Data Structures and Algorithms

String Algorithms

InverseHackermann

Definitions

Character: one letter of the alphabet. Often appears in single quotes.

String: sequence of characters. Often appears in double quotes. Length denoted as |S|

For s="onedragon", S[0]='o', S[1]='n', S[2]='e', etc.

Substring: A range of characters in a string. S[3 .. 6]="drag".

Prefix: substring starting at the beginning of the string.

Suffix: substring ending at the end of the string.

Border: both a prefix and a suffix of the string. "on" is a border of "onedragon".

Proper: not the entire string. A string is always a border of itself, but not a proper border.

String matching

Goal: look for occurrences of pattern P within S.

Directly comparing at each index can take |S|*|P| comparisons.

Index:	0	1	2	3	4	5	6	7	8	9	10
S:	R	А	W	R	Α	R	Α	R	W	R	Α
P:	R	Α	R	W							
No match:				R	Α	R	W				
Match:						R	Α	R	W		

Knuth-Morris-Pratt (KMP) string matching

Solution: Combine string into T = P#S('#') is a unique character).

Store LPB[i] = length of longest proper border for each prefix T[0 .. i].

To compute LPB[i], extend a border of T[0 .. i] by one character.

LPB[i - 1] is longest proper border of T[0 .. i]; how to get next shorter border?

Index:	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
T:	R	Α	R	W	#	R	А	W	R	А	R	Α	R	W	R	А
lpb[i]:	0	0	1	0	0	1	2	0	1	2	3	2	3	4	1	2
lpb[10]=3:																
lpb[2]=1:																

Z-function

Again, combine string into T = P#S('#') is a unique character). Store Z[i] = length of longest matching prefix starting at index i. (skip <math>i=0)

To compute Z[i], compute estimate for match length, then scan forwards.

Track r = furthest scanned index; When r increases, set 1 to match index, so S[1 ... r] = S[0 ... r - 1 + 1]. Use min(Z[i - 1], r - 1 + 1) as estimate for Z[i].

Index:	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Т:	R	А	R	W	#	R	Α	W	R	А	R	А	R	W	R	А
z[i]:		0	1	0	0	2	0	0	3	0	4	0	1	0	2	0
z[8]=3																
Estimates										0	1					

Suffix Array

Array of sorted suffixes of a string. (Often last character is unique smallest char \$)

Lexicographic order: compare first characters, break ties by next character.

Helps to think about **text** indices vs **lex** indices: regular string order is text index, lexicographic order is lex index. (I'll try keep i as text indices and j as lex indices.)

Suffix array: text indices sorted in lex order. Maps lex indices to text indices.

Similar substrings end up in nearby locations in suffix array.

Suffix Array

Text	S
0	R
1	Α
2	W
3	R
4	Α
5	R
6	Α
7	R
8	W
9	R
10	Α
11	\$

Lex	Text												
0	11	\$											
1	10	Α	\$										
2	4	Α	R	Α	R	W	R	Α	\$				
3	6	Α	R	W	R	Α	\$						
4	1	Α	W	R	Α	R	Α	R	W	R	Α	\$	
5	9	R	Α	\$									
6	3	R	Α	R	Α	R	W	R	Α	\$			
7	5	R	Α	R	W	R	Α	\$					
8	0	R	Α	W	R	Α	R	Α	R	W	R	Α	\$
9	7	R	W	R	Α	\$							
10	8	W	R	Α	\$								
11	2	W	R	Α	R	Α	R	W	R	Α	\$		

Burrows–Wheeler Transform (BWT)

Closely related to suffix array: start with sorted cyclic shifts of string.

Take last characters of each cyclic shift.

Due to sorted structure, similar substring appearances are grouped into ranges of equal characters, so run-length encoding after BWT can be used for compression.

"LF mapping": ith occurrence of character c in first character of shifts is at the same text index as ith occurrence of character c in last character of shifts.

Use LF property to compute inverse BWT.

Burrows–Wheeler Transform (BWT)

Text	S
0	R
1	Α
2	W
3	R
4	Α
5	R
6	Α
7	R
8	W
9	R
10	Α
11	\$

Lex	Text	F											L
0	11	\$	R	Α	W	R	Α	R	Α	R	W	R	Α
1	10	Α	\$	R	Α	W	R	Α	R	Α	R	W	R
2	4	Α	R	Α	R	W	R	Α	\$	R	Α	W	R
3	6	Α	R	W	R	Α	\$	R	Α	W	R	Α	R
4	1	Α	W	R	Α	R	Α	R	W	R	Α	\$	R
5	9	R	Α	\$	R	Α	W	R	Α	R	Α	R	W
6	3	R	Α	R	Α	R	W	R	Α	\$	R	Α	W
7	5	R	Α	R	W	R	Α	\$	R	Α	W	R	Α
8	0	R	Α	W	R	Α	R	Α	R	W	R	Α	\$
9	7	R	W	R	Α	\$	R	Α	W	R	Α	R	Α
10	8	W	R	Α	\$	R	Α	W	R	Α	R	Α	R
11	2	W	R	Α	R	Α	R	W	R	Α	\$	R	Α

Longest Common Prefix (LCP) Array and Longest Common Extension (LCE) Queries

Given suffix array SA, LCP[j] = longest common prefix of SA[j] and SA[j - 1]. (Skip j = 0).

Then LCE(ia, ib) = most matching characters starting at text indices ia and ib = min(LCP[SA[ia] + 1 .. SA[ib]]). (May need to swap so SA[ia] < SA[ib].)

Let SAINV[SA[j]] = j. (inverse permutation, maps text indices to lex indices.)

Value of LCP[j] - 1 can estimate next text suffix LCP[SAINV[1 + SA[j]]].

LCP Array and LCE Queries

Text	S
0	R
1	Α
2	W
3	R
4	Α
5	R
6	Α
7	R
8	W
9	R
10	Α
11	\$

Lex	Text	LCP												
0	11	N/A	\$											
1	10	0	Α	\$										
2	4	1	Α	R	Α	R	W	R	Α	\$				
3	6	2	Α	R	W	R	Α	\$						
4	1	1	Α	W	R	Α	R	Α	R	W	R	Α	\$	
5	9	0	R	Α	\$									
6	3	2	R	Α	R	Α	R	W	R	Α	\$			
7	5	3	R	Α	R	W	R	Α	\$					
8	0	2	R	Α	W	R	Α	R	Α	R	W	R	Α	\$
9	7	1	R	W	R	Α	\$							
10	8	0	W	R	Α	\$								
11	2	3	W	R	Α	R	Α	R	W	R	Α	\$		

Range Minimum Queries (RMQ)

Precompute power-of-two-length RMQs.

Break general length RMQ into min of two power-of-two-length RMQ.

Some elements are counted twice, but counting twice in a min is fine.

Index	0	1	2	3	4	5	6	7	8	9
Original RMQ										
Left power of two										
Right power of two										