return (Q, \Sigma_{sol}, \delta, s_{sol}, F)
24:
```

#### 3.1.4 Creating equivalent state pairs does not change D

To prove this statement, we will prove two minor propositions first. In this context we will call a word w distinguishing word of p, q, iff  $d_A(w, p, q)$  whereas

$$d_A(w, p, q) \Leftrightarrow (\delta^*(p, w) \in F \Leftrightarrow \delta^*(q, w) \notin F)$$

The following lemma and its proof are in parts inspired by [26, ch. 4 p. 18].

**Lemma 1.** In the context of ComEquivPairs the following is true: If and only if  $(p,q) \in$ 

m(n), the shortest distinguishing word of p, q has length n. Formally:

$$(p,q) \in m(n) \iff \exists w \in \Sigma^* \colon (|w| = n \land d_A(w,p,q))$$
  
  $\land \nexists v \in \Sigma^* \colon (|v| < |w| \land d_A(v,p,q))$ 

*Proof.* Per induction on the number of ComEquivPairs-iterations n.

 $n=0, "\Leftrightarrow".$ 

$$(p,q) \in m(0) = \{(p,q), (q,p) \mid p \in F, q \notin F\} \text{ (see 0, line 2)})$$
  
 $\Leftrightarrow$  one of  $p,q$  in  $F$ , one not  
 $\Leftrightarrow$  one of  $\delta^*(p,\varepsilon), \delta^*(q,\varepsilon)$  in  $F$ , one not  
 $\Leftrightarrow \exists w \in \Sigma^* \colon (|w| = 0 \land \text{one of } \delta^*(p,w), \delta^*(q,w) \text{ in } F, \text{ one not)}$   
 $\Leftrightarrow \exists w \in \Sigma^* \colon (|w| = 0 \land d_A(w,p,q))$   
and there cannot be no shorter such word  $\checkmark$ 

 $0 \dots n-1 \to n$ , " $\Rightarrow$ ". Then the following holds for some states p, q (see 0, line 5):

$$(p,q) \in m(n) = \{(p,q), (q,p) \mid (p,q) \notin \bigcup m(\cdot) \land \exists \sigma \in \Sigma \colon (\delta(p,\sigma), \delta(q,\sigma)) \in m(n-1)\}$$
(3.1)

So in particular there exists a symbol  $\sigma$  such that  $(\delta(p,\sigma), \delta(q,\sigma)) \in m(n-1)$ . Let  $(p',q') = (\delta(p,\sigma), \delta(q,\sigma))$ . Per induction there exists a word w', |w'| = n-1 to p',q' such that one of  $\delta^*(p',w'), \delta^*(q',w')$  in F, one not and there is no shorter word.

Thus one of  $\delta^*(p, \sigma w')$ ,  $\delta^*(q, \sigma w')$  is in F, one not, which makes  $\sigma w'$  a distinguishing word of length n for p, q. Since (p, q) is not in any m(i), i < n (recall eq. 3.1), there is per induction no shorter distinguishing word.  $\checkmark$ 

 $0 \dots n-1 \to n$ , " $\Leftarrow$ ". Then the following holds for some states p,q:

$$\exists w \in \Sigma^* \colon (|w| = n \land d_A(w, p, q))$$
  
 
$$\land \nexists v \in \Sigma^* \colon (|v| < |w| \land d_A(v, p, q))$$

So there exists a word w with |w| = n > 0 such that one of  $\delta^*(p, w), \delta^*(q, w)$  is in F, one not and there is no shorter word fulfilling this property.

Consequently there exists a symbol  $\sigma$  such that  $w = \sigma w'$ . Let  $(\delta(p, \sigma), \delta(q, \sigma)) = (p', q')$ . Thus, if one of  $\delta^*(p, w), \delta^*(q, w)$  is in F and one not, then the same must hold for  $\delta^*(p', w'), \delta^*(q', w)'$ . The word w' is also the shortest fulfilling this property, because, if there existed a shorter word v', |v'| < |w'|, then  $\sigma v'$  would be a distinguishing word shorter than w for p, q which is contradictory.

Per induction we can deduce, that  $(p',q') \in m(n-1)$ . The pair (p,q) is not in any m(i), i < n, since else-wise per induction the shortest distinguishing word would be shorter than w and thus not w, which is contradictory. Since  $(p',q') \in m(n-1)$  and  $(\delta(p,\sigma),\delta(q,\sigma)) = (p',q')$ , we can then deduce, that  $(p,q) \in m(n)$ .  $\checkmark$ 

**Lemma 2.** If Comequiverairs has done D(A) iterations and terminated, then the longest word w, that is a shortest distinguishing word for any state pair, has length D(A) - 1.

*Proof.* Via direct proof. Assume m-ComEquivPairs(A) has done n iterations (so D(A) = n). We then know, that

1. 
$$\forall i \in [0, n-1] : m(i) \neq \emptyset$$

$$2. m(n) = \emptyset$$

3. 
$$\forall i > n : m(i) = \bot$$

m-ComEquivPairs(A) terminates iff  $m(i) = \emptyset$ . If the first point would not hold, then the algorithm would have stopped before. Since the algorithm did n iterations, the internal variable i must be n at the end of the last iteration. The terminating condition is  $m(i) \neq \emptyset$  and no more iterations are done thereafter; thus follow the second and third point.

We prove that then there exists a shortest distinguishing word of length n-1, but the existence of a longer distinguishing word that is shortest for some state pair leads to a contradiction. Recall lemma 1:

$$(p,q) \in m(n) \iff \exists w \in \Sigma^* \colon (|w| = n \land d_A(w,p,q))$$
  
  $\land \nexists v \in \Sigma^* \colon (|v| < |w| \land d_A(v,p,q))$ 

Following this lemma and point 1, we can deduce that there exists at least one shortest distinguishing word w with |w| = n - 1 = D(A) - 1 for some  $p, q \in Q$ .

There cannot be any shortest distinguishing word w' with |w'| = k > n - 1 for any two states  $p', q' \in Q$  fulfilling this property. Following the lemma again, m(k) for some k > n - 1 would be defined and non-empty, which is contradictory to point 2 and 3.  $\square$ 

**Lemma 3.** If w is shortest distinguishing word for p, q and  $q \sim_A q'$ , then w is a shortest distinguishing word for p, q'.

*Proof.* Via direct proof. The relation  $q \sim_A q'$  says that for all words z,  $\delta^*(q, z)$  and  $\delta^*(q', z)$  are both in F or both not in F. Consequently asking whether the following is true

$$\delta^*(p,z) \in F \Leftrightarrow \delta^*(q,z) \in F$$

is the same as asking whether this is true:

$$\delta^*(p,z) \in F \Leftrightarrow \delta^*(q',z) \in F$$

This implies, that q and q' have exactly the same distinguishing words with other states (see definition of distinguishing words). As a consequence they too have the same shortest distinguishing words with other states.

**Theorem 4.** Creating equivalent state pairs in a minimal DFA A does not increase the number of iterations when the ComequivPairs-algorithm is applied on it.

*Proof.* Per contradiction. Let us assume there were n states  $r_0, \ldots, r_n$  added to a given minimal DFA  $A = (Q, \Sigma, \delta, s, F)$  resulting in a DFA  $A' = (Q', \Sigma, \delta', s, F')$  such that:

- $\bullet \ Q' = Q \cup \{r_1, \dots, r_n\}$
- $\forall i \in [1, n]: \exists q \in Q: r_i \sim_{A'} q$
- D(A) < D(A')

According to lemma 2 the longest shortest distinguishing word w of A has length |w| = D(A) - 1, while its counterpart w' in A' has length D(A') - 1. Consequently |w| < |w'|.

There exist states  $p', q' \in Q'$  such that w' distinguishes p', q'. We differentiate between three cases regarding the belonging of p', q'.

Both p', q' in Q: Since the longest shortest distinguishing word any state pair in Q can have is guaranteed shorter than w', we may conclude that  $p', q' \in Q$  and p', q' have w' as shortest distinguishing word is a contradiction.

One of p', q' in Q, one in  $Q' \setminus Q$ : W.l.o.g.  $p' \in Q$ ,  $q' \in Q' \setminus Q$ . We then know, that  $q' \in \{r_1, \ldots, r_n\}$ . This implies that  $q' \sim_{A'} q$  for an  $q \in Q$ .

By lemma 3 w' is a shortest distinguishing word for p', q. But this is contradictory to  $p', q \in Q$ . 4

**None of** p', q' **in** Q: Since p', q' both have to be in  $Q' \setminus Q = \{r_1, \ldots, r_n\}$ , we can find states  $p, q \in Q$  equivalent to p', q'. For these states w' would be a shortest distinguishing word, which is contradictory again. 4

#### 3.2 Adding unreachable states

From step 1 of the minimization algorithm we can deduce how to add unreachable states. These can easily be added to a DFA by adding non-start states with no ingoing transitions (see def. 1). Number and nature of outgoing transitions may be arbitrary.

```
1: function AddUnreachableStates (A, \mathcal{Q}_{unr}, c)
 2:
          for Q_{unr} times do
               q \leftarrow \text{unused state label}
 3:
 4:
               Q \leftarrow Q \cup \{q\}
               outSymbols \leftarrow c = 1 ? \Sigma: random subset of \Sigma
 5:
               R \leftarrow \text{random chosen sample of } |outSymbols| \text{ states from } Q \setminus \{q\}
 6:
               for \sigma in outSymbols do
 7:
                    q' \in R
 8:
                    R \leftarrow R \setminus \{q'\}
 9:
                    \delta \leftarrow \delta \cup \{((q, \sigma), q')\}
10:
          return A
11:
```

# Chapter 4

# Conclusion

What happens, if we change start and accepting states? What happens, if we add transitions only?

dfa specific planarity test? use planarity test information for better drawing?

### Appendix A

# **Bibliography**

- [1] André Almeida, Marco Almeida, José Alves, Nelma Moreira, and Rogério Reis. Fado and guitar. volume 5642, pages 65–74, 07 2009. doi:10.1007/978-3-642-02979-0\_10.
- [2] Marco Almeida, Nelma Moreira, and Rogério Reis. Aspects of enumeration and generation with a string automata representation. *Theoretical Computer Science*, 387:93–102, 06 2009. doi:10.1016/j.tcs.2007.07.029.
- [3] Frederique Bassino and Andrea Sportiello. Linear-time generation of specifiable combinatorial structures: general theory and first examples, 2013. arXiv:1307.1728.
- [4] Frédérique Bassino, Julien David, and Cyril Nicaud. Regal: A library to randomly and exhaustively generate automata. 07 2007. doi:10.1007/978-3-540-76336-9\_28.
- [5] Frédérique Bassino, Julien David, and Andrea Sportiello. Asymptotic enumeration of minimal automata. *Leibniz International Proceedings in Informatics*, *LIPIcs*, 14, 09 2011. doi:10.4230/LIPIcs.STACS.2012.88.
- [6] Frédérique Bassino and Cyril Nicaud. Enumeration and random generation of accessible automata. *Theoretical Computer Science*, 381:86–104, 08 2007. doi: 10.1016/j.tcs.2007.04.001.
- [7] Daniel Berend and Aryeh Kontorovich. The state complexity of random dfas. *Theoretical Computer Science*, 652, 07 2013. doi:10.1016/j.tcs.2016.09.012.
- [8] Arnaud Carayol and Cyril Nicaud. Distribution of the number of accessible states in a random deterministic automaton. volume 14, pages 194–205, 02 2012. doi: 10.4230/LIPIcs.STACS.2012.194.
- [9] Jean-Marc Champarnaud and Thomas Paranthoën. Random generation of dfas. *Theoretical Computer Science*, 330:221–235, 02 2005. doi:10.1016/j.tcs.2004.03.072.
- [10] Michael Domaratzki, Derek Kisman, and Jeffrey Shallit. On the number of distinct languages accepted by finite automata with n states. *J. Autom. Lang. Comb.*, 7(4):469–486, September 2002.
- [11] Philippe Duchon, Philippe Flajolet, Guy Louchard, and Gilles Schaeffer. Boltzmann samplers for the random generation of combinatorial structures. Comb. Probab. Comput., 13(4-5):577-625, July 2004. URL: https://doi.org/10.1017/S0963548304006315, doi:10.1017/S0963548304006315.

- [12] Gesellschaft für Informatik. Empfehlungen für bachelor- und masterprogramme im studienfach informatik an hochschulen (juli 2016), 2016. Accessed: 2020-02-12. URL: https://dl.gi.de/bitstream/handle/20.500.12116/2351/58-GI-Empfehlungen\_Bachelor-Master-Informatik2016.pdf.
- [13] John Hopcroft and Robert Tarjan. Efficient planarity testing. J. ACM, 21(4):549–568, October 1974.
- [14] John E. Hopcroft. An n log n algorithm for minimizing states in a finite automaton. Technical report, Stanford, CA, USA, 1971.
- [15] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to automata theory, languages, and computation (2. ed.)*. Addison-Wesley series in computer science. Addison-Wesley-Longman, 2001.
- [16] Pierre-Cyrille Héam and Jean-Luc Joly. On the uniform random generation of deterministic partially ordered automata using monte carlo techniques. 12 2014.
- [17] Johanna Högberg and Lars Larsson. Dfa minimisation using the myhill-nerode theorem. Accessed: 2020-02-09. URL: http://www8.cs.umu.se/kurser/TDBC92/VT06/final/1.pdf.
- [18] William Kocay. The hopcroft-tarjan planarity algorithm. October 1993.
- [19] A. Korshunov. Enumeration of finite automata. *Problemy Kibernetiki*, page 34:5–82, 1959. In russian.
- [20] C. Nicaud. Étude du comportement en moyenne des automates finis et des langages rationnels. PhD thesis, Université Paris 7, 2000.
- [21] Cyril Nicaud. Random deterministic automata. pages 5-23, 08 2014. doi:10.1007/978-3-662-44522-8\_2.
- [22] Albert Nijenhuis and Herbert S. Wilf. *Combinatorial Algorithms*. Academic Press, 1978.
- [23] Rogério Reis, Nelma Moreira, and Marco Almeida. On the representation of finite automata. 7th International Workshop on Descriptional Complexity of Formal Systems, DCFS 2005 Proceedings, 06 2005.
- [24] Uwe Schöning. Theoretische Informatik kurzgefasst. Spektrum, Akad. Verl, Heidelberg Berlin, 2001.
- [25] V. Vyssotsky. A counting problem for finite automata. Technical report, 1959. Bell Telephon Laboratories.
- [26] Thomas Schwentick Wim Martens. Theoretische informatik i ss18. Lecture notes, 2018.

# Appendix B

# Erklärung

Hiermit versichere ich, Gregor Hans Christian Sönnichsen, dass ich die vorliegende Arbeit selbständig verfasst habe, keine anderen als die von mir angegebenen Quellen und Hilfsmittel benutzt habe und die Arbeit nicht bereits zur Erlangung eines akademischen Grades eingereicht habe.

Bayreuth, den 8. Februar 2020.

Gregor Hans Christian Sönnichsen