# Lecture 8

# Pattern Matching and Regular Expressions

Here are some interesting and important questions:

- How hard is it to determine whether a given string $x$ matches a given pattern $\alpha$? This is an important practical question. There are very efficient algorithms, as we will see.

- Is every set represented by some pattern? Answer: no. For example, the set

$$\{a^n b^n \mid n \geq 0\}$$

  is not represented by any pattern. We'll prove this later.

- Patterns $\alpha$ and $\beta$ are *equivalent* if $L(\alpha) = L(\beta)$. How do you tell whether $\alpha$ and $\beta$ are equivalent? Sometimes it is obvious and sometimes not.

- Which operators are redundant? For example, we can get rid of $\epsilon$ since it is equivalent to $\sim(\#@)$ and also to $\emptyset^*$. We can get rid of @ since it is equivalent to $\#^*$. We can get rid of unary $^+$ since $\alpha^+$ is equivalent to $\alpha\alpha^*$. We can get rid of $\#$, since if $\Sigma = \{a_1, \ldots, a_n\}$ then $\#$ is equivalent to the pattern

$$a_1 + a_2 + \cdots + a_n.$$

The operator $\cap$ is also redundant, by one of the De Morgan laws:

$\alpha \cap \beta$ is equivalent to $\sim(\sim\alpha + \sim\beta)$.

Redundancy is an important question. From a user's point of view, we would like to have a lot of operators since this lets us write more succinct patterns; but from a programmer's point of view, we would like to have as few as possible since there is less code to write. Also, from a theoretical point of view, fewer operators mean fewer cases we have to treat in giving formal semantics and proofs of correctness.

An amazing and difficult-to-prove fact is that the operator $\sim$ is redundant. Thus every pattern is equivalent to one using only atomic patterns $a \in \Sigma$, $\epsilon$, $\emptyset$, and operators $+$, $\cdot$, and $^*$. Patterns using only these symbols are called *regular expressions*. Actually, as we have observed, even $\epsilon$ is redundant, but we include it in the definition of regular expressions because it occurs so often.

Our goal for this lecture and the next will be to show that the family of subsets of $\Sigma^*$ represented by patterns is exactly the family of regular sets. Thus as a way of describing subsets of $\Sigma^*$, finite automata, patterns, and regular expressions are equally expressive.

## Some Notational Conveniences

Since the binary operators $+$ and $\cdot$ are associative, that is,

$$L(\alpha + (\beta + \gamma)) = L((\alpha + \beta) + \gamma),$$
$$L(\alpha(\beta\gamma)) = L((\alpha\beta)\gamma),$$

we can write

$\alpha + \beta + \gamma$   and   $\alpha\beta\gamma$

without ambiguity. To resolve ambiguity in other situations, we assign precedence to operators. For example,

$\alpha + \beta\gamma$

could be interpreted as either

$\alpha + (\beta\gamma)$   or   $(\alpha + \beta)\gamma$,

which are not equivalent. We adopt the convention that the concatenation operator $\cdot$ has higher precedence than $+$, so that we would prefer the former interpretation. Similarly, we assign $^*$ higher precedence than $+$ or $\cdot$, so that

$\alpha + \beta^*$

is interpreted as

$$\alpha + (\beta^*)$$

and not as

$$(\alpha + \beta)^*.$$

All else failing, use parentheses.

## Equivalence of Patterns, Regular Expressions, and Finite Automata

Patterns, regular expressions (patterns built from atomic patterns $a \in \Sigma$, $\epsilon$, $\emptyset$, and operators $+$, $^*$, and $\cdot$ only), and finite automata are all equivalent in expressive power: they all represent the regular sets.

Theorem 8.1    *Let $A \subseteq \Sigma^*$. The following three statements are equivalent:*

*(i) $A$ is regular; that is, $A = L(M)$ for some finite automaton $M$;*

*(ii) $A = L(\alpha)$ for some pattern $\alpha$;*

*(iii) $A = L(\alpha)$ for some regular expression $\alpha$.*

*Proof.* The implication (iii) $\Rightarrow$ (ii) is trivial, since every regular expression is a pattern. We prove (ii) $\Rightarrow$ (i) here and (i) $\Rightarrow$ (iii) in Lecture 9.

The heart of the proof (ii) $\Rightarrow$ (i) involves showing that certain basic sets (corresponding to atomic patterns) are regular, and the regular sets are closed under certain closure operations corresponding to the operators used to build patterns. Note that

- the singleton set $\{a\}$ is regular, $a \in \Sigma$,

- the singleton set $\{\epsilon\}$ is regular, and

- the empty set $\emptyset$ is regular,

since each of these sets is the set accepted by some automaton. Here are nondeterministic automata for these three sets, respectively:



Also, we have previously shown that the regular sets are closed under the set operations $\cup$, $\cap$, $\sim$, $\cdot$, $^*$, and $^+$; that is, if $A$ and $B$ are regular sets, then so are $A \cup B$, $A \cap B$, $\sim A = \Sigma^* - A$, $AB$, $A^*$, and $A^+$.

These facts can be used to prove inductively that (ii) $\Rightarrow$ (i). Let $\alpha$ be a given pattern. We wish to show that $L(\alpha)$ is a regular set. We proceed by

induction on the structure of $\alpha$. The pattern $\alpha$ is of one of the following forms:

| | | | |
|---|---|---|---|
| (i) | $a$, where $a \in \Sigma$; | (vi) | $\beta + \gamma$; |
| (ii) | $\epsilon$; | (vii) | $\beta \cap \gamma$; |
| (iii) | $\emptyset$; | (viii) | $\beta\gamma$; |
| (iv) | #; | (ix) | $\sim\beta$; |
| (v) | @; | (x) | $\beta^*$; |
| | | (xi) | $\beta^+$. |

There are five base cases (i) through (v) corresponding to the atomic patterns and six induction cases (vi) through (xi) corresponding to compound patterns. Each of these cases uses a closure property of the regular sets previously observed.

For (i), (ii), and (iii), we have $L(a) = \{a\}$ for $a \in \Sigma$, $L(\epsilon) = \{\epsilon\}$, and $L(\emptyset) = \varnothing$, and these are regular sets.

For (iv), (v), and (xi), we observed earlier that the operators #, @, and $^+$ were redundant, so we may disregard these cases since they are already covered by the other cases.

For (vi), recall that $L(\beta + \gamma) = L(\beta) \cup L(\gamma)$ by definition of the + operator. By the induction hypothesis, $L(\beta)$ and $L(\gamma)$ are regular. Since the regular sets are closed under union, $L(\beta + \gamma) = L(\beta) \cup L(\gamma)$ is also regular.

The arguments for the remaining cases (vii) through (x) are similar to the argument for (vi). Each of these cases uses a closure property of the regular sets that we have observed previously in Lectures 4 and 6.    □
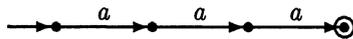
**Example 8.2**    Let's convert the regular expression
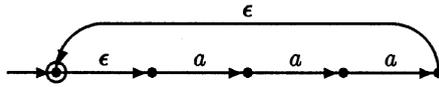
$$(aaa)^* + (aaaaa)^*$$

for the set

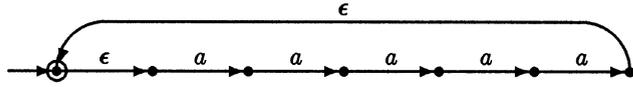$$\{x \in \{a\}^* \mid |x| \text{ is divisible by either 3 or 5}\}$$

to an equivalent NFA. First we show how to construct an automaton for $(aaa)^*$. We take an automaton accepting only the string $aaa$, say
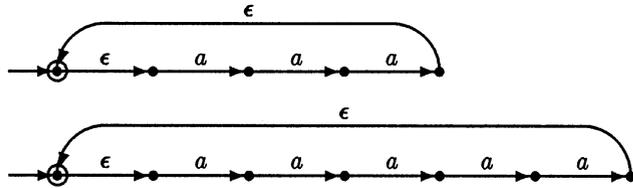


Applying the construction of Lecture 6, we add a new start state and $\epsilon$-transitions from the new start state to all the old start states and from all the old accept states to the new start state. We let the new start state be the only accept state of the new automaton. This gives

The construction for $(aaaaa)^*$ is similar, giving



To get an NFA for $(aaa)^* + (aaaaa)^*$, we can simply take the disjoint union of these two automata: