

Exact Matching, Part 1

photophorescent



Steven Bedrick

CS/EE 5/655, 11/12/14

Plan for today:

String-matching problems

Exact matching: naïve algorithm

Exact matching: linear-time algorithm

The exact string-matching problem:

Given a short pattern P ...

... and a longer text T ...

... identify the locations in T where instances of P occur.

Often, we have a large number of patterns...

... and a large amount of text...

... and we want to find instances efficiently.

Worst-case: $O(nm)$

Ideally: sub-linear

Approximate matching:

Like exact matching, but allowing “small differences”:

Functionally similar amino acid substitutions/indels?

Vowel shifts, systematic consonant dropping, morphological variants, etc.?

Review: distance metrics:

1. Non-negative property:

$$d(\mathbf{a}, \mathbf{b}) \geq 0 \quad \text{for all } \mathbf{a} \text{ and } \mathbf{b}$$

2. Zero property:

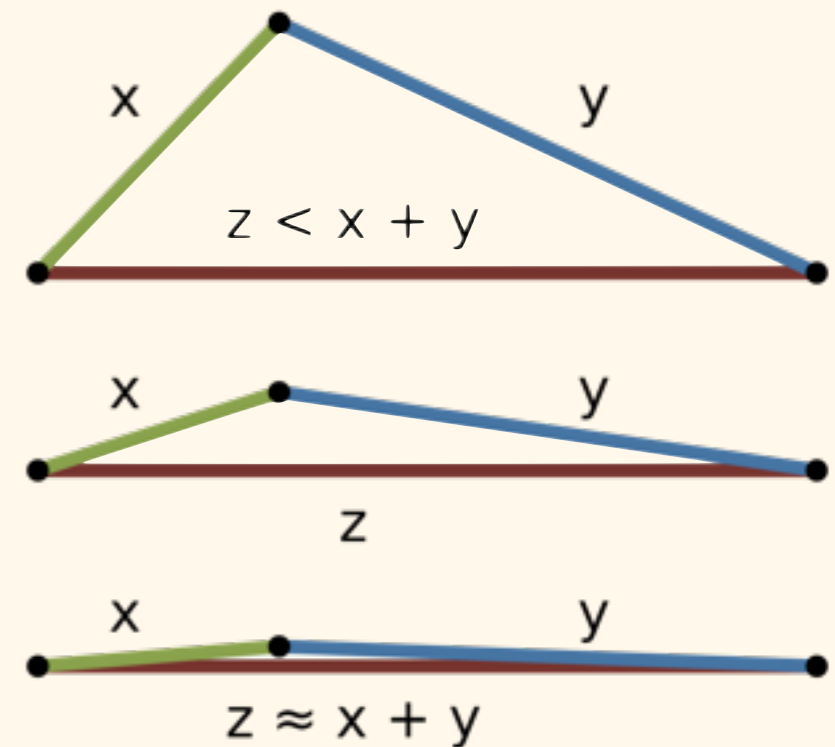
$$d(\mathbf{a}, \mathbf{b}) = 0 \quad \text{iff } \mathbf{a} = \mathbf{b}$$

3. Symmetry:

$$d(\mathbf{a}, \mathbf{b}) = d(\mathbf{b}, \mathbf{a}) \quad \text{for all } \mathbf{a} \text{ and } \mathbf{b}$$

4. Triangle inequality:

$$d(\mathbf{a}, \mathbf{b}) + d(\mathbf{b}, \mathbf{c}) \geq d(\mathbf{a}, \mathbf{c}) \quad \text{for all } \mathbf{a}, \mathbf{b} \text{ and } \mathbf{c}$$



Approximate matching:

Precise formulation of “distance” will is heavily application-dependent...

... but the algorithms don't care!

Today, though, we're talking about exact.

$P = \text{"xtpxtd"}$

$T = \text{"xluxtpxtdqwtdxtpxtsyxtpxtdy"}$

Occurrences of P in T :

xluxtpxtdqwtdxtpxtsyxtpxtdy

Demo

Performance:

If P has length m , and T length n , worst-case is $O(mn)$.

- suppose P is *aaaaab*,
and T is *aaaaaaaaaaaaaaaaaaaa*

Don't use the naïve algorithm!

Performance:

How might we improve?

$P = \text{"xtpxtd"}$

$T = \text{"xluxtpxtdqwt dxtpxtsyxtpxtdy"}$

Does our pattern have any internal structure?

Demo

Improvements to exact match almost always involve pre-processing either P or T .

The goal: save time (comparisons) by identifying repetitive elements.

The choice of which to focus on is application dependent:

Mapping new DNA sample? Probably transform P , since you'll be working with sequence-tagged sites.

Of course, there are probably errors in your STS library, so...

The choice of which to focus on is application dependent:

Searching for arbitrary words in a large text? Focus on transforming the text.

Today's algorithm will focus on patterns.

Definitions: for a string S :

$S[i,j]$ = contiguous substring starting at i and ending at j . $S(i) = S[i,i]$

$S = \text{aardvark}$

$S[2,4] = \text{ard}$ $S(4) = \text{d}$

For $i > 1$, $Z_i(S)$ is the length of the longest prefix of $S[i,|S|]$ that is also a prefix of S .

For $i > 1$, $Z_i(S)$ is the length of the longest prefix of $S[i, |S|]$ that is also a prefix of S .

$$S = \text{xtpxtd}$$

$$S[4, |S|] = \text{xttd}$$

$$Z_4(S) = 2$$

For $i > 1$, $Z_i(S)$ is the length of the longest prefix of $S[i, |S|]$ that is also a prefix of S .

$S = \underline{x}tpxtd$

$S[4, |S|] = \underline{x}td$

$Z_4(S) = 2$

In this string, for all other positions k ,
 $Z_k(P) = 0$.

$P = \text{aardvark}$

$$Z_2(P) = 1$$

$P = \underline{a}ardvark$



$$Z_2(P) = 1$$

$P = \underline{a}ardvark$



$$Z_2(P) = 1$$

$$Z_6(P) = 1$$

$P = \text{alfalfa}$



$$Z_4(P) = 4$$

P = aardvark



$$Z_2(P) = 1$$

P = alfalfa



$$Z_4(P) = 4$$

$$Z_6(P) = 1$$

P = photophosphorescent



$$Z_6(P) = Z_{10}(P) = 3$$

P = aardvark



$$Z_2(P) = 1$$

P = alfalfa



$$Z_4(P) = 4$$

$$Z_6(P) = 1$$

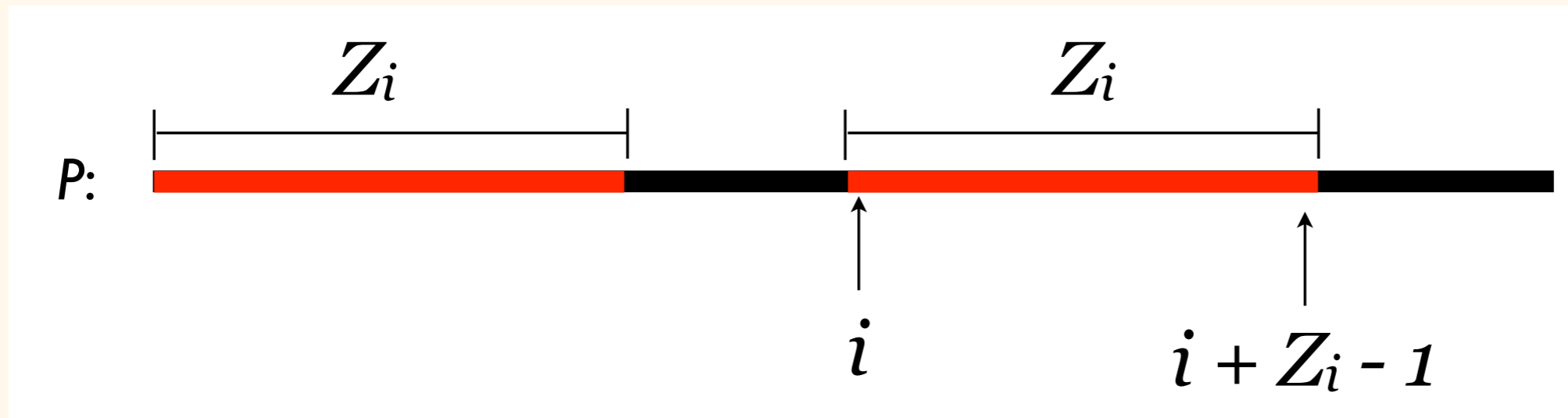
P = photophosphorescent



$$Z_6(P) = Z_{10}(P) = 3$$

These regions of prefix-overlap are called *z-boxes*.

P = photophosphorescent



Why do we care?

Let's stick P and T together with a sentinel character:

$P\$T = \text{pho\$photophosphorescent}$

p	h	o	\$	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
0	0	0	0	3	0	0	0	0	3	0	0	0	3	0	0	0	0	0	0	0	0	0

Any point whose $Z_i = |P|$ matches!

How to calculate Z_i ?

Naïve way: $O(nm)$

There exists a linear-time solution!

Next steps:

Knuth-Morris-Pratt:

A linear-time extension to what we've seen, with some clever preprocessing to allow larger shifts.

Boyer-Moore:

A totally different approach, moving right-to-left and often achieving sub-linear performance.