

CDM

Fast Minimization

Klaus Sutner
Carnegie Mellon University

24-minimization 2017/12/15 23:19



- 1 Partition Refinement
- 2 Hopcroft's Algorithm
- 3 Valmari-Lehtinen
- 4 Equivalence Testing
- 5 Characterizations of Regularity

Battleplan

3

Minimization is a good example where computation forces one to think more carefully than in math alone.

Mathematical Thinking: behavioral equivalence. Once the concept of behavior is clear, there is a straightforward brute-force algorithm for minimization. And, it's even polynomial time.

Algorithmic Thinking: refinement of equivalence relations. A better algorithm is obtained by thinking clearly about computing with equivalence relations (Moore). The reward is a clean, quadratic time algorithm (which is often much better than quadratic).

Smart Algo Thinking: baby-steps vs. giant-steps. Now things get tricky: all sub-quadratic algorithms require a much more careful argument and deeper algorithmic methods. A bit of creative insight is required to get down to log-linear. And doing things elegantly and efficiently is quite difficult.

Aside: Partitions versus Equivalence Relations

4

We will switch back and forth between two natural representations of the same concept.

Equivalence Relations

A relation $\rho \subseteq A \times A$ that is reflexive, symmetric and transitive.

Partition

A collection B_1, B_2, \dots, B_k of pairwise disjoint, non-empty subsets of A such that $\bigcup B_i = A$ (the blocks of the partition).

As always, we need to worry about appropriate data structures and algorithms that operate on these data structures.

Total Recall

5

Recall our abstract scenario: we have an equivalence relation $\rho \subseteq Q \times Q$ and an endofunction $f : Q \rightarrow Q$.

We want to find the coarsest refinement $\hat{\rho}$ of ρ that is compatible with f :

$$p \hat{\rho} q \Rightarrow f(p) \hat{\rho} f(q)$$

This is accomplished by repeated application of a refinement operator R_f :

$$p \rho_f q \Leftrightarrow f(p) \rho f(q)$$

$$R_f(\rho) = \rho \sqcap \rho_f$$

In other words: $\hat{\rho}$ is the fixed point of ρ under R_f .

Moore's Algorithm

6

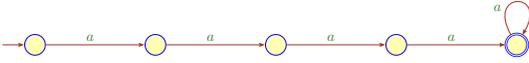
The refinement operator in Moore's algorithm works by representing all relations as canonical selector functions (aka int arrays) and scans the whole arrays for each step.

So each refinement step is $\Theta(n)$ (with good constants but still linear in n).

Quite often the algorithm uses fewer than n rounds, so the total time complexity may well be sub-quadratic.

Alas, there are cases when Moore requires $\Theta(n)$ rounds, producing quadratic running time overall.

Here is the standard example that demonstrates that Moore's algorithm may be quadratic: the minimal DFA for $\{a\}^{\geq k}$.



For this automaton, a single Moore round will split off only one state from the right end of block $D = \{0, 1, \dots, r\}$, at a cost of $\Theta(n)$ steps.

The split occurs only because of block $B = \{r + 1\}$, nothing else matters.

Critical Idea: Maybe we could get mileage out of trying to guide the refinement by single blocks, instead of blindly hitting the whole carrier set.

Suppose ρ is a partition of Q , and consider two blocks D and B . Let $f : Q \rightarrow Q$ be some endofunction.

We say that B splits D if

$$D \cap f^{-1}(B) \neq \emptyset \quad \text{and} \quad D - f^{-1}(B) \neq \emptyset.$$

In other words, $f(D)$ intersects both B and $Q - B$ and is not f -compatible yet. At this point we cannot factor f to get a function on Q/ρ (in the application to finite state machines this would result in nondeterminism).

Let's define a new, more complicated refinement operator $\rho' = R_f(\rho, D, B)$ as follows:

$$p \rho' q \iff (p, q \notin D \wedge p \rho q) \vee (p, q \in D \wedge (f(p) \in B \iff f(q) \in B))$$

In other words: outside of block D we keep the old ρ . Inside of D we check for B -equivalence of children.

So $R_f(\rho, D, B)$ is indeed a refinement of ρ : block D is split in two.

Proposition

- $R_f(\rho) \sqsubseteq R_f(\rho, D, B)$.
- $R_f(\rho) \neq \rho$ implies that $R_f(\rho, D, B) \neq \rho$ for some D and B .

In other words, we make no mistakes and we can't get stuck.

Proof.

$R_f(\rho)$ is $\prod_{B,D} R_f(\rho, D, B)$ and thus finer than each part.

If $R_f(\rho) \neq \rho$ there must be some block D and $p, q \in D$ such that $\neg(f(p) \rho f(q))$.

Let B be the block containing $f(p)$, done. □

Of course, from a complexity perspective this may not sound too promising: we are breaking one giant step into multiple baby steps. It is not unreasonable to suspect that this might even increase running time.

But: The baby steps provide much better control over the selection of the next refinement step: we can choose the blocks involved at will.

With a little effort this feature can be exploited sufficiently to speed up the whole process.

- Partition Refinement
- Hopcroft's Algorithm
- Valmari-Lehtinen
- Equivalence Testing
- Characterizations of Regularity

Suppose we have an endofunction $f : Q \rightarrow Q$ and a coloring χ of Q .

The algorithm maintains two data structures:

- a **partition** P of Q , representing the equivalence relation,
- a **split list** S with entries some of the blocks in the partition.

Both partitions are initialized by the coloring.

We refer to the blocks C in S as **active**: we will refine P by $f^{-1}(C)$.

The algorithm extracts an active block from the split list and tries to refine the blocks in the partition accordingly.

It then updates the split list in a clever way, and stops when the split list becomes empty.

initialize partitions P and S to the given coloring

```

while  $S$  not empty do
  extract  $C$  from  $S$ 
  compute  $\hat{C} = f^{-1}(C)$ 
  foreach block  $B$  split by  $C$  do
     $B^+ = B \cap \hat{C}$ 
     $B^- = B - B^+$ 
    replace  $B$  by  $B^+$  and  $B^-$  in  $P$ 
    if  $B$  is in  $S$ 
      then replace  $B$  by  $B^+$ ,  $B^-$  in  $S$ 
      else add smaller of  $B^+$ ,  $B^-$  to  $S$ 
  end

```

At first glance it may seem like we are not doing enough work: in the last case it feels like both B^+ and B^- , the parts of B obtained by splitting wrto the critical block C , should be added to the split list.

Otherwise we might miss out on splitting some block C that is not split by B but split by B^+ or B^- . The algorithm might stop without having produced an f -compatible relation.

Fortunately, this cannot happen: if B^+ splits block D then so does B^- , and conversely.

Call a set $Z \subseteq Q$ of states **safe** for partition P (or simply P -safe) if $f^{-1}(Z)$ does not split any block in P .

We will show that the following assertion is a loop invariant:

$$\forall X \in P - S \exists A \subseteq S (X \cup \bigcup A \text{ } P\text{-safe})$$

Clearly, this assertion holds before the loop ever executes.

As is customary, we indicate the value of a variable after one more execution of the loop-body by attaching a prime: so P' is the partition after one more round.

In the following we argue about the state of affairs at the end of a round.

Case 1: Block X is old: $X \in P$

Case 1.1: $X \in S$

The only way this can happen is $C = X$. But then safety is a direct result of the construction.

Case 1.2: $X \notin S$

So we have $Z = X \cup \bigcup A$ is P -safe for some $A \subseteq S$.

Case 1.2.1: $C \notin A$

In this case we can replace split blocks in A : replace X by X^+ and X^- (which are both active) and we have $Z = X \cup \bigcup A'$ is P -safe. It is easy to check that Z is also P' -safe.

Case 1.2.2: $C \in A$

Again, we can replace split blocks in A ; alas, we lose $C \notin S'$ to get $X \cup \bigcup A_0$. But no block $Z \in P'$ can be affected by this: it was already split wrto $f^{-1}(C)$ during the round.

Case 2: Block X is new: $X \in P' - P$

In this case there is a block $Z \in P$ such that, say, $X = Z^+$ and $X \cup \bigcup A$ is P -safe for some $A \subseteq S$.

As before we can handle split blocks in A .

Case 2.1: $C \notin A$

Then $X \cup (\bigcup A' \cup Z^-)$ is P' -safe as in Case 1.2.1. Note that indeed $Z^- \in S'$ by construction.

Case 2.2: $C \in A$

In this case, $X \cup (\bigcup A_0 \cup Z^-)$ is P' -safe as in Case 1.2.2.

Each block in P is represented by a doubly-linked list.

We maintain an array of pointers to these lists for P and similarly for S . We also keep track of the cardinality of each block.

Furthermore, we have an array of pointers so that $\text{pos}[p]$ points to the list node containing p , plus information about the current block containing p .

The key part of each round is the computation of $\hat{D} = f^{-1}(D)$. We may assume that $f^{-1}(p)$ has been precomputed for each state p . We can traverse D in time linear in $|D|$.

When a block B is hit for the first time, we start splitting it into two lists B^+ and B^- . If, in the end, $B^- = \emptyset$ we simply replace B by B^+ .

All this can be handled in time $O(|f^{-1}(D)|)$.

Let us say that a state p is **active** if $p \in \bigcup A$, **inactive** otherwise.

At level 0, all states are active.

Each state in the critical block C becomes inactive, but maybe reactivated later. Hence we can naturally assign activation levels 0, 1, 2, ... to all active states.

Recall that we only reactivate at most half the states in a block, so no state can be activated more than $\log n$ times. But then the total work computing pre-images of active states during the whole execution is just $O(n \log n)$.

Hence the running time of the whole algorithm is bounded by $O(n \log n)$.

We modify the split list in the algorithm to contain entries

$$(a, B) \in S$$

where B is a block and $a \in \Sigma$: the intent is that we later refine via $\delta_a^{-1}(C)$.

Of course, (a, X) is **smaller** than (a, Y) if $|X| \leq |Y|$.

Note that we need to add (a, X) for all $a \in \Sigma$, which produces a running time of $O(kn \log n)$.

The following example uses a machine over alphabet $\{a, b\}$ with 15 states. The transition matrix is

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
a	2	4	6	8	10	12	14	8	10	12	14	8	10	12	14
b	3	5	7	9	11	13	15	9	11	13	15	9	11	13	15

The final states are $\{12, 13, 14, 15\}$.

The following table shows the element extracted from the split list in the first column and the blocks in the second.

Note that split list entry (a, i) means: use the i th block in the current partition with respect to $f = \delta_a$.

split	partition
–	((12, 13, 14, 15), (1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11))
(a, 1)	((14, 15), (1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11), (12, 13))
(a, 1)	((14, 15), (7, 11), (12, 13), (1, 2, 3, 4, 5, 6, 8, 9, 10))
(a, 1)	((15), (7, 11), (12, 13), (1, 2, 3, 4, 5, 6, 8, 9, 10), (14))
(a, 4)	((15), (7, 11), (12, 13), (1, 2, 3, 4, 5, 6, 8, 9, 10), (14))
(a, 4)	((15), (7, 11), (12, 13), (1, 2, 3, 4, 5, 8, 9), (14), (6, 10))
(a, 5)	((15), (7, 11), (12, 13), (1, 2, 3, 4, 5, 8, 9), (14), (6, 10))
(a, 5)	((15), (7, 11), (12, 13), (1, 2, 3, 4, 5, 8, 9), (14), (6, 10))
(a, 6)	((15), (7, 11), (13), (1, 2, 3, 4, 5, 8, 9), (14), (6, 10), (12))
(a, 6)	((15), (7, 11), (13), (3, 5, 9), (14), (6, 10), (12), (1, 2, 4, 8))

Nothing changes beyond this point but there are several more steps since all of $(a, 7)$, $(a, 8)$, $(b, 1)$, $(b, 3)$, $(b, 2)$, $(b, 5)$, $(b, 6)$, $(b, 7)$, $(b, 8)$ are still on the split list.

The algorithm is quite messy to implement correctly, as can be seen from the following papers:

D. Gries
Describing an algorithm by Hopcroft
Acta Informatica, 2 (1973) 97–109.

T. Knuutila
Re-describing an algorithm by Hopcroft
Theoretical Computer Science, 250 (2001) 333–363.

Exercise

Implement Hopcroft's algorithm, correctly.

Theorem (Hopcroft 1971)

Hopcroft's algorithm minimizes a DFA in $O(kn \log n)$ steps, where n is the state complexity of the DFA and k the size of the alphabet.

Unsurprisingly, the running time of Hopcroft's algorithm depends on the size k of the alphabet.

Knuutila pointed out in 2001, one can produce cubic (in n) running time when $k = n/2$ and a poor method of choosing the "smaller" block is used.

Note that Hopcroft's algorithm is nondeterministic in several ways.

- We can extract any element from the split list (e.g., could use a stack, queue, ...).
- Likewise we can place the new entries anywhere in the split list.
- When B^+ and B^- have the same size, we can pick either one.

None of these choices effect correctness, but they may well influence running time.

It should be noted that the algorithm often takes far fewer than $n \log n$ steps.

Given a reasonable implementation, every round is linear. It turns out to be quite difficult to construct inputs where the algorithm requires $\log n$ many rounds.

The best result known today is that there are some DFAs such that the algorithm takes $n \log n$ steps for a certain choice of active blocks in the main loop.

Alas, for these machines a different choice of active blocks results in linear running time.

- When is the running time $\Omega(n \log n)$ regardless of the chosen split list protocol? (A unary example is known where the execution sequence is essentially unique and reaches the log-lin bound.)
- What is the average complexity of Hopcroft's algorithm (average with respect to input automaton and/or split list protocol)?
- Is Hopcroft faster than Moore on average? For the uniform distribution, Moore has expected behavior $O(n \log n)$ and it may be that the constants are smaller (Bassino, David, Nicaud 2009).

- Partition Refinement

- Hopcroft's Algorithm

- Valmari-Lehtinen

- Equivalence Testing

- Characterizations of Regularity

There are many examples where the number of transitions m in a partial DFA is much smaller than $k \cdot n$, k the size of the alphabet and n the number of states. This leads naturally to an algorithmic question:

Is there a $O(m \log n)$ minimization algorithm that deals directly with partial transition functions?

This would also nicely encapsulate problems with alphabet size in a parameter that really reflects the size of the data structure: unused symbols do not inflate the transition function.

The high-level logic is similar as in Hopcroft's algorithm: one maintains a partition of Q and tries to refine the partition by splitting with sets of the form

$$\delta_a^{-1}(B)$$

The choice of B is a bit more complicated, though.

Special care is taken to avoid unnecessary computation when preimages of blocks are computed. To this end, the algorithm also maintains and refines a second partition of the transitions.

As an implementation detail: the partition data structure is all array-based (unlike Hopcroft's algorithm).

Write $[n]_0$ for $\{0, 1, \dots, n-1\}$.

Suppose we wish to maintain a partition of $[n]_0$. Keep track of r , the number of blocks, and maintain two maps (arrays)

$$\text{lo, hi}: [r]_0 \longrightarrow [n]_0$$

plus an array $P[\mathbb{N}]$ such that

B_d , the block number d , is located in $P[\text{lo}[d], \text{hi}[d] - 1]$

We also have location and block-number maps

$$\text{loc}: [n]_0 \longrightarrow [n]_0$$

$$\text{bnum}: [n]_0 \longrightarrow [r]_0$$

such that $P[\text{loc}[p]] = p$, $\text{loc}[P[i]] = i$ and $p \in B_{\text{bnum}[p]}$.

It is convenient to subdivide the splitting process into three phases:

- Pre-splitting: initialize offset pointers $\text{mrk}[d] = \text{lo}[d]$ for all blocks, create empty hit list.
- Splitting: process a sequence of elements, swap each to the "marked" part $P[\text{lo}[d], \text{mrk}[d] - 1]$ of their respective blocks. If block B_d is encountered for the first time, add to hit list.
- Post-splitting: walk through blocks in hit list and update to maintain invariants.

It is straightforward to arrange the post-splitting phase so that whenever B splits into B_1 and B_2 , the larger part replaces B and the smaller part receives the higher block index (the current value of r).

As already mentioned, the minimization algorithm uses two PRDS:

- P : represents a partition of Q ; initialized to $(F, Q - F)$, as usual.
- T : represents a partition of the transitions; initialized to blocks containing all transitions with the same label.

The main loop of the algorithm looks like this:

```
foreach T-block C do
  split blocks in P via C
  foreach P-block B do
    split blocks in T via B
```

New blocks are "appended" in both partitions, so the traversals end when all blocks have been processed.

```
pre-split P
foreach transition p  $\xrightarrow{a}$  q in block C do
  mark source p in P-partition
post-split P
```

```
pre-split T
foreach state p in block B do
  foreach transition t : q  $\xrightarrow{a}$  p do
    mark t in T-partition
post-split T
```

Note that for the transition splitting operation one needs to be able to traverse all transitions with a fixed target.

To this end one precomputes two arrays

$$\text{trn}: [m]_0 \longrightarrow [m]_0$$

$$\text{fst}: [m] \longrightarrow [m]_0$$

such that for each state p

the transitions with target p are located in $\text{trn}[\text{fst}[p], \text{fst}[p+1] - 1]$

This is easy via counting sort.

Recall that $\text{bnum}(p)$ is the block number of state p . By abuse of notation, $\text{bnum}(t)$ is the block number of transition t .

Proposition

- *States:* $\text{bnum}(p) \neq \text{bnum}(q)$ implies $\llbracket p \rrbracket \neq \llbracket q \rrbracket$
- *Transitions:* $\text{bnum}(s) \neq \text{bnum}(t)$ and $\text{lab}(s) = \text{lab}(t)$ implies $\llbracket \text{trg}(s) \rrbracket \neq \llbracket \text{trg}(t) \rrbracket$

Proposition

Let s, t be two transitions in the same T -block.

If $\text{src}(s)$ and $\text{src}(t)$ are in different P -blocks, then at least one of the blocks is unprocessed.

Proposition

Let p, q be two states in the same P -block and $s : p \xrightarrow{a} p', t : q \xrightarrow{a} q'$ two transitions in different T -blocks. Then at least one of these blocks is unprocessed.

If $s : p \xrightarrow{a} p'$ and there is no a -transition with source q , then s is in an unprocessed T -block.

Theorem

The algorithm correctly minimizes a trim partial DFA in $O(m \log n)$ steps.

Proof.

It follows from the propositions that two states have the same behavior iff, upon completion, they are in the same P -block.

For running time, note that since we are dealing with a trim automaton, we have $n \leq m + 1$.

As in Hopcroft's algorithm, each state can be active at most $\log n$ times and, likewise, a transition can be active and at most $\log m$ times.

So the total running time is $O(m \log n)$.

- Partition Refinement
- Hopcroft's Algorithm
- Valmari-Lehtinen
- Equivalence Testing
- Characterizations of Regularity

Let's return to our old problem of testing two DFAs for equivalence.

Problem: **Equivalence**
 Instance: Two DFAs M_1 and M_2 .
 Question: Are the two machines equivalent?

We now have two algorithms for this:

- Using Boolean closure algorithms to check that language inclusion holds in both directions.
- Minimizing both machines, and checking that they are isomorphic: the same up to renaming of states.

More precisely, two DFAs M_1 and M_2 are **isomorphic** if there is a bijection $f : Q_1 \rightarrow Q_2$ such that

$$\begin{aligned} f(q_{10}) &= q_{20} \\ f(\delta_1(p, a)) &= \delta_2(f(p), a) \\ f(F_1) &= F_2 \end{aligned}$$

As usual, there is an associated decision problem.

Problem: **Isomorphism of DFA**
 Instance: Two DFAs M_1 and M_2 .
 Question: Are M_1 and M_2 isomorphic?

To check whether two DFAs are isomorphic one can use a variant of depth-first-search. We may safely assume the machines have the same size.

- Set $f(q_{01}) = q_{02}$.
- Extend the domain of f according to δ_1 , and the range according to δ_2 .
- Stop with failure if there ever is a clash (new point on one side, old point on the other; two different old points).
- Check if f maps the final states properly.
If so, return Yes, otherwise return No.

Exercise

Figure out the details of this isomorphism testing algorithm.

Theorem

Given two DFAs one can test whether they are equivalent in log-linear time.

Proof.

Given M_1 and M_2 we can compute the corresponding minimal automata M'_1 and M'_2 in log-linear time. Then we can check in linear time whether M'_1 and M'_2 are isomorphic. \square

Note that even using Moore's algorithm this will often produce sub-quadratic running time.

Since isomorphism testing part of our algorithm is very fast, the real cost is incurred in minimization.

Wild Idea: Can we somehow avoid minimization?

Of course, we also want to avoid a product machine construction.

Suppose the two given DFAs have disjoint state sets Q_1 and Q_2 and set $Q = Q_1 \cup Q_2$.

We could try to define an equivalence relation E on Q that will tell us (among other things) whether the two initial states have the same behavior. And to it fast.

```

set E = identity;
set active = ( (q01,q02) );

while( active != empty )
  (p,q) = active.extract();
  if( ! p E q ) {
    set p E q;
    foreach a in S do
      set (delta(p,a),delta(q,a)) active;
  }

```

The algorithm returns false whenever two states $p \in Q_1$ and $q \in Q_2$ are defined to be equivalent under E but

$$p \in F_1 \text{ xor } q \in F_2.$$

Otherwise we return true.

We can use a simple stack for the active list, the interesting question is how to maintain the equivalence relation E .

Note that this is the dynamical situation: initially E is just the identity relation, but as we go along we discover more related pairs.

So really we need to compute the **equivalential closure** of all the pairs (p, q) we have discovered. Union/Find is perfect for this.

By initialization $q_{01} E q_{02}$.

If $p \in Q_1$ and $q \in Q_2$ such that $p E q$ then $\delta(p, a) E \delta(q, a)$ for each $a \in \Sigma$.

By induction for any word x

$$\delta(q_{01}, x) E \delta(q_{02}, x).$$

Using a loop invariant one can show that in fact

$$E \cap Q_1 \times Q_2 = \{ (\delta(q_{01}, x), \delta(q_{02}, x)) \mid x \in \Sigma^* \}$$

Exercise

Give a detailed proof for the correctness of this algorithm.

We have seen $O(n \log n)$ minimization algorithms. However, no better algorithm is known to date.

Exercise

Why can't this equivalence testing algorithm (or rather: some slight modification thereof) be used to compute the minimal automaton in nearly linear time?

Note that we do not claim that equivalence is polynomial time decidable for nondeterministic machines. In fact, one can show that equivalence testing is NP-hard (even PSPACE-hard) for nondeterministic machines, even if one of the two machines simply accepts all words over the input alphabet.

Likewise it is computationally hard to find minimal nondeterministic machines for a given language. Moreover, the minimal nondeterministic machine is not unique in general.

- Partition Refinement
- Hopcroft's Algorithm
- Valmari-Lehtinen
- Equivalence Testing
- ⑤ Characterizations of Regularity

Recall that our original definition of regularity was based on DFAs: a language is regular if it is accepted by some DFA.

As Rabin-Scott showed, one can naturally generalize the underlying machines:

Theorem

A language is regular iff it is accepted by an NFA (or even an NFAE).

We have already seen the following.

Theorem

A language $L \subseteq \Sigma^*$ is regular iff there is a finite monoid M , $M_0 \subseteq M$ and a monoid homomorphism $f: \Sigma^* \rightarrow M$ such that $L = f^{-1}(M_0)$.

A similar result is the following.

Theorem

A language $L \subseteq \Sigma^*$ is regular iff there is a congruence on Σ^* of finite index that saturates L .

The last theorem also holds when congruence is replaced by right congruence.

Also, our quotient approach to minimization provides as a pleasant side effect yet another characterization.

Lemma

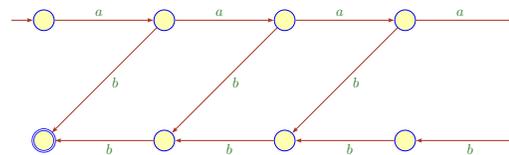
A language is regular iff it has finitely many quotients.

This follows immediately from the preceding results.

Interesting point: the quotient automaton makes perfect sense even for non-regular languages, it just has infinitely many states in this case.

If the language is simple, the structure of the diagram will be very regular. This leads naturally to machines more powerful than finite state machines yet not as strong as Turing machines (such as pushdown automata).

Part of the quotient machine for the non-regular language $L = \{a^i b^i \mid i \geq 0\}$ looks like this (we omit the pesky sink).



The picture very much suggests the augmentation necessary for a plain DFA to deal with L : a single counter will do (a stack with just one stack symbol).

Exercise

Figure out how to extend a DFA so it can deal with L , keeping the additions as simple as possible.

Alas, making sure that there are indeed infinitely many quotients can be difficult, a less challenging test would be helpful.

Lemma (Pumping Lemma)

For every regular language L there is a constant n such that for all words $x \in L$ with $|x| \geq n$ we have $x = uvw$ where $v \neq \varepsilon$, $|uv| \leq n$ and $uv^t w \in L$ for all $t \geq 0$.

Proof.

Consider the minimal DFA M for L and let n be the number of states of M . Then any word in L of length at least n must trace a loop in the diagram of M . The claim follows. \square

The Pumping lemma is useless to establish regularity but often the weapon of choice to refute it.

$L = \{a^i b^i \mid i \geq 0\}$ fails to be regular.

Assume otherwise.

Let n be as in the PL and consider $x = a^n b^n \in L$.

Then $x = uvw$ and $v = a^i$ for some $i > 0$.

But then $uv^t w \notin L$ for all $t \neq 1$, contradiction.

It follows that the language P of balanced parentheses is also non-regular.

For otherwise $P \cap a^* b^* = L$ would also be regular, which assertion we already know to be false.

$L = \{zz \mid z \in \Sigma^*\}$ fails to be regular.

Assume otherwise.

Let n be as in the PL and consider $x = ab^n ab^n \in L$.

Then $x = uvw$ and $|uv| \leq n$. If $v = a$ we get a contradiction with $t = 0$. If $v = b^i$ for some $i > 0$ we also get a contradiction with $t = 0$.

The problem here really is that a DFA cannot remember an arbitrarily long prefix z which is needed to check the remainder of the input.

$L = \{zz^r \mid z \in \Sigma^*\}$ fails to be regular.

Assume otherwise.

Let n be as in the PL and consider $x = (ab)^n (ba)^n \in L$.

Then $x = uvw$ and $|uv| \leq n$.

A straightforward but tedious argument shows that $uv^t w$ cannot be a palindrome for any $t \neq 1$.

As in the last example, a DFA cannot remember an arbitrarily long prefix z which is needed to check the remainder of the input.

Incidentally, between

$$L_1 = \{zz \mid z \in \Sigma^*\}$$

and

$$L_2 = \{zz^r \mid z \in \Sigma^*\}$$

the second one (even length palindromes) is much simpler:

To recognize these words one only needs to attach a stack to a FSM (context free language).

For L_1 , a simple stack is not sufficient (context sensitive language).

Back to quotients. Suppose we want to check divisibility by 5 for numbers given in **reverse binary** notation where the MSD comes last. In this case the value function ν^R looks like this:

$$\nu^R(x_0 x_1 \dots x_k) = \nu(x_k x_{k-1} \dots x_0) = \sum_{i \leq k} x_i 2^i$$

Appending a new digit changes the numerical value in a more complicated way here than in ordinary binary:

$$\nu^R(xa) = \nu^R(x) + a2^{|x|}$$

Cop-out: We could push the digits on a stack and then use our old DFA for standard binary notation, but that requires extra memory.

So how do we construct a DFA, a memory-less device, for reverse binary directly?

Note that the real issue here is that we read input strings

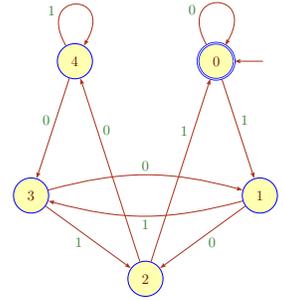
$$x_0x_1x_2 \dots x_{n-1}x_n$$

from left to right; this convention is hard-wired into our definition of finite state machine. Of course, the choice of direction is completely arbitrary (more precisely, a western cultural artifact).

Of course, we have a theorem that shows that L^{op} is regular whenever L is. However, the result relies on nondeterminism and determinization; we would really like a structural understanding of the machine for L^{op} .

In ordinary base 2 the state set is $Q = \{0, \dots, 4\}$ and the transition function of the canonical Horner automaton is determined by

$$\nu(xa) = 2 \cdot \nu(x) + a \pmod{5}$$



It is tempting to use the same state set $Q = \{0, \dots, 4\}$ for reverse binary. But the transition function now is given by

$$\nu^R(xa) = \nu^R(x) + a \cdot 2^{|x|} \pmod{5}$$

We can represent $\nu^R(x)$ as a state, but not the $|x|$.

There are two basic solutions:

- Make the state set more complicated so we can keep track of the missing information. Quite possible, but leads to a larger machine.
- Try to find another way of describing a transition function that works for Q directly (assuming that such a thing really exists).

What state set do we need to implement the transition function

$$\nu^R(xa) = \nu^R(x) + a \cdot 2^{|x|} \pmod{5}$$

We want to keep track of

- The current value $\nu^R(x) \pmod{5}$, and
- the current multiplier $2^{|x|} \pmod{5}$.

Hence we have a state set $Q = \{0, 1, 2, 3, 4\} \times \{1, 2, 3, 4\}$ and transitions

$$\delta((p, q), a) = (p + aq, 2q) \pmod{5}$$

This is the solution anyone understanding the basics should be able to come up with.

Here is a more elegant line of attack. Write

$$L_r = \{x \in 2^* \mid \nu^R(x) = r \pmod{m}\}$$

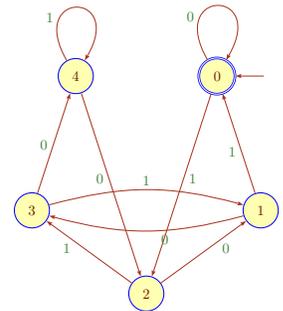
for $r = 0, \dots, 4$, so that L_0 is the language we want to recognize.

Can we perhaps compute its quotients?

$$\begin{aligned} a^{-1}L_r &= \{x \in 2^* \mid ax \in L_r\} \\ &= \{x \in 2^* \mid \nu^R(ax) = r \pmod{5}\} \\ &= \{x \in 2^* \mid a + 2\nu^R(x) = r \pmod{5}\} \\ &= \{x \in 2^* \mid \nu^R(x) = 3r + 2a \pmod{5}\} \\ &= L_{3r+2a \pmod{5}} \end{aligned}$$

So we can actually construct the quotient automaton directly in this case.

$$\nu^R(xa) = 3 \cdot \nu^R(x) + 2 \cdot a \pmod{5}$$



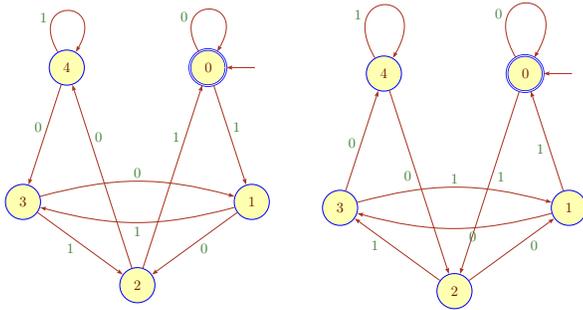
Exercise

Give a more careful description of the automaton using the Cartesian product state set for arbitrary m and reverse base B .

Exercise

Carry out the quotient construction for some other values of m and B .

It is worth drawing the machines right next to each other.



Observations?

We could have done this much cheaper ...

Here is a good test for implementations: determine the state complexity of

$$L_k = \{uv \in \{a, b\}^* \mid |u| = |v| = k, u \neq v\}.$$

Note that L_k is finite, so these languages are trivially regular.

Question: What is the state complexity of L_k ?

k	1	2	3	4	5	6
sc	5	12	25	50	99	196

Exercise

Conjecture the state complexity of L_k from the last table and prove your conjecture.

- DFA minimization can be handled in time $n \log n$.
- Moore's quadratic time algorithm may be competitive on average.
- Equivalence testing of DFAs can be handled in nearly linear time without minimization.
- Quotients can be useful to describe machines for regular languages.