

# LPS & KMP Algorithm

## 0. Preview

⇒ 'bem'이 'benji'의 부분 문자열인지 알고 싶다면?

bem의 길이    benji 길이

한 문자씩 대조해 가면서 알아보는 방법이 있는데, 시간복잡도가 최악의 경우

$O(N \times m)$   
비효율적!!!

## 1. LPS (Longest prefix suffix)

### 1-1) 접두사(Prefix)와 접미사(Suffix)

⇒ LPS를 위해서 문자열 맨 앞에 있는 접두사와 맨 뒤에 있는 접미사가 같은 경우를 먼저 찾아야 한다.

ex) **ABA****AB** ⇒ 이처럼 LPS는 서로 동일한 접두사 & 접미사의 가장 긴 길이를 의미한다.  
Prefix                  suffix

### 1-2) 파이 배열

⇒ 위 문자열 'ABAAB'에서, 인덱스 별로 부분 문자열을 만들어 가며 살펴본다.

ex) [index] 인덱스 별로 같은 문자들의 부분 문자열 LPS를 구한다. (앞↔끝만 봄)

0	A	B	A	A	B	[LPS] 0
1	A	B	A	A	B	0
2	A	B	A	A	B	1
3	A	B	A	A	B	1
4	A	B	A	A	B	2

이 LPS 값을 배열로 정리하면,  
[0, 0, 1, 1, 2]가 나오는데  
위 배열을 **파이 배열** 이라고 한다!

또 다른 예제로, 'ABCABDAB'의 파이 배열을 구해 보자.

[index]	[pattern]	[LPS]
0	ABCABDAB	0
1	ABCABDAB	0
2	ABCABDAB	0
3	ABCABDAB	1
4	ABCABDAB	2
5	ABCABDAB	0
6	ABCABDAB	1
7	ABCABDAB	2

$\therefore LPS = [0, 0, 0, 1, 2, 0, 1, 2]$

여기서 살펴볼 수 있는 점은, 이전 index의 LPS 값을 활용해서 현재 index의 LPS 값을 계산할 수 있다.

⇒ 앵..? 그렇다면 index 5번은 왜..? (이건 밑에 여제에서!)

[Example] 이때, 'pat = ABCABDAB' 라 가정!

- ① 위에서 본 index == 3 일 때, ABCABDAB 로 LPS 값이 1이었다. (동일한 접두사 & 접미사)
- ② 그 다음 index == 4 일 때, ABCABDAB 인데... 이전 LPS 값이 1이므로 B가 추가되기전 접두사 & 접미사가 1개로 이미 일치한다. 즉,  $pat[0] == pat[3]$ 이다.  
 때문에, 이전 LPS에서 확인된 것을 생략하고 그 다음 문자만 비교하면 된다.  
 즉!  $pat[1] == pat[4]$  만 확인하면 된다. → 만약 같다면, 이전 LPS 값에 (+1)을 하면 된다.
- ③ 그렇다면, index 4와 5로 예시로 한번 살펴보자.  
 index == 5 일 때, <sup>0 1 2 3 4 5</sup> ABCABDAB 가 된다. index == 4 일 때 LPS가 2였다. ( $pat[0] == pat[4]$ )  
 따라서  $pat[1] == pat[5]$ 를 검사하는데.. 일치하지 않는다! 이런 때는..  $pat[2]$ 가 아닌  $pat[0]$ 이다.

LPS 배열에서 LPS[최대 LPS값-1]로 한다. 즉 LPS[1] 값이 pat[10]과 pat[5]를 비교한다.

Q. 왜 이렇게 비교하는가? (맞은 평가는 아님!)

→ 이전 index 1개의 LPS값

A. ABCABAB 앞뒤 가정하면.. pat[10]==pat[5]가 될수 있는 case가 존재하므로..!

## 2. KMP Algorithm

⇒ KMP는 Knuth, Morris, Pratt 라는 세 명의 개발자 이름을 따온 명칭이다.

앞서 pi 배열을 사용하여 불필요한 검사를 건너 뛴다.

### 2-1) 동작과정

원본 문자열과 부분 문자열을 비교한다.

A	B	C	D	A	B	C	D	A	B	E
---	---	---	---	---	---	---	---	---	---	---

A	B	C	D	A	B	E
---	---	---	---	---	---	---

① 순서 좋아 'ABCDAB'까지는 일치했다.

그러나, 마지막 문자 E가 일치하지 않았다.

② 우리는 'ABCDAB'까지, 즉 index가 5까지 일치했다는 사실을 알고 있다.

이 사실과 LPS값을 사용해서 불필요한 계산을 줄인다.

A	B	C	D	A	B	C	D	A	B	E
---	---	---	---	---	---	---	---	---	---	---

A	B	C	D	A	B	E
---	---	---	---	---	---	---

LPS: 0 0 0 0 1 2

⇒ index 5에서의 LPS 값은 2이다. 즉, 앞의 두 문자와 뒤의 두 문자가 일치한다는 의미이다.

③ 뒤 AB 자리에 앞 AB가 오도록 패턴을 오른쪽으로 당길수 있다.

A	B	C	D	A	B	C	D	A	B	E
---	---	---	---	---	---	---	---	---	---	---

A	B	C	D	A	B	E
---	---	---	---	---	---	---

⇒ 새롭게 위치한 상태에서 첫 'AB'가 일치한다는 것을 알고 있으므로 (ps 만큼) 뒤로가서 C에서  
검사를 새롭게 하면 된다!

∴ 시간 복잡도  $O(M+N)$   
↓  
전체 문자열      부분 문자열