

# Lecture 12

# 스트링 알고리즘(1)

컴퓨터과학과 | 김진욱 교수



학습목차

- 1 | 기본 개념
- 2 | 라빈-카프 알고리즘
- 3 | KMP 알고리즘



01. 기본 개념

# 스트링 알고리즘 관련 기본 개념

### **스트링 string?**

- 문자가 연속적으로 나열된 문자열
  - ✓ (예) ATATCGCCCACGTAT
  - √ (예) 001011010001110101

### 의 말파벳 alphabet **Σ**

- 스트링에 사용되는 문자들의 집합
  - ✓ DNA 서열 →  $\Sigma$ ={A, C, G, T}
  - **✓** 이진 데이터 → Σ={0, 1}

# 스트림 알고리즘?

### ▶ 스트링에 대한 다양한 문제를 해결하는 알고리즘을 통칭

- 스트링 매칭
- 스트링 압축
- 최장 공통 부분 수열
- 최장 반복 서브스트링
- 최장 회문 서브스트링
- 접미부 트리
- 접미부 배열

...

# 스트링 매칭?

### ■ 텍스트에서 패턴이 나라나는 위치를 찾는 것

- 텍스트  $T=t_0t_1t_2\cdots t_{n-1}$  → 긴 스트링, 길이 n
- 패턴  $P=p_0p_1p_2\cdots p_{m-1}$  → 짧은 스트링, 길이 m

01

# 브루트-포스 스트링 매칭 알고리즘

- **▶** Brute-force algorithm **또는** Naïve algorithm
  - 텍스트의 각 위치에서부터 패턴의 길이만큼 문자를 비교하며 매치를 찾는 방법

```
위치 012345678
T = aabaabaaa
P = aabaa
```

- **>>** Brute-force algorithm **또는** Naïve algorithm
  - 텍스트의 각 위치에서부터 패턴의 길이만큼 문자를 비교하며 매치를 찾는 방법

- **▶** Brute-force algorithm **또는** Naïve algorithm
  - 텍스트의 각 위치에서부터 패턴의 길이만큼 문자를 비교하며 매치를 찾는 방법

01

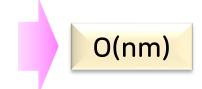
# 브루트-포스 스트링 매칭 알고리즘

- **▶** Brute-force algorithm **또는** Naïve algorithm
  - 텍스트의 각 위치에서부터 패턴의 길이만큼 문자를 비교하며 매치를 찾는 방법

- - 텍스트의 각 위치에서부터 패턴의 길이만큼 문자를 비교하며 매치를 찾는 방법

# 브루트-포스 스트링 매칭 알고리즘

```
BruteForce (n, T[], m, P[])
                         O(n)
 for (i = 0; i <= n-m; i++) { // 텍스트의 가능한 모든 위치에서
  flag = true;
                       0(m)
  for (j = 0; j < m; j++) // 패턴의 길이만큼
   if (T[i+j] != P[j]) { // 텍스트와 패턴의 문자를 비교
     flag = false; // 일치하지 않으면 표시 지움
     break;
  if (flag) 위치 i 출력; // 표시가 남아있으면 매치
```



01 기본 개념

### ▶ 패턴을 전처리

- 라빈-카프 알고리즘
- KMP 알고리즘
- 보이어-무어 알고리즘

### ▶ 텍스트를 전처리

- 접미부 트리
- 접미부 배열



02. 라빈-카프 알고리즘



- 패턴의 해시값으로 매치의 후보를 찾고, 후보에 대해서만 문자별로 비교해서 매치를 찾는 방법
  - 해시 함수: 문자열을 위한 가중 합  $h(S) = \left(\sum_{i=0}^{m-1} {s'}_i |\Sigma|^{m-1-i}\right) \operatorname{mod} M$ 
    - ✓ (예) Σ={a, b, c, d, e, f, g, ···, x, y, z}, M=101

P = a a b a a

h(aabaa)

$$= (0 \times 26^4 + 0 \times 26^3 + 1 \times 26^2 + 0 \times 26^1 + 0 \times 26^0) \mod 101$$

- 패턴의 해시값으로 매치의 후보를 찾고, 후보에 대해서만 문자별로 비교해서 매치를 찾는 방법
  - 해시 함수: 문자열을 위한 가중 합  $h(S) = \left(\sum_{i=0}^{m-1} s'_i |\Sigma|^{m-1-i}\right) \bmod M$

- 패턴의 해시값으로 매치의 후보를 찾고, 후보에 대해서만 문자별로 비교해서 매치를 찾는 방법
  - 해시 함수: 문자열을 위한 가중 합  $h(S) = \left(\sum_{i=0}^{m-1} s'_i |\Sigma|^{m-1-i}\right) \bmod M$

- 패턴의 해시값으로 매치의 후보를 찾고, 후보에 대해서만 문자별로 비교해서 매치를 찾는 방법
  - 해시 함수: 문자열을 위한 가중 합  $h(S) = \left(\sum_{i=0}^{m-1} s'_i |\Sigma|^{m-1-i}\right) \bmod M$

# 텍스트 위치별 해시값 계산

02 | 라빈-카프 알고리즘

의치 0: 0(m)

$$h(aabaa) = (0 \times 26^4 + 0 \times 26^3 + 1 \times 26^2 + 0 \times 26^1 + 0 \times 26^0) \mod 101 = 70$$

× 26

- 의치 1~n-m: 0(1)
  - 해시 함수의 특징을 이용

$$h(abaab) = (0 \times 26^4 + 1 \times 26^3 + 0 \times 26^2 + 0 \times 26^1 + 1 \times 26^0) \mod 101$$

$$= (26 \times h(aabaa) - 0 \times 26^5 + 1 \times 26^0) \mod 101$$

$$= (26 \times 70 - 0 + 1) \mod 101 = 3$$



# 라빈-카프 Rabin-Karp 알고리즘

### 02 | 라빈-카프 알고리즘

```
RabinKarp (n, T[], m, P[])
 int hp = 0, ht = 0;
 int dm = pow(26, m) \% M;
 for (j = 0; j < m; j++) {
 • hp = (hp * 26 + P[j] - 97) % M;
--> ht = (ht * 26 + T[j] - 97) % M;
```

` 패턴의 해시값
` 텍스트 위치 0 해시값

```
for (i = 0; i \le n-m; i++) {
 if (hp == ht) {
                      // 해시값 일치하는 위치에 대해
  flag = true;
  for (j = 0; j < m; j++) // 패턴의 길이만큼
    if (T[i+j] != P[j]) { // 텍스트와 패턴의 문자를 비교
     flag = false; // 일치하지 않으면 표시 지움
     break:
  if (flag) 위치 i 출력; // 표시가 남아있으면 매치
 if (i < n-m)
                     // 다음 위치 해시값 계산
  ht = (ht*26 - dm*(T[i]-97) + (T[i+m]-97)) % M;
                                     한극방송통신대학교
```

# ■ 텍스트 T=10011100에서 패턴 P=0011이 매치되는 모든 위치?

•  $\Sigma = \{0, 1\}, M = 11$ 

$$h(1001) = (1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0) \mod 11 = 9$$

$$h(0011)$$
=  $(2 \times 9 - 1 \times 2^4 + 1 \times 2^0) \mod 11 = 3$ 

$$h(0111)$$
=  $(2 \times 3 - 0 \times 2^4 + 1 \times 2^0) \mod 11 = 7$ 

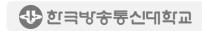
$$h(1110)$$
=  $(2 \times 7 - 0 \times 2^4 + 0 \times 2^0) \mod 11 = 3$ 

$$h(0011) = (0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0) \mod 11 = 3$$

- 성능 → 0(n+km)
  - 전처리 → 0(m)
  - 텍스트에서 해시값 계산 → 0(n)
  - 후보 위치는 문자 직접 비교. 매치 개수 k → 0(km)
- 매치 개수가 삼수 → O(n)
- 모든 위치에서 매치 → O(nm)



# 03. KMP 알고리즘



- ▶ Knuth-Morris-Pratt 알고리즘
  - 패턴 내의 문자들의 관계를 이용하여 매칭 시 중복된 비교를 줄임
  - 텍스트의 첫 위치에서 패턴의 앞부분부터 문자 비교

```
위치 012345678
T = aabaa baaa
P = aabaa
aabaa
aabaa
aabaa
aabaa
aabaa
aabaa
```

03 KMP 알고리즘

- ▶ Knuth-Morris-Pratt 알고리즘
  - 패턴 내의 문자들의 관계를 이용하여 매칭 시 중복된 비교를 줄임
  - 텍스트의 첫 위치에서 패턴의 앞부분부터 문자 비교

```
위치 012345678
T = aababbbaaa
P = aabaa
aabaa
aabaa
aabaa
aabaa
aabaa
```

### ▶ 일치한 서브스트링에 대한 접두부와 접미부의 최대 일치 정보

- $f_i$  → 패턴의 서브스트링  $p_0p_1\cdots p_i$ 에서 최대 일치인 접두부의 끝 문자 위치
  - ✓ 접두부와 접미부의 최대 일치가 없으면  $-1 \rightarrow f_0 = -1$

위치 01234  
P = aabaa  
aabaa  
1 
$$f_4 = 1$$

위치 01234  
P=aabaa  
aabaa  
aabaa  
1 
$$f_4 = 1$$
  
P=aabaa  
aabaa  
-1

```
PreKMP (m, P[])
                         패턴의 위치 0부터 m-1까지 차례대로 F[i], 즉 최대 일치 정보 f_i 구함
 int F[m], idx = -1;
 F[0] = -1;
 for (i=1; i<m; i++) {
   while (idx \geq 0 \&\& P[i] != P[idx+1])
     idx = F[idx]; // 최대 접두부 찾기
   if (P[i] == P[idx+1])
     idx++;
                     // 마지막 문자 일치
   F[i] = idx;
                     // i에서의 최대 접두부 설정
 return (F);
```

```
idx
P = a a b a a
F -10-101
```

# KMP 알고리즘

```
KMP (n, T[], m, P[])
 F[0..m-1] = PreKMP (m, P); // 패턴의 전처리
 \inf j = -1;
 for (i=0; i < n; i++) {
   while (j \ge 0 \&\& T[i] != P[j+1]) j = F[j];
   if (T[i] == P[j+1]) j++;
   if (j == m-1) { // 패턴의 마지막 문자까지 일치하면
     위치 i-j 출력; // 매치 발견
     j = F[j];
```

### 03 KMP 알고리즘

■ 텍스트 T=aabaabaaa에서 패턴 P=aabaa가 매치되는

```
무는 위치?
의치 012345678
T = aabaabaa
P = aabaa
```

## KMP 알고리즘\_예\_2

03 | KMP 알고리즘

● 텍스트 T=ATATATGATATGAA에서 패턴 P=ATATGAT가 매치되는

```
무는 위치?
의치 012345678910111213

T = ATATATGATATGAA

P = ATATGAT
```

# 성능과 특징

### 03 KMP 알고리즘

- 전처리 → O(m)
  - idx 값은 최대 m-1 증가
  - while 문은 최대 m-1만큼 수행
- 마침 → O(n)
  - j 값은 최대 n 증가
  - while 문은 최대 n만큼 수행
- n≥m → 전체 성능 O(n)

```
for (i=1; i<m; i++) {
  while (idx >= 0 && P[i] != P[idx+1]) idx = F[idx];
  if (P[i] == P[idx+1]) idx++;
  F[i] = idx;
}
```

```
for (i=0; i < n; i++) {
  while (j >= 0 && T[i] != P[j+1]) j = F[j];
  if (T[i] == P[j+1]) j++;
  if (j == m-1) {
    위치 i-j 출력;
    j = F[j];
  }
}
```



#### 1. 기본 개념

- 스트링 매칭-텍스트(길이 n)에서 패턴(길이 m)이 나타나는 위치를 찾는 문제
- 브루트-포스 스트링 매칭 알고리즘-성능 O(nm)

### 2. 라빈-카프 알고리즘

- 패턴의 해시값으로 매치의 후보를 찾고, 후보에 대해서만 문자별로 비교
- 성능-매칭 0(n+km) (k: 매치의 개수)
- 매치의 개수에 따라 최선 0(n), 최악 0(nm)

#### 3. KMP 알고리즘

- 패턴 내의 문자들의 관계를 이용하여 매칭 시 중복된 비교를 줄임
- 텍스트의 첫 위치에서 패턴의 앞부분부터 문자 비교
- 전처리-패턴의 각 위치별로 접두부와 접미부의 최대 일치 정보를 구함
- 성능 0(n)-전처리 0(m), 매칭 (n) (단, n≥m)

**P** ALGORITHM 말고리즘

다음시간에는

Lecture 13

스트링 알고리즘 (2)

컴퓨터과학과 | 김진욱 교수

