เอกสารประกอบการอบรม
เรื่อง String Matching Algorithms
อาจารย์ ดร.พิชญะ สิทธิอมร
นำมาจาก Slide วิชา Algorithms
โดย รองศาสตราจารย์ ดร.สมชาย ประสิทธิ์จูตระกูล

การจับคู่สตริง

้นำมาจาก slides ของ อ. สมชาย ประสิทธิ์จูตระกูล

พิชญะ สิทธิอมร

หัวข้อ

นิยามปัญหา **Brute-force (naïve)** Rabin-Karp **Finite-automaton Knuth-Morris-Pratt Boyer-Moore** Horspool

Exact String Matching

You are the fairest of your sex,

Let me be your hero;

I love you as one over x,

As x approaches zero.

text

Positively. You

pattern

Input

You are the fairest of **you**r sex,

Let me be **you**r hero;

I love **you** as one over x,

As *x* approaches zero.

Positively.

Output

Approximate String Matching

You are the fairest of your sex,

Let me be your hero;

I love you as one over x.

As x approaches zero.

Positively. heero

text pattern

Input

You are the fairest of your sex,

Let me be your hero;

I love you as one over x,

As x approaches **zero**.

Positively.

Output

ปัญหา String Matching

* Finite alphabet :
$$\Sigma$$

- $⋄ Σ = {a, b, c, ..., z}$
- $\diamond \Sigma = \{0, 1\}$
 - $\star \Sigma = \{A,C,G,T\}, \dots$
- * Text : T[1..n]
- Pattern : P[1..m]

ต้องการหา all valid shifts first valid shift last valid shift

 $0 \le s \le n - m$

s เป็น <u>valid shift</u> เมื่อ P[1..m] = T[s+1 .. s+m]

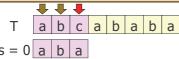
P[1..m]

อัลกอริทึมการเทียบสตริงแบบแม่นตรง

- √ ♦ Brute Force algorithm
- **Deterministic Finite Automaton** algorithm
- √ ♦ Rabin-Karp algorithm
 - Shift Or algorithm Morris-Pratt algorithm
- ✓ Knuth-Morris-Pratt algorithm
- Simon algorithm
- Colussi algorithm
- Galil-Giancarlo algorithm
- **Apostolico-Crochemore algorithm**
- Boyer-Moore algorithm
- Turbo BM algorithm Apostolico-Giancarlo algorithm
- Reverse Colussi algorithm
- ✓ ♦ Horspool algorithm Quick Search algorithm
- **Tuned Boyer-Moore algorithm**

- Zhu-Takaoka algorithm
- Berry-Ravindran algorithm Smith algorithm
- Raita algorithm
- Reverse Factor algorithm
- Turbo Reverse Factor algorithm
- Forward Dawg Matching
- algorithm **Backward Nondeterministic Dawg**
- Matching algorithm **Backward Oracle Matching**
- algorithm
- Galil-Seiferas algorithm
- Two Way algorithm
- String Matching on Ordered Alphabets algorithm
- Optimal Mismatch algorithm
- **Maximal Shift algorithm**
- Skip Search algorithm
 - KMP Skip Search algorithm
- Alpha Skip Search algorithm

Brute Force (Naive) Algorithm



T a b c a b a b a s = 1 a b a

T a b c a b a b a s = 2 a b a

T a b c a b a b a

T a b c a b a b a s = 4 a b a

Brute Force (Naïve) Algorithm

ตัวที่ 1 ไม่เหมือน
$$\left(\frac{d-1}{d}\right)$$
 \times 1 ตัวแรกเหมือน $\left(\frac{1}{d}\right)^{1}\left(\frac{d-1}{d}\right)$ \times 2 สองตัวแรกเหมือน $\left(\frac{1}{d}\right)^{2}\left(\frac{d-1}{d}\right)$ \times 3 $\sum_{k=1}^{m}\frac{k}{d^{k-1}}\left(\frac{d-1}{d}\right)$... m $-$ 1 ตัวแรกเหมือน $\left(\frac{1}{d}\right)^{m-1}\left(\frac{d-1}{d}\right)$ \times m $+$ เหมือนกันทั้ง m ตัว $\left(\frac{1}{d}\right)^{m}$ \times m $+$ $\frac{m}{d}$

$$\frac{m}{d^m} + \sum_{k=1}^m \frac{k}{d^{k-1}} \left(\frac{d-1}{d} \right) = \frac{m}{d^m} + \left(1 - \frac{1}{d} \right) \sum_{k=1}^m \frac{k}{d^{k-1}} = ?$$

$$\sum_{k=1}^{m} x^{k} = \frac{1 - x^{m+1}}{1 - x}$$

$$\sum_{k=1}^{m} kx^{k-1} = \frac{-(1 - x)(m+1)x^{m} + (1 - x^{m+1})}{(1 - x)^{2}}$$

$$= \frac{1 - (m+1)x^{m} + mx^{m+1}}{(1 - x)^{2}}$$

$$= \frac{1 - x^{m} + mx^{m}(x-1)}{(1 - x)^{2}}$$

$$\frac{mx^m(x-1)}{(x-1)^2}$$

$$\frac{m}{d^m} + \left(1 - \frac{1}{d}\right) \sum_{k=1}^m \frac{k}{d^{k-1}} = \frac{1 - d^{-m}}{1 - d^{-1}}$$

$$\sum_{k=1}^m k x^{k-1} = \frac{1 - x^m + m x^m (x - 1)}{(1 - x)^2}$$

$$m x^m + (1 - x) \sum_{k=1}^m k x^{k-1} = m x^m + (1 - x) \left(\frac{1 - x^m + m x^m (x - 1)}{(1 - x)^2}\right)$$

$$= m x^m + \left(\frac{1 - x^m + m x^m (x - 1)}{(1 - x)}\right)$$

$$= \frac{1 - x^m}{1 - x} \longrightarrow \text{TM} \quad x = d^{-1}$$

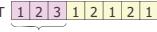
$$\frac{1-d^{-m}}{1-d^{-1}}(n-m+1) < 2(n-m+1)$$

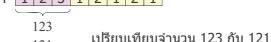
$$\frac{1-d^{-m}}{1-d^{-1}} < \frac{1}{1-d^{-1}}, \ d \ge 2$$

$$\le \frac{1}{1-2^{-1}} = 2$$

Rabin-Karp Matcher

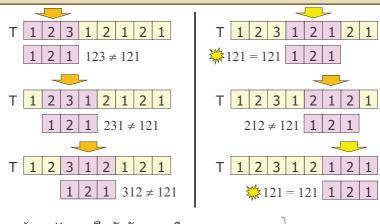






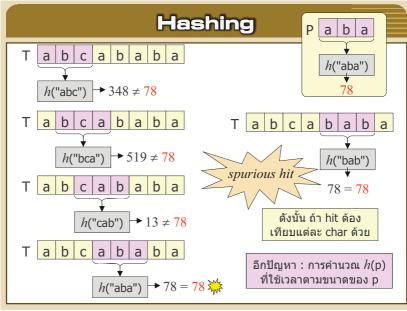


Rabin-Karp Algorithm



ถ้า pattern เป็นตัวอักษร หรือ ถ้าเป็นตัวเลข แต่มีขนาดใหญ่เกิน int เกิน long

ong } จะทำอย่างไร ?



การคำนวณค่า h : incremental

```
❖ h(s) = (มอง s เป็นจำนวน ในระบบเลขฐาน)
h("abc") = (1 \times 37^2 + 2 \times 37^1 + 3 \times 37^0) = 1446
h("bce") = (2\times37^2 + 3\times37^1 + 5\times37^0) = 2853
                 = 37 \times (1 \times 37^{2} + 2 \times 37^{1} + 3 \times 37^{0}) + 5 \times 37^{0}
                 = 37 \times (h("abc") - 1 \times 37^2) + 5 \times 37^0
           เลื่อนตำแหน่ง "bc"
                                             ลบ "a"
                                                           เพิ่ม "e"
               ไปทางซ้าย
         h("...e") = 37 \times (h("a...") - 1 \times 37^{m-1}) + 5 \times 37^{0}
```

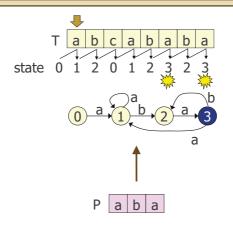
Rabin-Karp Algorithm

```
Rabin-Karp-Matching(T[1..n], P[1..m], d) {
  ht = 0; hp = 0
  for (i = 1 \text{ to } m) {
                                                     d = |\Sigma|
    ht = d*ht + T[i]
    hp = d*hp + P[i]
  for (s = 0 \text{ to } n - m) {
    if (ht == hp) {
       for (i = 1 \text{ to } m)
         if (T[s+i] \neq P[i]) break
       if (i > m) print(s)
                                        ถ้า m มีค่ามาก ขนาด
                                        ของ ht และ hp ก็ใหญ่
    if (s < n-m) {
       ht = d*(ht - T[s+1]*d^{m-1}) + T[s+m+1]
       h("...e") = 37 \times (h("a...") - 1 \times 37^{m-1}) + 5 \times 37^{0}
```

Rabin-Karp Algorithm

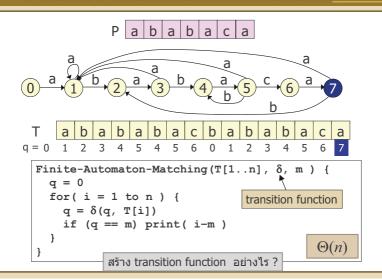
```
Rabin-Karp-Matching( T[1..n], P[1..m], d, q ) {
  ht = 0; hp = 0; dm1 = d^{m-1} % q
  for (i = 1 \text{ to } m) {
    ht = (d*ht + T[i]) % q
                                      worst case \Theta(nm)
    hp = (d*hp + P[i]) % q
  for (s = 0 to n - m) {
                                      avg case O(n + m)
    if (ht == hp) {
      for (i = 1 \text{ to } m)
         if (T[s+i] \neq P[i]) break
      if (i > m) print(s)
    if (s < n-m) {
      ht = (d*(ht - T[s+1]*dm1) + T[s+m+1]) % q
       h("...e") = 37 \times (h("a...") - 1 \times 37^{m-1}) + 5 \times 37^{0}
```

Finite Automaton Matcher

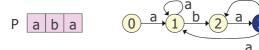


สร้าง state transition diagram จาก pattern

Finite Automaton Matcher



Transition Function



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

state a b c

0 1 0 0

m+1 1 1 2 0

2 3 0 0

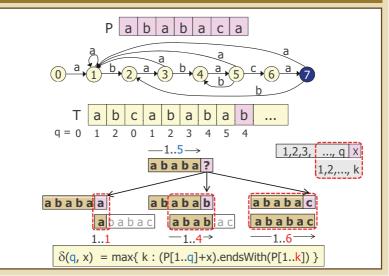
3 1 2 0

input

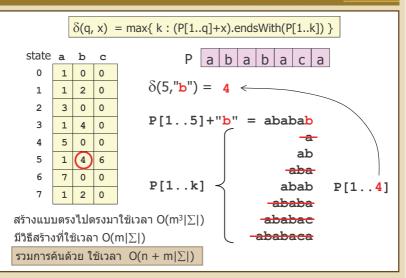
 δ (state, input)

ตารางมีขนาด $(m+1)|\Sigma|$

Transition Function



ตัวอย่าง Transition Function



Knuth-Morris-Pratt Matcher

❖ ของเดิม (transition function)

```
\delta(q, x) = \max\{k : (P[1..q]+x).\text{endsWith}(P[1..k])\}
```

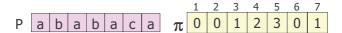
- ❖ ตารางมีขนาด (m+1)| Σ | 1,2,3, ..., q X
- ♦ ใช้เวลาสร้าง $O(m|\Sigma|)$

1,2,..., k

ของใหม่ (prefix function)

- $\star \pi(q) = \max\{k : k < q \text{ and } (P[1..q]).\text{endsWith}(P[1..k])\}$
- match มาได้ g ตัว แต่ mismatch ตัวที่ g+1 ให้ทำต่อที่ P[π(g)+1]
 - ๓ารางมีขนาด m
- 1,2,3, ⋄ ใช้เวลาสร้าง Θ(m)
- ❖ สร้าง π และคัน ใช้เวลารวม Θ(n+m)

การใช้ Prefix Function



Tababaabababaca

c a $\pi[5] = 3$

Tabababab

| C | a | $\pi[3] = 1$

 $\pi[1] = 0$

 $\pi[5] = 3$

Tababaabababaca

KMP-Matcher

```
KMP-Matcher(T[1..n], P[1..m]) {
  \pi = \text{Compute-Prefix-Function}(P) // \Theta(m)
  for ( i = 1 to n ) {
    while (q > 0 \text{ AND } P[q+1] \neq T[i]) {
      q = \pi[q]
    if (P[q+1] == T[i]) q++
    if (q == m) {
                                      \Theta(n+m)
      print( i - m )
       q = \pi[q]
                      h
```

 $\pi[4] = 2$

```
\pi(q) = \max\{ k : k < q \text{ and } (P[1..q]).endsWith(P[1..k]) \}
                      q = 1 \mid a
                                                          \pi[\ 1\ ] = 0
                                                          \pi[2] = 0
                  q = 2
                          a b
                q = 3 \mid a \mid b
                                                          \pi[3] = 1
            q = 4
                              b
                                                          \pi[4] = 2
                           а
                           а
                                 a
         q = 5
                                                          \pi[5] = 3
                a
                                 b
                              a
                                                          \pi[6] = 0
      q = 6
                    а
                           а
                                 а
                                     b
                                           b
                                        а
   q = 7 a b
                    b
                        a c
                                                          \pi[7] = 1
                                           а
                   P[1..q]
                                      P[1..k], k < q
```

```
\pi(q) = \max\{ k : k < q \text{ and } (P[1..q]).endsWith(P[1..k]) \}
                                                        \pi[1] = 0
                                                        \pi[2] = 0
а
                                                        \pi[3] = 0
      а
         Compute-Prefix-Function (P[1..m]) {
            \pi[1] = 0; k = 0
            for (q = 2 \text{ to } m) {
              if (P[k+1] == P[q]) k++
              \pi[q] = k
            return π
```

```
\pi(q) = \max\{ k : k < q \text{ and } (P[1..q]).endsWith(P[1..k]) \}
                                                         \pi[4] = 1
                                                         \pi[5] = 0
                                                        \pi[6] = 1
         Compute-Prefix-Function (P[1..m]) {
            \pi[1] = 0; k = 0
            for (q = 2 \text{ to } m) {
              while (k > 0 \text{ and } P[k+1] \neq P[q]) {
                k = \pi[k]
              if (P[k+1] == P[q]) k++
              \pi[q] = k
            return π
```

```
\pi(q) = \max\{ k : k < q \text{ and } (P[1..q]).endsWith(P[1..k]) \}
                                                        \pi[7] = 2
                                                        \pi[8] = 3
                                                        \pi[9] = 4
         Compute-Prefix-Function (P[1..m]) {
            \pi[1] = 0; k = 0
            for (q = 2 \text{ to } m) {
              while (k > 0 \text{ and } P[k+1] \neq P[q]) {
                 k = \pi[k]
              if (P[k+1] == P[q]) k++
              \pi[q] = k
            return π
```

```
\pi(q) = \max\{ k : k < q \text{ and } (P[1..q]).endsWith(P[1..k]) \}
     \pi[9] = 4 → P[1..4] = P[6..9] → ทดสอบ P[4+1] = P[10]?
     \pi[4] = 1 → P[1..1] = P[4..4] → ทดสอบ P[1+1] = P[10]?
                                                      \pi[10] = 2
     clalclaltlclal
         Compute-Prefix-Function (P[1..m]) {
           \pi[1] = 0; k = 0
            for (q = 2 \text{ to } m) {
              while (k > 0 \text{ and } P[k+1] \neq P[q]) {
                k = \pi[k]
              if (P[k+1] == P[q]) k++
              \pi[q] = k
            return π
```

```
\pi(q) = \max\{ k : k < q \text{ and } (P[1..q]).endsWith(P[1..k]) \}
                                                        \pi[11] = 3
                        а
                                                        \pi[12] = 4
              а
                        а
         Compute-Prefix-Function (P[1..m]) {
            \pi[1] = 0; k = 0
            for (q = 2 \text{ to } m) {
              while (k > 0 \text{ and } P[k+1] \neq P[q]) {
                k = \pi[k]
              if (P[k+1] == P[q]) k++
              \pi[q] = k
                                                  \Theta(m)
            return π
```

Boyer-Moore Matcher

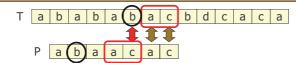
Knuth-Morris-Pratt (left-to-right)

A STRING SEARCHING EXAMPLE CONSISTING OF ... STING

Boyer Moore (right-to-left)

A STRING SEARCHING EXAMPLE CONSISTING OF ... STING

Boyer-Moore: Two heuristics



Bad character

bad-character แนะให้ เลื่อน 3 ตำแหน่ง

Tababa(b)acbdcaca

Pabaacac

เลือกตัวมาก

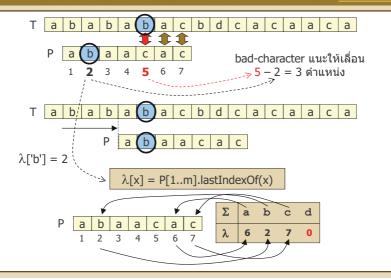
Good suffix

good-suffix แนะให้ เลื่อน 2 ตำแหน่ง

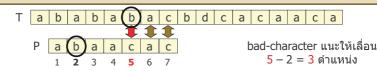
T a b a b a b a c b d c a c a

Pabaacac

Bad-Character Heuristics



Bad-Character Heuristics



T a b a b a b a c b d c a c a a c a
เลื่อน 3 ช่อง p
a b a a c a c
1 2 3 4 5 6 7

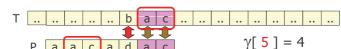
bad-character แนะให้เลื่อน
6-0=6

Tabababacbdcacacac
เลื่อน 6 ช่อง Pabaacac

Σ a b c d λ 6 2 7 0

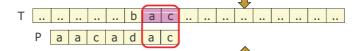
bad-character แนะให้เลื่อน 3 - 7 = -4 (ignored)

Good-Suffix Heuristics

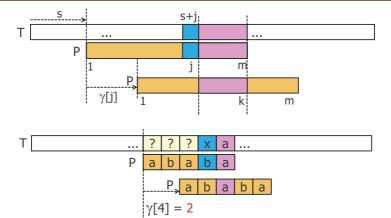


1 2 3 4 5 6 7

mismatch ตำแหน่งที่ <mark>5</mark> good-suffix แนะให้เลื่อน 4 ตำแหน่ง

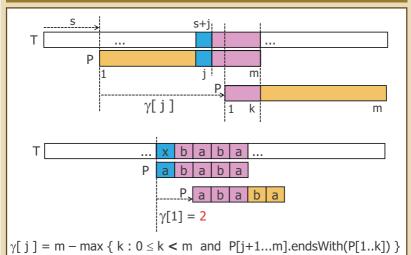


Good-Suffix Heuristics



 $\gamma[\ j\] = m - \max \left\{\ k: 0 \le k < m \ \text{ and } \ P[1..k]. ends With(P[j+1...m])\ \right\}$

Good-Suffix Heuristics



```
\gamma[j] = m - \max \{ k : 0 \le k < m \text{ and } P[j+1...m] \sim P[1..k] \}
                                            \gamma [7] = 1
        AGAGGAG
  j=6
                    A G G A G
                                            \gamma[6] = 7 - 5 = 2
                 G
  j=5
                    G A G G A G
                                            \gamma[5] = 7 - 4 = 3
  j=4
                   G A G G A G
                                            \gamma [4] = 7 - 4 = 3
  j=3
                 GGA
                       A G A G G A G \gamma[3] = 7 - 2 = 5
  j=2
                 G
                            A G G A G \gamma[2] = 7 - 2 = 5
  j=1
          G
                 GG
                            A G G A G \gamma[1] = 7 - 2 = 5
  j=0
           G
              AGGA
                               G G A G \gamma [0] = 7 - 2 = 5
```

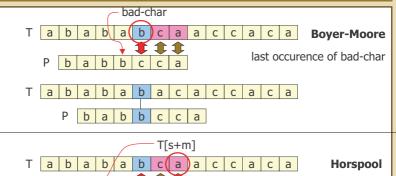
Boyer-Moore Matcher

```
Boyer-Moore-Matcher (T[1..n], P[1..m], \Sigma)
  \lambda = LAST-OCCURRENCE(P, m, \Sigma) // O(|\Sigma|+m)
  \gamma = GOOD-SUFFIX(P, m) // O(m)
  s = 0
  while (s \le n - m) {
    \dot{I} = M
    while (i > 0 \text{ and } P[i] = T[s+i]) i--
    if (j==0) {
      print(s)
       s = s + \gamma[0]
     } else {
       s = s + max(\gamma[j], j - \lambda[T[s+j]])
 worst: O((n-m+1)m+|\Sigma|)
                               best : O(n/m + m + |\Sigma|)
```

Boyer-Moore Matcher

```
Bover Moore-Matcher (T[1..n], P[1..m], \Sigma)
 \lambda = \text{LAST OCCURRENCE}(P, m, \Sigma)
     - GOOD SUFFIX (P, m)
  s = 0
  while (s \le n - m) {
     i = m
     while (i > 0 \text{ and } P[i] = T[s+i]) i--
     if (j==0) {
       print(s)
        s = s + \frac{\gamma(0)}{\gamma} 1
     } else {
        s = s + \frac{\max(\gamma[j], j - \lambda[T[s+j]])}{1}
```

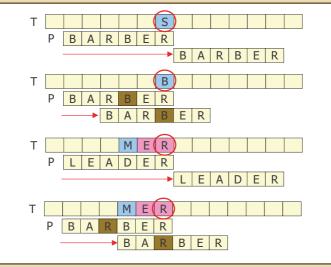
Horspool Matcher



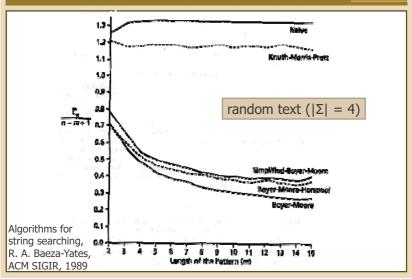
last occurence of T[s+m]

Tababacaacca

Horspool Matcher



ประสิทธิภาพการทำงาน



ประสิทธิภาพการทำงาน

