

Generic ϵ -Removal and Input ϵ -Normalization Algorithms for Weighted Transducers

Mehryar Mohri
`mohri@research.att.com`
AT&T Labs – Research
180 Park Avenue, Room E135
Florham Park, NJ 07932, USA

Received October 30, 2000
Revised October 31, 2000
Communicated by Sheng Yu

ABSTRACT

We present a new generic ϵ -removal algorithm for weighted automata and transducers defined over a semiring. The algorithm can be used with any semiring covered by our framework and works with any queue discipline adopted. It can be used in particular in the case of unweighted automata and transducers and weighted automata and transducers defined over the tropical semiring. It is based on a general shortest-distance algorithm that we briefly describe. We give a full description of the algorithm including its pseudocode and its running time complexity, discuss the more efficient case of acyclic automata, an on-the-fly implementation of the algorithm and an approximation algorithm in the case of the semirings not covered by our framework. We illustrate the use of the algorithm with several semirings. We also describe an input ϵ -normalization algorithm for weighted transducers based on the general shortest-distance algorithm. The algorithm, which works with all semirings covered by our framework, admits an on-the-fly implementation.

Keywords: finite automata, finite-state transducers, rational power series, semirings, ϵ -removal, shortest-paths algorithms

1. Introduction

Weighted automata are general and efficient devices used in many applications such as text, speech and image processing [7, 3, 2]. The automata created in such applications are often the result of various complex operations, some of them introducing the empty string ϵ . For example, the classical method of Thompson can be used to construct in linear time and space a non-deterministic automaton with ϵ -transitions representing a regular expression [10]. The result of rational operations (closure, concatenation, union) also contains in general ϵ -transitions.

For the most efficient use of an automaton, it is preferable to *remove* empty strings since in general they induce a delay in their use. An algorithm constructing

an automaton B with no empty string equivalent to an input automaton A with empty strings is called ϵ -removal.

Textbooks do not present ϵ -removal of (unweighted) automata as an independent algorithm deserving a specific study. Instead, the algorithm is often mixed with other optimization algorithms such as determinization [1]. This usually makes the presentation of determinization more complex and the underlying ϵ -removal process obscure. Since ϵ -removal is not presented as an independent algorithm, it is usually not analyzed and its running time complexity not clearly determined.

We present a new generic ϵ -removal algorithm for weighted automata and transducers defined over a semiring. The algorithm can be used with any semiring covered by our framework and works with any queue discipline adopted. It can be used in particular in the case of unweighted automata and transducers and weighted automata and transducers defined over the tropical semiring. It is based on a general shortest-distance algorithm that we briefly describe. We give a full description of the algorithm including its pseudocode and its running time complexity, discuss the more efficient case of acyclic automata, an on-the-fly implementation of the algorithm and an approximation algorithm in the case of the semirings not covered by our framework. We illustrate the use of the algorithm with several semirings. We also describe an input ϵ -normalization algorithm for weighted transducers which is based on the general shortest-distance algorithm. The algorithm, which works with all semirings covered by our framework, admits an on-the-fly implementation.

2. Preliminaries

Weighted automata are automata in which the transitions are labeled with weights in addition to the usual alphabet symbols. For various operations to be well-defined, the weight set needs to have the algebraic structure of a semiring [6].

Definition 1 A system $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ is a right semiring if:

1. $(\mathbb{K}, \oplus, \bar{0})$ is a commutative monoid with $\bar{0}$ as the identity element for \oplus ,
2. $(\mathbb{K}, \otimes, \bar{1})$ is a monoid with $\bar{1}$ as the identity element for \otimes ,
3. \otimes right distributes over \oplus : $\forall a, b, c \in \mathbb{K}, (a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c)$,
4. $\bar{0}$ is an annihilator for \otimes : $\forall a \in \mathbb{K}, a \otimes \bar{0} = \bar{0} \otimes a = \bar{0}$.

Left semirings are defined in a similar way by replacing right distributivity with left distributivity. $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ is a *semiring* if both left and right distributivity hold. Thus, more informally, a semiring is a ring that may lack negation. As an example, $(\mathbb{N}, +, \cdot, 0, 1)$ is a semiring defined on the set of nonnegative integers \mathbb{N} .

A semiring $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ is said to be *idempotent* if for any $a \in \mathbb{K}$, $a \oplus a = a$. The boolean semiring $\mathcal{B} = (\{0, 1\}, \vee, \wedge, 0, 1)$ and the tropical semiring $\mathcal{T} = (\mathbb{R}_+ \cup \{\infty\}, \min, +, \infty, 0)$ are idempotent, but $(\mathbb{N}, +, \cdot, 0, 1)$ is not.

Definition 2 A weighted automaton $A = (\Sigma, Q, I, F, E, \lambda, \rho)$ over the semiring \mathbb{K} is a 7-tuple where Σ is the finite alphabet of the automaton, Q is a finite set of states, $I \subseteq Q$ the set of initial states, $F \subseteq Q$ the set of final states, $E \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times \mathbb{K} \times Q$

a finite set of transitions, $\lambda : I \rightarrow \mathbb{K}$ the initial weight function mapping I to \mathbb{K} , and $\rho : F \rightarrow \mathbb{K}$ the final weight function mapping F to \mathbb{K} .

Given a transition $e \in E$, we denote by $i[e]$ its input label, $w[e]$ its weight, $p[e]$ its origin or previous state and $n[e]$ its destination state or next state. Given a state $q \in Q$, we denote by $E[q]$ the set of transitions leaving q , and by $E^R[q]$ the set of transitions entering q .

A path $\pi = e_1 \cdots e_k$ in A is an element of E^* with consecutive transitions: $n[e_{i-1}] = p[e_i]$, $i = 2, \dots, k$. We extend n and p to the set of paths by setting: $n[\pi] = n[e_k]$ and $p[\pi] = p[e_1]$. We denote by $P(q, q')$ the set of paths from q to q' . P can be extended to subsets $R \subseteq Q$ $R' \subseteq Q$, by:

$$P(R, R') = \bigcup_{q \in R, q' \in R'} P(q, q')$$

The labeling function i and the weight function w can also be extended to paths by defining the label of a path as the concatenation of the labels of its constituent transitions, and the weight of a path as the \otimes -product of the weights of its constituent transitions:

$$\begin{aligned} i[\pi] &= i[e_1] \cdots i[e_k] \\ w[\pi] &= w[e_1] \otimes \cdots \otimes w[e_k] \end{aligned}$$

Given a string $x \in \Sigma^*$, we denote by $P(x)$ the set of paths from I to F labeled with x :

$$P(x) = \{\pi \in P(I, F) : i[\pi] = x\}$$

Since the additive operation is commutative, the output weight associated by A to an input string $x \in \Sigma^*$ can be defined by:

$$\llbracket A \rrbracket(x) = \bigoplus_{\pi \in P(x)} \lambda(p[\pi]) \otimes w[\pi] \otimes \rho(n[\pi])$$

If $P(x) = \emptyset$, $\llbracket A \rrbracket(x)$ is defined to be $\bar{0}$. Note that weighted automata over the boolean semiring are equivalent to the classical unweighted finite automata.

These definitions can be easily generalized to cover the case of weighted automata with ϵ -transitions. An ϵ -removal algorithm computes for any input weighted automaton A with ϵ -transitions an equivalent weighted automaton B with no ϵ -transition, that is such that:

$$\forall x \in \Sigma^*, \llbracket A \rrbracket(x) = \llbracket B \rrbracket(x)$$

Weighted finite-state transducers are defined in a similar way.

Definition 3 A weighted transducer $T = (\Sigma, \Omega, Q, I, F, E, \lambda, \rho)$ over the semiring \mathbb{K} is a 8-tuple where Σ is the finite input alphabet of the transducer, Ω is the finite output alphabet of T , Q is a finite set of states, $I \subseteq Q$ the set of initial states, $F \subseteq Q$ the set of final states, $E \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Omega \cup \{\epsilon\}) \times \mathbb{K} \times Q$ a finite set of transitions, $\lambda : I \rightarrow \mathbb{K}$ the initial weight function mapping I to \mathbb{K} , and $\rho : F \rightarrow \mathbb{K}$ the final weight function mapping F to \mathbb{K} .

The output weight associated by T to an input string $x \in \Sigma^*$ and output string $y \in \Omega$ is:

$$\llbracket T \rrbracket(x, y) = \bigoplus_{\pi \in P(x, y)} \lambda(p[\pi]) \otimes w[\pi] \otimes \rho(n[\pi])$$

where $P(x, y)$ is the set of successful paths with input label x and output label y .

The following definitions will help us define the framework for our generic ϵ -removal algorithm [8].

Definition 4 Let $k \geq 0$ be an integer. A commutative semiring $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ is k -closed if:

$$\forall a \in \mathbb{K}, \bigoplus_{n=0}^{k+1} a^n = \bigoplus_{n=0}^k a^n$$

When $k = 0$, the previous expression can be rewritten as:

$$\forall a \in \mathbb{K}, \bar{1} \oplus a = \bar{1}$$

and \mathbb{K} is then said to be *bounded*. Semirings such as the boolean semiring $\mathcal{B} = (\{0, 1\}, \vee, \wedge, 0, 1)$ and the tropical semiring $\mathcal{T} = (\mathbb{R}_+ \cup \{\infty\}, \min, +, \infty, 0)$ are bounded.

Definition 5 Let $k \geq 0$ be an integer, $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ a commutative semiring, and A a weighted automaton over \mathbb{K} . $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ is right k -closed for A if for any cycle π of A :

$$\bigoplus_{n=0}^{k+1} w[\pi]^n = \bigoplus_{n=0}^k w[\pi]^n$$

By definition, if \mathbb{K} is k -closed, then it is k -closed for any automaton over \mathbb{K} .

3. Algorithm

Let $A = (\Sigma, Q, I, F, E, \lambda, \rho)$ be a weighted automaton over the semiring \mathbb{K} with ϵ -transitions. We denote by A_ϵ the automaton obtained from A by removing all transitions not labeled with ϵ . We present a general ϵ -removal algorithm for weighted automata based on a generic shortest-distance algorithm which works for any semiring \mathbb{K} k -closed for A_ϵ [8]. In what follows, we assume that \mathbb{K} has this property. This class of semirings includes in particular the boolean semiring $(\{0, 1\}, \vee, \wedge, 0, 1)$, the tropical semiring $(\mathbb{R}_+ \cup \{\infty\}, \min, +, \infty, 0)$ and other semirings not necessarily idempotent. ^a

We present the algorithm in the case of weighted automata. The case of weighted transducers can be straightforwardly derived from the automata case by viewing a transducer as an automaton over the alphabet $\Sigma \cup \{\epsilon\} \times \Sigma \cup \{\epsilon\}$. An ϵ -transition of a transducer is then a transition labeled with (ϵ, ϵ) .

For p, q in Q , the ϵ -distance from p to q in the automaton A is denoted by $d[p, q]$ and defined as:

$$d[p, q] = \bigoplus_{\pi \in P(p, q), i[\pi] = \epsilon} w[\pi]$$

^aThe class of semirings with which our generic shortest-distance algorithm works can in fact be extended to that of right semirings *right k -closed* for A_ϵ [8], but for reasons of space we do not describe this more general framework here.

This distance is well-defined for any pair of states (p, q) of A when \mathbb{K} is a semiring k -closed for A_ϵ [8]. By definition, for any $p \in Q$, $d[p, p] = \bar{1}$. $d[p, q]$ is the distance from p to q in A_ϵ .

3.1. Description and Proof

The algorithm works in two steps. The first step consists of computing for each state p of the input automaton A its ϵ -closure denoted by $C[p]$:

$$C[p] = \{(q, w) : q \in \epsilon[p], d[p, q] = w \in \mathbb{K} - \{\bar{0}\}\}$$

where $\epsilon[p]$ represents the set of states reachable from p via a path labeled with ϵ . We describe this step in more detail later.

The second step consists of modifying the outgoing transitions of each state p by removing those labeled with ϵ and by adding to $E[p]$ non- ϵ -transitions leaving each state $q \in \epsilon[p]$ with their weights pre- \otimes -multiplied by $d[p, q]$. The following is the pseudocode of the second step of the algorithm.

ϵ -removal(A)

```

1  for each  $p \in Q$ 
2      do  $E[p] \leftarrow \{e \in E[p] : i[e] \neq \epsilon\}$ 
3      for each  $(q, w) \in C[p]$ 
4          do  $E[p] \leftarrow E[p] \cup \{(p, a, w \otimes w', r) : (q, a, w', r) \in E[q], a \neq \epsilon\}$ 
5          if  $q \in F$ 
6              then if  $p \notin F$ 
7                  then  $F \leftarrow F \cup \{p\}$ 
8                   $\rho[p] \leftarrow \rho[p] \oplus (w \otimes \rho[q])$ 
```

State p is a final state if some state $q \in \epsilon[p]$ is final and the final weight $\rho[p]$ is then:

$$\rho[p] = \bigoplus_{q \in \epsilon[p] \cap F} d[p, q] \otimes \rho[q]$$

Theorem 1 *Let $A = (\Sigma, Q, I, F, E, \lambda, \rho)$ be a weighted automaton over the semiring \mathbb{K} right k -closed for A . Then the weighted automaton B result of the ϵ -removal algorithm just described is equivalent to A .*

Proof. We show that the function defined by the weighted automaton A is not modified by the application of one step of the loop of the pseudocode above. Let $A = (\Sigma, Q, I, F, E, \lambda, \rho)$ be the automaton just before the removal of the ϵ -transitions leaving state $p \in Q$ (lines 2-8). For any successful path π in A , define $s[\pi]$ by:

$$s[\pi] = \lambda(p[\pi]) \otimes w[\pi] \otimes \rho(n[\pi])$$

Let $x \in \Sigma^*$, and let $Q(p)$ denote the set of successful paths labeled with x passing through p and either ending at p or following an ϵ -transition of p . Define S_1 and S_2 by:

$$S_1 = \bigoplus_{\pi \in P(x) - Q(p)} s[\pi] \quad S_2 = \bigoplus_{\pi \in Q(p)} s[\pi]$$

By commutativity of \oplus , we have:

$$\llbracket A \rrbracket(x) = S_1 \oplus S_2$$

The ϵ -removal at state p does not affect the paths not in $Q(p)$, thus we can limit our attention to the second term S_2 . A path in $Q(p)$ can be factored in: $\pi = \pi_1 \pi_\epsilon \pi_2$, where π_ϵ is a portion of the path from p to $n[\pi_\epsilon] = q$ labeled with ϵ . The distributivity of \otimes over \oplus gives us:

$$S_2 = \left(\bigoplus_{n[\pi_1]=p} \lambda(p[\pi_1]) \otimes w[\pi_1] \right) \otimes S_{22}$$

with:

$$\begin{aligned} S_{22} &= \left(\bigoplus_{p[\pi_\epsilon]=p, n[\pi_\epsilon]=q} w[\pi_\epsilon] \right) \otimes \bigoplus_{p[\pi_2]=q} w[\pi_2] \otimes \rho(n[\pi_2]) \\ &= d[p, q] \otimes \bigoplus_{p[\pi_2]=q} w[\pi_2] \otimes \rho(n[\pi_2]) \\ &= \bigoplus_{p[\pi_2]=q} (d[p, q] \otimes w[\pi_2]) \otimes \rho(n[\pi_2]) \end{aligned}$$

For paths π such that $\pi_2 = \epsilon$:

$$S_{22} = \bigoplus_{p[\pi_2]=q, q \in \epsilon[p]} d[p, q] \otimes \rho(q)$$

which is exactly the final weight associated to p by ϵ -removal at p . Otherwise, by definition of π_ϵ , π_2 does not start with an ϵ -transition. S_{22} can thus be rewritten as:

$$S_{22} = \bigoplus_{e \in E[q], i[e] \neq \epsilon, \pi_2 = e\pi'_2} (d[p, q] \otimes w[e]) \otimes w[\pi'_2] \otimes \rho(n[\pi_2])$$

The ϵ -removal step at state p follows exactly that decomposition. It adds non ϵ -transitions e leaving q with weights $(d[p, q] \otimes w[e])$ to the transitions leaving p . This ends the proof of the theorem. \square

After removing ϵ 's at each state p , some states may become inaccessible if they could only be reached by ϵ -transitions originally. Those states can be removed from the result in time linear in the size of the resulting machine using for example a depth-first search of the automaton.

Figures 1 (a)-(c) illustrate the use of the algorithm in the specific case of the tropical semiring. ϵ -transitions can be *removed* in at least two ways: in the way described above, or the reverse. The latter is equivalent to applying ϵ -removal to the reverse of the automaton and re-reversing the result. The two methods may lead to results of very different sizes as illustrated by the figures. This is due to two factors that are independent of the semiring:

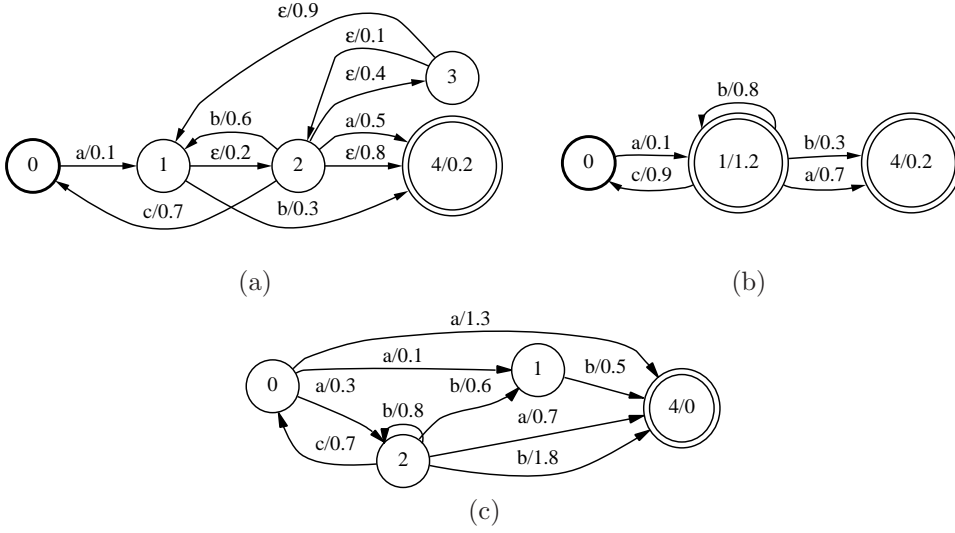
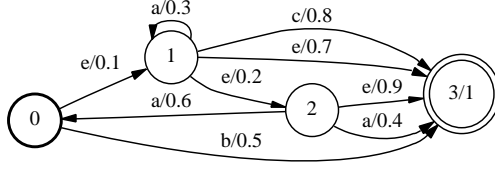


Figure 1: ϵ -removal in the tropical semiring. (a) Weighted automaton A with ϵ -transitions. (b) Weighted automaton B equivalent to A result of the ϵ -removal algorithm. (c) Weighted automaton C equivalent to A obtained by application of ϵ -removal to the reverse of A .

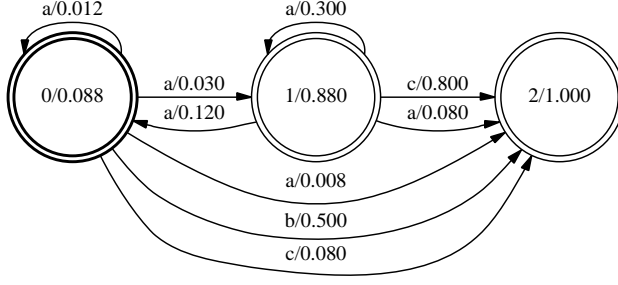
- the number of states in the original automaton whose incoming transitions (or outgoing transitions in the reverse case) are all labeled with the empty string. As mentioned before, those states can be removed from the result. For example, state 3 of the automaton of figure 1 (a) can only be reached by ϵ -transitions and admits only outgoing transitions labeled with ϵ . Thus, that state does not appear in the result in both methods (figures 1 (b)-(c)). The incoming transitions of state 2 are all labeled with ϵ and thus it does not appear in the result of the ϵ -removal with the first method, but it does in the reverse method because the outgoing transitions of state 2 are not all labeled with ϵ ;
- the total number of non ϵ -transitions of the states that can be reached from each state q in A_ϵ (the reverse of A_ϵ in the reverse case). This corresponds to the number of outgoing transitions of q in the result of ϵ -removal.

In practice, one can use some heuristics to reduce the number of states and transitions of the resulting machine although this will not affect the worst case complexity of the algorithm. One can for instance remove some ϵ -transitions in the reverse way when that creates less transitions and others in the way corresponding to the first method when that helps reducing the resulting size.

Figures 2 (a)-(b) illustrate the algorithm in the case of another semiring, the semiring of real numbers. Our general algorithm applies in this case since A_ϵ is acyclic.



(a)



(b)

Figure 2: ϵ -removal in the real semiring $(\mathbb{R}, +, *, 0, 1)$. (a) Weighted automaton A . A_ϵ is k -closed for $(\mathbb{R}, +, *, 0, 1)$. (b) Weighted automaton B equivalent to A output of the ϵ -removal algorithm.

3.2. Computation of ϵ -closures

As noticed before, the computation of ϵ -closures is equivalent to that of all-pairs shortest-distances over the semiring \mathbb{K} in A_ϵ . There exists a generalization of the algorithm of Floyd-Warshall [4, 12] for computing the all-pairs shortest-distances over a semiring \mathbb{K} under some general conditions [8].^b However, the running time complexity of that algorithm is cubic:

$$O(|Q|^3(T_\oplus + T_\otimes + T_*))$$

where T_\oplus , T_\otimes , and T_* denote the cost of \oplus , \otimes , and closure operations in the semiring considered. The algorithm can be improved by first decomposing A_ϵ into its strongly connected components, and then computing all-pairs shortest distances in each component visited in reverse topological order. However, it is still impractical for large automata when A_ϵ has large cycles of several thousand states. The quadratic space complexity $O(|Q|^2)$ of the algorithm also makes it prohibitive for such large automata. Another problem with this generalization of the algorithm of Floyd-Warshall is that it does not exploit the sparseness of the input automaton.

There exists a generic single-source shortest-distance algorithm that works with any semiring covered by our framework [8]. The algorithm is a generalization of the

^bThat algorithm works in particular with the semiring $(\mathbb{R}, +, *, 0, 1)$ when the weight of each cycle of A_ϵ admits a well-defined closure.

classical shortest-paths algorithms to the case of the semirings of this framework.

This generalization is not trivial and does not require the semiring to be idempotent. In particular, a straightforward extension of the classical algorithms based on a relaxation technique would not produce the correct result in general [8]. The algorithm is also generic in the sense that it works with any queue discipline. The following is the pseudocode of the algorithm.

GENERIC-SINGLE-SOURCE-SHORTEST-DISTANCE (B, s)

```

1  for each  $p \in Q$ 
2      do  $d[p] \leftarrow r[p] \leftarrow \bar{0}$ 
3   $d[s] \leftarrow r[s] \leftarrow \bar{1}$ 
4   $S \leftarrow \{s\}$ 
5  while  $S \neq \emptyset$ 
6      do  $q \leftarrow \text{head}(S)$ 
7          DEQUEUE( $S$ )
8           $R \leftarrow r[q]$ 
9           $r[q] \leftarrow \bar{0}$ 
10         for each  $e \in E[q]$ 
11             do if  $d[n[e]] \neq d[n[e]] \oplus (R \otimes w[e])$ 
12                 then  $d[n[e]] \leftarrow d[n[e]] \oplus (R \otimes w[e])$ 
13                      $r[n[e]] \leftarrow r[n[e]] \oplus (R \otimes w[e])$ 
14                     if  $n[e] \notin S$ 
15                         then ENQUEUE( $S, n[e]$ )
16  $d[s] \leftarrow \bar{1}$ 

```

Figure 3: Generic single-source shortest-distance algorithm.

A queue S is used to maintain the set of states whose leaving transitions are to be relaxed. S is initialized to $\{s\}$ (line 4). For each state $q \in Q$, two attributes are maintained: $d[q] \in \mathbb{K}$ an estimate of the shortest distance from s to q , and $r[q] \in \mathbb{K}$ the total weight added to $d[q]$ since the last time q was extracted from S . Lines 1 – 3 initialize arrays d and R . After initialization, $d[q] = r[q] = \bar{0}$ for $q \in Q - \{s\}$, and $d[s] = r[s] = \bar{1}$.

Given a state $q \in Q$ and an transition $e \in E[q]$, a relaxation step on e is performed by lines 11-13 of the pseudocode, where R is the value of $r[q]$ just after the latest extraction of q from S if q has ever been extracted from S , its initialization value otherwise.

Each time through the **while** loop of lines 5-15, a state q is extracted from S (lines 6-7). The value of $r[q]$ just after extraction of q is stored in R , and then $r[q]$ is set to $\bar{0}$ (lines 8-9). Lines 11-13 relax each transition leaving q . If the tentative shortest distance $d[n[e]]$ is updated during the relaxation and if $n[e]$ is not already in S , the state $n[e]$ is inserted in S so that its leaving transitions be later relaxed (lines 14-15). $r[n[e]]$ is updated whenever $d[n[e]]$ is. $r[n[e]]$ stores the total weight \oplus -added to $d[n[e]]$ since $n[e]$ was last extracted from S or since the time after initialization if $n[e]$ has never been extracted from S . Finally, line 16 resets the value of $d[s]$ to $\bar{1}$.

In the general case, the complexity of the algorithm depends on the semiring

considered and the queue discipline chosen for S :

$$O(|Q| + (T_{\oplus} + T_{\otimes} + C(A))|E| \max_{q \in Q} N(q) + (C(I) + C(X)) \sum_{q \in Q} N(q))$$

where $N(q)$ denotes the number of times state q is extracted from S , $C(X)$ the worst cost of extracting a state q from the queue S , $C(I)$ that of inserting q in S , and $C(A)$ the cost of an assignment.^c

In the case of the tropical semiring $(\mathbb{R}_+ \cup \{\infty\}, \min, +, \infty, 0)$, the algorithm coincides with classical single-source shortest-paths algorithms. In particular, it coincides with Bellman-Ford's algorithm when a FIFO queue discipline is used and with Dijkstra's algorithm when a shortest-first queue discipline is used. Using Fibonacci heaps [5], the complexity of Dijkstra's algorithm in the tropical semiring is:

$$O(|E| + |Q| \log |Q|)$$

The complexity of the algorithm is linear in the case of an acyclic automaton with a topological order queue discipline:

$$O(|Q| + (T_{\oplus} + T_{\otimes})|E|)$$

Note that the topological order queue discipline can be generalized to the case of non-acyclic automata A_{ϵ} by decomposing A_{ϵ} into its strongly connected components. Any queue discipline can then be used to compute the all-pairs shortest distances within each strongly connected component [8].

The all-pairs shortest-distances of A_{ϵ} can be computed by running $|Q|$ times the generic single-source shortest-distance algorithm. Thus, when A_{ϵ} is acyclic, that is when A admits no ϵ -cycle, then the all-pairs shortest distances can be computed in quadratic time:

$$O(|Q|^2 + (T_{\oplus} + T_{\otimes})|Q| \cdot |E|)$$

When A_{ϵ} is acyclic, the complexity of the computation of the all-pairs shortest distances can be substantially improved if the states of A_{ϵ} are visited in reverse topological order and if the shortest-distance algorithm is interleaved with the actual removal of ϵ 's. Indeed, one can proceed in the following way. For each state p of A_{ϵ} visited in reverse topological order:

1. run a single-source shortest-distance algorithm with source p to compute the distance from p to each state q reachable from p by ϵ 's;
2. remove the ϵ -transitions leaving q as described in the previous section.

The reverse topological order guarantees that the ϵ -paths leaving p are reduced to the ϵ -transitions leaving p . Thus, the cost of the shortest-distance algorithm run from p only depends on the number of ϵ -transitions leaving p and the total cost of the computation of the shortest-distances is linear:

$$O(|Q| + (T_{\oplus} + T_{\otimes})|E|)$$

^cThis includes the potential cost of reorganization of the queue to perform this assignment.

In the case of the tropical semiring and using Fibonacci heaps, the complexity of the first stage of the algorithm is:

$$O(|Q| \cdot |E| + |Q|^2 \log |Q|)$$

In the worst case, in the second stage of the algorithm each state q belongs to the ϵ -closure of each state p , and the removal of ϵ 's can create in the order of $|E|$ transitions at each state. Hence, the complexity of the second stage of the algorithm is:

$$O(|Q|^2 + |Q| \cdot |E|)$$

Thus, the total complexity of the algorithm in the case of an acyclic automaton A_ϵ is:

$$O(|Q|^2 + (T_\oplus + T_\otimes)|Q| \cdot |E|)$$

In the case of the tropical semiring and with a non acyclic automaton A_ϵ , the total complexity of the algorithm is:

$$O(|Q| \cdot |E| + |Q|^2 \log |Q|)$$

4. Remarks and experiments

We have fully implemented the ϵ -removal algorithm presented here and incorporated it in recent versions of the FSM library [9]. For some automata of about a thousand states with large ϵ -cycles, our implementation is up to 600 times faster than the previous implementation based on a generalization of the Floyd-Warshall algorithm.

An important feature of our algorithm is that it admits a natural on-the-fly implementation. Indeed, the outgoing transitions of state q of the output automaton can be computed directly using the ϵ -closure of q . However, with an on-the-fly implementation, a topological order cannot be used for the queue S even if A_ϵ is acyclic since this is not known ahead of the time. Thus, we have implemented both an off-line and an on-the-fly version of the algorithm.

The algorithms presented in textbooks often combine the classical powerset construction, or determinization, and ϵ -removal [1]. Each state of the resulting automaton corresponds to the ϵ -closure of a subset of states constructed as in determinization.

If one wishes to apply determinization after ϵ -removal, then the integration of ϵ -removal within determinization may often be more efficient. This is because to compute the ϵ -closure of a subset created by determinization, one can run a single shortest-distance algorithm from the states of that subset rather than a distinct one for each state of the subset. Since some states can be reached by several elements of the subset, the first method provides more sharing. This is also corroborated by experiments carried out on various automata in the context of natural language processing applications [11].

This integration of determinization and ϵ -removal can be extended to the weighted case if one replaces the classical unweighted determinization by weighted determinization [7] and if one uses the weighted ϵ -closure presented in the previous sections. It is however limited to the cases where the weighted automaton resulting from ϵ -removal is determinizable since in the weighted case this is not always guaranteed.

Notice, however, that in the integrated algorithm presented in textbooks, the computation of the ϵ -closure is limited to the use of a FIFO queue discipline [1]. Thus, it is a special case of our general algorithm which may be less efficient both in terms of complexity and in practice than using a shortest-first queue discipline as in Dijkstra's algorithm, or a topological order in the case of acyclic automata. Our experiments confirm that in the case of acyclic automata, the ϵ -removal based on a topological order queue discipline can be many times faster than the one based on a shortest-first queue discipline, which itself is faster than the one based on a FIFO queue discipline.

The algorithm we presented can be straightforwardly modified to remove transitions with a label a different from ϵ . This can be done for example by replacing ϵ by a new label and a by ϵ , applying ϵ -removal and then restoring original ϵ 's.

The shortest-distance algorithm presented in the previous sections admits an approximation version where the equality of line 11 in the pseudocode of figure 3 is replaced by an approximate equality modulo some predefined constant δ [8]. This can be used to remove ϵ -transitions of a weighted automaton A over a semiring \mathbb{K} such as $(\mathbb{R}, +, *, 0, 1)$, even when \mathbb{K} is not k -closed for A_ϵ . Although the transition weights are then only approximations of the correct results, this may be satisfactory for many practical purposes such as speech processing applications. Furthermore, one can arbitrarily improve the quality of the approximation by reducing δ . As mentioned before, one can also use a generalization of the Floyd-Warshall algorithm to compute the result in an exact way in the case of machines A_ϵ with relatively small strongly connected components and when the weight of each cycle of A_ϵ admits a well-defined closure.

5. Input ϵ -normalization

This section presents an algorithm closely related to the ϵ -removal algorithm described in the previous sections. Let T_1 be a weighted transducer over a k -closed semiring \mathbb{K} . T_1 may have transitions with input label ϵ . In some applications, one may wish to construct a new transducer T_2 equivalent to T_1 such that T_2 has no ϵ -transition and such that along any successful path of T_2 , no transition with input label different from ϵ is preceded with a transition with input label ϵ .

A transducer such as T_2 is then said to be *input-normalized* and an algorithm that constructs T_2 from T_1 is called *input ϵ -normalization*.^d The advantage of the use of T_2 is that matching with its input in algorithms such as composition does not require any unnecessary delay due to the presence of input ϵ 's. Figure 4 illustrates

^dThere are other criteria for the definition of an input-normalize transducer. Here, we are considering a specific normalization criterion.

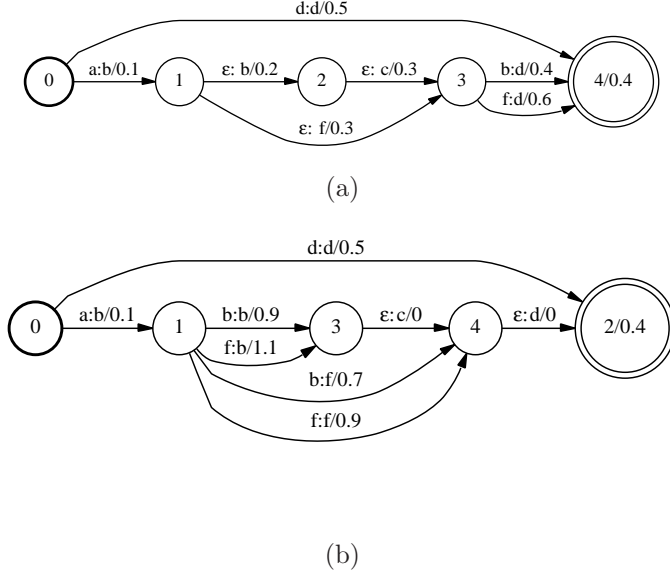


Figure 4: Input ϵ -normalization in the tropical semiring. (a) Weighted transducer T_1 . (b) Weighted transducer T_2 equivalent to T_1 output of the ϵ -normalization algorithm.

the application of the algorithm in the case of a weighted transducer over the tropical semiring.

Note that the system $(\Sigma^*, \cup, \cap, \emptyset, \Sigma^*)$ defines an idempotent semiring \mathbb{K}' over the set of strings Σ^* . A weighted transducer T_1 over the semiring \mathbb{K} can be viewed as a weighted automaton A over the cross-product semiring $\mathbb{S} = \mathbb{K}' \times \mathbb{K}$. An element of \mathbb{S} is a pair (x, w) where $x \in \Sigma^*$ is a string and $w \in \mathbb{K}_1$ the original *weight* of the transducer.

Our input ϵ -normalization algorithm applies only to weighted transducers T_1 over an arbitrary semiring \mathbb{K} that do not admit any cycle with input label ϵ .^e Since T_1 does not admit any cycle with input label ϵ , A_ϵ does not admit any ϵ -cycle and the semiring \mathbb{S} is k -closed for A_ϵ . Thus the shortest-distance algorithm presented in the previous sections can be used with A_ϵ and more generally the ϵ -removal algorithm presented in the previous section applies to A .

Our input ϵ -normalization algorithm is very close to that ϵ -removal algorithm in the specific case of semirings of the type of \mathbb{S} . The result of the ϵ -removal algorithm is an acceptor with no ϵ but with weights in \mathbb{S} . Figure 5 (a) shows the result of that algorithm when applied to the transducer T_1 of figure 4 (a). Our representation of transducers does not allow output labels to be strings, but we can replace transitions with output strings by a set of consecutive transitions with outputs in $\Omega \cup \{\epsilon\}$ in a trivial fashion. The result of that transformation is of interest and is useful in

^eNote that cycles with both input and output ϵ can be removed using for example the ϵ -removal algorithm described in previous sections.

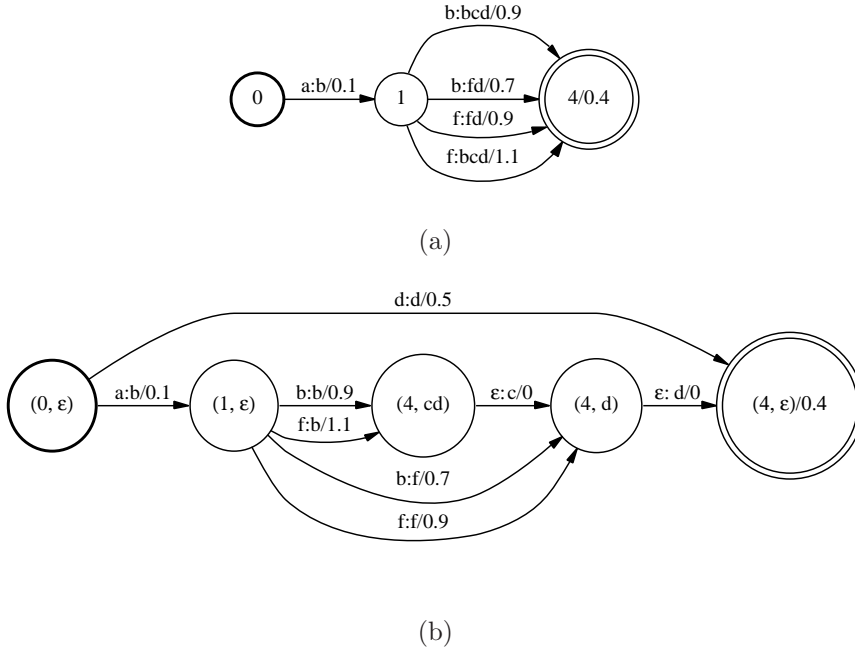


Figure 5: Input ϵ -normalization in the tropical semiring. (a) Weighted transducer with output string labels result of ϵ -removal of T_1 . (b) Construction of the weighted transducer T_2 output of the ϵ -normalization algorithm.

many contexts. However, the resulting transducer does not have the normalization property described above. To ensure that input ϵ 's are not found on a path before an input non- ϵ label, we need to normalize the result of ϵ -removal in a different way.

To each state p , we associate a residual output string y that needs to be output from that state. States in the result of ϵ -normalization correspond to pairs (p, y) where p is a state of the original machine and y a residual output string. The residual string associated to an initial state is just ϵ . The outgoing transitions of a state (p, y_0) are determined as follows. For each transition from state p to q with input label x , output string y_1 and weight w_1 in the result of the ϵ -removal algorithm, a transition is created from state (p, y_0) to state $(q, a^{-1}y_0y_1)$, with input x , output a and weight w_1 , where a is the first letter of y_0y_1 , $a = \epsilon$ if $y_0y_1 = \epsilon$. A state (p, y_0) is final iff p is a final state of the original machine and if $y_0 = \epsilon$. When p is final and $y_0 \neq \epsilon$, then a transition is created from (p, y_0) to (p, y_1) with input ϵ and output a , with $y_0 = ay_1$.

Figure 5 (b) illustrates this construction for the particular case of the transducer of figure 4 (a). As in the case of ϵ -removal, input ϵ -normalization admits a natural on-the-fly implementation since the outgoing transitions of state (p, y_0) of the output automaton can be computed directly using only y_0 and the input ϵ -closure of p .

6. Conclusion

A generic algorithm for ϵ -removal of weighted automata was given. The algorithm works with any semiring covered by a general framework. This framework includes in particular the tropical semiring and the boolean semiring and does not impose idempotent semirings. The algorithm was shown to be more efficient than the existing algorithm both in space and time complexity. Experiments confirm this improvement of efficiency in practice. The algorithm also admits a natural on-the-fly implementation which can be combined with other on-the-fly algorithms such as weighted determinization or composition. An input ϵ -normalization algorithm for weighted transducers was also presented. Composition with input- ϵ -normalized transducers is more efficient since it is not affected by any delay induced by the presence of input ϵ 's.

Acknowledgements

I thank Michael Riley for various discussions about this work and the referees for suggestions which helped improve the initial draft of this paper.

References

1. A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers, Principles, Techniques and Tools*. Addison Wesley: Reading, MA, 1986.
2. C. Cortes and M. Mohri. Context-Free Recognition with Weighted Automata. *Grammars*, 3(2-3), 2000.
3. K. Culik II and J. Kari. Digital images and formal languages. In G. Rozenberg and A. Salomaa, editors, *Handbook of Formal Languages*, pages 599–616. Springer, 1997.
4. R. W. Floyd. Algorithm 97 (shortest path). *Communications of the ACM*, 18, 1968.
5. M. L. Fredman and R. E. Tarjan. Fibonacci heaps and their uses in improved optimization problems. *Communications of the ACM*, 34, 1987.
6. W. Kuich and A. Salomaa. *Semirings, Automata, Languages*. Number 5 in EATCS Monographs on Theoretical Computer Science. Springer-Verlag, Berlin-New York, 1986.
7. M. Mohri. Finite-State Transducers in Language and Speech Processing. *Computational Linguistics*, 23:2, 1997.
8. M. Mohri. General Algebraic Frameworks and Algorithms for Shortest-Distance Problems. Technical Memorandum 981210-10TM, AT&T Labs - Research, 62 pages, 1998.
9. M. Mohri, F. C. N. Pereira, and M. Riley. The design principles of a weighted finite-state transducer library. *Theoretical Computer Science*, 231:17–32, January 2000.
10. K. Thompson. Regular expression search algorithm. *Comm. ACM*, 11, 1968.
11. G. van Noord. Treatment of epsilon moves in subset construction. *Computational Linguistics*, 26:1, 2000.
12. S. Warshall. A theorem on boolean matrices. *Journal of the ACM*, 9(1):11–12, 1962.