

MA/CSSE 473

Day 29

Boyer-Moore



Boyer-Moore Recap 1

- Based on two main ideas:
- compare pattern characters to text characters from right to left
- precompute the shift amounts in **two** tables
 - **bad-symbol table** indicates how much to shift based on the text's character that causes a mismatch
 - **good-suffix table** indicates how much to shift based on matched part (suffix) of the pattern



Boyer-Moore Recap 2

n	length of text
m	length of pattern
i	position in text that we are trying to match with rightmost pattern character
k	number of characters (from the right) successfully matched before a mismatch

After successfully matching $0 \leq k < m$ characters, the algorithm shifts the pattern right by

$$d = \max \{d_1, d_2\}$$

where $d_1 = \max\{t_1[c] - k, 1\}$ is the *bad-symbol* shift
($t_1[c]$ is from Horspool table)

$d_2[k]$ is the *good-suffix* shift

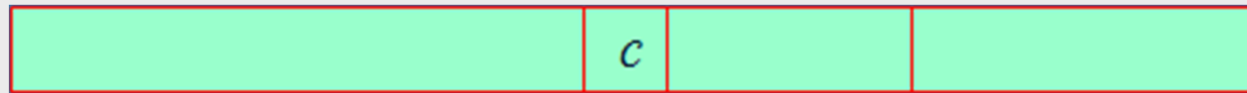
(next we explore how to compute it)



Bad-symbol shift in Boyer-Moore

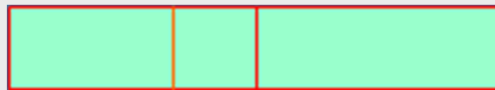
- If the rightmost character of the pattern does not match, Boyer-Moore algorithm acts much like Horspool's
- If the rightmost character of the pattern does match, BM compares preceding characters right to left until either
 - all pattern's characters match, or
 - a mismatch on text's character c is encountered after $k > 0$ matches

text



\neq k matches

pattern



bad-symbol shift: How much should we shift by?

$d_1 = \max\{t_1(c) - k, 1\}$,
where $t_1(c)$ is the value from the Horspool shift table.



Good-suffix Shift in Boyer-Moore

- Good-suffix shift d_2 is applied after the k last characters of the pattern are successfully matched
 - $0 < k < m$
- How can we take advantage of this?
- As in the bad suffix table, we want to pre-compute some information based on the characters in the suffix.
- We create a **good suffix table** whose indices are $k = 1 \dots m-1$, and whose values are how far we can shift after matching a k -character suffix (from the right).
- Example patterns: CABABA AWOWWOW
 WOWWOW ABRACADABRA



Boyer-Moore Example

pattern = abracadabra

text =

abracadabtabradabracadabcbcadaxbrabbbacabadabraxxxxxxabracadabracadabra

m = 11, n = 67

badCharacterTable: a3 b2 r1 a3 c6 x11

GoodSuffixTable: (1,3) (2,10) (3,10) (4,7) (5,7) (6,7) (7,7) (8,7)
(9,7) (10,7)

abracadabtabradabracadabcbcadaxbrabbbacabadabraxxxxxxabracadabracadabra
abracadabra

i = 10 k = 1 t1 = 11 d1 = 10 d2 = 3

abracadabtabradabracadabcbcadaxbrabbbacabadabraxxxxxxabracadabracadabra
abracadabra

i = 20 k = 1 t1 = 6 d1 = 5 d2 = 3

abracadabtabradabracadabcbcadaxbrabbbacabadabraxxxxxxabracadabracadabra
abracadabra

i = 25 k = 1 t1 = 6 d1 = 5 d2 = 3

abracadabtabradabracadabcbcadaxbrabbbacabadabraxxxxxxabracadabracadabra
abracadabra

i = 30 k = 0 t1 = 1 d1 = 1



Boyer-Moore Example

First step is a repeat from the previous slide

abracadabtabradabracadabcbadaxbrabbracadabraxxxxxxabracadabracadabra
abracadabra

i = 30 k = 0 t1 = 1 d1 = 1

abracadabtabradabracadabcbadaxbrabbracadabraxxxxxxabracadabracadabra
abracadabra

i = 31 k = 3 t1 = 11 d1 = 8 d2 = 10

abracadabtabradabracadabcbadaxbrabbracadabraxxxxxxabracadabracadabra
abracadabra

i = 41 k = 0 t1 = 1 d1 = 1

abracadabtabradabracadabcbadaxbrabbracadabraxxxxxxabracadabracadabra
abracadabra

i = 42 k = 10 t1 = 2 d1 = 1 d2 = 7

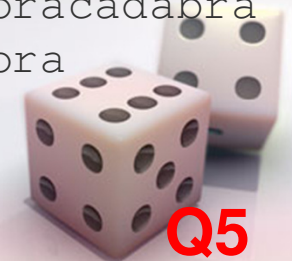
abracadabtabradabracadabcbadaxbrabbracadabraxxxxxxabracadabracadabra
abracadabra

i = 49 k = 1 t1 = 11 d1 = 10 d2 = 3

abracadabtabradabracadabcbadaxbrabbracadabraxxxxxxabracadabracadabra
abracadabra

49

Brute force took 50 times through the outer loop; Horspool took 13; Boyer-Moore 9 times.



Q5

Boyer-Moore Example

- On Moore's home page
- <http://www.cs.utexas.edu/users/moore/best-ideas/string-searching/fstrpos-example.html>

