

Patterns, Regular Expressions and Finite Automata

(LECTURE 5)

Patterns and their defined languages

- S : a finite alphabet
- A pattern is a string of symbols representing a set of strings in S^* .
- The set of all patterns is defined inductively as follows:
 1. atomic patterns:
 $a \in S, e, \emptyset, \#, @$.
 2. compound patterns: if a and b are patterns, then so are: $a + b, a \cap b, a^*, a^+, \sim a$ and $a \cdot b$.
- For each pattern a , $L(a)$ is the language represented by a and is defined inductively as follows:
 1. $L(a) = \{a\}, L(e) = \{e\}, L(\emptyset) = \{\}, L(\#) = S, L(@) = S^*$.
 2. If $L(a)$ and $L(b)$ have been defined, then
$$L(a + b) = L(a) \cup L(b), \quad L(a \cap b) = L(a) \cap L(b).$$
$$L(a^+) = L(a)^+, \quad L(a^*) = L(a)^*,$$
$$L(\sim a) = S^* - L(a), \quad L(a \cdot b) = L(a) \cdot L(b).$$

More on patterns

- We say that a string x matches a pattern a iff $x \in L(a)$.
- Some examples:
 1. $S^* = L(@) = L(\#^*)$
 2. $L(x) = \{x\}$ for any $x \in S^*$
 3. for any x_1, \dots, x_n in S^* , $L(x_1 + x_2 + \dots + x_n) = \{x_1, x_2, \dots, x_n\}$.
 4. $\{x \mid x \text{ contains at least 3 a's}\} = L(@a@a@a@)$
 5. $S\{a\} = \# \cap \sim a$
 6. $\{x \mid x \text{ does not contain } a\} = (\# \cap \sim a)^*$
 7. $\{x \mid \text{every 'a' in } x \text{ is followed sometime later by a 'b'}\} =$
 $= \{x \mid \text{either no 'a' in } x \text{ or } \$ \text{ 'b' in } x \text{ followed no 'a'}\}$
 $= (\# \cap \sim a)^* + @b(\# \cap \sim a)^*$

More on pattern matching

- Some interesting and important questions:
 1. How hard is it to determine if a given input string x matches a given pattern a ?
==> efficient algorithm exists
 2. Can every set be represented by a pattern ?
==> no! the set $\{a^n b^n \mid n > 0\}$ cannot be represented by any pattern.
 3. How to determine if two given patterns a and b are equivalent ?
(I.e., $L(a) = L(b)$) --- an exercise !
 4. Which operations are redundant ?
 - $e = \sim(\#^+ \cap @) = \emptyset^*$; $a^+ = a \cdot a^*$
 - $\# = a_1 + a_2 + \dots + a_n$ if $S = \{a_1, \dots, a_n\}$
 - $a + b = \sim(\sim a \cap \sim b)$; $a \cap b = \sim(\sim a + \sim b)$
 - It can be shown that \sim is redundant.

Equivalence of patterns, regular expr. & FAs

- Recall that regular expressions are those patterns that can be built from: $a \in S, \epsilon, \emptyset, +, \cdot$ and $*$.
- Notational conventions:
 - $a + br$ means $a + (br)$
 - $a + b^*$ means $a + (b^*)$
 - $a b^*$ means $a (b^*)$

Theorem 8: Let $A \subseteq S^*$. Then the followings are equivalent:

1. A is regular (i.e., $A = L(M)$ for some FA M),
2. $A = L(a)$ for some pattern a ,
3. $A = L(b)$ for some regular expression b .

pf: Trivial part: (3) \Rightarrow (2).

(2) \Rightarrow (1) to be proved now!

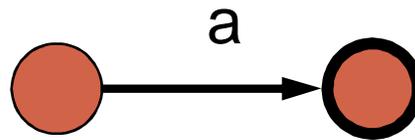
(1) \Rightarrow (3) later.

(2) \Rightarrow (1) : Every set represented by a pattern is regular

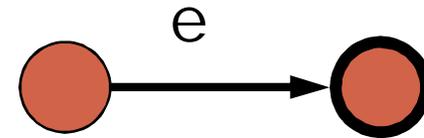
Pf: By induction on the structure of pattern a .

Basis: a is atomic: (by construction!)

1. $a = a$:



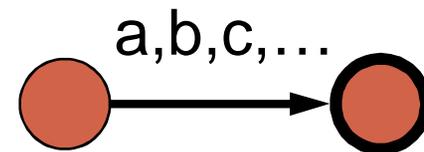
2. $a = \epsilon$:



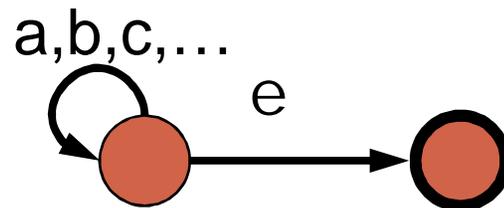
3. $a = \emptyset$:



4. $a = \#$:



5. $a @ = \#^*$:



Inductive cases: Let M_1 and M_2 be any FAs accepting $L(b)$ and $L(g)$, respectively.

6. $a \approx b g : \Rightarrow L(a) = L(M_1 \cdot M_2)$

7. $a \approx b^* : \Rightarrow L(a) = L(M_1^*)$

8. $a \approx b + g$, $a \approx \sim b$ or $a \approx b \cap g$: By ind. hyp. b and g are regular. Hence by closure properties of regular languages, a is regular, too.

9. $a \approx b^+ = b b^*$: Similar to case 8.

Some examples patterns & their equivalent FAs

1. $(aaa)^* + (aaaaa)^*$

(1) \Rightarrow (3): Regular languages can be represented by reg. expr.

$M = (Q, S, d, S, F)$: a NFA; $X \subseteq Q$: a set of states; $m, n \in Q$: two states

- $p^X(m, n) =_{\text{def}} \{y \in S^* \mid \exists \text{ a path from } m \text{ to } n \text{ labeled } y \text{ and all intermediate states } \in X \}$.
- Note: $L(M) = ?$
- $p^X(m, n)$ can be shown to be representable by a regular expr, by induction as follows:

Let $D(m, n) = \{ a \mid (m \xrightarrow{a} n) \in d \} = \{a_1, \dots, a_k\}$ ($k \geq 0$)

= the set of symbols by which we can reach from m to n , then

Basic case: $X = \emptyset$:

1.1 if $m \neq n$: $p^{\emptyset}(m, n) = \{a_1, a_2, \dots, a_k\} = L(a_1 + a_2 + \dots + a_k)$ if $k > 0$,
= $\{\}$ = $L(\emptyset)$ if $k = 0$.

1.2 if $m = n$: $p^{\emptyset}(m, n) = \{a_1, a_2, \dots, a_k, e\} = L(a_1 + a_2 + \dots + a_k + e)$ if $k > 0$,
= $\{e\}$ = $L(e)$ if $k = 0$.

3. For nonempty X , let q be any state in X , then :

$$p^X(m, n) = p^{X-\{q\}}(m, n) \cup p^{X-\{q\}}(m, q) (p^{X-\{q\}}(q, q))^* p^{X-\{q\}}(q, n).$$

By Ind.hyp.(why?), there are regular expressions a, b, g with
 $L([a, b, g]) = \{ p^{X-\{q\}}(m, n), p^{X-\{q\}}(m, q), (p^{X-\{q\}}(q, q)), p^{X-\{q\}}(q, n) \}$

$$\begin{aligned} \text{Hence } p^X(m, n) &= L(a) \cup L(b) \cup L(g)^* L(r), \\ &= L(a + b g^* r) \end{aligned}$$

and can be represented as a reg. expr.

- Finally, $L(M) = \{x \mid s \xrightarrow{x} f, s \in S, f \in F\}$
 $= \bigcup_{s \in S, f \in F} p^0(s, f)$, is representable by a regular expression.

Some examples

Example (9.3): M :

- $L(M) = p^{\{p,q,r\}}(p,p) = p^{\{p,r\}}(p,p) + p^{\{p,r\}}(p,q) (p^{\{p,r\}}(q,q))^* p^{\{p,r\}}(q,p)$
- $p^{\{p,r\}}(p,p) = ?$
- $p^{\{p,r\}}(p,q) = ?$
- $p^{\{p,r\}}(q,q) = ?$
- $p^{\{p,r\}}(q,p) = ?$

Hence $L(M) = ?$

	0	1
>pF	{p}	{q}
q	{r}	{}
r	{p}	{q}

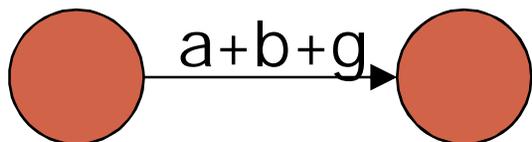
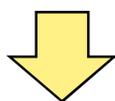
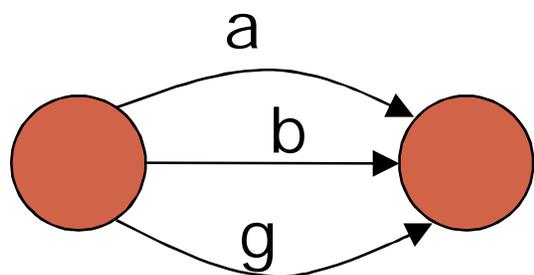
Another approach

- The previous method
 - easy to prove,
 - easy for computer implementation, but
 - **hard for human computation.**
- The strategy of the new method:
 - reduce the number of states in the target FA and
 - encodes path information by regular expressions on the edges.
 - until there is one or two states : one is the start state and one is the final state.

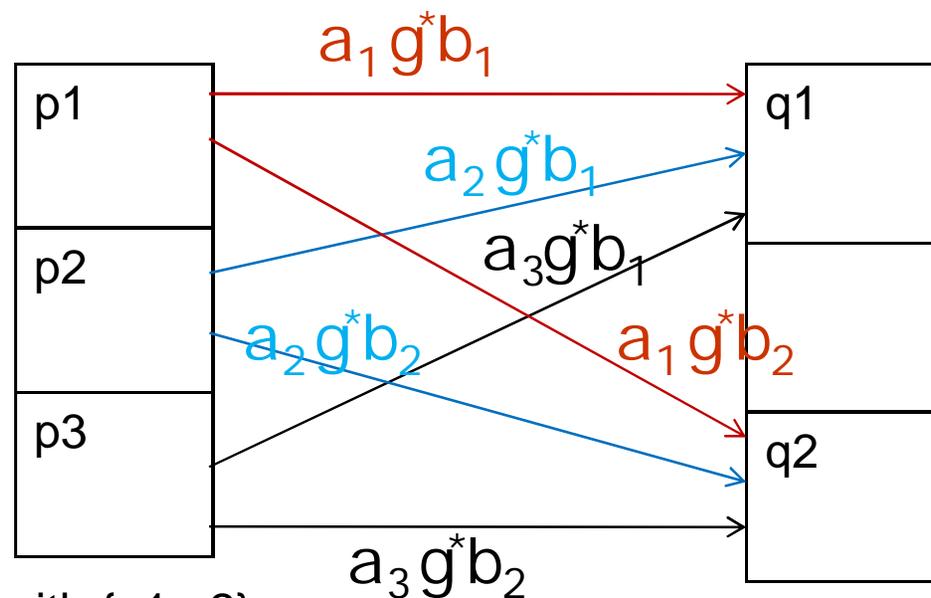
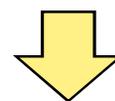
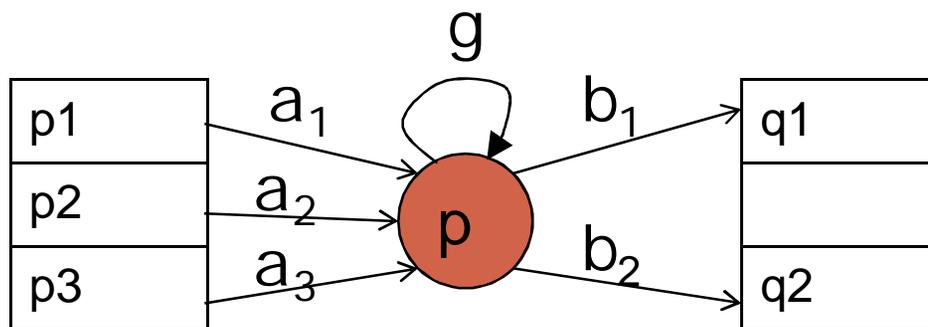
Steps

0. Assume the machine M has only one start state and one final state. Both may probably be identical.
1. While there exists a third state p that is neither start nor final:
 - 1.1 (**Merge edges**) For each pair of states (q,r) that has more than 1 edges with labels t_1, t_2, \dots, t_n , respectively, then merge these edges by a new one with regular expression $t = t_1 + t_2 \dots + t_n$.
 - 1.2 (**Replace state p by edges; remove state**)
Let $(p_1, a_1, p), \dots, (p_n, a_n, p)$ where $p_j \neq p$ be the collection of all edges in M with p as the destination state, and
 $(p, b_1, q_1), \dots, (p, b_m, q_m)$ where $q_j \neq p$ be the collection of all edges with p as the start state. Now the state p together with all its connecting edges can be removed and replaced by a set of $m \times n$ new edges :
 $\{ (p_i, a_i t^* b_j, q_j) \mid i \text{ in } [1,n] \text{ and } j \text{ in } [1,m] \}$.
The new machine is equivalent to the old one.

- Merge Edges :



- Replace state by Edges



Note: {p1,p2,p3} may intersect with {q1,q2}.

2. perform 1.1 once again (merge edges)

// There are one or two states now

3 Two cases to consider:

3.1 The final machine has only one state, that is both start and final. Then if there is an edge labeled t on the state,

then t^* is the result, other the result is ϵ .

3.2 The machine has one start state s and one final state f .

Let $(s, s \rightarrow s, s)$, $(f, f \rightarrow f, f)$, $(s, s \rightarrow f, f)$ and $(f, f \rightarrow s, s)$ be the collection of all edges in the machine, where $(s \rightarrow f)$ means the regular expression or label on the edge from s to f .

The result then is

Example

	0	1
>p	{p,r}	{q,r}
q	{r}	{p,q,r}
rF	{p,q}	{q,r}

1. another representation

	p	q	r
p	0	1	0,1
q	1	1	0,1
r	0	0,1	1

Merge edges

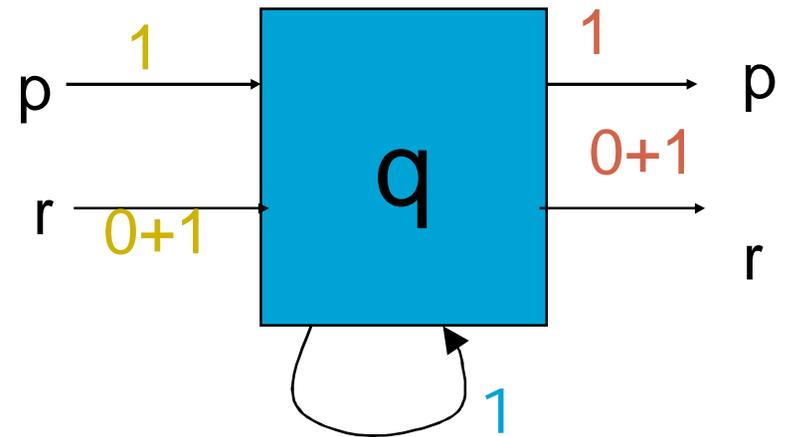
	p	q	r
p	0	1	0,1
q	1	1	0,1
r	0	0,1	1

	p	q	r
p	0	1	0+1
q	1	1	0+1
r	0	0+1	1

remove q

	p	q	r
p	0, 11*1	1	0+1, 11*(0+1)
q	1	1,	0+1
r	0, (0+1)1*1	0+1	1, (0+1)1*(0+1)

	p	q	r
p	0	1	0+1
q	1	1	0+1
r	0	0+1	1



Form the final result

	p	r
$p \rightarrow p$	$0 + 11^*1$	$0 + 1 + 11^*(0 + 1)$
$r \rightarrow r$	$0 + (0 + 1)1^*1$	$1 + (0 + 1)1^*(0 + 1)$

Final result := $[p \rightarrow p + (p \rightarrow r) (r \rightarrow r)^* (r \rightarrow p)]^* (p \rightarrow r) (r \rightarrow r)^*$

$[(0 + 11^*1) + (0 + 1 + 11^*(0 + 1)) (1 + (0 + 1)1^*(0 + 1))^* (0 + (0 + 1)1^*1)]^*$
 $(0 + 1 + 1^*(0 + 1))(1 + (0 + 1)(0 + 1)^*)$