

CSE 417: Algorithms and Computational Complexity

Winter 2001
Lecture 15
Instructor: Paul Beame

1

Pattern Matching

- Given
 - a string, s , of n characters
 - a pattern, p , of m characters
 - usually $m \ll n$
- Find
 - all occurrences of the pattern p in the string s
- Obvious algorithm:
 - try to see if p matches at each of the positions in s , stopping at a failed match

2

String $s = xyxyxyxyxyxyxyxyxyxyxy$
Pattern $p = xyxyxyxyxx$

3

String $s = xyxyxyxyxyxyxyxyxyxyxy$
 $xyxyxyxyxx$

4

String $s = xyxyxyxyxyxyxyxyxyxyxy$
 $xyxy$
 $xyxyxyxyxx$

5

String $s = xyxyxyxyxyxyxyxyxyxyxy$
 $xyxy$
 x
 $xyxyxyxyxx$

6

```
String s = xyxxyxyxyxyxyxyxyxyxyxyxyx
  xyxy
  x
  xy
  xyxyxyxyxyxx
```

7

```
String s = xyxxyxyxyxyxyxyxyxyxyxyxyx
  xyxy
  x
  xy
  xyxyy
  xyxyxyxyxyxx
```

8

```
String s = xyxxyxyxyxyxyxyxyxyxyxyxyx
  xyxy
  x
  xy
  xyxyy
  x
  xyxyxyxyxyxx
```

9

```
String s = xyxxyxyxyxyxyxyxyxyxyxyxyx
  xyxy
  x
  xy
  xyxyy
  x
  xyxyxyxyxyxx
  xyxyxyxyxyxx
```

10

```
String s = xyxxyxyxyxyxyxyxyxyxyxyxyx
  xyxy
  x
  xy
  xyxyy
  x
  xyxyxyxyxyxx
  x
  xyxyxyxyxyxx
```

11

```
String s = xyxxyxyxyxyxyxyxyxyxyxyxyx
  xyxy
  x
  xy
  xyxyy
  x
  xyxyxyxyxyxx
  x
  xyx
  xyxyxyxyxyxx
```

12

String $s = xyxxyxyxyxyxyxyxyxyxyxyx$

13

String $s = xyxxyxyxyxyxyxyxyxyxyxyx$

14

String $s = xyxxyxyxyxyxyxyxyxyxyxyx$

15

String $s = xyxxyxyxyxyxyxyxyxyxyxyx$

Worst-case time $O(mn)$

16

String $s = xyxxyxyxyxyxyxyxyxyxyxyx$

Suppose we knew the pattern well

Lots of wasted work

Since we know earlier agreement of the string with the pattern, these can't possibly match

17

Preprocess the Pattern

- After each character in the pattern figure out ahead of time what the next useful work would be if it failed to match there
- i.e. how much can one shift over the pattern for the next match

18

String $s = xyxyxyxyxyxyxyxyxyxyxy$
 $xyxy$

$xyxyy$

$xyxyxyxyxyxx$

$xyxyxyxyxyxx$

19

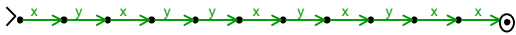
Preprocessing the pattern

- At each mismatch
 - Look at the last part that matched plus extra mismatched character
 - Try to fit pattern as far to the left in this as possible
 - i.e. look for the longest prefix of the pattern that matches the end of the sequence so far.

20

Preprocessing the pattern

Pattern $p = xyxyxyxyxyxx$

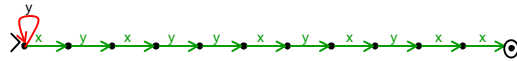


Each dot represents how far in the pattern things are matched

21

Preprocessing the pattern

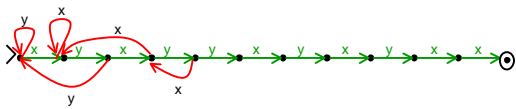
Pattern $p = xyxyxyxyxyxx$



22

Preprocessing the pattern

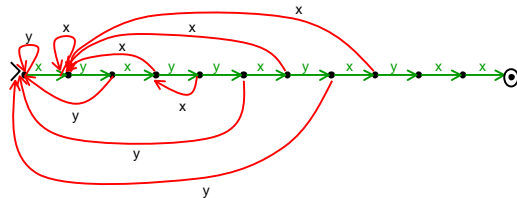
Pattern $p = xyxyxyxyxyxx$



23

Preprocessing the pattern

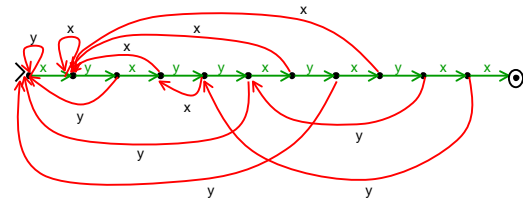
Pattern $p = xyxyxyxyxyxx$



24

Preprocessing the pattern

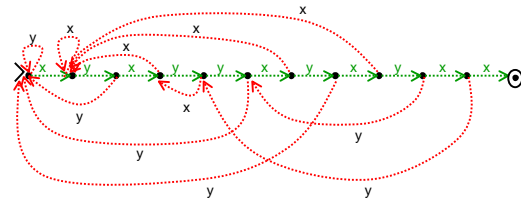
Pattern $p = xyxyxyxyxx$



25

Running on the string

String $s = xyxyxyxyxyxyxyxyxyxyxy$



26

Knuth-Morris-Pratt Algorithm

- Once the preprocessing is done there are only n steps on any string of size n
 - just follow your nose
- Obvious algorithm for doing preprocessing is $O(m^2)$ steps
 - still usually good since $m \ll n$
- Knuth-Morris-Pratt Algorithm can do the pre-processing in $O(m)$ steps
 - Total $O(m+n)$ time

27

Finite State Machines

- The diagram we built is a special case of a **finite automaton**
 - start state
 - goal or accepting state(s)
 - an arc out of each state labeled by each of the possible characters
- Finite automata take strings of characters as input and decide what to do with them based on the paths they follow

28

Finite State Machines

- Many communication protocols, cache-coherency protocols, VLSI circuits, user-interfaces, even adventure games are designed by making finite state machines first.
 - The "strings" that are the input to the machines can be
 - a sequence of actions of the user
 - the bits that arrive on particular ports of the chip
 - a series of values on a bus

29

Finite State Machines

- Can search for arbitrary combinations of patterns not just a single pattern
 - Given two finite automata can build a single new one that accepts strings that either of the original ones accepted
- Typical text searches are based on finite automata designs
 - Perl builds this in as a first-class component of the programming language.

30

Next time

- We start the computability and complexity portion of the course.
- We discuss **Turing machines** which are similar in style to finite-state machines but much more powerful
 - as powerful as any programming language

31