

形式言語とオートマトン

Note3 正規表現

2020.4.03作成

2020.4.27update 5.18update 5.23update 5.25 update

5.30 update

2021.4.26update 5.11 5.23 5.24

中野眞一 群馬大学

正規表現

正規表現は、語の集合の代数的な表現のひとつである。
(代数とは、集合とその上の演算のことである。)

アルファベット Σ 上の正規表現とは、 Σ 中の記号、および、6つの記号

"(" と ")" と "*" と "+" (もしくは "|") と " ϵ " と " ϕ " からなる有限の長さの語のうち、

次のように再帰的に定義されるものである。

正規表現 α が表す語の集合を $\|\alpha\|$ と書く。

正規表現(つづき)

(1) Σ の要素である記号は、 Σ 上の正規表現である。(長さ1のもの)

$$a \in \Sigma \text{ のとき } \|a\| = \{a\}$$

(2) ε と ϕ はそれぞれ Σ 上の正規表現である。(長さ0の文字列と空集合)

$$\|\varepsilon\| = \{\varepsilon\}$$

$$\|\phi\| = \phi$$

(3) (長さ2以上のもの)

α と β をそれぞれ Σ 上の正規表現とする。

(a) $(\alpha + \beta)$ は Σ 上の正規表現である。(和)

$$\|(\alpha + \beta)\| := \|\alpha\| \cup \|\beta\|$$

(b) $(\alpha \beta)$ は Σ 上の正規表現である。(接続)

$$\|(\alpha \beta)\| := \|\alpha\| \cdot \|\beta\| \quad (\text{言語の接続}) \quad (\text{記号} \cdot \text{は略してもよい})$$

(c) (α^*) は Σ 上の正規表現である。(閉包)

$$\|(\alpha^*)\| := \|\alpha\|^* \quad (\text{言語} L \text{の閉包})$$

(d) $(\sim \alpha)$ は Σ 上の正規表現である。

$$\|(\sim \alpha)\| := \Sigma^* - \|\alpha\|$$

正規表現 (例)

例) $\Sigma = \{ 1, 2, 3 \}$

2 は Σ 上の正規表現である。 $\| 2 \| = \{ 2 \}$ である。

ε は Σ 上の正規表現である。 $\| \varepsilon \| = \{ \varepsilon \}$ である。

ϕ は Σ 上の正規表現である。 $\| \phi \| = \phi$ である。

$(1+2)$ は Σ 上の正規表現である。 $\| (1+2) \| = \| 1 \| \cup \| 2 \| = \{ 1, 2 \}$ である。

$(2+3)$ は Σ 上の正規表現である。 $\| (2+3) \| = \{ 2, 3 \}$ である。

$((1+2)3)$ は Σ 上の正規表現である。 $\| ((1+2)3) \| = \{ 1, 2 \} \cdot \{ 3 \} = \{ 13, 23 \}$ である。

$((1+2)(2+3))$ は Σ 上の正規表現である。 $\| ((1+2)(2+3)) \| = \{ 12, 13, 22, 23 \}$

(2^*) は Σ 上の正規表現である。 $\| (2^*) \| = \{ \varepsilon, 2, 22, 222, \dots \}$

$((1+2)^*)$ は Σ 上の正規表現である。

$\| ((1+2)^*) \| = \{ \varepsilon, 1, 2, 11, 12, 21, 22, 111, 112, \dots \}$

正規表現

ただし、冗長なカッコは削除してもよい。
(和、接続、閉包の順に優先順位が強くなるとする。)

例1) $a + bc^*d$ は
($a + ((bc^*)d)$)の意味です。

例2) $(hot + cold)(apple + blueberry)(pie + tart)$
が表す語の集合は何かな？要素はいくつある？

正規表現

正規表現によって表すことのできる言語を正規言語 (もしくは正規集合 正則言語) という。

正規言語は、3つの演算、(1)和、(2)連接、(3)閉包について **閉じている**。すなわち、任意の正規言語に、3つの演算をしても、正規言語のままである。

正規表現の性質

α と β をそれぞれ Σ 上の正規表現とする。次が成り立つ。
(両辺とも言語を表すことに注意しよう。)(人間がわかりやすく変形できる!)

$$(1) \alpha + \beta = \beta + \alpha$$

$$(2) \alpha + (\beta + \gamma) = (\alpha + \beta) + \gamma$$

$$(3) \alpha + \phi = \alpha$$

$$(4) \alpha + \alpha = \alpha$$

$$(5) \alpha^* \alpha^* = \alpha^*$$

$$(6) (\alpha^*)^* = \alpha^*$$

$$(7) (\varepsilon + \alpha) = (\alpha + \varepsilon) \neq \alpha \text{ (最後の不等号が成立するときは、どんなときかな?)}$$

$$(8) \beta \beta^* \neq \beta^*$$

$$(9) (\alpha^* + \beta^*)^* = (\alpha^* \beta^*)^* = (\alpha + \beta)^*$$

正規表現はいろいろなシステムで利用できます

UNIXのawk, grep, lex等のコマンドで正規表現が使用できる。

Python言語には正規表現モジュールが用意されている。

Ruby言語には正規表現クラスが用意されている。

Javascriptでも正規表現を扱うしくみがある。

UNIXにおける正規表現の拡張

UNIXでの拡張表現の例。

[aghinostw]のように "["と"]"で囲むと[]内の文字の集合を表す。

[a-z]で範囲内のどれか1文字の集合を表す。

^で行のはじめを示し、\$で行の終りを示す。

. は任意の1文字を示す。(wild card)(don't care)

aの0回以上の繰り返しをa*で示す。

1回以上の繰り返しを+で示す。 a+ = aa* である。

aの0回もしくは1回の出現を a?で示す。

¥* 特殊文字"*"を示す。

pythonでの拡張の例

{m} 直前の正規表現をちょうど m 回繰り返したもの 例え、 a{6} は 6 個ちょうどの 'a' 文字

{m,n}? 前にある正規表現を、 m 回から n 回まで繰り返したものにマッチし、できるだけ少なく繰り返したもの

その他たくさんあります。。。

/usr/share/dict/words でいろいろな語を検索してみよう。

```
%egrep man$ /usr/share/dict/words (2291語)
```

```
%egrep woman$ /usr/share/dict/words (176語)
```

```
%egrep ".*man$|.*son$" /usr/share/dict/words (7932語)
```

```
%egrep "^a.*a$" /usr/share/dict/words (1433語)
```

正規表現とFAの等価性

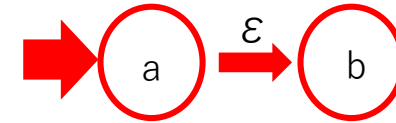
正規表現とNFAは互いに変換できることを示す。

正規表現 α が与えられたとき、 $\|\alpha\| = L(M)$ となるオートマトン M が構成できる。

オートマトン M が与えられたとき、 $L(M) = \|\alpha\|$ となる正規表現 α が構成できる。

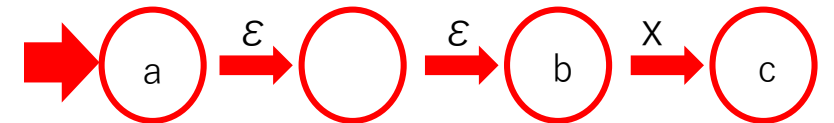
正規表現 = > NFA

まず、NFAに、 ϵ 遷移を導入する。
状態aから状態bへ入力がなくとも遷移できるとする。



ϵ 遷移付きのNFA Mは、次のようにして、(ϵ 遷移を用いない)通常のNFA Nに変換できる。

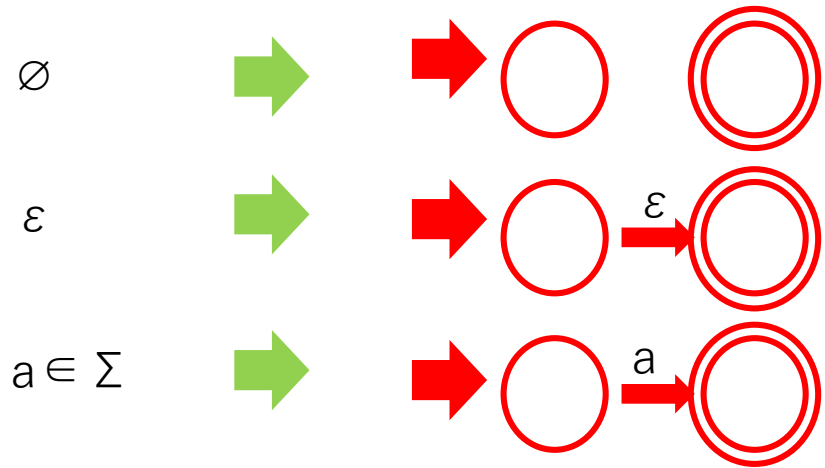
Mにおいて、状態aから ϵ 遷移のみを用いて状態bに遷移でき、
かつ、状態bから通常の遷移1回で(入力xとする)状態cに遷移できるとき、
Nにおいて、 $\delta(a,x)=c$ とする。



つまり、正規表現を、 ϵ 遷移付きのNFAに変換できれば、
(ϵ 遷移付きのNFAは、通常のNFAに変換できるので)
正規表現を通常のNFAに変換できることがわかる。

正規表現 = \Rightarrow NFA

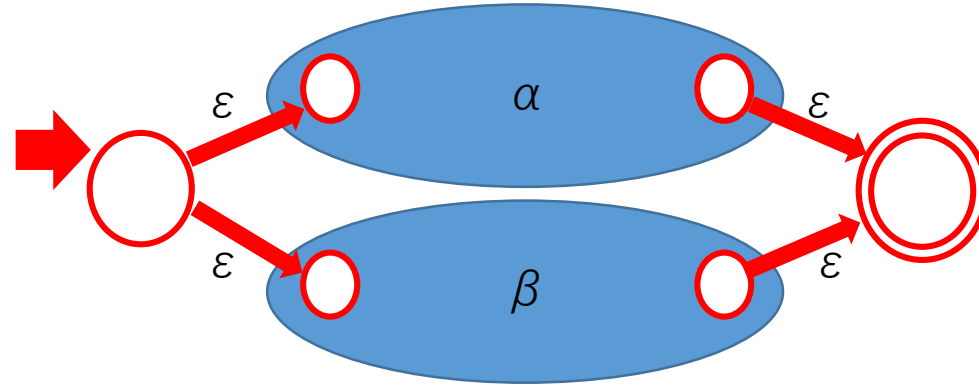
正規表現の定義に基づき、構成的にNFAを作成する方法を説明する。
定義のそれぞれについて構成方法を示す。



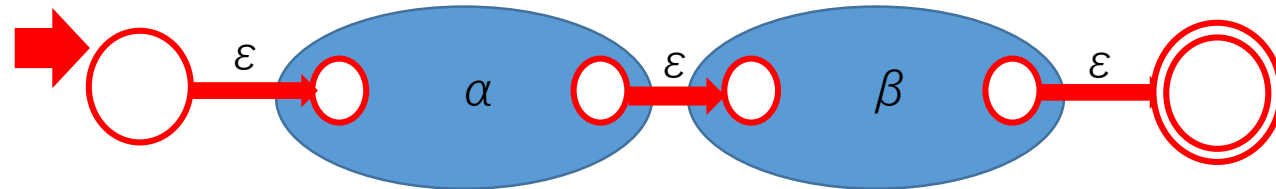
正規表現 \Rightarrow NFA (つづき)

(最後にまとめて、 ϵ 遷移を解消する。)

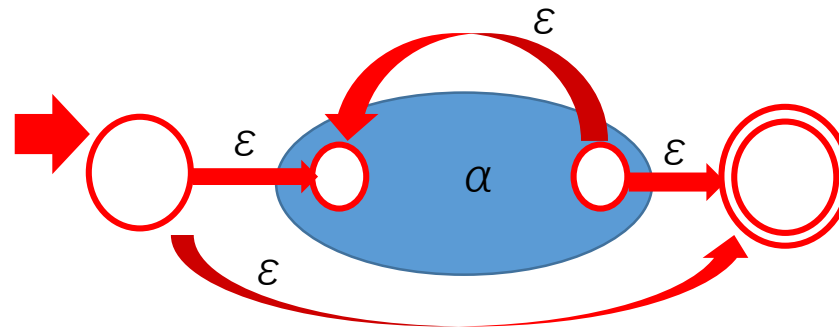
$\alpha + \beta$



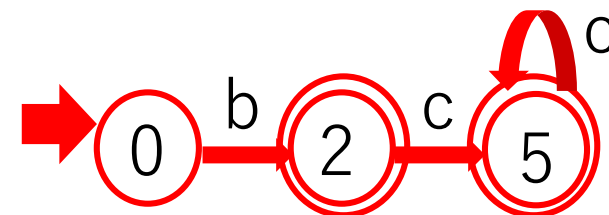
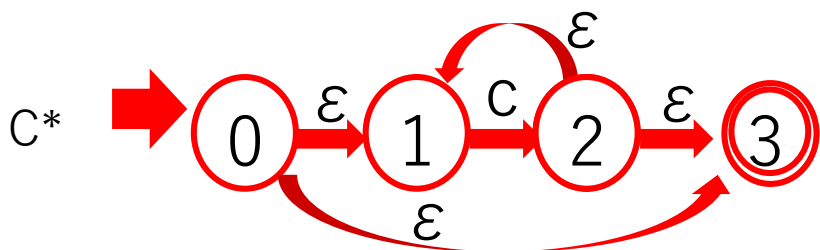
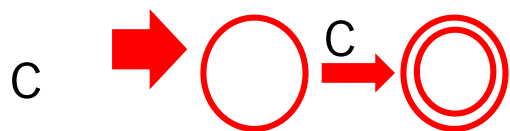
$\alpha \beta$



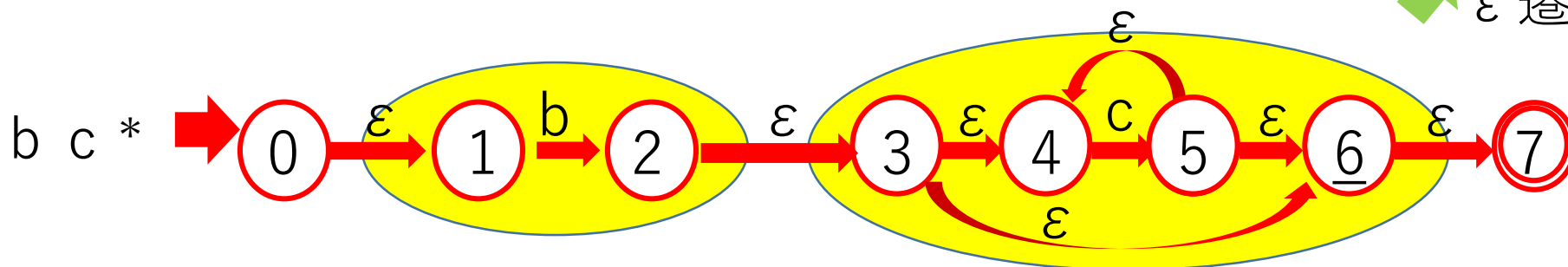
α^*



bc* に対応するオートマトン



ε 遷移を削除



NFA => 正規表現

NFAの初期状態点以外の状態点を1つずつ消去する。

その際にNFAの辺のラベルを、次のように、順次、正規表現に置き換える。

点の個数が1個や2個になれば、正規表現がわかる。(ゴールが多数あるときは、ゴールごとに正規表現を考える。)

(base) 各辺のラベルを、長さ1の正規表現とみなす。

(induction) 点uを消すとする。

uに向かう辺を、 (s_1, u) 、 (s_2, u) 、...、 (s_n, u) とする。

ラベルを、それぞれ、 S_1 、 S_2 、...、 S_n とする。

uから出ていく辺を、 (u, t_1) 、 (u, t_2) 、...、 (u, t_m) とする。

ラベルを、それぞれ、 T_1 、 T_2 、...、 T_n とする。

uからuへの辺(自己ループ辺)のラベルをUとする。

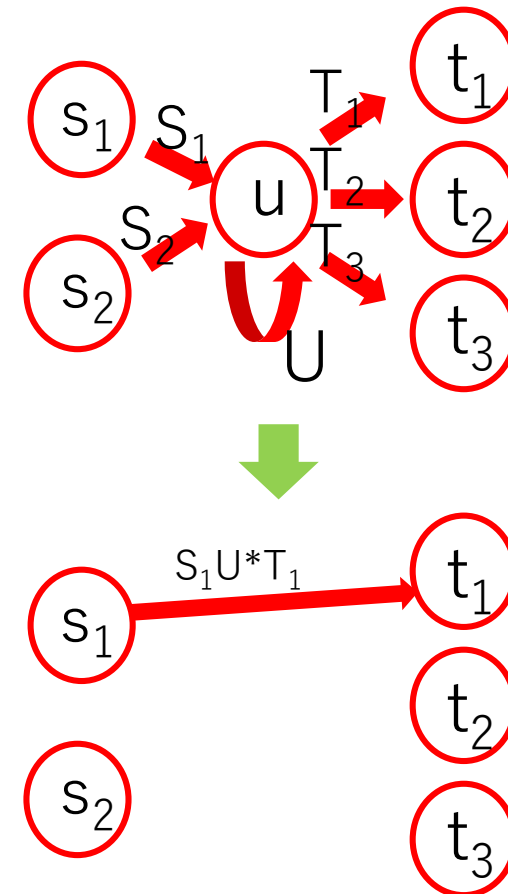
(もし、uからuへの辺がなければ、そのラベルを ϕ とする。)

各辺 (s_i, t_j) の"現在の"ラベルを R_{ij} としよう。

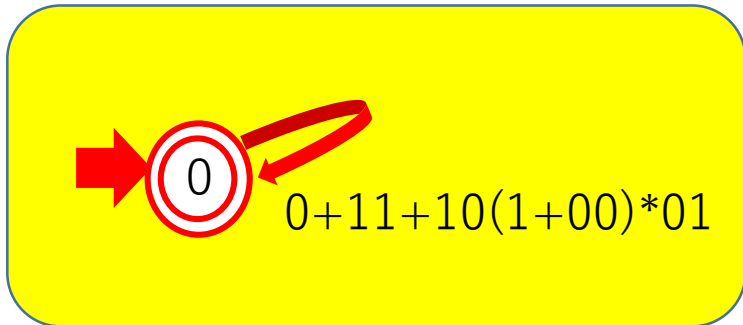
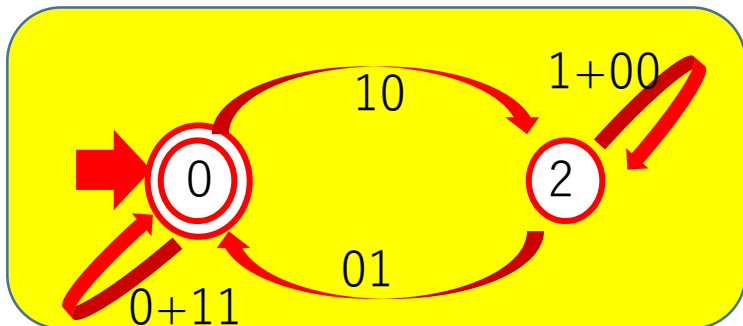
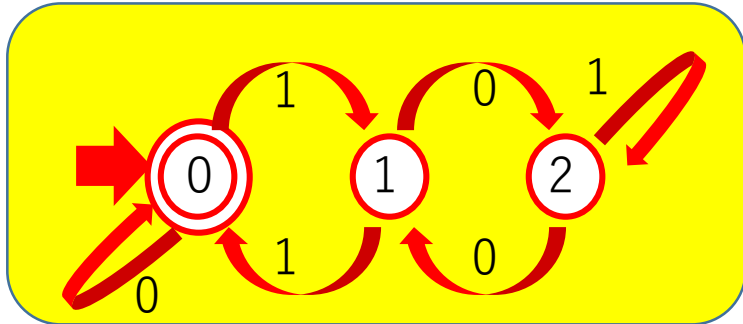
各辺 (s_i, t_j) のラベルを $R_{ij} + S_i U^* T_j$ に変更する。

(もし、 $U = \phi$ のときには $R_{ij} + S_i T_j$ に変更するのと同じである。)

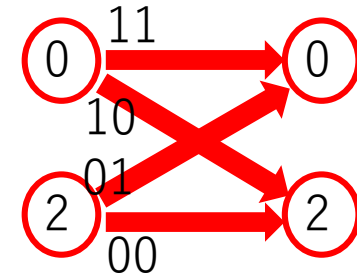
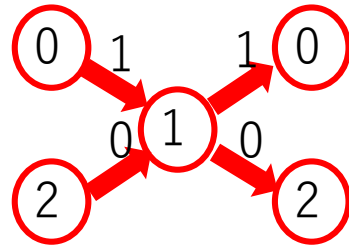
点uを削除しても、同じすごろくのまま。(ゴールするパスの集合は同じ!)



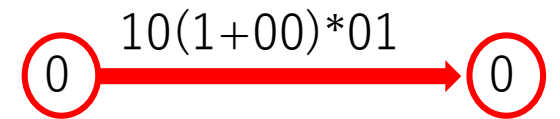
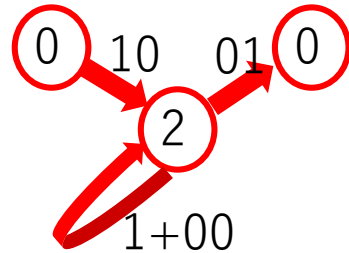
NFA => 正規表現 の例1



状態1を消す



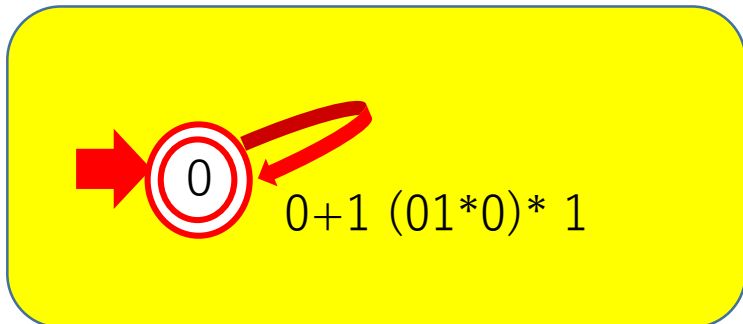
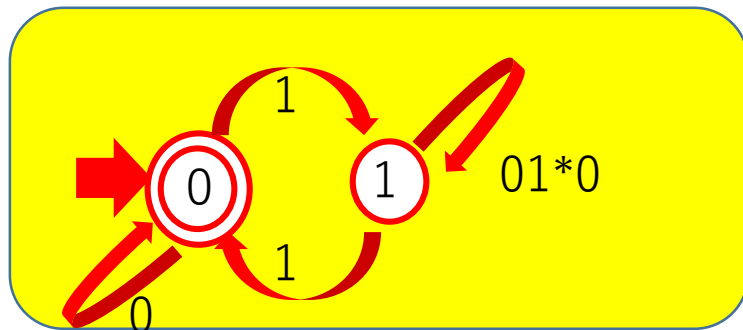
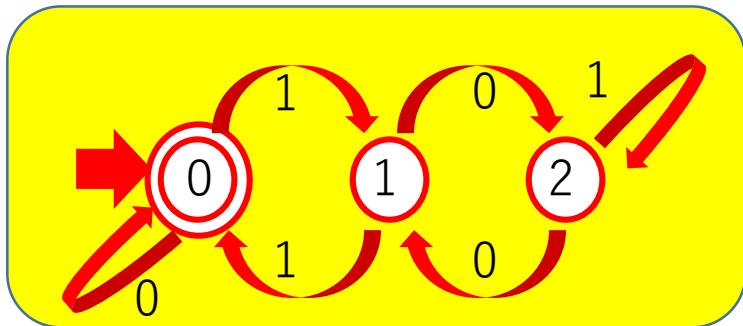
状態2を消す



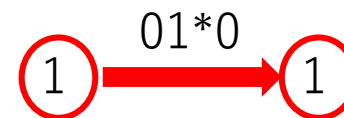
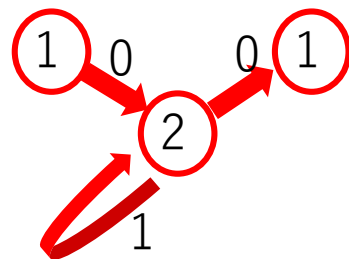
0, 11, 110, 1001, 1111,
10010, 10101, 11000, ...

正規表現は $(0 + 11 + 10(1+00)^*01)^*$

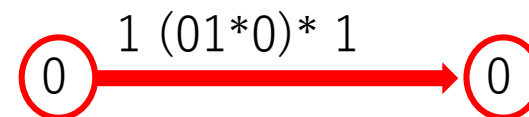
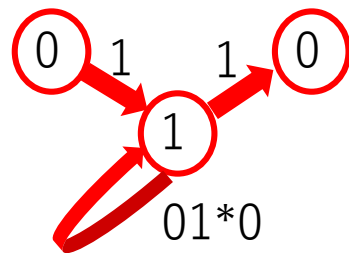
NFA => 正規表現 の例2



状態2を消す



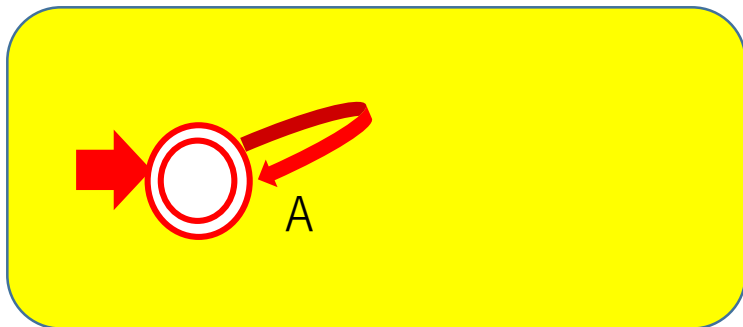
状態1を消す



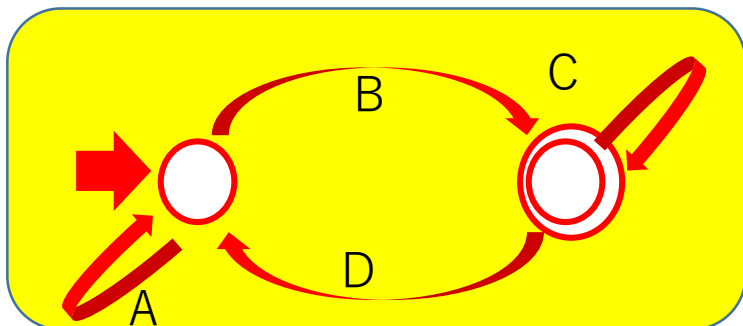
0, 11, 110, 1001, 1111,
10010, 10101, 11000, ...

正規表現は $(0 + 1(01^*0)^*1)^*$

NFA => 正規表現 最後の变换



正規表現は $(A)^*$



正規表現は $A^* B (C + DA^*B)^*$

初ゴール

再ゴール

ポンプ定理

正規言語は必ず次の条件を満たす。(ワンパターンである！)

(ポンプ定理)

L が正規言語(正規集合)ならば、次の条件を満たす正整数 n が存在する。

任意の $|x| \geq n$ なる $x \in L$ は $x = uvw$ と書け、任意の整数 $i \geq 0$ について

$uv^i w \in L$ である。

ただし、 $|uv| \leq n$ かつ $|v| \geq 1$ である。

ポンプ定理の意味

Pumping lemma for regular languages, the fact that all sufficiently **long** strings in such a language have a substring that can be **repeated arbitrarily many times**, usually used to prove that certain languages are not reg

(https://en.wikipedia.org/wiki/Pumping_lemma)

ポンプの定理 (証明)

任意の正規言語 L に対して、 $L(M) = L$ となるオートマトン M が存在する。この M の状態数を n とする。

長さが n より大きな任意の語 $x \in L$ を選ぼう。 $x = a_1a_2\dots a_m$ とする。

$n < m$ なので、この入力に対するパスには、2回以上現れる状態がある。

このように2回以上現れる状態のうち、2回目が最も早く現れる状態を q としよう。 s 文字目および t 文字目を読んだときに状態 q とする

すなわち、 $a_1a_2\dots a_s$ まで入力を読んだとき状態 q となり、その後 $a_{s+1}a_{s+2}\dots a_t$ まで入力を読んだとき、再び状態 q となるでしょう。

q の選び方より、 a_1, a_2, \dots, a_{t-1} はすべて異なる。

$$u = a_1a_2\dots a_s$$

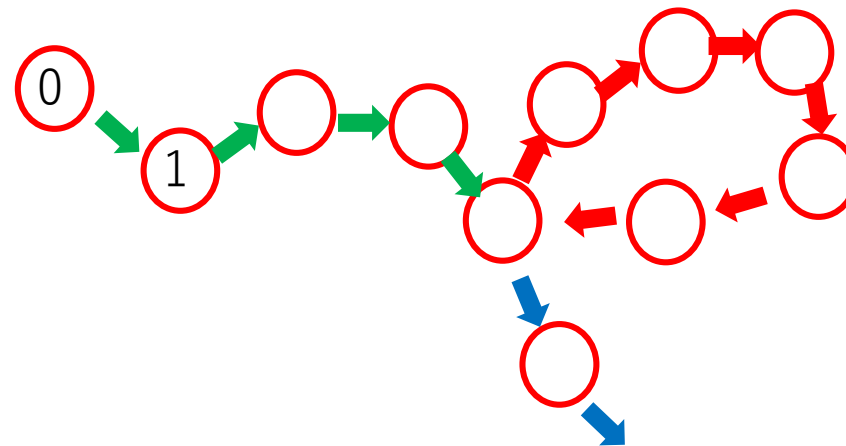
$$v = a_{s+1}a_{s+2}\dots a_t$$

$$w = a_{t+1}a_{t+2}\dots a_m \quad \text{とする。}$$

(オートマトンの略図を書いてみよう。)

このとき、任意の整数 $i \geq 0$ について $uv^iw \in L$ である。

また、 $|uv| \leq n$ かつ $|v| \geq 1$ である。(証明終)



正規表現の限界

言語 $L = \{ a^k b^k \mid k \geq 0 \} = \{ \varepsilon, ab, aabb, aaabbb, \dots \}$

は正規言語ではない。

(証明)

もし L が正規言語だとすると、ポンプの定理が成立するような整数 n がある。

語 $x = a^{n+1}b^{n+1}$ とすると $|x| = 2n+2 > n$ なので、ポンプの定理より

$x = uvw$ と書け、(ただし、 $|uv| \leq n$ かつ $|v| \geq 1$ である)

任意の整数 $i \geq 0$ について、 $uv^i w \in L$ である。

$|uv| \leq n$ より uv は a のみからなる語である。

また、 w は、1個以上の a と、 $n+1$ 個の b を含む。

このとき、 $y = uv^2w \in L$ とは **ならない**。

なぜならば、 y に含まれる b の個数は x と同じであるが、 a の個数は増えるからである。

よって **矛盾** する。(証明終)

正規表現の限界

対応のとれたカッコからなる言語も、
正規表現では記述できない。

(これまで、いくつ余分に左カッコが現れたかを、
有限の状態では記憶できないので。)

(もし、スタックが使えるれば、記憶できるのですが。。。)

理解確認クイズ

- (1) $\Sigma = \{a, b, c\}$ とする。 Σ 上の次の言語を正規表現で表せ。
- (a) 少なくとも3個のaを含む語からなる言語
 - (b) 高々3個のaを含む語からなる言語
 - (c) 最後から3文字目はaである語からなる言語
 - (d) 最後から3文字目はaでない語からなる言語
 - (e) 長さが5の語からなる言語
 - (f) aの次の文字は必ずbである語からなる言語
- (2) $\Sigma = \{0, 1\}$ とする。 Σ 上の次の言語を正規表現で表せ。
- (a) 偶数個の1を含む語からなる言語
 - (b) 4個以上1が連続する部分をもつ語からなる言語
- (3) 正規表現 $((a+b)(a+b+d))^*$ をオートマトンに変換せよ。

理解確認クイズ(1)

(1) $\Sigma = \{a, b, c\}$ とする。 Σ 上の次の言語を正規表現で表せ。

(a) 少なくとも3個のaを含む語からなる言語

$$(b+c)^* a (b+c)^* a (b+c)^* a (a+b+c)^*$$

(b) 高々3個のaを含む語からなる言語

$$(b+c)^* + (b+c)^* a (b+c)^* + (b+c)^* a (b+c)^* a (b+c)^* + (b+c)^* a (b+c)^* a (b+c)^* a (b+c)^*$$

(c) 最後から3文字目はaである語からなる言語

$$(a+b+c)^* a (a+b+c) (a+b+c)$$

(d) 最後から3文字目はaでない語からなる言語

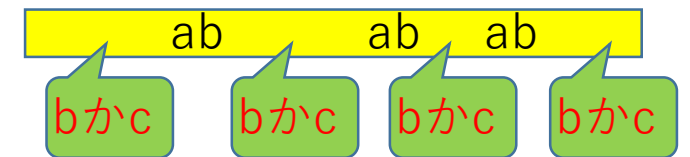
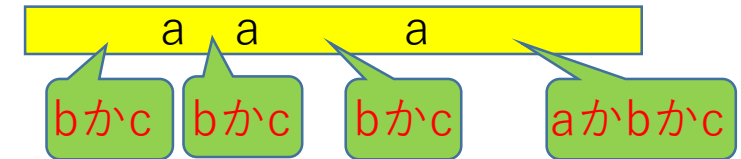
$$(a+b+c)^* (b+c) (a+b+c) (a+b+c)$$

(e) 長さが5の語からなる言語

$$(a+b+c) (a+b+c) (a+b+c) (a+b+c) (a+b+c)$$

(f) aの次の文字は必ずbである語からなる言語

$$((b+c)^* ab)^* (b+c)^*$$

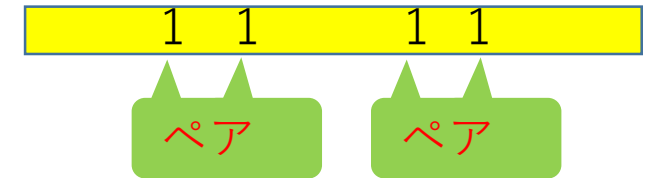


理解確認クイズ(2)

(2) $\Sigma = \{0,1\}$ とする。 Σ 上の次の言語を正規表現で表せ。

(a) 偶数個の1を含む語からなる言語

$$(0^*10^*1)^* 0^*$$



(b) 4個以上1が連続する部分をもつ語からなる言語

$$(0+1)^*1111(0+1)^*$$



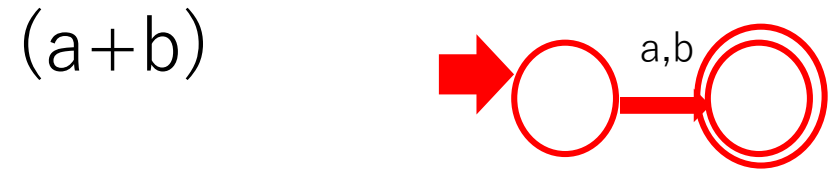
理解確認クイズ(3)

(3) 正規表現 $((a+b)(a+b+d))^*$ をオートマトンに変換せよ。

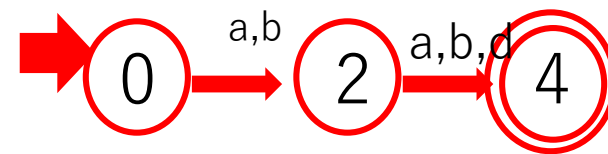
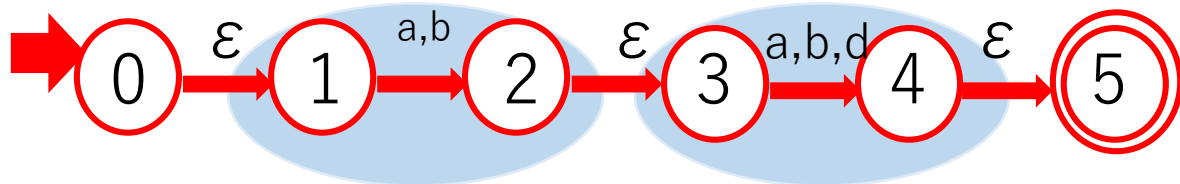
実際に挑戦してみましよう！

理解確認クイズ(3)つづき

$((a+b)(a+b+d))^*$



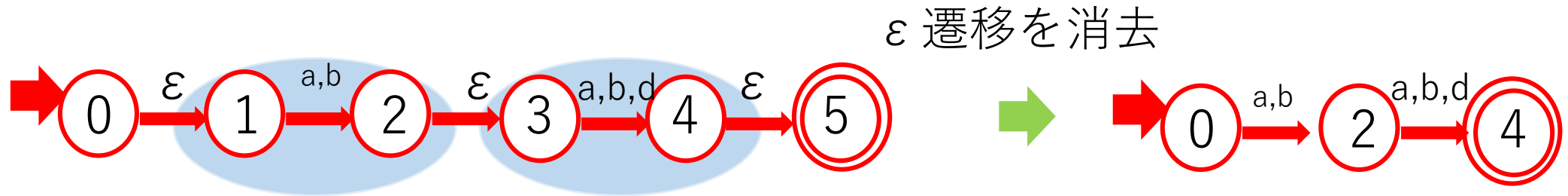
$(a+b)(a+b+d)$



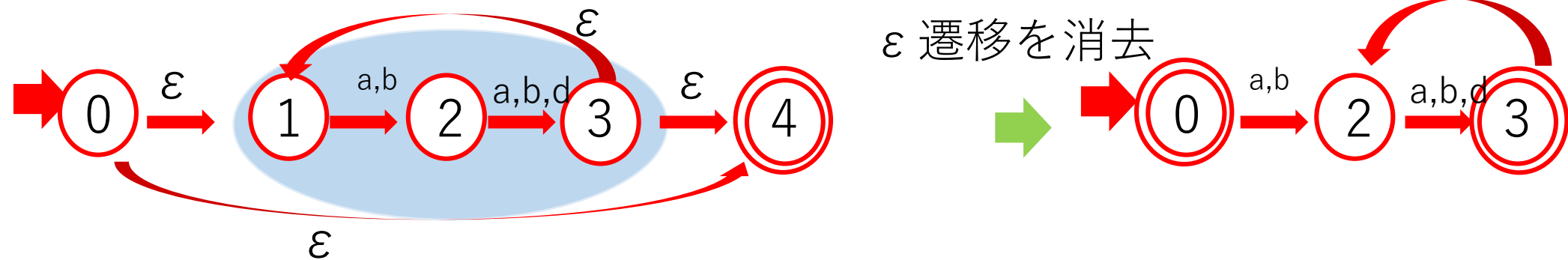
理解確認クイズ(3)つづき

$((a+b)(a+b+d))^*$

$(a+b)(a+b+d)$



$((a+b)(a+b+d))^*$



理解確認クイズ

(4) $\Sigma = \{0,1\}$ とする。 Σ 上の語で、二進数として解釈すると値が $4n$, $n=0,1,2,\dots$ になるものからなる言語を正規表現で表せ。

(5) $\Sigma = \{0,1\}$ とする。 Σ 上の語で、二進数として解釈すると3の倍数になるものからなる言語を正規表現で表せ。

(ヒント

オートマトンを設計してから正規表現に変換します。)

理解確認クイズ

(4) $\Sigma = \{0,1\}$ とする。 Σ 上の語で、二進数として解釈すると値が $4n$, $n=0,1,2,\dots$ になるものからなる言語を正規表現で表せ。

0	0
4	100
8	1000
12	1100
16	10000
20	10100
24	11000

$$(0+1)^*00 + 0$$

理解確認クイズ

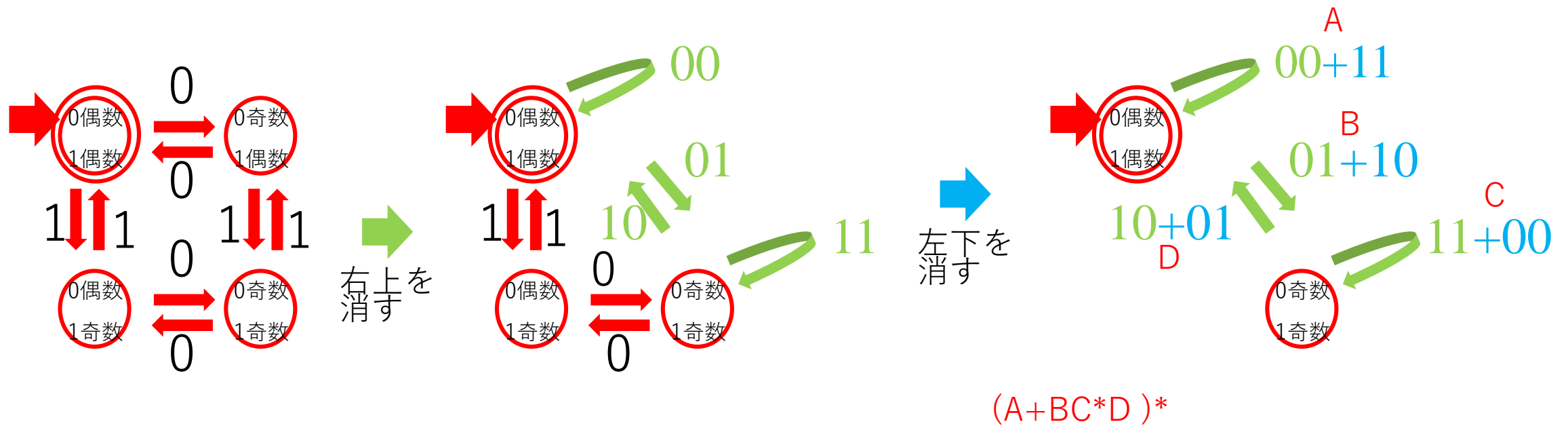
(6) $\Sigma = \{0,1\}$ とする。 Σ 上の語で、0の個数も1の個数も偶数個であるものからなる言語を正規表現で表せ。

ヒント

まずオートマトン、それを、正規表現に変換しましょう！

理解確認クイズ

(6) $\Sigma = \{0,1\}$ とする。 Σ 上の語で、0の個数も1の個数も偶数個であるものからなる言語を正規表現で表せ。



$$(00+11 + (01+10) (11+00)^* (10+01))^*$$