

# Formal Languages

Wednesday, September 19, 2007

Reading: Stoughton 2.1, 2.2, end of 2.3, beginning of 3.1

---

## CS235 Languages and Automata

Department of Computer Science  
Wellesley College

## Alphabets, Strings, and Languages

An **alphabet** is a set of symbols.

E.g.:  $\Sigma_1 = \{0,1\}$ ;  $\Sigma_2 = \{-,0,+ \}$   $\Sigma_3 = \{a,b, \dots, y, z\}$ ;  $\Sigma_4 = \{\odot, \Rightarrow, \boxed{a}, \boxed{aa}\}$

A **string over  $\Sigma$**  is a sequence of symbols from  $\Sigma$ .

The empty string is traditionally written  $\epsilon$ , but Stoughton uses  $\%$ .

$\Sigma^*$  denotes all strings over  $\Sigma$ . E.g.:

- $\Sigma_1^*$  contains  $\%, 0, 1, 00, 01, 10, 11, 000, \dots$
- $\Sigma_2^*$  contains  $\%, -, 0, +, --, -0, -+, 0-, 00, 0+, +- , +0, ++, ---, \dots$
- $\Sigma_3^*$  contains  $\%, a, b, \dots, aa, ab, \dots, bar, baz, foo, wellesley, \dots$
- $\Sigma_4^*$  contains  $\%, \odot, \Rightarrow, \boxed{a}, \boxed{aa}, \dots, \boxed{a} \Rightarrow \odot, \dots, \boxed{a \odot a}, \boxed{aa a}, \dots$

A **language over  $\Sigma$**  (Stoughton's  $\Sigma$ -language) is any subset of  $\Sigma^*$ .

I.e., it's a (possibly countably infinite) set of strings over  $\Sigma$ . E.g.:

- $L_1$  over  $\Sigma_1$  is all sequences of 1s and all sequences of 10s.
- $L_2$  over  $\Sigma_2$  is all strings with equal numbers of  $-$ ,  $0$ , and  $+$ .
- $L_3$  over  $\Sigma_3$  is all lowercase words in the OED.
- $L_4$  over  $\Sigma_4$  is  $\{\odot, \odot \Rightarrow \odot, \boxed{a \odot a}\}$ .

## String Operations

**Length:**  $|s|$  is the length of a string  $s$ . E.g.:

$$|\%| = 0, |\text{foo}| = 3, |\text{☺} \Rightarrow \boxed{\text{a|aa}}| = 4$$

**Concatenation:** If  $x, y$  in  $\Sigma^*$ , then  $x@y$  in  $\Sigma^*$  is the string consisting of all symbols in  $x$  followed by all symbols in  $y$ . (Concatenation is traditionally written  $xy$ , and Stoughton uses this as an abbreviation.)  
E.g.  $\text{baz@quux} = \text{bazquux}$

*Concatenation Properties:*

- $(x@y)@z = xyz = x@(y@z)$  (Associativity)
  - $x@% = x = %@x$  (Identity)
  - $|x@y| = |x| + |y|$
- Monoid**

*Other Definitions:*

- $x$  is a **prefix** of  $y$  iff  $y = xv$  for some  $v$
- $x$  is a **suffix** of  $y$  iff  $y = ux$  for some  $u$
- $x$  is a **substring** of  $y$  iff  $y = uxv$  for some  $u$  and  $v$

There are **proper** versions of these, too.

What are all prefixes, suffixes, substrings of **bar**?

Formal Languages 7-3

## More String Operations

**String Powers:** Suppose  $x$  is a string.

- $x^0 = \%$
- $x^n = x@x^{n-1}$ , abbreviated  $xx^{n-1} = x(x^{n-1})$

*Power Properties:*

- $x^{a+b} = x^a x^b$
- $|x^n| = n \cdot |x|$

**String Reversal:** Suppose  $a$  is a symbol and  $x$  is a string.

- $\%^R = \%$
- $(a@x)^R = x^R@a$

*Reversal Properties:*

- $(x@y)^R = y^R@x^R$
- $(x^R)^R = x$
- $|x^R| = |x|$

**Alphabet:** Suppose  $a$  is a symbol and  $x$  is a string.

- $\text{alphabet}(\%) = \{ \}$
- $\text{alphabet}(ax) = \{a\} \cup \text{alphabet}(x)$

What are some alphabet properties?

Formal Languages 7-4

## String Induction

Suppose  $P(w)$  is a property of strings  $w$  in  $\Sigma^*$ . Can prove  $P(w)$  by *natural induction (or strong induction) on  $|w|$* . Equivalently:

### Right String Induction:

Suppose that

1. (basis step)  $P(\epsilon)$  holds.

2. (inductive step) For all  $a \in \Sigma$  and  $x$  in  $\Sigma^*$ ,  $P(x) \Rightarrow P(ax)$ .

Then  $P(w)$  holds for all  $w \in \Sigma^*$ .

the inductive  
hypothesis (IH)

$P(x)$

### Left String Induction:

(inductive step) For all  $a \in \Sigma$  and  $x$  in  $\Sigma^*$ ,  $P(x) \Rightarrow P(ax)$ .

### Strong String Induction:

(inductive step) For all  $w \in \Sigma^*$ ,  $(\bigwedge_{x \in \Sigma^* \text{ s.t. } |x| < |w|} P(x)) \Rightarrow P(w)$ .

Formal Languages 7-5

## String Induction Example: Reversal

Prove that  $(x@y)^R = y^R@x^R$

Which form of induction to use?

(basis step)

(inductive step)

What is I.H.?

Formal Languages 7-6

## Language Operations

Suppose  $L_1$  and  $L_2$  are  $\Sigma$ -languages.

**Set Operations:** The following are all  $\Sigma$ -languages:

$$L_1 \cup L_2, L_1 \cap L_2, L_1 - L_2, \overline{L_1} (= \Sigma^* - L_1)$$

**Language Concatenation:**

$$L_1 @ L_2 = \{x @ y \mid x \text{ in } L_1 \text{ and } y \text{ in } L_2\} \text{ (abbreviated } L_1L_2)$$

$$\text{E.g. } \{CS, PHYS\} @ \{110, 111, 115\} =$$

*Concatenation Properties:*

- $(L_1 @ L_2) @ L_3 = L_1 @ (L_2 @ L_3)$  (Associativity)
- $\{\%\} @ L = L = L @ \{\%\}$  (Identity)
- $\emptyset @ L = \emptyset = L @ \emptyset$  (Zero)
- $|L_1 @ L_2| = |L_1| \cdot |L_2|$  for finite  $L_1, L_2$

Formal Languages 7-7

## More Language Operations

**Language Powers:**

- $L^0 = \{\%\}$
- $L^n = L @ L^{n-1}$

$$\text{E.g., } \{0,1\}^2 =$$

*Power Properties:*

- $L^{a+b} = L^a @ L^b$
- $|L^n| = |L|^n$  for finite  $L$
- $\{x\}^n = \{x^n\}$
- $\{\%\}^n = \{\%\}$

**Kleene Closure (Kleene Star):** (pronounced "clay knee")

$$L^* = \{L^n \mid n \text{ in Nat}\}$$

$$\text{E.g. } \{0,1\}^* =$$

(This is consistent with notation  $\Sigma^*$ .)

Formal Languages 7-8

## Stoughton's Conventions

He fixes the set of symbols to be **Sym**, which contains:

- digits 0, 1, ..., 9
- upper case letters A, B, ..., Z
- lower case letters a, b, ..., z
- < and > delimiting any sequence of ASCII characters in which < and > are properly nested. E.g., <plus>, <<a>, <b, c>>

Symbols are ordered first by length, and then lexicographically:

0 < ... < 9 < A < ... < Z < a ... < z < <0> < ... < z> < <00> ... < <zz> < ...

He fixes the universe of strings to be **Str** = all strings over **Sym**.

He uses % for the empty string.

Strings are ordered first by length, and then lexicographically. E.g.:

% < a < <ab> < ab < ac < a<bbb> < <bb><aa> < <bbb>a < aa<bb> < <bb>aa

Can't write  $(a + b) * c$ . Instead, write <LPAREN>a<PLUS>b<RPAREN><TIMES>c