

データ構造とアルゴリズム (第13回)

文字列のアルゴリズム

文字列のアルゴリズム

- 7.1 文字列照合問題とは
- 7.2 カマかせ法
- 7.3 ラビン-カーブ法
- 7.4 クヌース-モリス-ブラッツ法
- 7.5 ボイヤー-ムーア法 (概要)

この章の学習目標

- 文字列照合問題とは何か, 応用例を用いて説明できる
- カマかせ法とは何か, また, その時間計算量を説明できる
- ハッシュを利用した文字列照合アルゴリズムのアイデアとその時間計算量を説明できる
- 線形時間の文字列照合アルゴリズムのアイデアとその時間計算量を説明できる
- 文字列照合アルゴリズムのプログラムを書ける

Part1:文字列照合問題,
BF法, Rabin-Karp法

文字列照合問題(パターンマッチング)

- 利用分野の例
 - ▣ テキストエディタ, 検索エンジン, 遺伝子解析...

- **文字列照合問題**

- ▣ **入力:**

- テキスト (長さ n の文字列)
 - パターン (長さ m の文字列)

テキスト: $T[0]T[1]...T[n-1]$
パターン: $P[0]P[1]...P[m-1]$

- ▣ **出力:**

- テキスト中にパターンが現れるかどうか判定
 - テキスト中でパターンの位置を求める

- ▣ 多くの場合, $n \gg m$

文字列照合アルゴリズム

$T[0]T[1]T[2]T[3]T[4]T[5]T[6]T[7]T[8]T[9]T[10]T[11]T[12]T[13]$

この中からパターンを見つける

パターン

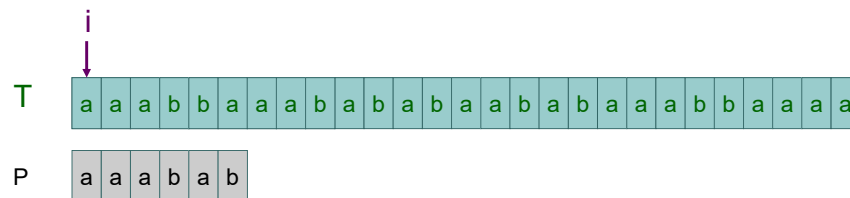
$P[0]P[1]P[2]P[3]$

$n=14, m=4$

→ 答え: $i=5$

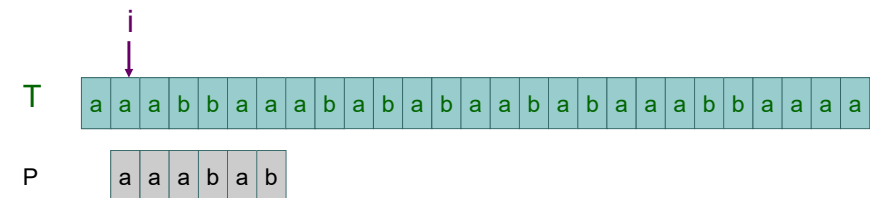
アルゴリズム	計算時間
力まかせ(Brute force)法	$O(nm)$
Rabin-Karp法	平均 $O(n)$ 最悪 $O(nm)$
Knuth-Morris-Pratt法	最悪 $O(n + m)$
Boyer-Moore法	最悪 $O(n + m)$

Brute Force 法



一致せず!

Brute Froce 法



一致せず!

これを繰り返せばよい
1回の比較: 最悪 m
比較開始位置: 最悪 $n-(m-1)$

Brute Force法アルゴリズム

```
// テキスト T[0]T[1] ... T[n-1]
// パターン P[0]P[1] ... P[m-1]
for ( i = 0; i <= n-m; i++) {
  for ( j = 0; j < m; j++)
    if (T[i+j] != P[j]) break;
  if (j==m) return i;
}
return -1;
```

ここはフルには回らない

最悪 $O(nm)$
(実用的にはそこまで
悪くならないことも多い)

最悪を実現する場合

$T = \text{aaaaa...aaaaab}$
 $P = \text{aaaaab}$

Rabin-Karp法

Brute Force + ハッシュ法

ハッシュ関数を用意する

- 長さ m の文字列を引数に取って、整数値を返す
- (m は入力パターンのパターン長)

ハッシュ値は偽陽性(False-positive性)を持つ

- 2つの同じ入力文字列が与えられると、必ず同じ値
- 異なる2つの文字列を与えても(たまに)同じ値(衝突)
 - どれくらいの頻度で起きるかはハッシュ値の値域による
 - (基本的には、値域が広いほうが衝突が起きにくい)

Rabin-Karp法

ハッシュ値

6 5 3 0 6 5 3 6 ...

T a a a b b a a a b a b a b a a b a b a a a a

P a a a b a b

$h(\text{a a a b a b}) = 5$

$h(\text{a a a b b a}) = 6$

$h(\text{a a b b a a}) = 5$

$h(\text{a b b a a a}) = 3$

ハッシュ値が異なるときは
確実に異なるパターン



ハッシュ値が等しいとき
だけBrute-force法と
同様にチェックする
(衝突があるのでチェックは必要)

Rabin-Karp法アルゴリズム

```
// テキスト T[0]T[1] ... T[n-1]
// パターン P[0]P[1] ... P[m-1]
// hp (パターンのハッシュ値) と
// ht (テキストのハッシュ値) の初期値(=h(T[0]...T[m-1]))
```

```
for ( i = 0; i <= n-m; i++) {
  if (hp==ht) { // ハッシュ値が等しい時だけチェック
    for ( j = 0; j < m; j++)
      if (T[i+j] != P[j]) break;
    if (j==m) return i;
  }
  // htの更新(ht = h(T[i+1]...T[i+m]))
}
return -1;
```

動画で $T[0] \dots T[m-1] \rightarrow T[1] \dots T[m]$ と訂正していますが、嘘です。
このままで正しいです。すいません $m(_)m$

ハッシュ関数

□ ハッシュ関数 h

- 長さ m の文字列 \rightarrow ハッシュ値 (整数値)

d : 文字種数
 q : 素数

- $h(a_1 a_2 \dots a_m) = (a_1 d^{m-1} + a_2 d^{m-2} + \dots + a_m) \bmod q$

各文字は $0 \sim d-1$ の数字で符号化
されているとする
(e.g., アルファベット = {a,b} なら $a=1, b=0$ のように)

d 進数として数値化

$$\begin{aligned} h_i^T &= h(T[i]T[i+1] \dots T[i+m-1]) \\ &= (T[i]d^{m-1} + T[i+1]d^{m-2} + \dots + T[i+m-1]) \bmod q \end{aligned}$$

$$\begin{aligned} h_{i+1}^T &= h(T[i+1]T[i+2] \dots T[i+m]) \\ &= (T[i+1]d^{m-1} + T[i+2]d^{m-2} + \dots + T[i+m]) \bmod q \\ &= ((h_i^T - T[i]d^{m-1})d + T[i+m]) \bmod q \end{aligned}$$

- h_{i+1}^T は h_i^T から定数時間で計算可能

Rabin-Karp法アルゴリズム(完全版)

```
// テキスト T[0]T[1] ... T[n-1]
// パターン P[0]P[1] ... P[m-1]
dm=1; hp=0; ht=0;
for(i=0; i<m-1; i++) dm=(d*dm) % q; // d^{m-1}を事前に計算
for(i=0; i<m; i++) {
    hp=(hp*d+P[i]) % q; ht=(ht*d+T[i]) % q;
}
for ( i = 0; i <= n-m; i++) {
    if (hp==ht) { // ハッシュ値が等しい時だけチェック
        for (j = 0; j < m; j++)
            if (T[i+j] != P[j]) break;
        if (j==m) return i;
    }
    ht=((ht-((T[i]*dm) % q) + q)*d+T[i+m]) % q;
}
return -1;
```

Part2 : Knuth-Morris-Pratt法, Boyer-Moore法

Rabin-Karp法

- Rabin-Karp法はハッシュ値の衝突が毎回生じるとパフォーマンスが低下

- 平均時間 $O(n+m)$
- 最悪時間 $O(nm)$

最悪時間 $O(n+m)$ のアルゴリズムは？

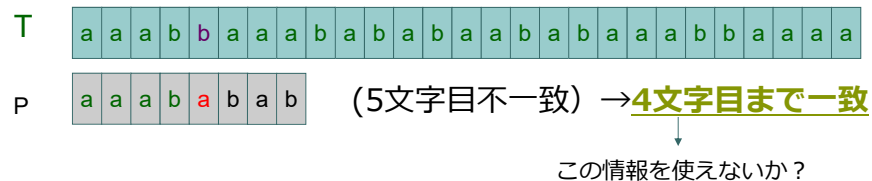


テキストに対してパターンをひとつずつずらすのが、
効率を悪くしている。

じゃあ、どうすればいいか？

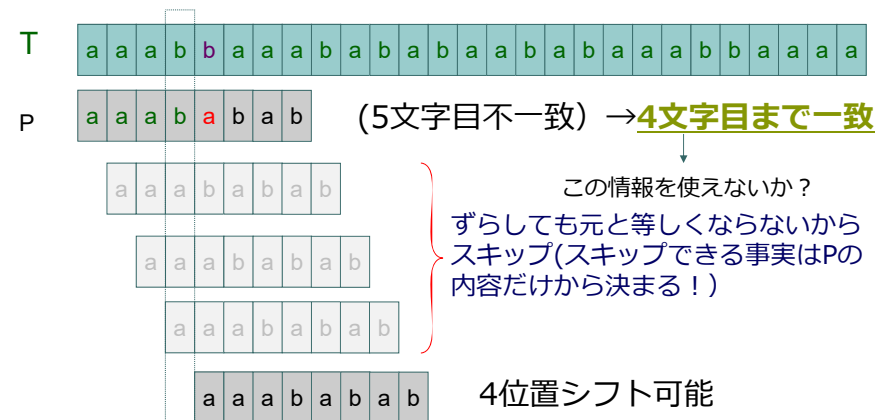
Knuth-Morris-Pratt法

Knuth-Morris-Pratt法 = 照合位置スキップ



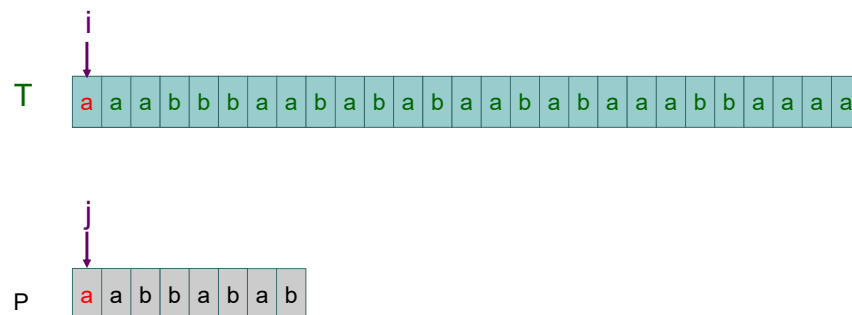
Knuth-Morris-Pratt法

Knuth-Morris-Pratt法 = 照合位置スキップ



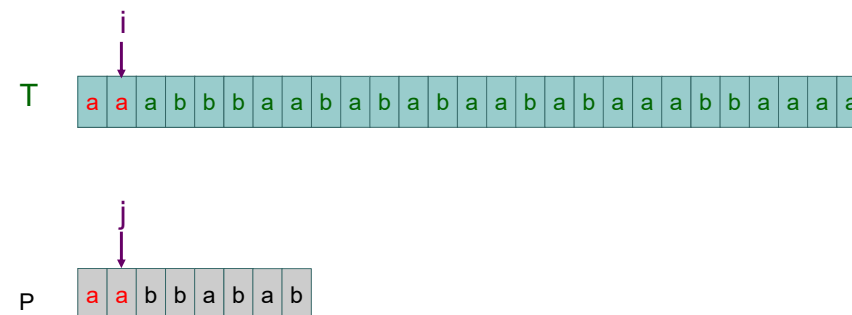
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



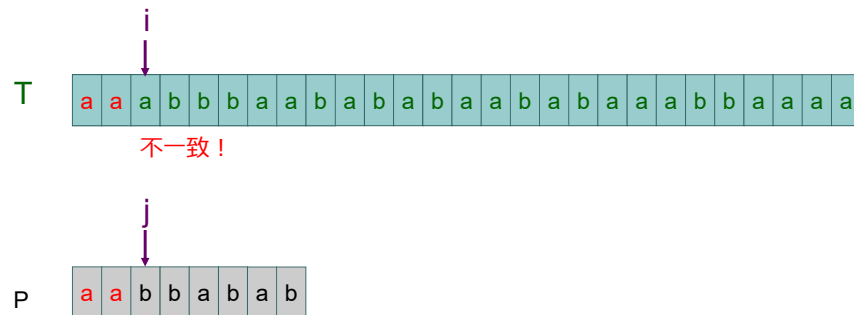
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



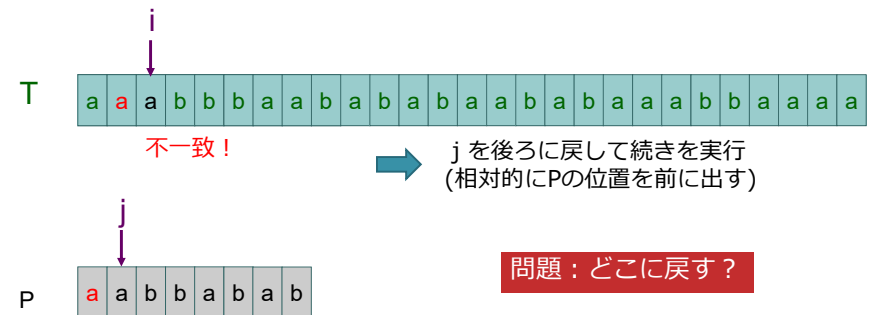
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



Knuth-Morris-Pratt法

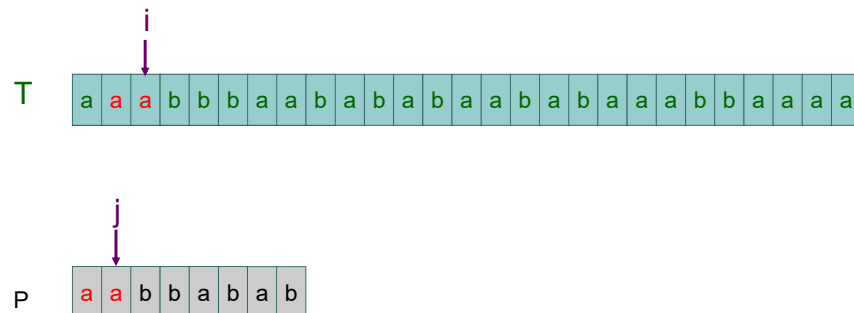
- 実際にはP側の添字位置を「巻き戻す」



$next[k]$ ($-1 \leq k \leq m-1$): $j=k$ で不一致が起きたときの巻戻し位置
これを前もって計算しておけばよい

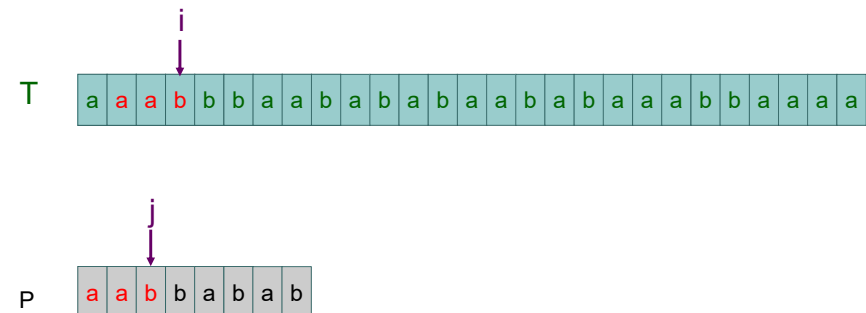
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



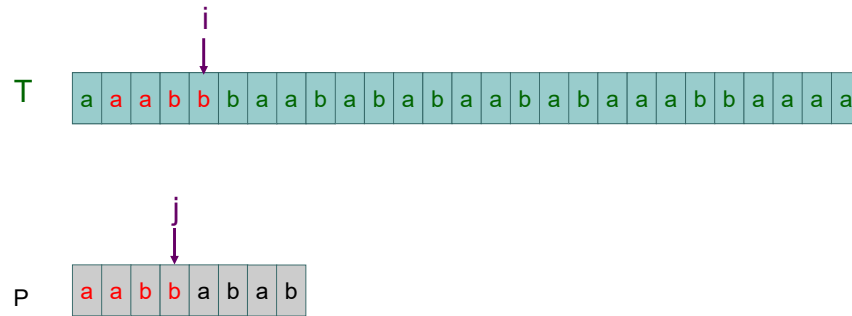
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



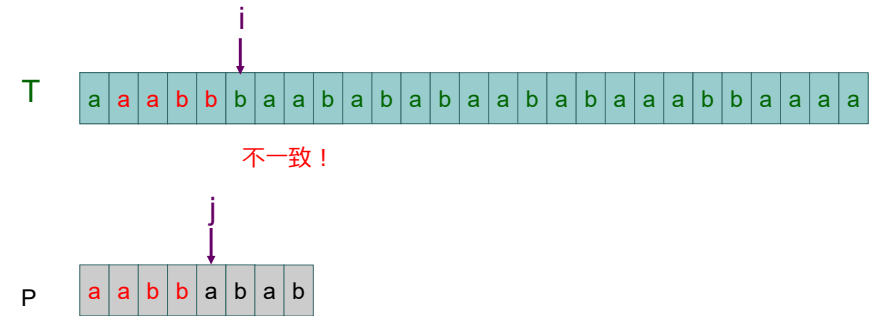
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



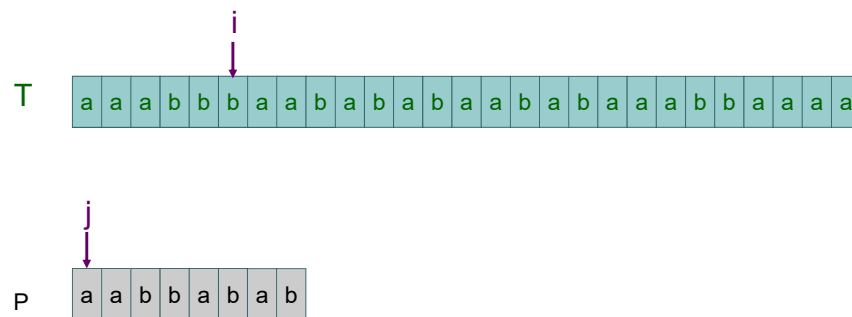
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



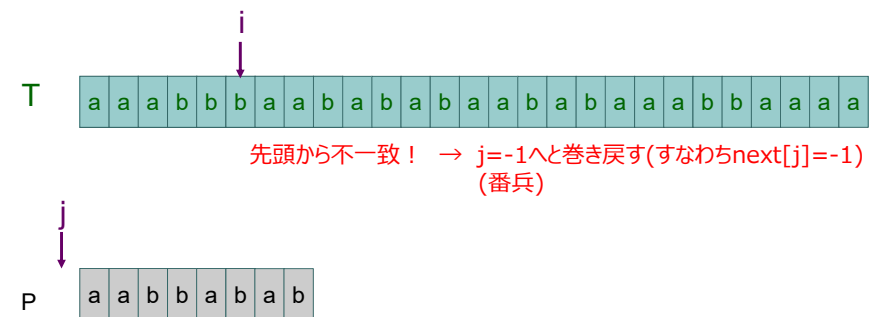
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



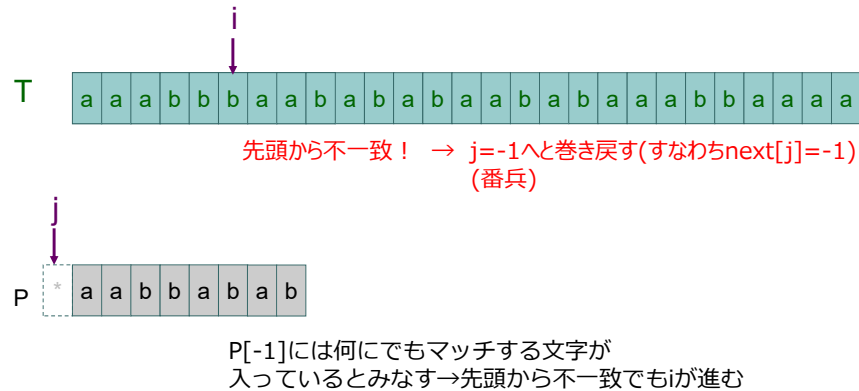
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



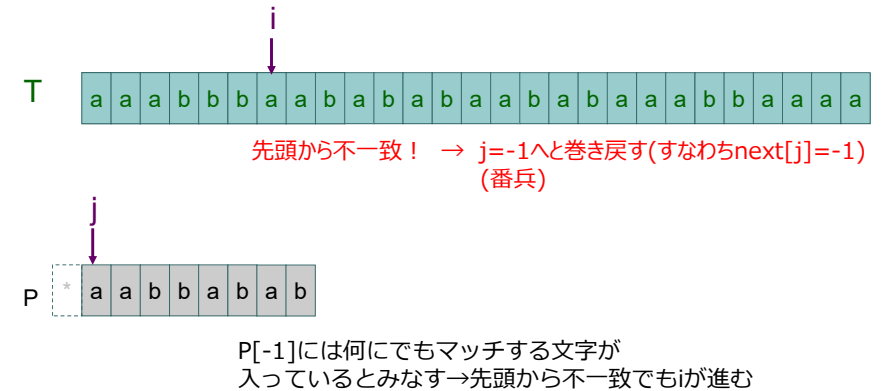
Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



Knuth-Morris-Pratt法

- 実際にはP側の添字位置を「巻き戻す」



Knuth-Morris-Pratt法

```
// テキスト T[0]T[1] ... T[n-1]
// パターン P[0]P[1] ... P[m-1]
// 再開位置 : next[0]=-1, next[1]~next[m-1]
j=0;
for ( i = 0; i < n; i++) {
    while ((j>=0) && (T[i]!=P[j]))
        j = next[j];
    if (j == m-1) return i-m+1;
    j++;
}
return -1;
```

i: テキストの何文字目か
j: パターンの何文字目か

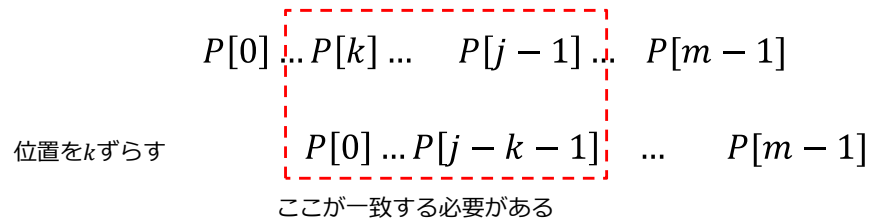
KMP法の計算時間

- KMP法は悪くとも以下の状況では終了する
 - iがnに到達したとき
 - パターンを右までシフトしきったとき (実際には、それに伴ってiの値も大きくなっている)
- While文の中のループ実行1回で、パターンは少なくとも右に1個ずれる
 - パターンが左にずれることはない
→Whileループ内部は高々n回実行

nextの構成を除いて $O(n)$ 時間

nextの計算法

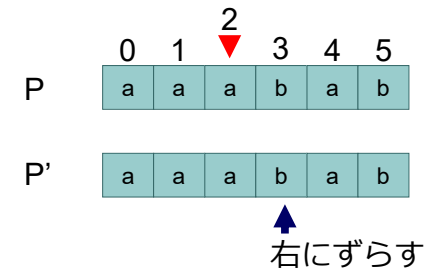
- 巻き戻し位置はどこか?
- パターンの j 文字目で不一致が発生



- $P[k]P[k+1] \dots P[j-1] = P[0]P[1] \dots P[j-k-1]$ なる最小の k に対して, $\text{next}[j]=j-k$ となる

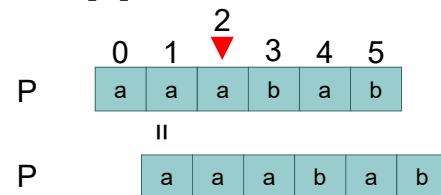
nextの計算法

- パターンとパターンを文字列照合する
- 例: $\text{next}[2]$ を計算してみる



nextの計算法

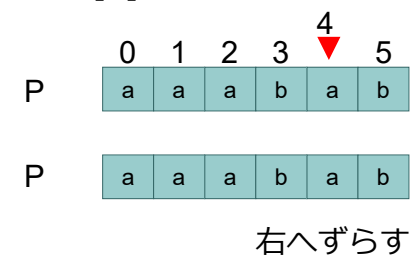
- パターンとパターンを文字列照合する
- 例: $\text{next}[2]$ を計算してみる



$\text{next}[2]=1$

nextの計算法

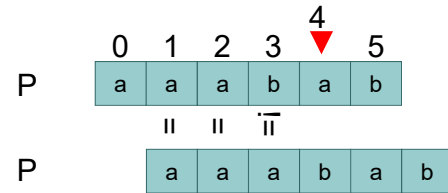
- パターンとパターンを文字列照合する
- 例: $\text{next}[4]$ を計算してみる



nextの計算法

- パターンとパターンを文字列照合する

- 例：next[4]を計算してみる

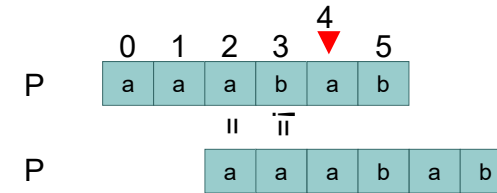


マッチしないのでもう一度

nextの計算法

- パターンとパターンを文字列照合する

- 例：next[4]を計算してみる

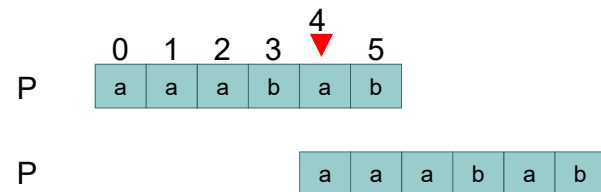


マッチしないのでもう一度
で、これを繰り返して…

nextの計算法

- パターンとパターンを文字列照合する

- 例：next[4]を計算してみる



マッチしないのでもう一度
で、これを繰り返して…

next[4]=0

nextの計算

- 力まかせでやると $O(m^2)$,
- 工夫してやると $O(m)$ ：nextの計算自体にKMPを使う

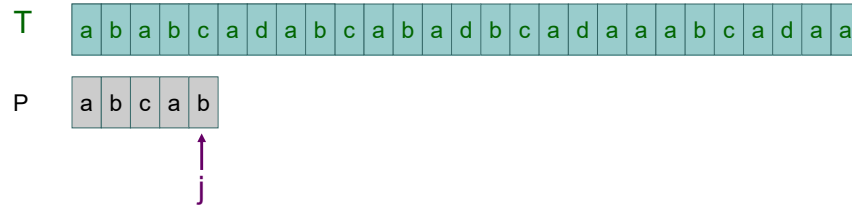
```
// パターン P[0]P[1] ... P[m-1]  
// 再開位置：next[i] の計算  
j=-1;  
for (i=0; i<m; i++) {  
    next[i]=j;  
    while ((j>=0) && (P[i] !=P[j]))  
        j=next[j];  
    j++;  
};
```

- 詳細は省略するが、これでnextを正しく計算できる
 - 実際にはiについての帰納法で証明する

Boyer-Moore法(素朴版)

□ 基本アイデア

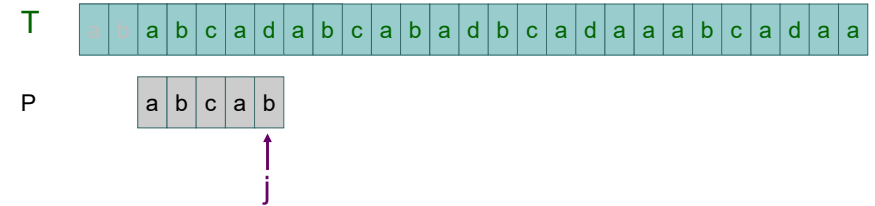
- テキストを前から、パターンを後ろから照合する



Boyer-Moore法(素朴版)

□ 基本アイデア

- テキストを前から、パターンを後ろから照合する

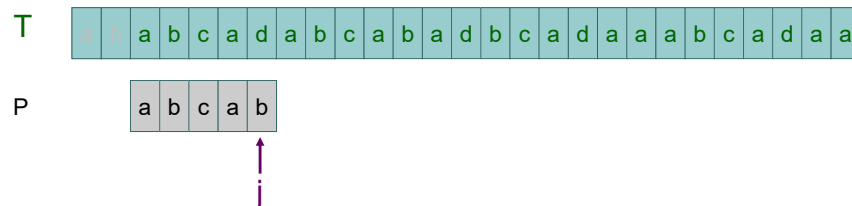


- $T[j]=c$ で不一致なので、右からチェックして最初にcが出てくるところまでPをシフトできる

Boyer-Moore法(素朴版)

□ 基本アイデア

- テキストを前から、パターンを後ろから照合する

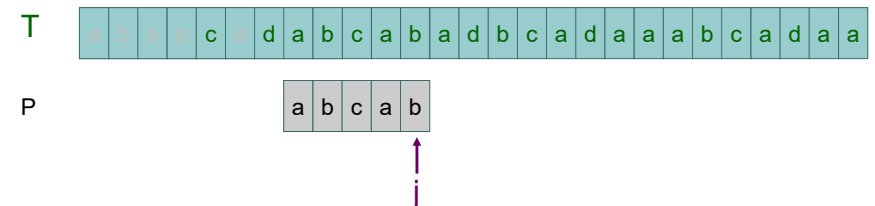


- $T[j]=b$ で不一致なので、右からチェックして最初にdが出てくるところまでPをシフトできる
→Pはdを含まないので端までシフトできる

Boyer-Moore法(素朴版)

□ 基本アイデア

- テキストを前から、パターンを後ろから照合する

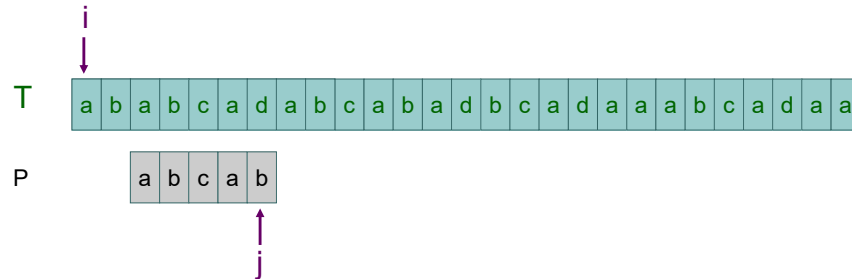


- $T[j]=b$ で不一致なので、右からチェックして最初にdが出てくるところまでPをシフトできる
→Pはdを含まないので端までシフトできる

Boyer-Moore法(素朴版)

基本アイデア

- テキストを前から、パターンを後ろから照合する

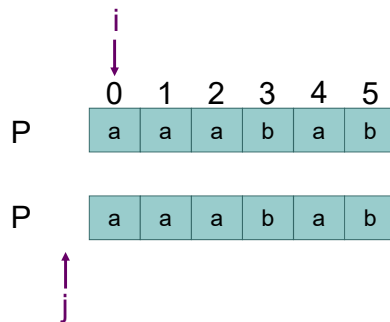


- $T[j]=c$ で不一致なので、右からチェックして最初にcが出てくるところまでPをシフトできる

Boyer-Mooreのご利益

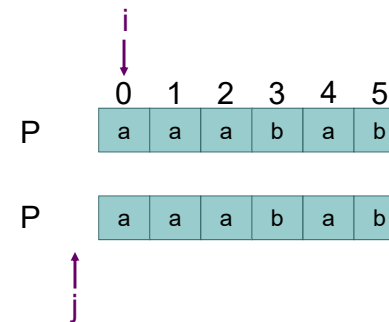
- Boyer-Mooreはテキストを全部見ないで済むことがある
 - うまくいくと $O(n/m)$ 時間で終わる（し、実用的にはそうなることも多い）
- 素朴版では最悪時は $O(mn)$ 時間
- KMPのように可能シフト量を事前計算する方法と組み合わせると、最悪時を $O(n + m)$ にできる(詳細は略)

参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置 : next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

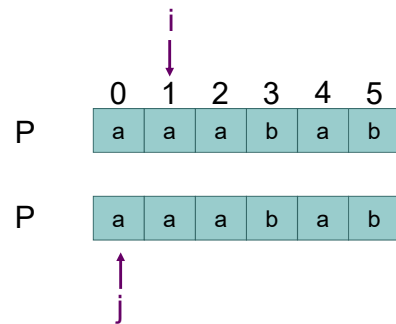
参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置 : next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

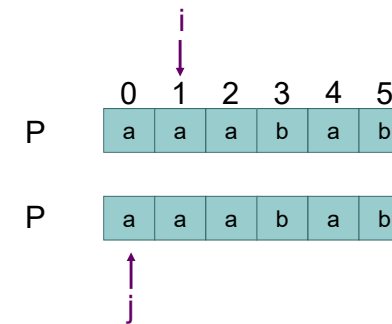
next[0] = -1

参考：KMP法でのnextの計算(例)



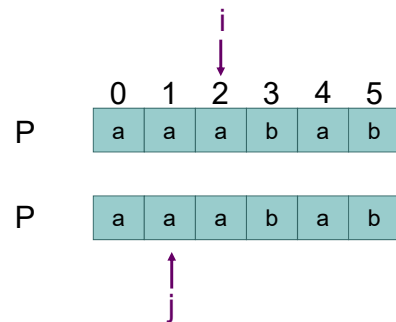
```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

参考：KMP法でのnextの計算(例)



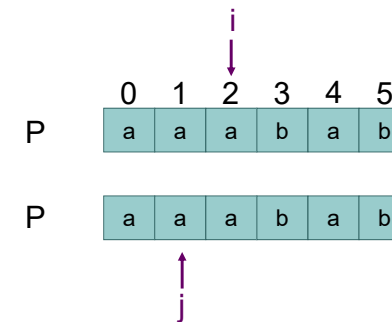
```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

参考：KMP法でのnextの計算(例)



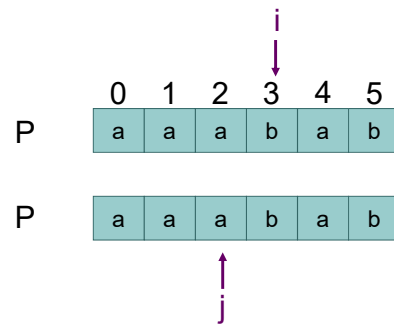
```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

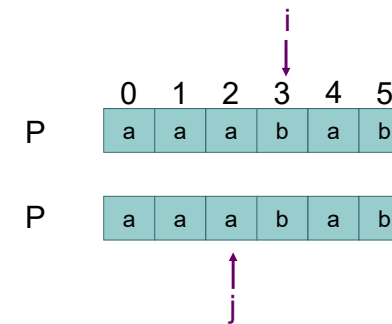
参考：KMP法でのnextの計算(例)



next[0] = -1, next[1] = 0, next[2] = 1

```
// パターン P[0]P[1] ... P[m-1]
// 再開位置 : next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

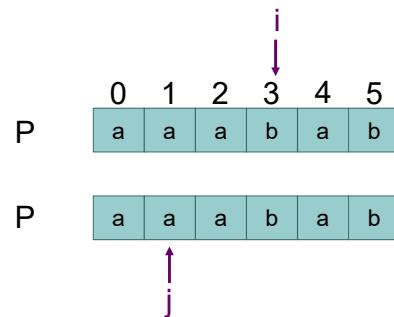
参考：KMP法でのnextの計算(例)



next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2

```
// パターン P[0]P[1] ... P[m-1]
// 再開位置 : next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

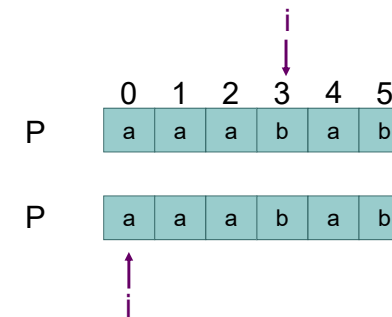
参考：KMP法でのnextの計算(例)



next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2

```
// パターン P[0]P[1] ... P[m-1]
// 再開位置 : next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

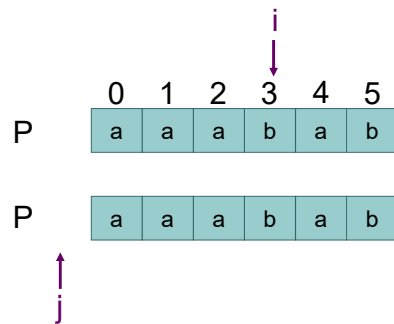
参考：KMP法でのnextの計算(例)



next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2

```
// パターン P[0]P[1] ... P[m-1]
// 再開位置 : next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

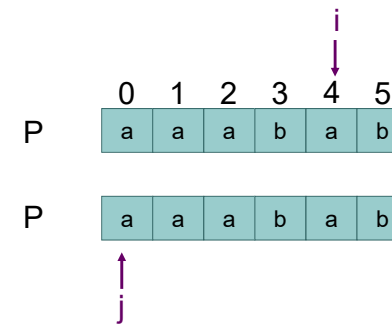
参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2

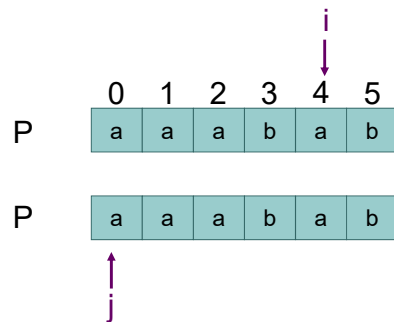
参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2

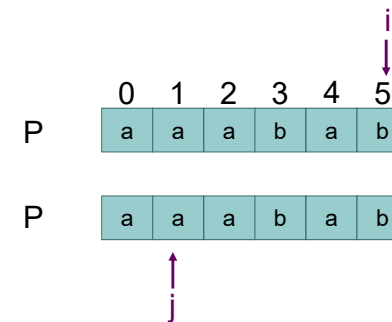
参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2, next[4]=0,

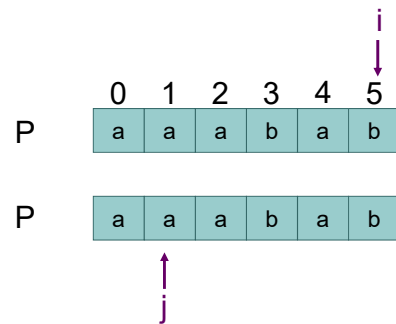
参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
    next[i]=j;
    while ((j>=0) && (P[i] !=P[j]))
        j=next[j];
    j++;
};
```

next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2, next[4]=0

参考：KMP法でのnextの計算(例)



```
// パターン P[0]P[1] ... P[m-1]
// 再開位置：next[i] の計算
j=-1;
for (i=0; i<m; i++) {
  next[i]=j;
  while ((j>=0) && (P[i] !=P[j]))
    j=next[j];
  j++;
};
```

next[0] = -1, next[1] = 0, next[2] = 1, next[3]=2, next[4]=0,
next[5] = 1