

# 形式言語とオートマトン

Note6 オートマトンの応用例

2020.5.12作成

2020.6.23update 7.6 2021.6.8 6.27 6.29

中野眞一 群馬大学

# 文字列のマッチング

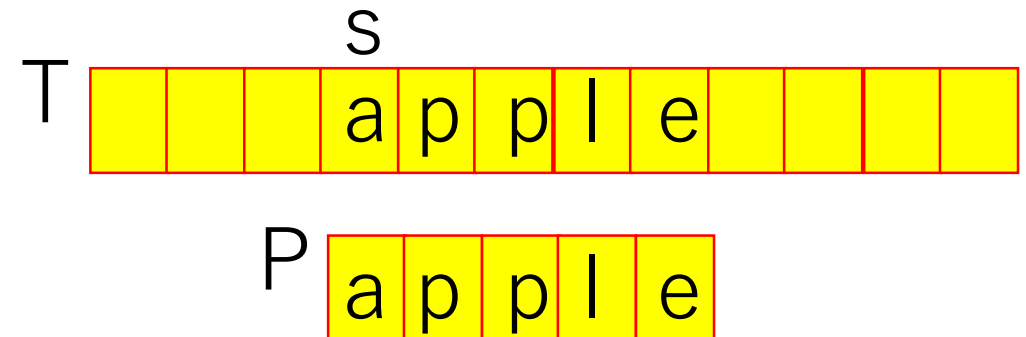
膨大な長さの文字列の中から、指定された文字列を高速に探索する問題を考えよう。  
(エディタや、DNA文字列の処理プログラム等で必要とされる。)

この問題を次のように定式化する。

## 文字列マッチング問題

膨大な文字列をテキストTと呼び、文字の配列  $T[0..n-1]$  に格納されているとする。  
探索したい文字列をパターンPと呼び、文字の配列  $P[0..m-1]$  に格納されているとする。  
 $m \leq n$  とする。TもPも、アルファベット  $\Sigma$  上の語(文字列)であるとする。

もし、各  $j=0,1,2,\dots,m-1$  について、  
 $T[s+j] = P[j]$  であるならば、  
パターンPは、位置sに現れたという。  
(matchするという。)



# 簡単なアルゴリズム (カブク法)

Brute-Force-String-Matcher(T,P)

for  $s = 0$  to  $n-m$

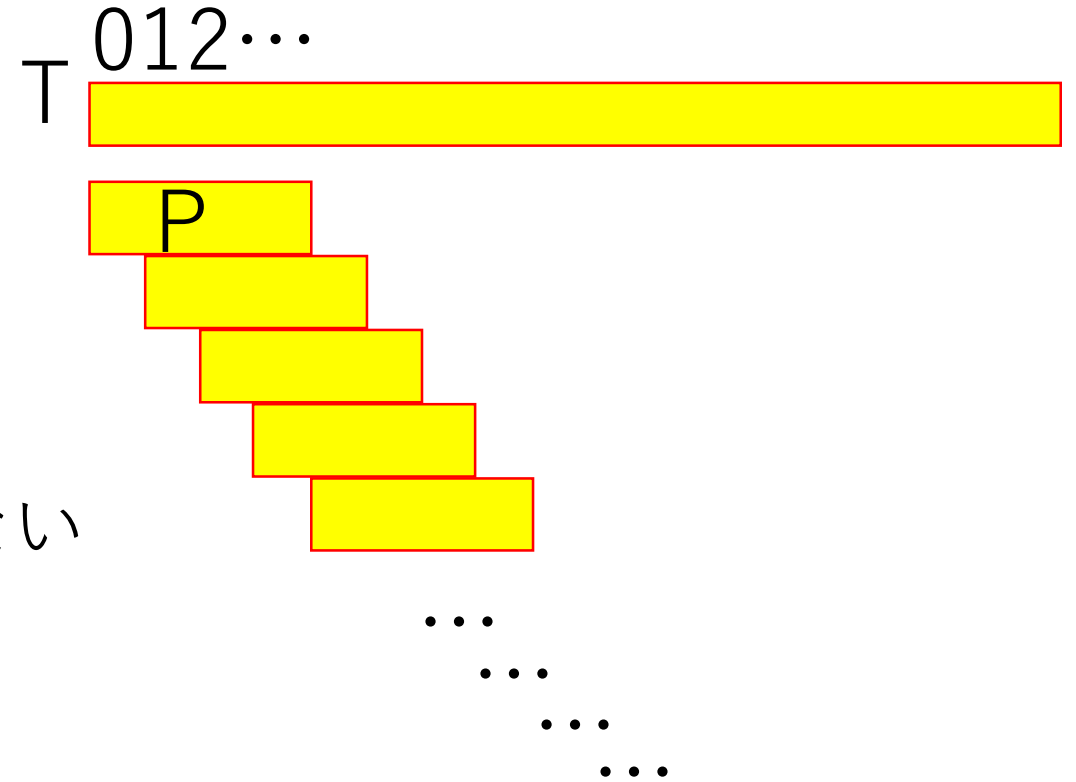
for  $i = 0$  to  $m-1$

if  $T[s+i]$  が  $P[i]$  と等しくない

then break

if  $i = m$

then  $s$  を出力



計算時間は  $O((n-m+1)m)$

# Rabin-Karp法

剰余(あまり)を上手に利用する方法である。

各文字に、 $0, 1, 2, \dots, |\Sigma| - 1$  の数字を割り当てる。

パターン $P$  を $m$ 桁の $|\Sigma|$ 進数とみなす。この値を $p$  とする。

テキスト $T$  の $s$ 文字目から連続する $m$ 文字を  $T_s$  とし、これを $m$ 桁の $|\Sigma|$ 進数とみなし、値を $t_s$  とする。

$T_s = P$  ならば、 $T$ 中の位置 $s$ に、 $P$ が出現したことになる。

このとき、 $t_s = p$  となっているはずである。

(問題1) しかし、 $t_s$  や  $p$  の値を求めるのに時間がかかる。

簡単な方法では、 $t_s$  を1つ計算するのに、 $O(m)$ 時間かかる。

例  $|\Sigma| = 10$  のとき、

$$p = P[m-1] + 10( P[m-2] + 10( P[m-3] + \dots + 10( P[1] + 10 P[0]) \dots )) \text{ である。}$$

(問題2) また、 $t_s$  や  $p$  の値は、非常に大きくなるかもしれない。(long型でも不足するかも。)

# Rabin-Karp法 (つづき)

問題1を解決しよう。各  $t_s$  を独立に計算することをやめ、

$t_s$  から、 $t_{s+1}$  を求めることにより、高速化しよう。

例  $|\Sigma| = 10$  のとき、

$$t_s = T[s+m-1] + 10(T[s+m-2] + 10(T[s+m-3] + \dots + 10(T[s+1] + 10 T[s]) \dots)) \text{ である。}$$

$$t_{s+1} = T[s+m] + 10(T[s+m-1] + 10(T[s+m-2] + 10(T[s+m-3] + \dots + 10(T[s+1]) \dots)) \text{ である。}$$

これより、 $t_{s+1} = 10(t_s - 10^{m-1}T[s]) + T[s+m]$  となる。この式は  $O(1)$  時間で計算できる。

(ただし、もし、 $10^{m-1}$  の値が計算してあれば。)

例えば、 $T = 3141592$ ,  $m = 5$  のとき、

$t_0 = 31415$  であり、

$t_1 = 10(31415 - 10^4 \cdot 3) + 9 = 14159$  である。

上記により、最初に、 $O(m)$  時間で、 $p$  と  $t_0$  を求める。

あとは、1個あたり  $O(1)$  時間で、 $t_1, t_2, \dots, t_{n-m}$  をそれぞれ求めることができる。

アルゴリズム全体で、 $O(m) + O(n-m+1) = O(n-m+1)$  時間となる。

# Rabin-Karp法 (つづき)

問題2を解決しよう。

適切に数 $q$ を選び、 $q$ の剰余のみの計算で、上記をすべて実行すればよい。  
 $q$ は、 $|\Sigma|q$ が、計算機の1 word におさまるような素数とするのが一般的である。

ただし、下記の注意が必要である。

もし、 $t_s \neq p \pmod{q}$  ならば、 $P$ は $T$ の位置 $s$ に出現しない。

しかし、 $t_s = p \pmod{q}$  ならば、 $P$ は $T$ の位置 $s$ に出現するとはいえない。

(まだどちらか不明です。)

(値ではなく、剰余で判定しているので。)

しかし、 $q$ を十分大きく選べば、 $t_s = p \pmod{q}$  なのに、  
 $P$ は $T$ の位置 $s$ に出現しないという場合は、ほとんどおこらないとみなせる。

# Rabin-Karp法 (つづき)

RK-String-Matcher( $T, P, d, q$ )    {  $d = |\Sigma|$  とする }

$h = d^{m-1} \bmod q$

$p = 0$

$t_0 = 0$

for  $i = 0$  to  $m-1$  {前処理}

$p = (d \cdot p + P[i]) \bmod q$

$t_0 = (d \cdot t_0 + T[i]) \bmod q$

for  $s = 0$  to  $n-m$

    if  $p = t_s$

        then {本当にmatchかどうか確認する}

            if  $P[0..m-1] = T[s..s+m-1]$  {ループが必要}

                then  $s$ を出力

    if  $s < n-m$

        then

$t_{s+1} = (d(t_s - T[s] h) + T[s+m]) \bmod q$

前処理時間は  $O(m)$  時間である

他の計算時間は  $O((n-m+1)m)$  時間である

最悪の場合は力づく法と同じ計算時間であるが、平均的・実用的には、こちらの方が早いことが多い。

# 有限オートマトン法

文字列x

appl .....

パターンP

appliancejuice

1234

有限オートマトンに基づくアルゴリズムを与える。

この方法は、各文字を、ちょうど**1回**づつしかチェックしない。

パタンの**前処理**のほかに、かかる時間は、わずかに  **$O(n)$ 時間**である。

(しかし、 $\Sigma$ が大きい時には、前処理にかかる時間も大きくなるので注意。)

各パターンPが与えられたとき、このパタンの出現を見つけるオートマトンを設計しよう。

パターンPが与えられたとき、suffix関数  $\sigma$  (シグマ) を次のように定義する。

$\sigma$  は、 $\Sigma^*$ から  $\{0,1,2,\dots,m\}$  への関数である。xは  $\Sigma$  上の語とする。

$\sigma(x) = \max\{k \mid P_k \text{ は } x \text{ の接尾語}\}$  と定義する。

語xの**最後**の方に、パタンの**最大**何文字目までが現れているかを表している。

ここで、 $P_k$ は、Pのk文字までの接頭語とする。

(**文字列xの最後とパターンPの最初を、最大、どのくらい重ねられるか。**)



# 有限オートマトン法(つづき)

オートマトンを設計する。

状態  $Q = \{0, 1, \dots, m\}$  とする。(状態  $i$  は、パタンの最初の  $i$  文字まで発見! の意味です.)

アルファベット  $\Sigma$  は文字の集合とする。

状態遷移関数を次のように定義する。  $\delta(q, a) = \sigma(P_q a)$

(右辺のかっこの中は、パターン  $P$  の  $q$  文字目までに文字  $a$  を追加した文字列である。)

初期状態  $q_0 = 0$  とする。

受理状態  $F = \{m\}$  とする。

このオートマトンは次の性質を持つ。

テキスト  $T$  を  $i$  文字目までチェックしたとき、状態が  $q$  であるとしよう。

このとき、 $q = \sigma(T_i)$  となる。

ここで、 $T_i$  は、 $T$  の  $i$  文字目までの接頭語とする。

(つまり、 $T$  の  $i$  文字目までチェックしたとき、 $P$  の  $q$  文字目までが出現している!)

もし、次の文字  $T[i+1]$  が、 $P$  の  $q+1$  文字目と同じならば、状態は  $q+1$  へ

もし、次の文字  $T[i+1]$  (これを  $a$  とする) が、 $P$  の  $q+1$  文字目と違うならば、状態は  $\sigma(P_q a)$  へ

(つまり、 $P$  の出現は無かったが、“次の”チャンス計算しようとしているのです。

次のチャンスまではスキップ! できる!)

(状態  $\sigma(P_q a)$  の計算に  $P$  は関係あるが、 $T$  は関係ないことに注意しよう。)

12345文字目

パターン apadd

パターン4文字目は違った

文字列x

apap

パターン3文字目まで発見!  
だけど4文字目はなかった

apad..

けどパターン2文字目  
までは発見

ap....

# 有限オートマトン法(つづき)

FA-String-Matcher(T,P)

{前処理で オートマトンを計算する}

q = 0

for i = 0 to n-1

q =  $\delta(q, T[i])$  /\* オートマトンを動かしているだけ\*/

if q = m

then i-m を出力

前処理時間は  $O(m |\Sigma|)$

matching 時間は  $O(n)$

# 理解確認クイズ

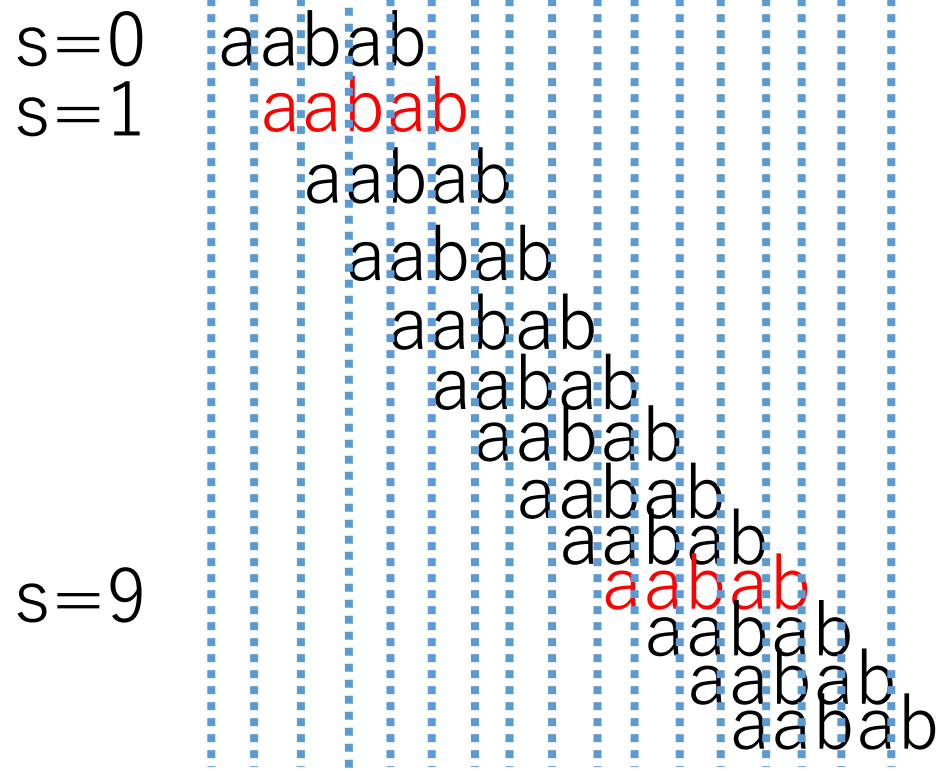
(1)

$P=abab$ ,  $T=aaababaabaabaab$ のとき,  
カづく法のアルゴリズムを説明せよ。

P=aabab, T=aaababaabaabaabのとき,  
カブク法のアルゴリズムを説明せよ。

01234567890123456

T=aaababaabaabaabのとき,



S=0,1,2,3,...の位置に  
Pがあるかをすべて調べる。  
S=1およびS=0の位置に  
Pがあるのでこれを出力する

# 理解確認クイズ

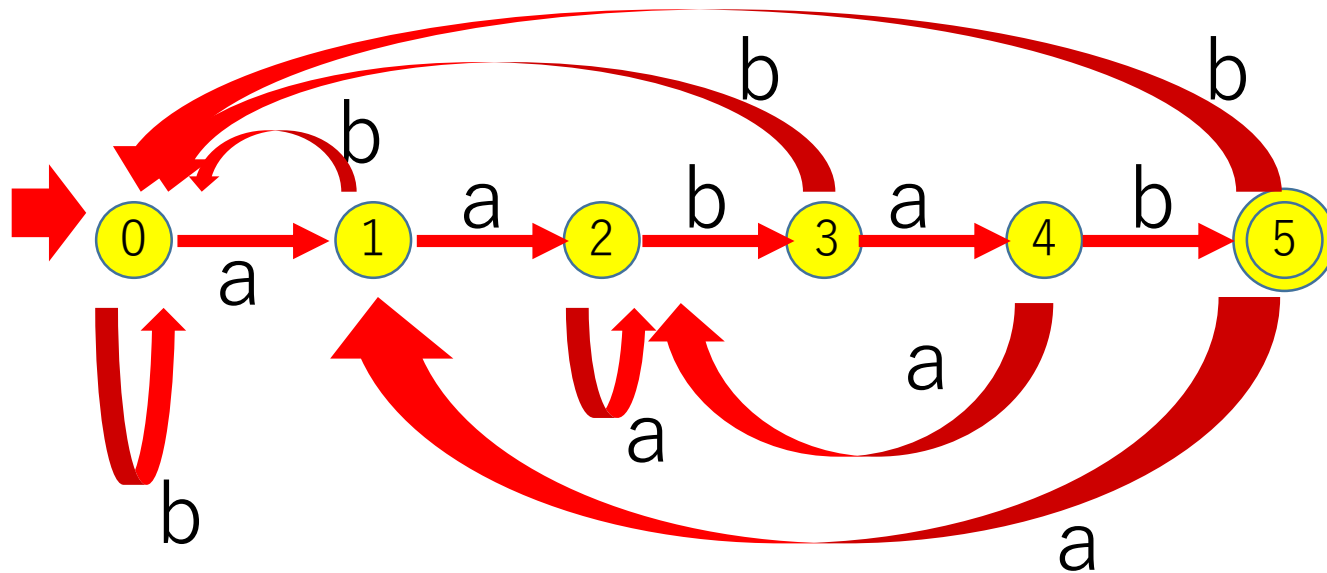
(2)

$P=aabab$ ,  $T=aaababaabaabaab$ のとき,  
有限オートマトン法のアルゴリズムを説明せよ。

# 理解認識クイズ

T    aababaabaabaab  
状態 012234**5**1234234**5**123

(2)  $P=aabab$ ,  $T=aaababaabaabaab$  のとき,  
有限オートマトン法のアルゴリズムを説明せよ。



状態2からの遷移

aaa 状態2へ

aab 状態3へ

# 理解確認クイズ

(3)

$P=ababaca$ ,  $T=abababacaba$ のとき,  
有限オートマトン法のアルゴリズムを説明せよ。

オートマトン設計してみましよう