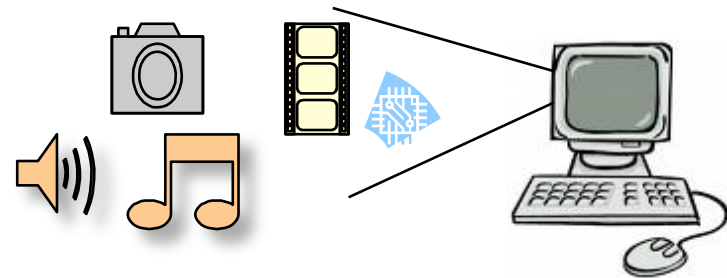


Words and Languages

Discrete Mathematics
Evgeny Skvortsov

Why Strings?

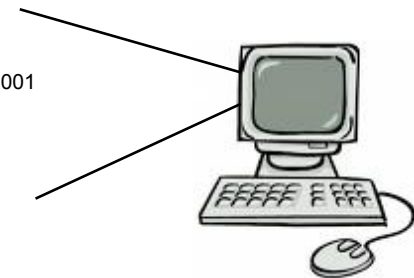
- Computer data is very diverse



- However, before it can be processed it must be converted into

- ... strings

```
00101101010101010100101010101010110101010010010101010101010101001
1100101101010101010100101010101010101101010100100101010101010101010
```



Alphabets and Strings

- An **alphabet** is any finite set. Σ

Its elements are called **symbols** or **letters**

$\{0,1\}$ a binary alphabet

$\{0,1,2,3,4,5,6,7,8,9\}$

$\{a,b,\dots,x,y,z\}$ Latin alphabet

$\{a, б, \dots, э, ю, я\}$ Cyrillic

.....

$\left\{ \begin{array}{l} \text{I 3 H T X - F X - I X F X} \\ \text{X 7 8 X X X X Y T 6 X 7 X} \end{array} \right\}$ the dancing men alphabet

- A **string over alphabet** Σ is a sequence of symbols from Σ connected by means of **juxtaposition** or **concatenation**

0010110101010101001010

Powers and Empty String

- Strings that are obtained by concatenation of the same number of symbols are grouped into **powers** of the alphabet.

- An induction definition:

$$\Sigma^1 = \Sigma$$

$\Sigma^{n+1} = \{xy \mid x \in \Sigma, y \in \Sigma^n\}$ where xy denotes the concatenation of x and y

- If $\Sigma = \{0,1\}$, then $\Sigma^2 = \{00,01,10,11\}$

- If $\Sigma = \{a,b,c,\dots,z\}$, then $\Sigma^3 = \{act, bad, cat, den, \dots\}$

Note that $fre, aat, lkj \in \Sigma^3$

- Empty string λ is the string containing no symbols.

λ is not the blank symbol, and λ cannot be a symbol in an alphabet

- $\Sigma^0 = \{\lambda\}$

More Powers

● If a string x is obtained by concatenation of n symbols, we say that it has **length** n , denoted $||x|| = n$

● $||\lambda|| = 0$

● $\Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \dots = \bigcup_{n=1}^{\infty} \Sigma^n$

● $\Sigma^* = \Sigma^+ \cup \{\lambda\}$ Kleene star

● Examples:

$$\{0,1\}^+ = \{0,1,00,01,10,11,000,001,010,100,011,101,110,111,\dots\}$$

$$\{a,b\}^* = \{\lambda, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$$

Equality and Concatenation

- Two strings $x = x_1x_2 \dots x_k$ and $y = y_1y_2 \dots y_l$ are **equal** if $k = l$ and $x_1 = y_1, x_2 = y_2, x_3 = y_3, \dots, x_k = y_k$
- **Concatenation** of strings $x = x_1x_2 \dots x_k$ and $y = y_1y_2 \dots y_l$ is the string $xy = x_1x_2 \dots x_ky_1y_2 \dots y_l$
- Observe that $\|xy\| = \|x\| + \|y\|$
- For the empty string we define

$$\begin{aligned}\lambda x &= \lambda x_1x_2 \dots x_k = x_1x_2 \dots x_k = x \\ x\lambda &= x_1x_2 \dots x_k\lambda = x_1x_2 \dots x_k = x \\ \lambda\lambda &= \lambda\end{aligned}$$

- Power of a string: $x^0 = \lambda, x^1 = x, x^2 = xx, x^3 = xxx, \dots$

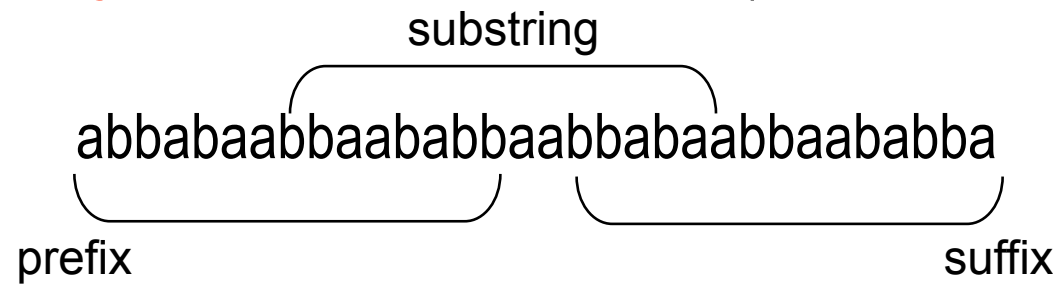
Prefixes and Suffixes

- For any strings x, w such that $w = xy$ for some string y , string x is called a **prefix** of w .
- For any strings y, w such that $w = xy$ for some string x , string y is called a **suffix** of w .
- Note that the empty string is a prefix and suffix of any string

Examples:

$x = \text{abbabaab}$ Then its prefixes are $\lambda, a, ab, abb, abba, abbab, \dots$
 and suffixes are $\lambda, b, ab, aab, baab, abaab, \dots$

A **substring** of w is a prefix of a suffix (or a suffix of a prefix)



Languages

- A (formal) **language** over alphabet Σ is a subset of Σ^*
The **empty language** \emptyset is also a language
- $\{001, 010, 100\}$ is a language over alphabet $\{0, 1\}$
- $\{0, 10, 100, 110, 1000, 1010, 1100, 1110, 10000, \dots\}$
the language of binary expansions of even numbers
- The English language = the set of grammatically correct English texts
over alphabet $\{a, \dots, z, A, \dots, Z, 0, \dots, 9, \dots, ;, :, !, ?, -, \dots\}$
- Thue-Morse language $\{a, b, ab, ba, abba, baab, abbabaab, baababba, \dots\}$
- The language of properly placed parenthesis
 $\{(), (()), ()(), ((())), (()()), (())(), ()()(), \dots\}$ Dyck language

Representing Languages: Language Operations

- Finite languages can be described by a list of their elements:
{if, then, while, for, to,...}

- Set theoretic operations:

Union of languages $A \cup B$

Intersection of languages $A \cap B$

Complement of a language $\bar{A} = \Sigma^* - A$

- Examples

A is the language of binary expansions of even numbers

B is the language of binary expansions of odd numbers

$A \cap B, A \cup B, \bar{A}$

Representing Languages: Concatenation

- Let A and B be languages over alphabet Σ . Then

$$AB = \{xy \mid x \in A, y \in B\}$$

- Let $A = \{a^k \mid k \in \mathbb{N}\}$
 $B = \{b^k \mid k \in \mathbb{N}\}$
 $C = \{c\}$

Then $ACB = \{a^k cb^l \mid k, l \in \mathbb{N}\}$

- Concatenation is not commutative

Let $A = \{a, ab, c\}$ and $B = \{\lambda, b\}$. Then $AB = \{a, ab, abb, c, cb\}$ and $BA = \{a, ab, c, ba, bab, bc\}$.

In particular, $|AB| = 5 \neq 6 = |BA|$

Properties of Concatenation

Theorem

Let Σ be an alphabet, let A, B, C be languages over Σ .

- | | |
|---|--|
| (1) $A\{\lambda\} = \{\lambda\}A = A,$ | (4) $(AB)C = A(BC)$ |
| (2) $A(B \cup C) = AB \cup AC,$ | (5) $(B \cup C)A = BA \cup CA$ |
| (3) $A(B \cap C) \subseteq AB \cap AC,$ | (6) $(B \cap C)A \subseteq BA \cap CA$ |

Proof

(4) Take $w \in (AB)C$. Then $w = xy$ such that $x \in AB$ and $y \in C$.

Then $x = uv$, where $u \in A$ and $v \in B$. If $u = u_1u_2 \dots u_k,$

$v = v_1v_2 \dots v_l,$ and $y = y_1y_2 \dots y_m,$ then

$$\begin{aligned} w &= (u_1u_2 \dots u_kv_1v_2 \dots v_l)y_1y_2 \dots y_m = u_1u_2 \dots u_kv_1v_2 \dots v_ly_1y_2 \dots y_m \\ &= u_1u_2 \dots u_k(v_1v_2 \dots v_ly_1y_2 \dots y_m) = u(vy) \in A(BC). \end{aligned}$$

The reverse inclusion is similar.

Properties of Concatenation (cntd)

(6) Take $w \in (B \cap C)A$. Then $w = xy$, where $x \in B \cap C$ and $y \in A$. Thus $x \in B$, and so $w = xy \in BA$. Similarly, $x \in C$, hence $w = xy \in CA$. We conclude $w \in BA \cap CA$.

Q. E. D.

● Example

Let $B = \{a\}$, $C = \{ab\}$, and $A = \{c, bc\}$. Then $BA = \{ac, abc\}$, $CA = \{abc, abbc\}$. Therefore $BA \cap CA = \{abc\}$, while $B \cap C = \emptyset$, and so $(B \cap C)A = \emptyset$.

Kleene Star

- For a language A over an alphabet Σ :

$$A^0 = \{\lambda\}, A^1 = A, \text{ and for } n \in \mathbb{N}, A^n = A^{n+1}A = \underbrace{AA \dots A}_{n \text{ times}}$$

$$A^+ = A \cup A^2 \cup A^3 \cup \dots = \bigcup_{n=1}^{\infty} A^n \quad \text{the positive closure of } A$$

$$A^* = \{\lambda\} \cup A^+ = \bigcup_{n=0}^{\infty} A^n \quad \text{the Kleene closure or Kleene star}$$

- Let $A = \{aa, ab, ba, bb\}$. Then A^* is the language of all strings of even length.
- If $B = \{a, b\}$ then BA is the language of all strings of odd length
- What are $\{a\}\{ba\}^*$ and $\{ab\}^*\{b\}$?

Regular Expressions

- An **atomic language** is a language that contains only one string, and this string has length 1. $\{a\}$
For short we denote such a language simply by a
- Every language that contains only one string can be represented as a concatenation of atomic languages. $A = \{abba\} = abba$
(Careful!!! The $abba$ in the parenthesis is a string, while the $abba$ in the end is a concatenation of languages.)
- Any finite language is a union of concatenations of atomic languages. $A = \{ab,ba,abba\} = ab \cup ba \cup abba$
- An expression constructed from atomic languages by means of concatenation, union, intersection, complementation, and Kleene star is called a **regular expression**

Regular Expressions: Examples

- What the languages a^*ba^*b , $(a \cup b)^*c^*$, $((ab \cup ba)^*c)^*$ are?
- Write a regular expression for the language over $\{a,b,c\}$ that contains string with exactly one occurrence of c
- with exactly two occurrences of c
- over $\{a,b,c,d\}$ with one occurrence of c and one occurrence of d
- with as many occurrences of c as you wish, but each such occurrence should be followed by an occurrence of d

Homework

Exercises from the Book:

No. 1, 7, 11, 13, 15 (page 317 – 318)