

アルゴリズムとデータ構造

第14回: 文字列照合

担当: 上原 隆平(uehara)

2014/05/29

文字列照合 (String matching)

- テキスト文字列: $\text{text} = t_1t_2\dots t_n$
- パターン文字列: $\text{patt} = p_1p_2\dots p_m$
- 問題: text の中に patt と等しいものが含まれるか. 含まれるときはその位置を出力せよ.
- 例:
 - $\text{text} = \text{ababcababcbab}$
 - $\text{patt} = \text{abcabac}$
 - 答: ababcababcabacbab (8)

Q. どうやって計算する?

素朴な方法: 逐次探索

- 入力
 - $t[]$: 1番目から n 番目
 - $p[]$: 1番目から m 番目
- 出力:
 - s if $t[s, s+m-1] = p[1, m]$
 - 0 otherwise

$a[s, e]$: s 番目から e 番
目までの部分列, e.g.,

$a_1 a_2 \dots \underbrace{a_i \dots a_j}_{a[i, j]} \dots a_n$

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++)
        if(t[s+i-1]!=p[i])
            break;
    if(i == m+1)
        return s;
}
return 0;
```

素朴な方法: 逐次探索

- 例:

– text=ababcababcbab

– patt=abcbac

不一致

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
}
return 0;
```

素朴な方法: 逐次探索

- 例:



– text=**a**babcababacbab

– patt=ab**c**abac



不一致

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
}
return 0;
```

素朴な方法: 逐次探索

- 例:

– text=ab[↓]abcababcabacbab

– patt=abcabac[↑]

不一致

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
}
return 0;
```

素朴な方法: 逐次探索

- 例:

– text=ababcbababcbab

– patt=abcbac

不一致

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
}
return 0;
```

素朴な方法: 逐次探索

- 例:

– text=ababcababcbabab

– patt=abcbac

不一致

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
}
return 0;
```


素朴な方法: 逐次探索

- 例:

– text=ababcababcbabab

– patt=abcbabac

不一致

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
}
return 0;
```

素朴な方法: 逐次探索

- 例:

– text=abab**ca**babcabacbab

– patt=ab**ca**bac

不一致

```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
}
return 0;
```

素朴な方法: 逐次探索

- 例:

– text=**ababcab**abcbacbab

– patt=abcbac



一致発見

```
for(s=1; s<=n-m+1; s++){  
    for(i=1; i<=m; i++){  
        if(t[s+i-1]!=p[i])  
            break;  
        if(i == m+1)  
            return s;  
    }  
    return 0;  
}
```

素朴な方法の計算時間

- 最悪の場合:

for $1 \leq s \leq n-m+1$,

$t[s, s+m-2] = p[1, m-1] \wedge t[s+m-1] \neq p[m]$

- 例:

– text = aaaaaaaaaa

– patt = aaab

- 計算量: $O(mn)$

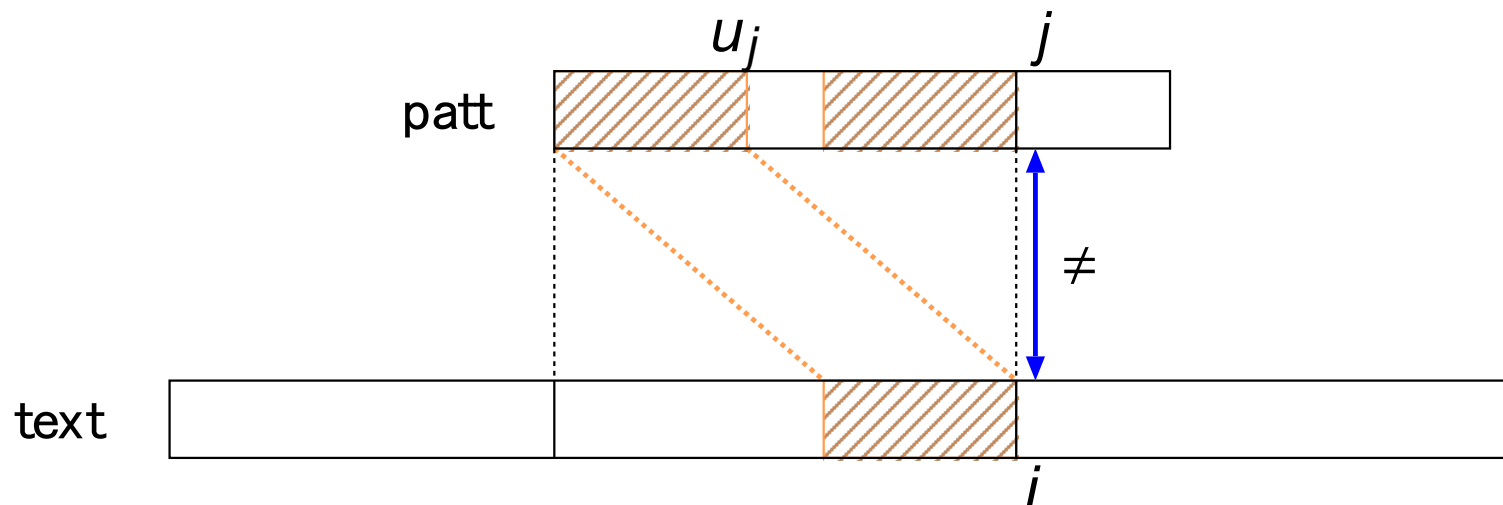
```
for(s=1; s<=n-m+1; s++){
    for(i=1; i<=m; i++){
        if(t[s+i-1]!=p[i])
            break;
        if(i == m+1)
            return s;
    }
    return 0;
}
```

素朴な方法が遅い理由

- テキストとパターンの不一致が起こるたびに、テキストを指すポインタが前に戻っている
- Q: 戻らないようにできるか？
- A: できる
 - $t[i]$ で不一致が起きた場合は、次の照合は $t[i]$ から始める（戻らない！）
 - $p[]$ の添字はうまくずらす

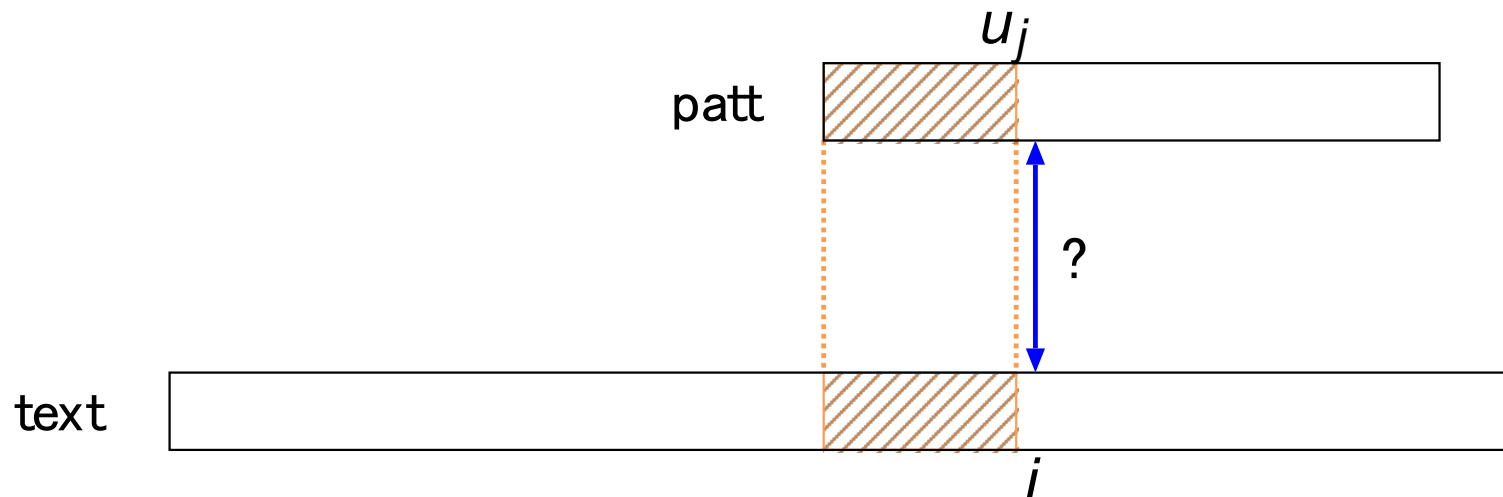
素朴な方法の改善

- $t[i-j+1, i-1] = p[1, j-1]$ かつ $t[i] \neq p[j]$ とする
- 次に比べるのは $t[i]$ と $p[u_j+1]$ から
 - u_j : $p[1, k] = p[j-k, j-1]$ となる最大の $k < j-1$
 - 1文字も一致しなければ 0
- $p[u_j+1]$: $t[i]$ とマッチする可能性のある文字で、最も右にあるもの



素朴な方法の改善

- $t[i-j+1, i-1] = p[1, j-1]$ かつ $t[i] \neq p[j]$ とする
- 次に比べるのは $t[i]$ と $p[u_j+1]$ から
 - $u_j: p[1, k] = p[j-k, j-1]$ となる最大の $k < j-1$
 - 1文字も一致しなければ 0
- $p[u_j+1]$: $t[i]$ とマッチする可能性のある文字で、最も右にあるもの





KMP (Knuth-Morris-Pratt) 法

- u_j の値が前もって計算されているとき
 - text[]: 1..n, patt[]: 1..m

```
i=1; j=1;
while (j<=m && i<=n)
  if(j==0 || text[i]==patt[j])
    i++; j++;
  else
    j=u[j];
if(j==m+1)
  パターンをテキストのi-m番目からの部分に見;
```


KMP法: u_j の計算

- u_j の計算でもパターン照合を行う
- patt の j 文字目まで調べたとき
 - u_k の値が $k = 1, \dots, j$ まで計算されていると仮定
 - 照合に必要なのは $k \leq j$ についての u_k

```
k=0; j=1; u[1]=0;
while(j < m){
  if(k==0 || patt[k]==patt[j])
    k++; j++; u[j] = k;
  else
    k = u[k];
}
```

ほぼ一緒

```
i=1; j=1;
while (j<=m && i<=n)
  if(j==0 || text[i]==patt[j])
    i++; j++;
  else
    j=u[j];
```

KMP法

KMP法: 計算時間

```
i=1; j=1;
while (j<=m && i<=n)
  if(j==0 || t[i]==p[j])
    i++; j++;
  else
    j=u[j];
```

- j の減少回数 $<$ j の増加回数
– 0でないときのみ減少する
- j の増加回数 = i の増加回数
- i の増加回数: n

- 全体の計算時間 $< 3n \in O(n)$

KMP法:

u_j の計算込みの全体の計算時間

- $u[]$ の構築はKMP法と変わらない
→ $O(m)$
 - m は $patt$ の長さ
- よって $u[]$ を構築して文字列照合を行うと
 $O(n + m)$
 - n は $text$ の長さ

≡二演習

- 以下の例に対してKMP法を実行せよ
 - text = aaaaaaaaaa
 - patt = aaab