

正規言語上の Abstract Numeration System の 文字列圧縮への応用

新屋 良磨

Abstract Numeration System (ANS) $S = (L, A, <)$ は, 全単射 $\text{val}_S : L \rightarrow \mathbb{N}$ とその逆写像 $\text{rep}_S : \mathbb{N} \rightarrow L$ によって, 数と文字列を一対一対応づける記数法である. $\text{val}_S(w) = n \Leftrightarrow \text{rep}_S(n) = w$ は, 長さ—辞書順において文字列 $w \in L$ が n 番目の要素であることを意味する. 本稿では, 正規言語上の $ANS S = (L, A, <), S' = (K, A', <')$ について, $P_{S \rightarrow S'} : L \rightarrow K, P_{S \rightarrow S'}(w) := \text{rep}_{S'}(\text{val}_S(w))$ を言語 L から K への基数変換として定義し, これが圧縮に応用できることを示す. また, 具体的なアルゴリズムとその計算量, および圧縮率に関する重要な定理を与える.

1 はじめに

Goldberg と Sipser [5] は, 与えられた言語について文字列の順番 — *Ranking* による圧縮を提案し, 文脈自由文法においては *Ranking* の計算が多項式時間で計算可能なことを示した. それ以降, 形式文法の文法クラスにおける *Ranking* の並列計算可能性^{†1}や効率の良いアルゴリズムの研究がいくつか行われている [6][3][9]. 文献 [8] §7.2 には *Ranking* の圧縮への応用についての情報が纏められている.

一方, *Ranking* とは独立に, 言語と数の全単射によって定められる記数法 *Abstract Numeration System (ANS)* が Lecomte と Rigo [7] によって提案され, その表現力や性質が研究されている [1].

本研究では, 正規言語上の *ANS* について, 任意の正規言語 L, K 上の基数変換 $P_{S \rightarrow S'} : L \rightarrow K$ を定義し, $P_{S \rightarrow S'}$ が圧縮の条件を満たす場合について考察を行なった. 本研究における貢献を挙げる.

- *ANS* に関するアルゴリズムの提案とその計算量の見積り (§ 3)
- 正規言語上の *ANS* による基数変換の提案 (§ 4.1)
- 圧縮率の定義とその計算方法 (§ 4.2)

本稿は本章を含んで 5 章から構成される. § 2 は本稿で使用する記法及び前提となる定義・定理について紹介する. § 3 にて, *ANS* のアルゴリズム (擬似コードを含む) とその計算量を示し, § 4 で圧縮への応用を提案する. § 5 を本稿のまとめとする.

2 準備

本稿における記法及び定義, いくつかの定理を紹介する. 読者には特に前提知識を仮定しない. 本章で扱う定義や定理の詳細は文献 [10][2][12][4] 等に依る.

2.1 正規言語とオートマトン

有限文字集合 A において, A^n を A 上の長さが n である文字列の集合として用いる. さらに, A 上の全ての文字列の集合を

$$A^* := \bigcup_{i=0}^{\infty} A^i$$

によって表す. A^* の部分集合 L を, A 上の言語と呼ぶ. 特に断らない限り, $w \in A^*, \sigma \in A$ の記号を用いる. さらに, 文字列においては, その長さを $|w|$ とし, $1 \leq i \leq |w|$ について $w_i \in A$ を i 番目の文字とする.

Text Compression using Abstract Numeration System on a Regular Language

Ryoma Sin'ya, 東京工業大学情報理工学研究所数理・計算科学専攻, Department of Mathematical and Computing Sciences, Graduate School of Information Science and Engineering, Tokyo Institute of Technology.

†1 NC-computability

定義 2.1. A 上の正規言語の集合は, 以下のように再帰的に定義される.

- 空言語 \emptyset , 空文字列の単集合 $\{\epsilon\}$, 及び A の元のみからなる単集合 $\forall \sigma \in A \ \{\sigma\}$ は全て正規言語である.
- 正規言語 K, L において, 下記の演算を適用した結果も正規言語である.
 - $L \cup M$ (和集合)
 - $KL := \{kl \in A^* \mid k \in K, l \in L\}$ (接続)
 - $L^* := \{\epsilon\} \cup L \cup LL \cup \dots$ (Kleene 閉包) \square

定義 2.2. 有限オートマトン \mathcal{A} は 5 つ組 $\mathcal{A} = (Q, A, \delta, I, F)$; 状態の有限集合 Q , 入力文字の有限集合 A , 遷移関数 $\delta: Q \times A \rightarrow \mathfrak{P}(Q)$, 初期状態集合 $I \subseteq Q$, 受理状態集合 $F \subseteq Q$ によって定義される. \square

オートマトン \mathcal{A} の大きさ $|A|$ を, その状態集合の元の個数 $\#Q$ とする. さらに, 説明のため以下に定義される拡張遷移関数 $\hat{\delta}: Q \times A^* \rightarrow \mathfrak{P}(Q)$ を用いる.

$$\hat{\delta}(q, \sigma w) := \bigcup_{q' \in \delta(q, \sigma)} \hat{\delta}(q', w)$$

$$\hat{\delta}(q, \epsilon) := \{q\}$$

ここで ϵ は空文字を表す. \mathcal{A} において, $p \in \hat{\delta}(q, w)$ となる q 及び文字列 w が存在する場合, この w による状態遷移を $q \xrightarrow{w} p$, 誤解のない場合は $q \xrightarrow{w} p$ と表記する.

オートマトン \mathcal{A} 上の計算 (computation) c は状態遷移列を意味し,

$$c := q_0 \xrightarrow{\sigma_1} q_1 \xrightarrow{\sigma_2} q_2 \cdots \xrightarrow{\sigma_n} q_n$$

と記述する. \mathcal{A} において q_0 が初期状態, q_n が受理状態である場合, 文字列 w を受理文字列と呼ぶ. \mathcal{A} の全受理文字列の集合を

$$L(\mathcal{A}) = \left\{ w \mid \exists q \in I, \exists p \in F \left[q \xrightarrow{w} p \right] \right\}$$

として記述する. オートマトン $\mathcal{A}, \mathcal{A}'$ が等価であるとは, $L(\mathcal{A}) = L(\mathcal{A}')$ を満たすことである. $L(\mathcal{A}) = L$ であるとき, \mathcal{A} は L を認識するという.

定義 2.3. 以下の条件をみたす場合, \mathcal{A} は無曖昧オートマトン (Unambiguous, UFA) であるという.

1. すべての状態対 $(q, p) \in Q \times Q$, 全ての文字列 $w \in A^*$ について, 計算 $c = q \xrightarrow{w} p$ がたかだか一つ存在する.

2. 全ての受理文字列 $w \in L(\mathcal{A})$ について, $q \xrightarrow{w} p$ を満たす初期状態 $q \in I$ と受理状態 $p \in F$ が一意に存在する.

特に, \mathcal{A} の初期状態集合及び遷移状態関数の像が単集合:

$$|I| = 1 \ \wedge \ \forall q \in Q, \forall \sigma \in A \ [|\delta(q, \sigma)| = 1]$$

を満たす場合は, \mathcal{A} は決定性オートマトン (Deterministic, DFA) であるという. DFA であることを強調したい場合, オートマトンを \mathcal{D} で記す \square

DFA においては遷移関数の像が集合である必要はないので, 簡約のため $\delta: Q \times A \rightarrow Q$ として扱う.

この節の最後に, オートマトンのベクトル及び行列による表現形式を定義する. オートマトンは有向グラフと見ることができると, その隣接行列を考えることでオートマトンから自然数上の正方行列を導出することができる.

定義 2.4. オートマトン \mathcal{A} について, その隣接行列 $M(\mathcal{A}) \in \mathbb{N}^{|A| \times |A|}$ を, 全ての q 行 p 列要素に対して

$$M(\mathcal{A})_{qp} := \# \left\{ \sigma \mid \exists q, p \in Q \left[q \xrightarrow{\sigma} p \right] \right\}$$

と定義する. また, 初期ベクトル (行ベクトル) 及び受理ベクトル (列ベクトル) $V_I(\mathcal{A}), V_F(\mathcal{A})^T \in \mathbb{B}^{|A|}$ を

$$V_I(\mathcal{A})_q := \text{if } q \in I \text{ then } 1 \text{ else } 0$$

$$V_F(\mathcal{A})_q := \text{if } q \in F \text{ then } 1 \text{ else } 0$$

と定義する. また, \mathcal{A} が固定されている場合は, それぞれ簡単に M, V_I, V_F と記述する. さらに, 隣接行列において M_q とした場合は, M における q 行目の行ベクトルを表すとす. \square

2.2 Combinatorial Complexity

形式言語理論における Combinatorial Complexity は, 言語 L における数え上げ関数 $\dagger^2 C_L: \mathbb{N} \rightarrow \mathbb{N}$ に関する複雑さとして定式化されている [12] [1] [3].

定義 2.5. 言語 L について, 長さが n 及び n 以下である文字列の個数をそれぞれ

$$C_L(n) := \# \{ w \in L \mid n = |w| \} = \#(L \cap A^n)$$

$$C_L^{\leq}(n) := \sum_{i=0}^n C_L(i)$$

で表す. \square

\dagger^2 Counting function 又は Complexity function とも呼ぶ

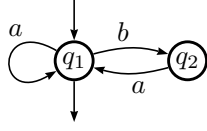


図 1 DFA \mathcal{D}_{fib} ; $L(\mathcal{D}_{fib}) = \{a, ba\}^*$

C_L の計算は、前述の UFA を用いて計算することができる。

補題 2.1. L を認識する UFA \mathcal{A} について、

$$C_L(n) = \mathbf{V}_I \mathbf{M}^n \mathbf{V}_F$$

が成り立つ。

証明. \mathcal{A} の隣接行列 \mathbf{M} について、 \mathbf{M}_{qp}^n は状態 q から p への文字列の長さが n の計算の総数にほかならない。これに初期ベクトルを掛けた結果、ベクトルの各要素 $(\mathbf{V}_I \mathbf{M}^n)_p$ は初期状態から状態 p への長さ n の計算の総数となる。さらに、受理ベクトルを掛けることで、補題の式から「初期状態から受理状態への長さ n の計算の総数」が得られる。

ここで \mathcal{A} は無曖昧であるので、計算の文字列において重複はない。よって補題が証明された。 \square

例 2.1. 図 1 の DFA \mathcal{D}_{fib} について、 C_L を考える。 $n \leq 4$ までを例示すると

$$\begin{aligned} C_L(0) &= \#\{\epsilon\} = 1, C_L(1) = \#\{a\} = 1 \\ C_L(2) &= \#\{aa, ba\} = 2 \\ C_L(3) &= \#\{aaa, aba, baa\} = 3 \\ C_L(4) &= \#\{aaaa, aaba, abaa, baaa, baba\} = 5 \end{aligned}$$

である。 \mathcal{D}_{fib} の隣接行列 及び初期ベクトル, 受理ベクトルはそれぞれ

$$\mathbf{M} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}, \mathbf{V}_I = \mathbf{V}_F^T = (1, 0)$$

となる。 \mathbf{M} の固有値 $\alpha = \frac{1+\sqrt{5}}{2}, \beta = \frac{1-\sqrt{5}}{2}$ から

$$\mathbf{M} = \mathbf{P} \mathbf{D} \mathbf{P}^{-1} = \frac{1}{\sqrt{5}} \begin{pmatrix} \alpha & \beta \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \alpha & 0 \\ 0 & \beta \end{pmatrix} \begin{pmatrix} 1 & -\beta \\ -1 & \alpha \end{pmatrix}$$

と対角化可能である。ここで、補題 2.1 を適用すると

$$\begin{aligned} C_L(n) &= \mathbf{V}_I \mathbf{M}^n \mathbf{V}_F = \mathbf{V}_I \mathbf{P} \mathbf{D}^n \mathbf{P}^{-1} \mathbf{V}_F \\ &= \frac{1}{\sqrt{5}} \mathbf{V}_I \begin{pmatrix} \alpha^{n+1} - \beta^{n+1} & \alpha^n - \beta^n \\ \alpha^n - \beta^n & \alpha^{n-1} - \beta^{n-1} \end{pmatrix} \mathbf{V}_F \\ &= \frac{1}{\sqrt{5}} \left\{ \left(\frac{1+\sqrt{5}}{2} \right)^{n+1} - \left(\frac{1-\sqrt{5}}{2} \right)^{n+1} \right\} \end{aligned}$$

となり、フィボナッチ数の $n+1$ 番目が得られる。 \square

次に、数え上げ関数の漸近に関する次の重要な定理を紹介する。この定理は § 4.2 における圧縮率の見積りで必要となる。

定理 2.1 (Shur [11][12]). DFA \mathcal{D} によって受理される任意の正規言語 $L = L(\mathcal{D})$ について、

$$C_L(n) = \Theta(n^m \alpha^n)$$

が成り立つ。ここで α は \mathcal{D} の隣接行列のフロベニウス根^{†3}, $m+1$ は \mathcal{D} の隣接行列において α を持つ、単一路で連結された強連結成分の最大個数である。

補足 2.1. ここでは定理 2.1 の証明は行わない。頂点数 n 及び辺の数が e である有向グラフについて、その隣接行列のフロベニウス根を $\Theta(n^3)$, 全ての強連結成分を求める強連結成分分解は $O(n+e) = O(n^2)$ の時間計算量で求められることを補足とする。 \square

例 2.2. 再び図 1 の \mathcal{D}_{fib} について考える。 \mathcal{D}_{fib} の隣接行列におけるフロベニウス根は $\alpha = \frac{1+\sqrt{5}}{2}$ である。さらに、 \mathcal{D} は唯一の強連結成分から構成されているので、 $m+1=1$ となる。例 2.1 で求めた結果から

$$C_L(n) = \Theta \left(\left(\frac{1+\sqrt{5}}{2} \right)^n \right)$$

を満たすことが分かる。 \square

2.3 Abstract Numeration System

記数法 (*Numeration System*) とは、数を文字列として表現するシステムを指す。その表現力や代数演算に関する性質など、それ自体 1 つの数学の研究分野として多く研究されている [1][7]。

ここでは、非常に柔軟な記数法である *Abstract Nu-*

^{†3} 非負の正方行列における最大の固有値。詳しくは [2] を参照。

meration System について説明を行う。基本となるアイデアは、順序付けられた文字集合上の言語と自然数間の全単射である。

定義 2.6. 全順序が定められた文字集合 $(A, <)$ において、長さ-辞書順^{†4} \prec を次のように定める。文字列 $u, v \in A^*$ において、 $|u| < |v|$ である場合 (長さ順) に $u \prec v$ である。さらに、 $|u| = |v|$ かつ適当な文字列 $p, q, r \in A^*$ 及び文字 $a, b \in A$ $a < b$ によって $u = paq, v = pbr$ の形に分解できる場合 (辞書順) も $u \prec v$ である。 \prec は A^* における全順序となる。□

定義 2.7. Abstract Numeration System (ANS) \mathcal{S} は 3 つ組み $(L, A, <)$; 文字の全順序集合 $(A, <)$ 上の言語 L によって定義される。 \mathbb{N} から言語 L への全単射 $\text{rep}_{\mathcal{S}} : \mathbb{N} \rightarrow L$ は、自然数 $n \in \mathbb{N}$ を、 L を長さ-辞書順によって整列した時の n 番目の文字列に写す写像である。 $\text{rep}_{\mathcal{S}}$ の逆写像を $\text{val}_{\mathcal{S}} : L \rightarrow \mathbb{N}$ で表す。□

本稿では正規言語上の ANS のみについて考える。ANS の仕組みは非常に単純であるが、記数法としての柔軟性は非常に高い。2 進法や 10 進法も ANS として表現することができる (一般には記数法は負数、さらには実数をも扱えるが、本稿では自然数との対応のみを扱う)。

例 2.3. ANS $\mathcal{S}_2 = (L, \{0, 1\}, <); L = \{0\} \cup \{1\}\{0, 1\}^*, 0 < 1$ について考える。 L を長さ-辞書順によって整列すると $\{0, 1, 10, 11, 100, \dots\}$ が得られ、 $\text{rep}_{\mathcal{S}_2}$ 及び $\text{val}_{\mathcal{S}_2}$ は 2 進法を定める。つまり

$$\begin{aligned} \text{rep}_{\mathcal{S}_2}(0) &= 0, \text{rep}_{\mathcal{S}_2}(1) = 1, \text{rep}_{\mathcal{S}_2}(2) = 10, \dots \\ \text{val}_{\mathcal{S}_2}(0) &= 0, \text{val}_{\mathcal{S}_2}(1) = 1, \text{val}_{\mathcal{S}_2}(10) = 2, \dots \end{aligned}$$

は $n \in \mathbb{N}$ とその 2 進表現を一対一に対応付ける。

同様に、任意の正整数 k について ANS $\mathcal{S}_k = (L, \{0, 1, \dots, \sigma_{k-1}\}, <); L = \{0\} \cup \{1, 2, \dots, \sigma_{k-1}\}\{0, 1, \dots, \sigma_{k-1}\}^*, 0 < 1, 1 < 2, \dots, \sigma_{k-1} < \sigma_k$ は k 進法を定める。□

3 DFA を用いた ANS のアルゴリズム

本章では、任意の DFA \mathcal{D} の認識する言語 $L = L(\mathcal{D})$ 上の ANS $\mathcal{S} = (L, A, <)$ について、 $\text{rep}_{\mathcal{S}}$ と $\text{val}_{\mathcal{S}}$ の

アルゴリズムを定める。アルゴリズムの中心となるのは、DFA \mathcal{D} の隣接行列とベクトルの演算である。

3.1 $\text{val}_{\mathcal{S}}$ の計算

$\text{val}_{\mathcal{S}}$ のアルゴリズムの説明に、次のベクトルの記法を用いる。

定義 3.1. ベクトル $\mathbf{V}_q^{<\sigma} \in \mathbb{N}^{|\mathcal{D}|}$ を、状態 q から σ 未満の入力文字による遷移の総数

$$(\mathbf{V}_q^{<\sigma})_p := \#\{\sigma' \in A \mid \sigma' < \sigma \wedge q \xrightarrow{\sigma'} p\}$$

と定義する。□

定理 3.1. DFA $\mathcal{D} = (Q, A, \delta, I, F); L = L(\mathcal{D})$ となる ANS $\mathcal{S} = (L, A, <)$ において

$$\begin{aligned} \text{val}_{\mathcal{S}}(w) &= \sum_{i=1}^{|w|} (\mathbf{V}_{q_i}^{<w_i} + \mathbf{V}_I) \mathbf{M}^{|w|-i} \mathbf{V}_F \\ \text{where } q_1 &\in I, q_i \xrightarrow{w_i} q_{i+1} \end{aligned}$$

が成り立つ。

証明. 定義より、

$$\begin{aligned} \text{val}_{\mathcal{S}}(w) &= \#\{f \in L \mid f \prec w\} \\ &= \#\{f \in L \cap A^{|w|} \mid f \prec w\} + C_L^{\leq}(|w| - 1) \end{aligned}$$

さらに、補題 2.1 より

$$C_L^{\leq}(|w| - 1) = \sum_{i=0}^{|w|-1} \mathbf{V}_I \mathbf{M}^i \mathbf{V}_F \quad (1)$$

が成り立つ。なので

$$\#\{f \in L \cap A^{|w|} \mid f \prec w\} \quad (2)$$

について、式を求めれば良い。これは、長さが $|w|$ かつ辞書順において w 未満の文字列による、 \mathcal{D} の全て受理計算にほかならない。よって

$$(2) = \sum_{i=1}^{|w|} (\mathbf{V}_{q_i}^{<w_i}) \mathbf{M}^{|w|-i} \mathbf{V}_F$$

が満たされる。(1), (2) より

$$\begin{aligned} \text{val}_{\mathcal{S}}(w) &= (1) + (2) \\ &= \sum_{i=0}^{|w|-1} \mathbf{V}_I \mathbf{M}^i \mathbf{V}_F + \sum_{i=1}^{|w|} (\mathbf{V}_{q_i}^{<w_i}) \mathbf{M}^{|w|-i} \mathbf{V}_F \\ &= \sum_{i=1}^{|w|} (\mathbf{V}_{q_i}^{<w_i} + \mathbf{V}_I) \mathbf{M}^{|w|-i} \mathbf{V}_F \end{aligned}$$

□

系 3.1. 式 1 のように簡単な変形を行うことで、 $\text{val}_{\mathcal{S}}(w)$ は

- $2|w| - 1$ 回のベクトル和
- $|w| - 1$ 回のベクトル—行列積
- 1 回のベクトル—ベクトル積

によって計算可能である。□

^{†4} length-lexicographic 又は radix order と呼ぶ

式 1 定理 3.1 の変形式 (Horner 法)

$$\begin{aligned}
\text{vals}(w) &= \sum_{i=1}^{|w|} (\mathbf{V}_{q_i}^{<w_i} + \mathbf{V}_I) M^{|w|-i} \mathbf{V}_F \\
&= (\mathbf{V}_{q_1}^{w_1} + \mathbf{V}_I) M^{|w|-1} \mathbf{V}_F + (\mathbf{V}_{q_2}^{w_2} + \mathbf{V}_I) M^{|w|-2} \mathbf{V}_F + \cdots + (\mathbf{V}_{q_{|w|}}^{w_{|w|}} + \mathbf{V}_I) M^0 \mathbf{V}_F \\
&= \left(\cdots \left(\left((\mathbf{V}_{q_1}^{w_1} + \mathbf{V}_I) M + (\mathbf{V}_{q_2}^{w_2} + \mathbf{V}_I) \right) M + (\mathbf{V}_{q_3}^{w_3} + \mathbf{V}_I) \right) M + \cdots + (\mathbf{V}_{q_{|w|}}^{w_{|w|}} + \mathbf{V}_I) \right) \mathbf{V}_F
\end{aligned}$$

3.2 rep_S の計算

基本的に rep_S(n) の計算は, val_S(w) の計算式 1 について各ベクトル $\mathbf{V}_{q_i}^{w_i}$ を満たす w_i を文字集合 A について線形探索をするだけである. ここで重要なのは, 入力 n から対応する w の長さ $|w|$ を求めることである. $|w|$ を同定しないかぎり, ベクトルの線形探索はできない. n から $|w|$ を求めるには,

$$\max_{\ell} \left(C_L^{\leq}(\ell - 1) \leq n \right)$$

を求めればよい. 補題 2.1 において, $C_L(n)$ の計算が冪行列とベクトルの積のみで計算できることを示した. この節では $C_L^{\leq}(n)$ も同様の計算で行えることを示す. そのために, 以下の性質を持つ正規言語を使う.

定義 3.2. 正規言語 $L \subseteq A^*$ について, その総和言語 $\tilde{L} \subseteq (A \cup \{\$\})^*$; $\$ \notin A$ を

$$\tilde{L} := L\{\$\}^*$$

と定義する. \square

補題 3.1. 任意の正規言語 L について,

$$C_L^{\leq}(n) = C_{\tilde{L}}^{\leq}(n)$$

が成り立つ.

証明. 全ての n について, 長さが L の n 以下の全ての文字列 w に対応する, 長さ n の文字列 $\tilde{w} \in \tilde{L}$ が存在し, かつそれ以外が \tilde{L} に含まれないことを示せば良い.

$\forall n \in \mathbb{N} \quad \forall w \in L \left[|w| \leq n \Rightarrow \tilde{w} = w \overbrace{\$\dots\$}^{|w|-n} \in \tilde{L} \right]$
 によって順方向を示せる. また, \tilde{w} は必ず $\$$ を含まない文字列 v の後ろに $\$$ が続く形になるので, 末尾の $\$$ を除いた文字列 v は必ず L の要素である. よって逆方向も示せた. \square

補題 3.2. 任意の UFA $\mathcal{A} = (Q, A, \delta, I, F)$ について, その総和言語を認識する UFA $\tilde{\mathcal{A}}; L(\tilde{\mathcal{A}}) = \tilde{L}(\mathcal{A})$ は,

ただだか 2 つの状態の追加で構成可能である.

証明. A に含まれない文字 $\$ \notin A$ について, $\tilde{\mathcal{A}} = (\tilde{Q}, A \cup \{\$\}, \tilde{\delta}, \tilde{I}, \tilde{F})$ を構成する. まず, \mathcal{A} が空文字を受け取らない場合 ($\epsilon \notin L(\mathcal{A})$)

$$\begin{aligned}
\tilde{Q} &:= Q \cup \{\tilde{q}\} \\
\tilde{\delta}(q \in Q, \sigma) &:= \begin{cases} \delta(q, \sigma) \cup \{\tilde{q}\} & \delta(q, \sigma) \cap F \neq \emptyset \\ \delta(q, \sigma) & \text{otherwise} \end{cases} \\
\tilde{\delta}(\tilde{q}, \$) &:= \{\tilde{q}\} \\
\tilde{I} &:= I \\
\tilde{F} &:= \{\tilde{q}\}
\end{aligned}$$

\mathcal{A} が空文字を受け取る場合は, 上の構成の後に, $\tilde{Q}, \tilde{I}, \tilde{F}$ にそれぞれ q_ϵ を追加する.

以上の構成から, \mathcal{A} が受理状態に遷移する場合は $\tilde{\mathcal{A}}$ も受理状態に遷移する. また, $\tilde{\mathcal{A}}$ は一度受理状態 \tilde{q} に遷移すると, 文字 $\$$ によってその後も (唯一) 自身に遷移することができる. よって $L(\tilde{\mathcal{A}}) = \tilde{L}(\mathcal{A})$ を満たす. \square

3.3 計算量

前節で説明したアルゴリズムを擬似コードとしてアルゴリズム 1, 2 に, 対応する計算量を表 1 に示す. 計算量については, 任意の桁数の整数乗算が定数時間で行える理想的な場合と, 実際の乗算アルゴリズムで知られている計算量の両方を記述した. n 桁の整数の乗算において, $O(n \log n \log \log n)$ の計算量を持つ Schnhage–Strassen アルゴリズム等が知られている. 長整数乗算の計算量の比較は文献 [4] §9.5.8 が詳しい.

4 ANS による文字列圧縮

本章では, ANS による正規言語の文字列と数の一対一対応を応用して, 正規言語間での文字列の一対一対応を提案し, 圧縮に関する考察を行う.

アルゴリズム 1 val_S	アルゴリズム 2 rep_S
入力: $w \in L$	入力: $0 \leq n < \#L$
結果: $n = \#\{f \in L \mid f \prec w\}$	結果: $w \in L$ s.t. $\#\{w' \in L \mid w' \prec w\} = n$
1: $q \leftarrow q_0 \in I$	1: $q \leftarrow q_0 \in I$
2: $V \leftarrow V_0 \in \mathbb{N}^{ w }$	2: $\ell \leftarrow \max_{\ell} (C_L^{\leq}(\ell - 1) \leq n)$
3: for $i = 1 \rightarrow w $ do	3: $m \leftarrow n - C_L^{\leq}(\ell - 1)$
4: $V \leftarrow V + V_q^{\sigma_i} + V_I$	4: $w \leftarrow \epsilon$
5: $q \leftarrow \delta(q, w_i)$	5: for $i = \ell \rightarrow 1$ do
6: if $i < w $ then	6: $\sigma \leftarrow \min A$
7: $V \leftarrow VM$	7: while $M_{\delta(q, \sigma)}^{\ell-i} {}^tV_F \leq m$ do
8: end if	8: $m \leftarrow m - M_{\delta(q, \sigma)}^{\ell-i} {}^tV_F$
9: end for	9: $\sigma \leftarrow \text{suc}(\sigma)$
10: $n \leftarrow V^t V_F$	10: end while
	11: $w \leftarrow w\sigma$
	12: $q \leftarrow \delta(q, \sigma)$
	13: end for

表 1 計算量

	時間計算量	行列積	行列-ベクトル積	ベクトル積
$\text{val}_S(w)$	$\Theta(w \times \mathcal{D} ^2) [O(w ^2 \log w \log \log w \times \mathcal{D} ^2)]$	0	$ w - 1$	1
$\text{rep}_S(n)$	$O(\ell \log \ell \times \mathcal{D} ^3) [O(\ell^2 \log^2 \ell \log \log \ell \times \mathcal{D} ^3)]$	$O(\ell \log \ell)$	0	$O(\ell \times A)$

$\ell = |\text{rep}_S(n)|$ □ 内は乗算コストを考慮した場合、行列積、行列-ベクトル積、ベクトル積は演算回数を表す。

4.1 ANS 間の基数変換

定義 4.1. ANS $S = (L, A, <)$, $S' = (K, A', <')$ において, S から S' への基数変換 $P_{S \rightarrow S'} : L \rightarrow K$ を

$$P_{S \rightarrow S'}(w) := \text{rep}_{S'}(\text{val}_S(w))$$

によって定める. □

前章において, val_S は文字列の長さに対して線形時間, rep_S は対応する文字列の長さに対して順線形時間で計算可能なことを示した. よって, 基数変換は両方の計算量の和で計算可能である.

補足 4.1. ANS $S = (L(\mathcal{D}), A, <)$, $S' = (L(\mathcal{D}'), A', <')$ において $P_{S \rightarrow S'}(w)$ の計算量は

$$\Theta(|w| \times |\mathcal{D}|^2) + O(\ell \log \ell \times |\mathcal{D}'|^3) \\ \text{where } \ell = |\text{rep}_{S'}(\text{val}_S(w))|$$

となる. □

4.2 圧縮率

定義 4.2. ANS $S = (L, A, <)$, $S' = (K, A', <')$ について, 基数変換 $P_{S \rightarrow S'} : L \rightarrow K$ の圧縮率 CR を

$$CR(P_{S \rightarrow S'}, n) := \frac{\log C_L^{\leq}(n)}{\log C_K^{\leq}(n)}$$

さらにその極限を

$$CR(P_{S \rightarrow S'}) := \lim_{n \rightarrow \infty} CR(P, n)$$

と定義する. 特に, $CR(P_{S \rightarrow S'}) < 1$ を満たすような $P_{S \rightarrow S'}$ を圧縮と呼ぶ. □

$CR(P_{S \rightarrow S'}, n)$ は長さ n の文字列についての情報量の比である. 単に圧縮率といった場合は, その極限を指す. 圧縮率に関しては, 明らかに次の系が成り立つ.

系 4.1. 任意の ANS $S = (L, A, <)$ に対して, ANS $S' = (K, A', <')$ への基数変換 $P_{S \rightarrow S'}$ は $K = A^*$ において最も圧縮率が高い. □

ここで、数え上げ関数に関する以下の補題を説明する。これは次の定理の証明に用いられる。

補題 4.1. 任意のオートマトン \mathcal{A} について、認識する言語を $L = L(\mathcal{A})$ 、その隣接行列のフロベニウス根を α とすると、

$$C_L^{\leq}(n) = \begin{cases} \Theta(n \times C_L(n)) = \Theta(n^{m+1}) & \alpha = 1 \\ \Theta(C_L(n)) = \Theta(n^m \alpha^n) & \alpha > 1 \end{cases}$$

が成り立つ。

証明. まず、 $\alpha = 1$ について証明する。定理 2.1 の簡単な式変形より、

$$\begin{aligned} C_L^{\leq}(n) &= \sum_{i=0}^n \Theta(i^m \alpha^i) = \Theta\left(\sum_{i=0}^n i^m\right) \\ &= \Theta(n^{m+1}) \end{aligned}$$

また、 $\alpha < 1$ のとき

$$C_L(n) = \Theta(n^m \alpha^n) \leq C_L^{\leq}(n) \quad (3)$$

さらに、

$$\begin{aligned} C_L^{\leq}(n) &= \sum_{i=0}^n \Theta(i^m \alpha^i) \\ &\leq n^m \sum_{i=0}^n \Theta(\alpha^i) = n^m \Theta(\alpha^n) \\ &= \Theta(n^m \alpha^n) \end{aligned} \quad (4)$$

(3), (4) より、上下界が漸近的に一致するため

$$C_L^{\leq}(n) = \Theta(C_L(n))$$

となる。□

準備が揃ったので、この説の最後に圧縮率に関する重要な定理を示す。

定理 4.1. 任意の UFA $\mathcal{A}, \mathcal{A}'$ の認識する言語上の、任意の ANS $\mathcal{S} = (L = L(\mathcal{A}), A, <), \mathcal{S}' = (K = L(\mathcal{A}'), A', <')$ において、 \mathcal{S} から \mathcal{S}' への基数変換 $P_{\mathcal{S} \rightarrow \mathcal{S}'} : L \rightarrow K$ の圧縮率は

$$CR(P_{\mathcal{S} \rightarrow \mathcal{S}'}) = \begin{cases} \frac{\log \alpha}{\log \alpha'} & \alpha > 1, \alpha' > 1 \\ \frac{m+1}{m'+1} & \alpha = 1, \alpha' = 1 \\ \infty & \alpha > 1, \alpha' = 1 \\ 0 & \alpha = 1, \alpha' > 1 \end{cases}$$

となる。ここで、 $\alpha(\alpha')$ は $\mathcal{A}(\mathcal{A}')$ の隣接行列におけるフロベニウス根、 $m(m')$ は $\mathcal{A}(\mathcal{A}')$ の $\alpha(\alpha')$ を持つ、単一路で連結された強連結成分の最大数である。

証明. $\alpha > 1, \alpha' > 1$ について考える。圧縮率の定義

と補題 4.1 より

$$\begin{aligned} CR(P_{\mathcal{S} \rightarrow \mathcal{S}'}) &= \lim_{n \rightarrow \infty} \frac{\log C_L^{\leq}(n)}{\log C_{K'}^{\leq}(n)} = \lim_{n \rightarrow \infty} \frac{\log \Theta(C_L(n))}{\log \Theta(C_{K'}(n))} \\ &= \lim_{n \rightarrow \infty} \frac{\log \Theta(n^m \alpha^n)}{\log \Theta(n^{m'} \alpha'^n)} \\ &= \lim_{n \rightarrow \infty} \frac{m \log n + n \log \alpha + \log c}{m' \log n + n \log \alpha' + \log c'} \\ &= \frac{\log \alpha}{\log \alpha'} \end{aligned}$$

が得られる。ここで c, c' は Θ における定数項である。

$\alpha = 1, \alpha' = 1$ や他の場合も、補題 4.1 から同様に計算することで各圧縮率が得られる。□

補足 4.2. 圧縮率 $CR(P)$ の計算は、 $\mathcal{A}, \mathcal{A}'$ のフロベニウス根の計算に帰着される。補足 2.1 で述べたように、その計算量は $O(|\mathcal{A}|^3), O(|\mathcal{A}'|^3)$ である。□

5 まとめ

本稿では ANS における圧縮を提案し、その圧縮率を与えられた任意の正規言語について計算できることを示した。また、ANS における基本的な計算 $\text{vals}, \text{rep}_{\mathcal{S}}$ 及び数え上げ関数 C_L, C_L^{\leq} のアルゴリズムの提案も行なった。

参考文献

- [1] Berth, V. and Rigo, M.: *Combinatorics, Automata and Number Theory*, Cambridge University Press, New York, NY, USA, 1st edition, 2010.
- [2] Brouwer, A. E. and Haemers, W. H.: *Spectra of graphs (to appear)*, Berlin: Springer, 2012.
- [3] Choffrut, C. and Goldwurm, M.: Rational Transductions and Complexity of Counting Problems, *Proceedings of the 17th International Symposium on Mathematical Foundations of Computer Science*, MFCS '92, London, UK, UK, Springer-Verlag, 1992, pp. 181–190.
- [4] Crandall, R., Pomerance, C., Crandall, R., and Pomerance, C.: *Prime numbers: a computational perspective. Second Edition*, 2005.
- [5] Goldberg, A. and Sipser, M.: Compression and ranking, *Proceedings of the seventeenth annual ACM symposium on Theory of computing*, STOC '85, New York, NY, USA, ACM, 1985, pp. 440–448.
- [6] Hemachandra, L. A. and Rudich, S.: On the Complexity of Ranking., *J. Comput. Syst. Sci.*, Vol. 41, No. 2(1990), pp. 251–271.

- [7] Lecomte, P. B. A. and Rigo, M.: Numeration systems on a regular language, *CoRR*, Vol. cs.OH/9903005(1999).
- [8] Li, M. and Vitnyi, P. M.: *An Introduction to Kolmogorov Complexity and Its Applications*, Springer Publishing Company, Incorporated, 3 edition, 2008.
- [9] Mkinen, E., Mkinen, E., and Mkinen, E.: Ranking and unranking left Szilard languages, Technical report, ISO/IEC JTC1/SC29/WG11/N2467, Atlantic City, 1997.
- [10] Sakarovitch, J.: *Elements of Automata Theory*, Cambridge University Press, New York, NY, USA, 2009.
- [11] Shur, A. M.: Combinatorial complexity of regular languages, *Proceedings of the 3rd international conference on Computer science: theory and applications*, CSR'08, Berlin, Heidelberg, Springer-Verlag, 2008, pp. 289–301.
- [12] Shur, A. M.: Combinatorial Characterization of Formal Languages, *CoRR*, Vol. abs/1010.5456(2010).