

CSCI-UA 480.4: APS

Algorithmic Problem Solving

Some String Algorithms

Instructor: Joanna Klukowska

created based on materials for this class by
Bowen Yu and materials shared by the authors of
the textbook Steven and Felix Halim

String Definitions

a **string** of length n consists of characters s_0, s_1, \dots, s_{n-1}

a **substring** is a sequence of consecutive characters in a string that starts at position i and ends at position j (inclusive on both ends)

a **prefix** is a substring for which $i = 0$

a **suffix** is a substring for which $j = n - 1$

a **subsequence** is any sequence of characters in a string in their original order (not necessarily consecutive)

Longest Common Subsequence

The **longest common subsequence** (lcs) of two strings is the longest string that appears as a subsequence in both strings.

Examples:

- "floor" and "door", the lcs is "oor"
- "caged" and "rage", the lcs is "age"
- "capsule" and "recaps", the lcs is "caps"

Longest Common Subsequence

The **longest common subsequence** (lcs) of two strings is the longest string that appears as a subsequence in both strings.

Examples:

- "floor" and "door", the lcs is "oor"
- "caged" and "rage", the lcs is "age"
- "capsule" and "recaps", the lcs is "caps"

Solution

given: two strings and

- function that returns length of the longest common subsequence of the prefixes
and

Longest Common Subsequence

The **longest common subsequence** (lcs) of two strings is the longest string that appears as a subsequence in both strings.

Examples:

- "floor" and "door", the lcs is "oor"
- "caged" and "rage", the lcs is "age"
- "capsule" and "recaps", the lcs is "caps"

Solution

given: two strings and

- function that returns length of the longest common subsequence of the prefixes
and

[Visualization of the algorithm](#)

Edit Distance

The **edit distance** between two strings is defined as the minimum number of editing operations that transform one string into the other.

The allowed operations may vary, but are often

- insert a character, "ABC" -> "ABCA"
- remove a character, "ABC" -> "AC"
- replace a character, "ABC" -> "ADC"
(this one can be thought of as two separate operations of remove followed by insert)

Edit Distance

The **edit distance** between two strings is defined as the minimum number of editing operations that transform one string into the other.

The allowed operations may vary, but are often

- insert a character, "ABC" -> "ABCA"
- remove a character, "ABC" -> "AC"
- replace a character, "ABC" -> "ADC"

(this one can be thought of as two separate operations of remove followed by insert)

Solution

given: two strings `s` and `t`

- function that returns the edit distance between the prefixes `s[0:i]` and `t[0:j]`

Edit Distance

The **edit distance** between two strings is defined as the minimum number of editing operations that transform one string into the other.

The allowed operations may vary, but are often

- insert a character, "ABC" -> "ABCA"
 - remove a character, "ABC" -> "AC"
 - replace a character, "ABC" -> "ADC"
- (this one can be thought of as two separate operations of remove followed by insert)

Solution

given: two strings s and t

- function that returns the edit distance between the prefixes $s[0..i]$ and $t[0..j]$

where $d[i][j]$ is the edit distance between $s[0..i]$ and $t[0..j]$ when $i > 0$ and $j > 0$, otherwise.

Edit Distance

The **edit distance** between two strings is defined as the minimum number of editing operations that transform one string into the other.

The allowed operations may vary, but are often

- insert a character, "ABC" -> "ABCA"
- remove a character, "ABC" -> "AC"
- replace a character, "ABC" -> "ADC"

(this one can be thought of as two separate operations of remove followed by insert)

Solution

given: two strings s and t

- function that returns the edit distance between the prefixes $s[0..i]$ and $t[0..j]$

where $d[i][j]$ is the edit distance between $s[0..i]$ and $t[0..j]$ when $i > 0$ and $j > 0$, otherwise.

Z-Array / Z-Algorithm

What is a Z-array?

Z-Array / Z-Algorithm

What is a Z-array?

The Z-array for a string of length is an array of length in which stores the length of the longest substring starting at that is also a prefix of .

Example:

Z-Array / Z-Algorithm

What is a Z-array?

The Z-array for a string of length is an array of length in which stores the length of the longest substring starting at that is also a prefix of .

Example:

What is a Z-Algorithm?

Z-Array / Z-Algorithm

What is a Z-array?

The Z-array for a string s of length n is an array of length n in which $Z[i]$ stores the length of the longest substring starting at $s[i]$ that is also a prefix of s .

Example:

What is a Z-Algorithm?

It's an algorithm that computes the Z-Array for a given string s .

Z-Algorithm

Idea: maintain a range $[l, r]$ such that l is a prefix of s , the value of r has been calculated and r is as large as possible.

Observation: l is the same as

Challenge: Pattern Matching

Find all locations of a pattern string in a given string .

Challenge: Pattern Matching

Find all locations of a pattern string `pattern` in a given string `text`.

Solution

- create a new string `text + pattern + '$'` in which `'$'` is a special character that does not occur in neither `text` nor `pattern`
- create the z-array for the new string
- the locations in the z-array for which the value is equal to the length of the pattern string `pattern` are the location of the pattern in `text` (adjust indexes by subtracting the length of `pattern`)

Challenge: Pattern Matching

Find all locations of a pattern string `pattern` in a given string `text`.

Solution

- create a new string `text + pattern + pattern` in which `pattern` is a special character that does not occur in neither `text` nor `pattern`
- create the z-array for the new string
- the locations in the z-array for which the value is equal to the length of the pattern string `pattern` are the location of the pattern in `text` (adjust indexes by subtracting the length of `pattern`)

[Visualization of the solution](#)

Challenge: Finding Borders

A **border** in a string is a substring that is both a prefix and a suffix of that string (but not the entire string, i.e., proper prefix and proper suffix).

Example:

the borders are

Challenge: Finding Borders

Solution

- create the z-array for
- borders are all suffixes , such that

Challenge: Finding Borders

Solution

- create the z-array for
- borders are all suffixes , such that

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

Knuth-Morris-Pratt's (KMP) Algorithm

Another algorithm for pattern matching in strings.

It uses the observation that when a mismatch occurs, the pattern itself contains sufficient information to determine where the next match could begin. This allows to skip re-examination of previously matched characters.

Example of the idea:

[Visualization of KMP](#)

[another visualization of KMP](#)

Source code implementation: [cpp](#), [java](#)

