

Bioinformatics Algorithms

(Fundamental Algorithms, module 2)

Zsuzsanna Lipták

Masters in Medical Bioinformatics
academic year 2018/19, II semester

Strings and Sequences in Computer Science

Some formalism on strings

- Σ a finite set called **alphabet**

Some formalism on strings

- Σ a finite set called **alphabet**
- its elements are called **characters** or **letters**

Some formalism on strings

- Σ a finite set called **alphabet**
- its elements are called **characters** or **letters**
- $|\Sigma|$ is the **size** of the alphabet (number of different characters)

Some formalism on strings

- Σ a finite set called **alphabet**
- its elements are called **characters** or **letters**
- $|\Sigma|$ is the **size** of the alphabet (number of different characters)
- a **string over Σ** is a finite sequence of characters from Σ

Some formalism on strings

- Σ a finite set called **alphabet**
- its elements are called **characters** or **letters**
- $|\Sigma|$ is the **size** of the alphabet (number of different characters)
- a **string over Σ** is a finite sequence of characters from Σ
- we write strings as $s = s_1 s_2 \dots s_n$ i.e. s_i is the i 'th character of s

Some formalism on strings

- Σ a finite set called **alphabet**
- its elements are called **characters** or **letters**
- $|\Sigma|$ is the **size** of the alphabet (number of different characters)
- a **string over Σ** is a finite sequence of characters from Σ
- we write strings as $s = s_1 s_2 \dots s_n$ i.e. s_i is the i 'th character of s

N.B.: We number strings from 1, not from 0

Some formalism on strings (cont.)

- $|s|$ is the **length** of string s

Some formalism on strings (cont.)

- $|s|$ is the **length** of string s
- ϵ is the **empty string**, the (unique) string of length 0

Some formalism on strings (cont.)

- $|s|$ is the **length** of string s
- ϵ is the **empty string**, the (unique) string of length 0
- Σ^n is the set of strings of length n

Some formalism on strings (cont.)

- $|s|$ is the **length** of string s
- ϵ is the **empty string**, the (unique) string of length 0
- Σ^n is the set of strings of length n
- $\Sigma^* = \bigcup_{n=0}^{\infty} \Sigma^n$

Some formalism on strings (cont.)

- $|s|$ is the **length** of string s
- ϵ is the **empty string**, the (unique) string of length 0
- Σ^n is the set of strings of length n
- $\Sigma^* = \bigcup_{n=0}^{\infty} \Sigma^n = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots$ is the set of all strings over Σ

Some formalism on strings: Examples

Examples

- **DNA:** $\Sigma = \{A, C, G, T\}$, alphabet size $|\Sigma| = 4$,
 $s = ACCTG$ is a string of length 5 of Σ , with
 $s_1 = A, s_2 = s_3 = C, s_4 = T, s_5 = G$.

Some formalism on strings: Examples

Examples

- **DNA:** $\Sigma = \{A, C, G, T\}$, alphabet size $|\Sigma| = 4$,
 $s = ACCTG$ is a string of length 5 of Σ , with
 $s_1 = A, s_2 = s_3 = C, s_4 = T, s_5 = G$.
- **RNA:** $\Sigma = \{A, C, G, U\}$, again alphabet size is 4

Some formalism on strings: Examples

Examples

- **DNA:** $\Sigma = \{A, C, G, T\}$, alphabet size $|\Sigma| = 4$,
 $s = ACCTG$ is a string of length 5 of Σ , with
 $s_1 = A, s_2 = s_3 = C, s_4 = T, s_5 = G$.
- **RNA:** $\Sigma = \{A, C, G, U\}$, again alphabet size is 4
- **protein:** $\Sigma = \{A, C, D, E, F, \dots, W, Y\}$, alphabet size is 20,
ANRFYWNL is a string over Σ of length 8

Some formalism on strings: Examples

Examples

- **DNA:** $\Sigma = \{A, C, G, T\}$, alphabet size $|\Sigma| = 4$,
 $s = ACCTG$ is a string of length 5 of Σ , with
 $s_1 = A, s_2 = s_3 = C, s_4 = T, s_5 = G$.
- **RNA:** $\Sigma = \{A, C, G, U\}$, again alphabet size is 4
- **protein:** $\Sigma = \{A, C, D, E, F, \dots, W, Y\}$, alphabet size is 20,
ANRFYWNL is a string over Σ of length 8
- **English alphabet:** $\Sigma = \{a, b, c, \dots, x, y, z\}$ of size 26,
alphabet is a string over Σ of length 8

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$ (i.e., a "contiguous piece" of s)

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...
- t is a **prefix** of s if $t = \epsilon$ or $t = s_1 \dots s_j$ for some $1 \leq j \leq n$
(i.e., a "beginning" of s)

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...
- t is a **prefix** of s if $t = \epsilon$ or $t = s_1 \dots s_j$ for some $1 \leq j \leq n$
(i.e., a "beginning" of s) AC, ACCTG, ...

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...
- t is a **prefix** of s if $t = \epsilon$ or $t = s_1 \dots s_j$ for some $1 \leq j \leq n$
(i.e., a "beginning" of s) AC, ACCTG, ...
- t is a **suffix** of s if $t = \epsilon$ or $t = s_i \dots s_n$ for some $1 \leq i \leq n$
(i.e., an "end" of s)

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...
- t is a **prefix** of s if $t = \epsilon$ or $t = s_1 \dots s_j$ for some $1 \leq j \leq n$
(i.e., a "beginning" of s) AC, ACCTG, ...
- t is a **suffix** of s if $t = \epsilon$ or $t = s_i \dots s_n$ for some $1 \leq i \leq n$
(i.e., an "end" of s) CCTG, G, ...

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...
- t is a **prefix** of s if $t = \epsilon$ or $t = s_1 \dots s_j$ for some $1 \leq j \leq n$
(i.e., a "beginning" of s) AC, ACCTG, ...
- t is a **suffix** of s if $t = \epsilon$ or $t = s_i \dots s_n$ for some $1 \leq i \leq n$
(i.e., an "end" of s) CCTG, G, ...
- t is a **subsequence** of s if t can be obtained from s by deleting some (possibly 0, possibly all) characters from s

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...
- t is a **prefix** of s if $t = \epsilon$ or $t = s_1 \dots s_j$ for some $1 \leq j \leq n$
(i.e., a "beginning" of s) AC, ACCTG, ...
- t is a **suffix** of s if $t = \epsilon$ or $t = s_i \dots s_n$ for some $1 \leq i \leq n$
(i.e., an "end" of s) CCTG, G, ...
- t is a **subsequence** of s if t can be obtained from s by deleting some
(possibly 0, possibly all) characters from s AT, CCT, ...

Some formalism on strings

Let $s = s_1 \dots s_n$ be a string over Σ .

ex. $s = \text{ACCTG}$

- t is a **substring** of s if $t = \epsilon$ or $t = s_i \dots s_j$ for some $1 \leq i \leq j \leq n$
(i.e., a "contiguous piece" of s) CCT, AC, ...
- t is a **prefix** of s if $t = \epsilon$ or $t = s_1 \dots s_j$ for some $1 \leq j \leq n$
(i.e., a "beginning" of s) AC, ACCTG, ...
- t is a **suffix** of s if $t = \epsilon$ or $t = s_i \dots s_n$ for some $1 \leq i \leq n$
(i.e., an "end" of s) CCTG, G, ...
- t is a **subsequence** of s if t can be obtained from s by deleting some
(possibly 0, possibly all) characters from s AT, CCT, ...

N.B.

string = sequence, but substring \neq subsequence!

Substrings etc.

N.B.

1. Every substring is a subsequence, but not every subsequence is a substring!

Substrings etc.

N.B.

1. Every substring is a subsequence, but not every subsequence is a substring!

Ex.: Let $s = \text{ACCTG}$, then ACT is a subsequence but not a substring.

Substrings etc.

N.B.

1. Every substring is a subsequence, but not every subsequence is a substring!
Ex.: Let $s = \text{ACCTG}$, then ACT is a subsequence but not a substring.
2. Every prefix and every suffix is a substring.

Substrings etc.

N.B.

1. Every substring is a subsequence, but not every subsequence is a substring!
Ex.: Let $s = \text{ACCTG}$, then ACT is a subsequence but not a substring.
2. Every prefix and every suffix is a substring.
3. t is substring of $s \iff t$ is prefix of a suffix of $s \iff t$ is suffix of a prefix of s

Counting substrings, subsequences etc.

Question

Given $s = s_1 \dots s_n$. How many

- prefixes,
- suffixes,
- substrings,
- subsequences

does s have (exactly, or at most, or at least)?