

Lecture 13

# 스트링 알고리즘 (2)

컴퓨터과학과 | 김진욱 교수

## 학습목차

**1 | 보이어-무어 알고리즘**

**2 | 데이터 압축 기본 개념**

**3 | RLE**

# 01. 보이어-무어 알고리즘

### ► Boyer-Moore 알고리즘

- 패턴 내의 문자들의 관계를 이용하여 매칭 시 중복된 비교를 줄임
- 텍스트의 첫 위치에서 **패턴의 뒷부분부터** 문자 비교
- 불일치 또는 매치 발생 → 패턴 이동
  - ✓ 불일치 문자 방법, 일치 접미부 방법 중 더 많이 이동시킬 수 있는 값 선택

위치 0 1 2 3 4 5 6 7 8 9 ...

T = a b a b a b c a b a b c a b a

P = a b a b c a b



### ▶ 일치 접미부 good suffix 방법

- 일치한 서브스트링(접미부)에 대한 접두부와 접미부의 최대 일치 정보 활용
- KMP 알고리즘과 유사하면서도 다름

위치 0 1 2 3 4 5 6 7 8 9 ...

T = a b a b a b c a b a b c a b a

P = a b a b c a b

```
BadChar (m, P[ ],  $\Sigma$ )  
{  
  for (알파벳의 각 문자 c)  
     $\delta_1[c] = -1$ ;           // -1로 배열 초기화(패턴에 없는 경우 대비)  
  for (i=0; i < m; i++)  
     $\delta_1[P[i]] = i$ ;        // 패턴에서 가장 마지막에 나타나는 위치만 남김  
  return ( $\delta_1$ );  
}
```

$\Sigma = \{a, b, c, d\}$

$P =$

0	1	2	3	4	5	6
a	b	a	b	c	a	b

	a	b	c	d
$\delta_1$	-1	6	4	-1

```
GoodSuf (m, P[ ])  
{  
    revP = P를 뒤집은 스트링;  
    revF[0..m-1] = PreKMP (m, revP); // 뒤집은 패턴의 최대 일치 정보 찾기  
    for (i=-1; i < m; i++)  
         $\delta_2[i] = m-1 - \text{revF}[m-1];$  // 전체 일치인 경우의 이동값으로 배열 초기화  
    for (k=m-1; k >= 0; k--)  
         $\delta_2[m-1-\text{revF}[k]-1] = k - \text{revF}[k];$   
    return ( $\delta_2$ );  
}
```

	-1	0	1	2	3	4	5	6
P =	a	b	a	b	c	a	b	
$\delta_2$	5	5	5	5	5	3	3	1

	0	1	2	3	4	5	6
revP =	b	a	c	b	a	b	a
revF	-1	-1	-1	0	1	0	1



BM (n, T[ ], m, P[ ],  $\Sigma$ )

{

$\delta_1 = \text{BadChar} (m, P, \Sigma);$

$\delta_2 = \text{GoodSuf} (m, P);$

$i = m-1;$

while ( $i < n$ ) {

$j = m-1;$                    // 패턴의 마지막 문자부터 시작

while ( $j \geq 0 \ \&\& \ P[j] == T[i]$ ) {  $i--; j--;$  }

if ( $j == -1$ ) {

    위치  $i+1$  출력;       // 매치 발견

$i += \delta_2[-1] + m;$    // 전체 일치인 경우로 이동

} else                   // 불일치 문자 정보와 일치 접미부 정보 중

$i += \max(j - \delta_1[T[i]], \delta_2[j]) + m-1-j;$    // 큰 값만큼 이동

}

}

위치 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14  
 $T =$  a b a b a b c a b a b c a b a  
 $P =$  a b a b c a b  
 0 1 2 3 4 5 6

```

i = m-1;
while (i < n) {
    j = m-1;
    while (j >= 0 && P[j] == T[i]) { i--; j--; }
    if (j == -1) {
        위치 i+1 출력;
        i +=  $\delta_2[-1] + m$ ;
    } else {
        i += max(j -  $\delta_1[T[i]]$ ,  $\delta_2[j]$ ) + m-1-j;
    }
}
    
```

	a	b	c	d
$\delta_1$	5	6	4	-1

	-1	0	1	2	3	4	5	6
$\delta_2$	5	5	5	5	5	3	3	1

### ▶ 전처리 $\rightarrow O(m)$

- 불일치 문자 방법  $\rightarrow O(m+|\Sigma|)$ 
  - ✓ 알파벳 크기( $|\Sigma|$ )와 패턴 길이( $m$ )만큼 반복
- 일치 접미부 방법
  - ✓ KMP 알고리즘의 전처리  $\rightarrow O(m)$
  - ✓ 패턴의 길이( $m$ )만큼 두 번 반복

```
for (알파벳의 각 문자 c)
     $\delta_1[c] = -1$ ;
for (i=0; i < m; i++)
     $\delta_1[P[i]] = i$ ;
```

```
revF[0..m-1] = PreKMP (m, revP);
for (i=-1; i < m; i++)
     $\delta_2[i] = m-1 - \text{revF}[m-1]$ ;
for (k=m-1; k >= 0; k--)
     $\delta_2[m-1-\text{revF}[k]-1] = k - \text{revF}[k]$ ;
```

# 성과 특징

01 | 보이어-무어 알고리즘

▶ 최악의 시간 복잡도  $\rightarrow O(nm)$

```
while (i < n) {  
    j = m-1;  
    while (j >= 0 && P[j] == T[i]) { i--; j--; }  
    if (j == -1) {  
        위치 i+1 출력;  
        i +=  $\delta_2[-1] + m$ ;  
    }  
    else  
        i += max(j -  $\delta_1[T[i]]$ ,  $\delta_2[j]$ ) + m-1-j;  
}
```

$O(n)$

$O(m)$

1 증가

T = aaaaaaaaaa...aaa

P = aaaaaa

# 성과 특징

01 | 보이어-무어 알고리즘

▶ 최선의 시간 복잡도  $\rightarrow O(n/m)$

```
while (i < n) {  
    j = m-1;  
    while (j >= 0 && P[j] == T[i]) { i--; j--; }  
    if (j == -1) {  
        위치 i+1 출력;  
        i +=  $\delta_2[-1] + m$ ;  
    }  
    else  
        i += max(j -  $\delta_1[T[i]]$ ,  $\delta_2[j]$ ) + m-1-j;  
}
```

$O(1)$

m 증가

T = aaaaaaaaaa...aaa

P = bbbbbb

## 02. 데이터 압축 기본 개념

### ▶ 주어진 데이터를 보다 적은 공간을 사용하여 표현하는 것

- 데이터 → 1차원(스트링), 2차원(이미지), 3차원(동영상) 등

```
uteo: ~/gzip Gnugi$  
uteo:~/gzip Gnugi$ ls -l  
  
-rw-rw-r-- 1 Gnugi 2609111 Aug 14 16:03 above_cys1_Se1,20000.txt  
uteo:~/gzip Gnugi$  
uteo:~/gzip Gnugi$  
uteo:~/gzip Gnugi$  
uteo:~/gzip Gnugi$ gzip above_cys1_Se1,20000.txt  
uteo:~/gzip Gnugi$  
uteo:~/gzip Gnugi$  
uteo:~/gzip Gnugi$  
uteo:~/gzip Gnugi$ ls -l  
  
-rw-rw-r-- 1 Gnugi 778760 Aug 14 16:03 above_cys1_Se1,20000.txt.gz  
uteo:~/gzip Gnugi$
```

### ▶ 무손실 압축

- 압축된 데이터로부터 원래의 데이터를 완벽하게 복원할 수 있는 압축 방법
- 데이터의 내용 하나하나가 모두 중요한 경우에 사용
- 종류 → RLE, 허프만 코딩, LZ77 등

### ▶ 손실 압축

- 압축된 데이터로부터 원래의 데이터를 완벽하게 복원할 수 없는 압축 방법
- 데이터의 내용이 약간 변형되어도 무방한 경우에 사용
- 종류 → JPEG 표준, MPEG 표준 등



### ▶ 인코딩 encoding

- 원래의 데이터를 압축된 데이터로 변환하는 것

### ▶ 디코딩 decoding

- 압축된 데이터를 압축되지 않은 데이터로 변환하는 것
- 손실 압축의 경우 디코딩 결과가 인코딩 이전 데이터와 동일하지 않을 수 있음

**03.**

**RLE**

## ▶ Run Length Encoding

- 스트링에서 **동일 문자**가 **연속**해서 나타나는 것을 그 문자와 반복 횟수로 압축하는 방법

bbbbbb



(b,5)

### ▶ 주어진 스트링을 차례로 보며 문자가 달라질 때까지 횟수 셈

```
RLE_enc (n, S[])
```

```
{
```

```
  int idx=0;
```

```
  for (i=0; i < n; i++) {
```

```
    count = 1;
```

```
    while (i+1 < n && S[i] == S[i+1]) {
```

```
      count++; i++;
```

```
    }
```

```
    C[idx++] = (S[i], count);
```

```
  }
```

```
  return (C[0..idx-1]);
```

```
}
```

aaabbbbbaaccccbaaaaaaa

↓ ↓ ↓ ↓ ↓ ↓  
(a,3)(b,5)(a,2)(c,4)(b,1)(a,7)

▶ 압축된 데이터를 차례로 보며 각 문자를 횟수만큼 반복

```
RLE_dec (m, C[])
```

```
{
```

```
    int i=0;
```

```
    for (idx=0; idx < m; idx++)
```

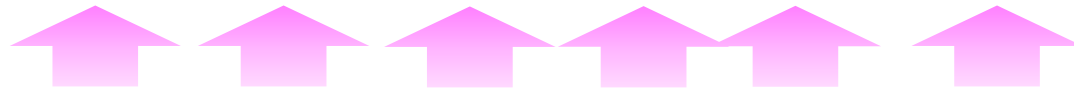
```
        for (j=0; j < C[idx].count; j++, i++)
```

```
            S[i] = C[idx].ch;
```

```
    return (S[0..i-1]);
```

```
}
```

aaabbbbbbbaaccccbaaaaaaa



(a,3)(b,5)(a,2)(c,4)(b,1)(a,7)

## ▶ 성능 → $O(n)$

- 인코딩과 디코딩 모두 이중 루프로 보이지만 실제로는 스트링의 길이( $n$ )만큼 수행

```
for (i=0; i < n; i++) {  
    count = 1;  
    while (i+1 < n && S[i] == S[i+1]) {  
        count++; i++;  
    }  
    C[idx++] = (S[i], count);  
}
```

i 증가

i 증가

i 증가

```
for (idx=0; idx < m; idx++)  
    for (j=0; j < C[idx].count; j++, i++)  
        S[i] = C[idx].ch;
```

## 1. 보이어-무어 알고리즘

- 패턴 내의 문자들의 관계를 이용하여 매칭 시 중복된 비교를 줄임
- 텍스트의 첫 위치에서 패턴의 뒷부분부터 문자 비교
- 불일치 문자 방법과 일치 접미부 방법 이용
- 전처리  $O(m)$ , 매칭-최악  $(nm)$ , 최선  $(n/m)$

## 2. 데이터 압축 기본 개념

- 주어진 데이터를 보다 적은 공간을 사용하여 표현하는 것
- 무손실 압축, 손실 압축
- 인코딩, 디코딩

## 3. RLE

- 스트링에서 동일 문자가 연속해서 나타나는 것을 그 문자와 반복 횟수로 압축
- 성능  $O(n)$

다음시간에는

Lecture **14**

**스트링 알고리즘 (3)**

컴퓨터과학과 | **김진욱** 교수