Contents



4

Knuth-Morris-Pratt algorithm

Boyer-Moore algorithm

String searching

2 Karp-Rabin fingerprint algorithm



Chittaranjan Mandal (IIT Kharagpur)

Section outline

- String search
- Brute force approach





String search

- Given a pattern string p, find first match in text t
- N: # characters in text
- M: # characters in pattern
- Length of pattern is small compared to the length of the text $(N \gg M)$
- Pattern can be pre-preprocessed
- Text cannot be pre-processed •



March 31, 2017

3/24

Chittaranian Mandal (IIT Kharagpur)

Brute force approach

- Check for pattern starting at every text position
- Running time depends on pattern and text
- Worst case: MN comparisons, in practice almost linear
- Slow if M and N are large, and have lots of repetition



Section outline

2 Karp-Rabin fingerprint algorithm

- String matching using hashing
- Efficient hash computation



String matching using hashing



- Match not possible unless computed hash of text substring under consideration matches hash of search string
- In case of match, actual comparison is needed to confirm match



6/24

Efficient hash computation

Example

• . . .

- Pre-compute: 10000 % 97 = 9
- First hash: 31415 % 97 = 84
- Previous hash: 41592 % 97 = 76
- Efficient next hash computation of 15926 (% 97):

```
= (41592 - (4 \times 10000)) \times 10 + 6
= (76 - (4 × 9)) × 10 + 6
= 406
= 18
```

7/24

Efficient hash computation

Example

• . . .

- Pre-compute: 10000 % 97 = 9
- First hash: 31415 % 97 = 84
- Previous hash: 41592 % 97 = 76
- Efficient next hash computation of 15926 (% 97):

$$= (41592 - (4 \times 10000)) \times 10 + 6$$

= (76 - (4 × 9)) × 10 + 6
= 406
= 18

- Choose modulus to be a large prime (q)
- Each window of M expected to be uniformly distributed in [0, q 1]
- Expected running time is $O\left(N + M\left(\frac{1}{q}(N M)\right)\right)$
- Worst case: Θ(MN) when?

ヨトィヨト

Efficient hash computation

Efficient hash computation

Example

• . . .

- Pre-compute: 10000 % 97 = 9
- First hash: 31415 % 97 = 84
- Previous hash: 41592 % 97 = 76
- Efficient next hash computation of 15926 (% 97):

 $= (41592 - (4 \times 10000)) \times 10 + 6$ = (76 - (4 × 9)) × 10 + 6 = 406 = 18

- Choose modulus to be a large prime (*q*)
- Each window of M expected to be uniformly distributed in [0, q 1]
- Expected running time is $O\left(N + M\left(\frac{1}{q}(N M)\right)\right)$
- Worst case: $\Theta(MN)$ when?
- Possible if the computed hash matches every time and confirmational matching is needed
- Published in 1987 as Efficient randomized pattern-matching algorithms

Section outline



Knuth-Morris-Pratt algorithm

- Optimised pattern matching with KMP
- KMP algorithm
- KMP failure function

computation

- KMP failure function algorithm
- Overall complexity of KMP
- Optimised failure function computation



	Se	arch	patte	ern	
1	2	3	4	5	6
а	а	С	а	а	b



- Suppose aacaa is received; current state will be 6 and b will be expected
- If b is not received, the pattern will have to be moved forward
- Instead of moving forward by one position (brute force approach), better to align the prefix

 aa with the suffix aa at the point of failure amounts to resuming comparison at state 3

Search pattern											
1	2	3	4	5	6						
а	а	С	а	а	b						



- Suppose aacaa is received; current state will be 6 and b will be expected
- If b is not received, the pattern will have to be moved forward
- Instead of moving forward by one position (brute force approach), better to align the prefix

 aa with the suffix aa at the point of failure amounts to resuming comparison at state 3
- We want the *longest prefix* (aa) that is a *suffix at the point of failure* (state 6)

Search pattern											
1	2	3	4	5	6						
а	а	С	а	а	b						



- Suppose aacaa is received; current state will be 6 and b will be expected
- If b is not received, the pattern will have to be moved forward
- Instead of moving forward by one position (brute force approach), better to align the prefix aa with the suffix aa at the point of failure – amounts to resuming comparison at state 3
- We want the *longest prefix* (aa) that is a *suffix at the point of failure* (state 6)
- Similarly, if after receiving aaca another a is not received; failure is at state 5; comparison may be resumed from state 2

	Se	arch	patte	ern	
1	2	3	4	5	6
а	а	С	а	а	b



- Suppose aacaa is received; current state will be 6 and b will be expected
- If b is not received, the pattern will have to be moved forward
- Instead of moving forward by one position (brute force approach), better to align the prefix aa with the suffix aa at the point of failure – amounts to resuming comparison at state 3
- We want the longest prefix (aa) that is a suffix at the point of failure (state 6)
- Similarly, if after receiving aaca another a is not received; failure is at state 5; comparison may be resumed from state 2
- Failure transitions are meant to step back in the pattern string, staying at the same place implies ∞-loop, so on failure at state 3 (on ¬c), matching is resumed at state 2

	Se	arch	patte	ern						
1	2 3 4 5 6 a c a b b b c a b b c c a b c c a b c									
а	а	С	a a b							
0	1	2	1	2	3					
	Fai	lure	funct	ion						



- Suppose aacaa is received; current state will be 6 and b will be expected
- If b is not received, the pattern will have to be moved forward
- Instead of moving forward by one position (brute force approach), better to align the prefix aa with the suffix aa at the point of failure – amounts to resuming comparison at state 3
- We want the longest prefix (aa) that is a suffix at the point of failure (state 6)
- Similarly, if after receiving aaca another a is not received; failure is at state 5; comparison may be resumed from state 2
- Failure transitions are meant to step back in the pattern string, staying at the same place implies ∞-loop, so on failure at state 3 (on ¬c), matching is resumed at state 2
- The point of resumption for failure at a certain point is the failure function

Algorithms

KMP algorithm

KMP algorithm

Example



- Use knowledge of search pattern
- Build automaton from pattern
- Run automaton on text
- On failure, go back to the longest proper prefix that is a suffix at the point of the last match – to avoid looping (at state 3, for example)



-

KMP algorithm

KMP algorithm

Example



- Use knowledge of search pattern
- Build automaton from pattern
- Run automaton on text
- On failure, go back to the longest proper prefix that is a suffix at the point of the last match – to avoid looping (at state 3, for example)

j ←1 // start of P

- ② for i←1 to N // span through T
- 3 while j > 0 and T[i] \neq P[j]
- i ← fail[j] // fail while no match
- if j=M return i-M+1
 - // terminate (success)
- j←j+1 // move forward in pattern
- 7 return NoMatch terminate (failure)



b

KMP failure function computation

Base case fail[1]=0 – need to start over again \checkmark

- Need to compute fail[*i*], assuming fail[*j*], j < i are available inductive cases, starting with k = 1 (when $\alpha' = \epsilon$)
- fail^k[i − 1] indicates the longest proper prefix (say α) that is a suffix at P[i − 1]; α = ε if fail^k[i − 1] = 0

Exit case where fail^{*k*}[i - 1] = 0 No prefix, so resume matching at the beginning, fail $[i] = 1 \checkmark$

Case
$$P[i-1] = P[fail^{k}[i-1]]$$



Thus, α · P[i − 1] is the longest proper prefix that is also a suffix at P[i], so fail[i] = fail^k[i − 1] + 1 √

KMP failure function computation (contd.)

Inductive cases, starting with k = 1 (when $\alpha' = \epsilon$, contd.)

fail^k[i − 1] indicates the longest proper prefix (say α) that is a suffix at P[i − 1]; α = ε if fail^k[i − 1] = 0

Case $P[i-1] \neq P[fail^{k}[i-1]]$

- Now, $\alpha \cdot \mathbf{P}[i-1]$ is not an admissible suffix at $\mathbf{P}[i]$, as $\mathbf{P}[i-1] \neq \mathbf{P}[\text{fail}^{k}[i-1]]$
- But, fail[fail^k[i 1]] = fail^{k+1}[i 1] indicates the longest proper prefix (say β) that is a suffix at **P**[fail^k[i 1]]



- Now, β is the longest proper prefix of **P** and also a suffix of α
- Thus, if $\mathbf{P}[\text{fail}^{k+1}[i-1]] = \mathbf{P}[i-1], \beta \cdot \mathbf{P}[i-1]$ is the longest proper prefix at $\mathbf{P}[i]$, so fail $[i] = \text{fail}^{k+1}[i-1] + 1 \checkmark$
- Continue induction with $k \leftarrow k + 1 \mathbb{Q}$

KMP failure function computation example

Example (Some steps of failure function computation for "aacaab")

- Consider failure at P[*i*] (say, P[6]=b)
- We would like to identify the longest proper prefix that is a suffix at P[i]
- The longest proper prefix that is a suffix at P[i − 1] is denoted by fail[i − 1]
- So, if P[i 1]=P[fail[i 1]], then fail[i]=fail[i 1]+1
- P[5]=a; P[fail[5]]=P[2]=a; so fail[6]=fail[5]+1=2+1=3

KMP failure function computation example

Example (Some steps of failure function computation for "aacaab")



- Consider failure at P[*i*] (say, P[6]=b)
- We would like to identify the longest proper prefix that is a suffix at P[i]
- The longest proper prefix that is a suffix at P[i − 1] is denoted by fail[i − 1]
- So, if P[i 1] = P[fail[i 1]], then fail[i] = fail[i 1] + 1
- P[5]=a; P[fail[5]]=P[2]=a; so fail[6]=fail[5]+1=2+1=3
- If P[i − 1]≠P[fail[i − 1]], then continue checking from P[fail[fail[i − 1]]]=P[fail²[i − 1]], and so on, but stopping at P[1]
- While computing fail[4], we find P[3]=c and P[fail[3](=2)]=a; P[3]≠P[fail[3](=2)], so go further back to fail²[3]=1 and stop there (at P[1]=a)

Chittaranjan Mandal (IIT Kharagpur)

KMP failure function computation algorithm



KMPCompFail(P[1..*M*])

Example (FF for "aacaab")

① j←0

- If or i←1 to M // span through M!
- I fail[i] ←j

// next prepare for fail[i+1]

- While (j>0 and P[i]≠P[j]) do
 - j←fail[j]
- one done

5

Ø j←j+1

		S	earch	earch pattern					
ı	1	2	3	4	5	6			
P [<i>i</i>]	а	а	С	а	а	b			
fail[1]	0	1	2	1	2	3			
			2			С			
P [<i>j</i> ₄]	-	а	a	а	а	а			
			a			а			
			1			2			
<i>Ĵ</i> 5	-	-	0	-	-	1			
			0			0			
Ĵ7	1	2	1	2	3	1			

Chittaranjan Mandal (IIT Kharagpur)

Algorithms

KMP failure function algorithm (contd.)



	KMPCompFail(P[1 <i>M</i>])	Exam	ole ((FF	for	"aa	caal	o")		
1	j←0				Se	earc	h pa	ttern		
2	for $i \leftarrow 1$ to M // span through M!	ı	1	2	3	4	5	6	7	8
3	fail[i]←j	P [<i>i</i>]	С	а	d	С	а	С	а	d
	<pre>// next prepare for fail[i+1]</pre>	fail[1]	0	1	1	1	2	3	2	3
4	while (j>0 and P[i] \neq P[j]) do	P [14]	_	С	С	С	а	d	а	а
5	j←fail[j]	• [J4]				Ŭ	~	С	~	~
6	done	15	_	0	0	_	_	1	_	_
7	j←j+1			_	_					
8	endfor	Ĵ7	1	1	1	2	3	2	3	4

Chittaranjan Mandal (IIT Kharagpur)

Algorithms

March 31, 2017

Complexity of computing fail[*i*] and running KMP

KMPCompFail(P[1..M])

- **①** j←0
- If i ←1 to M // span through M!
- I fail[i] ←j
- while (j>0 and P[i]≠P[j])
 - j←fail[j]
- one done

5

- **⊘** j←j+1
- endfor

-							
Exam	ple	(FF	for '	'aac	aab	")	
		Sea	arch	patt	ern		
	1	2	3	4	5	6	
	а	a	С	a	а	b	
	0	1	2	1	2	3	
		Fail	ure	func	tion		

- Note that fail[j]<j</p>
- In L5 j is decreased by at least 1
- Overall *j* can go back in L5 only as much as it has progressed in L7
- L7 is executed M times
- Complexity of KMPCompFail is *O*(*M*) (from 2*M*))



16/24

Complexity of computing fail[1] and running KMP

KMPCompFail(P[1..M])

- **①** j←0
- 2) for $i \leftarrow 1$ to M // span through M!
- I fail[i] ←j
- while (j>0 and P[i]≠P[j])
 - j←fail[j]
- one done

5

⊘ j←j+1

endfor

Example (FF for "aacaab")



- Note that fail[j]<j</p>
- In L5 *j* is decreased by at least 1
- Overall *j* can go back in L5 only as much as it has progressed in L7
- L7 is executed M times
- Complexity of KMPCompFail is *O*(*M*) (from 2*M*))
- Using similar reasoning complexity of the KMP algorithm is O(N) (from 2N)
- Overall complexity is O(M + N)
- Publication: Fast pattern matching in strings, D E Knuth, J H Morris, V R Pratt, SIAM JoC, v6, n2, June 1997

Optimised failure function computation



- Consider the failure at P[5]=a; fail[5]=2
- But P[2]=a, so after failing to match a at P[5], failure is guaranteed at P[2]
- This definite failure could be remedied by going all the way back to fail³[5]=0
- Function KMPOptFail does the required post-processisng – employing dynamic programming

KMPOptFail(P[1..M], fail[1..M])

- If for i←2 to M // bottom-up DP
- if P[i]=P[fail[i]] // definite failure
- o fail[i]←fail[fail[i]] // fail all
- endfor // way back via DP

Example (Opt FF for "aacaab")



Section outline



- Bad character shift rule
- Good suffix shift rule
- GSS computation



18/24

Key aspects

The Boyer-Moore algorithm is based on three ideas:

- Scanning the pattern pat from right to left: P[M], P[M 1], ...
- The "bad character shift rule": skips over parts of pattern where there is no possibility of matching the current character in the text
- The "good suffix shift rule": aligns only matching pattern characters against target characters already successfully matched
- These rules work independently, but are more effective together

19/24

< ロ > < 同 > < 回 > < 回 > < □ > <

March 31, 2017

									Text	stri	ng								
D	0		n	u	r	t	u	r	е		t	h	е	f	u	t	u	r	е
f	u	t	u	r	е														
S	Searc	h pa	ittern	i (r≠e	e)														

	Text string																			
D	0		n	u	r	t	u	r	e		t	h	е		f	u	t	u	r	е
f	u	t	u	r	е															
S	Searc	h pa																		
										stri	ng									
D	0		n	u	r	t	u	r	е		t	h	е		f	u	t	u	r	е
			+		r	<u> </u>														
	†	u	ι	u	11	6														







Bad character shift rule (contd)

- If the current character *c* of the text **T** does not match the corresponding character in the pattern, jump to the right most occurrence of *c* in **P** without shifting backwards
- If c does not occur ahead in P, just slide P all the way back

March 31, 2017

Bad character shift rule (contd)

- If the current character *c* of the text **T** does not match the corresponding character in the pattern, jump to the right most occurrence of *c* in **P** without shifting backwards
- If c does not occur ahead in P, just slide P all the way back

Definition ($R_i(x)$ **)**

 $R_i(x)$ is the position of the rightmost occurrence of character x before position i



21/24

Bad character shift rule (contd)

- If the current character *c* of the text **T** does not match the corresponding character in the pattern, jump to the right most occurrence of *c* in **P** without shifting backwards
- If c does not occur ahead in P, just slide P all the way back

Definition ($R_i(x)$ **)**

 $R_i(x)$ is the position of the rightmost occurrence of character x before position *i*

- $\forall c \in \Sigma, R_i(c) = \begin{cases} 0 & \text{if } \exists j < i | \mathbf{P}[j] = c \\ \max \{j < i | \mathbf{P}[j] = c\} & \text{otherwise} \end{cases}$
- $R_i(c)$ is computed in $O(|\Sigma| + M)$ time
- When a mismatch occurs at pattern position *i* in **P**, shift by *i* - *R_i*(**T**[*k*]) characters so that the next occurrence of **T**[*k*] in **P** is underneath position *k* in **T**
- Best case time complexity for BCS is O(N/M) sublinear!



Chittaranjan Mandal (IIT Kharagpur)

Algorithms



21/24

Good suffix shift rule

Example (Comparing shifts by BCS and GSS for $r \neq u_1$)





22/24

Chittaranjan Mandal (IIT Kharagpur)

Good suffix shift rule

Example (Comparing shifts by BCS and GSS for $r \neq u_1$)



- Let $s = \mathbf{P}[i..M] = \mathbf{T}[i + j..j + M]$ and $\mathbf{P}[i-1] \neq \mathbf{T}[i+j-1]$
- The GSS rule aligns the substring *s* with its rightmost occurrence in P (but not as a suffix) that is preceded by a character different from P[*i* - 1]
- Similar to the optimised KMP failure function from the right end



Algorithms

Example (P has the form: $\delta \cdot b \cdot \alpha \cdot \gamma \cdot a \cdot \alpha$, mismatch at $a \cdot \alpha$)



KMP: fail[6]=3, fail@P^r[..aab], resume@ P^r[aac]

BM: fail@P[6 - 3 + 1 = 4](caa), resume@ P[6 - 6 + 11](baa), gss[4]=1

Example (P has the form: $\delta \cdot b \cdot \alpha \cdot \gamma \cdot a \cdot \alpha$, mismatch at $a \cdot \alpha$)



- KMP: fail[6]=3, fail@P^r[..aab], resume@ P^r[aac]
- BM: fail@P[6 3 + 1 = 4](caa), resume@ P[6 6 + 11](baa), gss[4]=1
- KMP: fail[3]=2, fail@P'[..ac], resume@ P'[aa]
- BM: fail@P[6 2 + 1 = 5](aa), resume@ P[6 3 + 1 = 4](ca), gss[5]=4

Example (P has the form: $\delta \cdot b \cdot \alpha \cdot \gamma \cdot a \cdot \alpha$, mismatch at $a \cdot \alpha$)



- KMP: fail[6]=3, fail@P^r[..aab], resume@ P^r[aac]
- BM: fail@P[6 3 + 1 = 4](caa), resume@ P[6 6 + 11](baa), gss[4]=1
- KMP: fail[3]=2, fail@P'[..ac], resume@ P'[aa]
- BM: fail@P[6 2 + 1 = 5](aa), resume@ P[6 3 + 1 = 4](ca), gss[5]=4
- KMP: fail[1]=fail[2]=fail[4]=fail[5]=0
- BM: unfilled cells are marked with '-' indicating start over

Example (P has the form: $\delta \cdot b \cdot \alpha \cdot \gamma \cdot a \cdot \alpha$, mismatch at $a \cdot \alpha$)



- KMP: fail[6]=3, fail@P^r[..aab], resume@ P^r[aac]
- BM: fail@P[6-3+1 = 4](*caa*), resume@ P[6-6+11](*baa*), gss[4]=1
- KMP: fail[3]=2, fail@P'[..ac], resume@ P'[aa]
- BM: fail@P[6 2 + 1 = 5](aa), resume@ P[6 3 + 1 = 4](ca), gss[5]=4
- KMP: fail[1]=fail[2]=fail[4]=fail[5]=0
- BM: unfilled cells are marked with '-' indicating start over
- Here, KMP does not help for P[M], but BCS can, eg $R_6(c) = 4$
- On mismatch on *c* at $\mathbf{P}[i]$, can use max(GSS[i], $R_i(c)$)

GSS computation (contd.)

Example (P also has the form: (longest β) $\beta \cdot \gamma \cdot \beta$, mismatch at $a \cdot \alpha$)



Note the prospect of matching "fl" at the head of P, indicated by gss[8]=0

• gss[6]=-2, gss[5]=-3, gss[4]=-4, gss[3]=-5, gss[2]=-6, gss[1]=-7

• For
$$i < \max \{j | \mathbf{P}[j] = 0\} \land (fail[i] = 0 \lor fail[i] = 1), gss[i] = i - j$$