

## 2 - Approximate Matching

Sunday, 10 October 2021 17:06

BOYER-MOORE VERY POPULAR ALGORITHM

USED FOR QUERIES, WEB & DB SEARCH

MAIN IDEA: CAN WE LEARN FROM PAST COMPARISONS?

WE ARE GOING TO LEARN FROM OUR CHARACTER-COMPARISON AS MUCH AS WE CAN TO SKIP ALIGNMENTS THAT WILL NOT RESULT IN A MATCH

ALIGNMENTS IN LEFT-TO-RIGHT CHAR COMPARISON IN OPPOSITE ORDER (RIGHT-TO-LEFT)

IF WE SHIFT P BY  $\alpha$  WE DON'T CHECK  $\alpha-1$  ALIGNMENTS

THE BAD CHARACTER RULE

UPON MISMATCH SKIP ALIGNMENTS UNTIL:

- MISMATCH BECOMES A MATCH
- P MOVES PAST MISMATCHED CHARACTER

GOOD SUFFIX RULE

$E$  = SUBSTRING MATCHED BY INNER LOOP, SKIP UNTIL:

- THERE ARE NO MISMATCHES BETWEEN P AND E
- P MOVES PAST T

BOYER-MOORE

USE BAD-CHAR OR GOOD-SUFFIX WHICH EVER SKIPS MORE

PREPROCESSING USING TABLES

REAL GENOMES ARE NEARLY DIFFERENT FROM RANDOM ONES

REPETITIVE ELEMENTS?

EVOLUTIONARY PROCESS INTRODUCES PATTERNS

HUMAN GENOME IS EXTREMELY REPETITIVE

GENOMES GETS INVADED BY TRANSPOSONS

45% OF HUMAN GENOME IS COVERED BY TRANSPOSSABLE ELEMENTS

Alu, L1, SVA & OTHER

WHY SHOULD WE CARE?

THEY CREATE AMBIGUITY

LIKE A PUZZLE WHERE HALF OF THE PIECES ARE FEATURELESS BLUE SKY

WHERE DOES THE PIECE COME FROM?

REPETITIVE ELEMENTS CREATE PROBLEMS FOR OUR ALGORITHMS

BOYER-MOORE PRE-PROCESS PATTERN P IN ORDER TO BUILD LOOK-UP TABLES THAT HELP TO USE THE GOOD-SUFFIX & BAD-CHARACTER RULES

PREPROCESSING

IF WE WANT TO MATCH THE SAME PATTERN P AGAINST DIFFERENT TEXTS T, I CAN REUSE MY PREPROCESSING FOR THE PATTERN

THE COST OF PREPROCESSING CAN BE AMORTIZED OVER MANY PROBLEMS

OFFLINE VS ONLINE

PREPROCESSING THE TEXT T (ONLINE)

BOTH NAIIVE & BOYER-MOORE ARE 'ONLINE' ALGORITHMS WHEREAS WEB SEARCH ENGINE IS OFFLINE

WHEN WE SEARCH ON THE WEB, THE ENGINE ALREADY PREPROCESSED LOADS OF INFO. THIS MAKES QUERIES FASTER

READ ALIGNMENTS? WE CAN USE AN OFFLINE ONE IT WOULD BE HELPFUL FOR A REFERENCE GENOME

PREPROCESSING TO MATCH MANY READS QUICKLY

WHY REF-BASE ASSEMBLY IS OFFLINE?

REFERENCE GENOMES DO NOT CHANGE BETWEEN EXPERIMENTS, IT MAY BE UPDATED ONCE A YEAR

- INDEXING (LIKE IN COOKS)

LIST OF KEY-TERM & PLACES ON THE BOOK WHERE IT IS MENTIONED

- GROUPING (LIKE IN SUPERMARKETS)

HOW CAN WE ORGANIZE QUERIES IN ORDER TO MAKE THEM EASY TO RETRIEVE?

HOW CAN WE INDEX DNA?

TAKE EVERY SUBSTRING OF TEXT T OF A CERTAIN LENGTH, TREAT EACH OF THOSE SUBSTRINGS AS A WORD

K-MER SUBSTRING OF LENGTH K

WE CAN SOLVE THE EXACT-MATCHING PROBLEM

THESE K CHAR OF P (PATTERN) MATCH K CHAR OF T (TEXT) AT INDEX 'i'

FINDING THE INDEX-HITS

MULTI-MAP

MAPS EACH K-MER TO LIST OF ALL OCCURRENCES IN TEXT

(K-MER → OFFSET IN GENOME)

BINARY SEARCH  $O(\log n)$

FASTER QUERY SINCE MAP ORDERED ALPHABETICALLY

HASH-TABLE

USEFUL TO IMPLEMENT MULTI-MAPS (OFTEN AS LINKED LISTS)

PYTHON'S DICTIONARY IS AN IMPLEMENTATION OF A HASH TABLE

SUBSEQUENCE ≠ SUBSTRING (NOT CONTIGUOUS)

USING SUBSEQUENCE TENDS TO INCREASE THE SPECIFICITY OF THE INDEX

SUBSEQUENCES ARE SUBSEQUENCES BUT NOT THE OTHER WAY AROUND

GENOME INDEXES USED IN RESEARCH

INDEXING IS AT THE CORE IN A LOT OF METHODS TO ANALYZE READING-DATA

HOW CAN WE FIT GENOMES INTO A COMPACT INDEX?

VARIATIONS (DIFFERENT IMPLEMENTATIONS)

EXTRACT ONLY SOME OF THE K-MER (BOTH ODD & EVEN ODDS)

SUBSEQUENCE-INDEX (MORE SPECIFIC)

SUFFIX-INDEX (TREE, ANNAY, FM)

FM-INDEX (BASED ON BURROWS-WHEELER) ORIGINALLY DEVELOPED FOR COMPRESSION

BIO SEQUENCES ANALYSIS

APPROXIMATE MATCHING PROBLEM

ACCOUNTING FOR:

- SEQUENCING ERRORS
- NATURAL VARIANCE

DISTANCE

HOW 2 STRINGS DIFFER FROM EACH OTHER

- HAMMING DISTANCE

MINIMUM # OF SUBSTITUTIONS TO TURN ONE STRING INTO ANOTHER

$|x|=|y|$  (SAME LENGTH)

- EDIT DISTANCE

MINIMUM # OF EDITS TO TURN ONE STRING INTO ANOTHER

- SUBSTITUTION
- INSERTION
- DELETION

$|x|, |y|$  ANY LENGTH

PIGEONHOLE PRINCIPLE

IF WE HAVE 10 HOLES BUT 11 PIGEONS ONE HOLE WILL HAVE 2 PIGEONS.

THERE MUST EXIST SUCH A HOLE

THE EDIT-DISTANCE PROBLEM

DYNAMIC PROGRAMMING

METHOD TO SOLVE BIO-SEQ ANALYSIS PROBLEMS NEW FAMILY OF TOOLS

H-DIST EASY

E-DIST HARD

DISTANCE RELATIONSHIP

$$E(x,y) \leq H(x,y)$$

INSERTIONS & DELETION ARE VERY POWERFUL, WE CAN GET FROM X TO Y IN FEWER CHANGES

$$E(x,y) \geq ||x|-|y||$$

WE AVOID REPETITIONS IN DYNAMIC PROGRAMMING

Boyer-Moore: Good suffix rule

Like with the bad character rule, the number of skips possible using the good suffix rule can be precalculated into a few tables (Gusfield 2.24 and 2.25)

Rule on previous slide is the weak good suffix rule; there is also a strong good suffix rule (Gusfield 2.23)

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

T: CTTGCTTACCTTACT

P: CTTA**C**TAC

Weak ~~CTTACTTAC~~

Strong ~~CTTACTTAC~~

mismatch!

With the strong good suffix rule (and other minor modifications), Boyer-Moore is  $O(m)$  worst-case time. Gusfield discusses proof.

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

Boyer-Moore: Bad character rule

Upon mismatch, skip alignments until (a) mismatch becomes a match, or (b) P moves past mismatched character

Step 1: T: GCTT**G**CTACCTTACTTACT

P: CTTA**C**TAC

Step 2: T: GCTT**G**CTAC**C**CTTTGCGCGCGCGGGAA

P: T**C**TAC

With the strong good suffix rule (and other minor modifications), Boyer-Moore is  $O(m)$  worst-case time. Gusfield discusses proof.

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

Boyer-Moore: Bad character rule

Upon mismatch, skip alignments until (a) mismatch becomes a match, or (b) P moves past mismatched character

Step 1: T: GCTT**G**CTACCTTACTTACT

P: CTTA**C**TAC

Step 2: T: GCTT**G**CTAC**C**CTTTGCGCGCGCGGGAA

P: T**C**TAC

With the strong good suffix rule (and other minor modifications), Boyer-Moore is  $O(m)$  worst-case time. Gusfield discusses proof.

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

Boyer-Moore: Bad character rule

Upon mismatch, skip alignments until (a) mismatch becomes a match, or (b) P moves past mismatched character

Step 1: T: GCTT**G**CTACCTTACTTACT

P: CTTA**C**TAC

Step 2: T: GCTT**G**CTAC**C**CTTTGCGCGCGCGGGAA

P: T**C**TAC

With the strong good suffix rule (and other minor modifications), Boyer-Moore is  $O(m)$  worst-case time. Gusfield discusses proof.

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

Boyer-Moore: Bad character rule

Upon mismatch, skip alignments until (a) mismatch becomes a match, or (b) P moves past mismatched character

Step 1: T: GCTT**G**CTACCTTACTTACT

P: CTTA**C**TAC

Step 2: T: GCTT**G**CTAC**C**CTTTGCGCGCGCGGGAA

P: T**C**TAC

With the strong good suffix rule (and other minor modifications), Boyer-Moore is  $O(m)$  worst-case time. Gusfield discusses proof.

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

Boyer-Moore: Bad character rule

Upon mismatch, skip alignments until (a) mismatch becomes a match, or (b) P moves past mismatched character

Step 1: T: GCTT**G**CTACCTTACTTACT

P: CTTA**C**TAC

Step 2: T: GCTT**G**CTAC**C**CTTTGCGCGCGCGGGAA

P: T**C**TAC

With the strong good suffix rule (and other minor modifications), Boyer-Moore is  $O(m)$  worst-case time. Gusfield discusses proof.

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

Boyer-Moore: Bad character rule

Upon mismatch, skip alignments until (a) mismatch becomes a match, or (b) P moves past mismatched character

Step 1: T: GCTT**G**CTACCTTACTTACT

P: CTTA**C**TAC

Step 2: T: GCTT**G**CTAC**C**CTTTGCGCGCGCGGGAA

P: T**C**TAC

With the strong good suffix rule (and other minor modifications), Boyer-Moore is  $O(m)$  worst-case time. Gusfield discusses proof.

JOHNS HOPKINS

WRITING INCS OF ENGINEERING

Boyer-Moore: Bad character