

## 解説

## 原始語の集合と文脈自由言語について

新屋良磨 \*\*

On Primitive Words and Context-Free Languages

Ryoma Sin'ya\*\*

## Abstract

Primitive word is a word that can not be represented by any repetition of shorter words. Since every non-empty word is a repetition of the unique primitive word, primitive words play an important role in combinatorics on words. In this article, we explain a long-standing open problem called “primitive words conjecture” which has a deep connection with the theory of context-free languages.

## 1 はじめに

原始語とは「自身より短い語の繰り返し」では表されない語 (文字の有限列) のことである (正確な定義は次節にて行う). 任意の語はある原始語の繰り返しとして一意に分解することができ, そういった意味で原始語は語の世界における素数のような対象であり, 定義は単純であるが非常に奥深い性質を持っている.

形式言語理論においては「ある語の集合 (言語) が特定の言語族に属するかどうか?」といった種類の問題がしばしば興味の対象となる. ここでは「特定の言語族」とは計算論的あるいは代数的な特徴づけを持った言語族が主な興味の対象とされ, 典型的な例としては Chomsky の階層における正則言語 (regular languages), 文脈自由言語 (context-free languages), 文脈依存言語 (context-sensitive languages) などの族が挙げられる. 帰納的可算言語 (recursively enumerable language) まで行ってしまうと「ある語の集合が帰納的可算言語に属するかどうか? (すなわち計算可能であるか)」は (著者の感覚からすると) 形式言語理論というよりも計算論の問いになるものと思われる. その中でも特に文脈自由言語は, 代数における代数関数 (多項式系の解となる巻数) のある種の非可換化と捉えることができ, 代数的言語とも呼ばれている重要な言語族である. あらゆる性質が決定可能である正則言語とは異なり, 文脈自由言語は (等価性判定や普遍性判定を含む) 多くの問題が決定不能であり, 計算的には困難面も持ちながらも豊かな構造理論を持っている.

本稿では原始語全体の集合と文脈自由言語に対する有名な Dömösi-Horváth-Ito 予想について述べ, これまでの研究における予想へのいくつかのアプローチや今後の課題について解説を行う. Dömösi-Horváth-Ito 予想に対する既存のアプローチは多数あるものの, 本項では特に言語の母関数や測度に関連したアプローチについて取り上げる. 続く 2 章では原始語と文脈自由言語の定義を行い Dömösi-Horváth-Ito 予想を定式化する. 3 章にて原始語の集合と無曖昧文脈自由言語の関係について考察を行い, 4 章にて原始語の集合に対する著者なりの今後の課題について述べ本稿の結びとする.

## 2 原始語の集合と文脈自由言語

形式言語理論において, 文字集合とは単に空でない有限集合のことを差し, 文字集合の要素を文字と呼ぶ. 以降, 変数  $A$  は常に文字集合を表す. 文字集合  $A$  上の語とは  $A$  に属する文字  $a_i$  を有限個並べた列:  $a_1 a_2 \cdots a_n$

---

2018年8月10日受理

\*\* 秋田大学大学院理工学研究科数理・電気電子情報学専攻数理科学コース,  
Mathematical Science Course, Akita University Graduate School of Engineering Science.

である。語には接続  $\cdot$  という演算を考えることができ、2つの語  $u = a_1 \cdots a_i, v = b_1 \cdots b_j$  の接続  $u \cdot v$  は  $u$  と  $v$  を並べたものを表す：

$$u \cdot v \triangleq a_1 \cdots a_i b_1 \cdots b_j. \quad (1)$$

語  $w$  を  $n$  個接続した語を  $w^n$  で表す。例えば  $aaa = a^3$  で  $abab = (ab)^2$  である。語  $w = a_1 \cdots a_n$  の長さを  $|w|$  で表す： $|w| = n$ 。  $\varepsilon$  で長さが0の語(空語)を表す。 $A$  に属する文字から作られる語全ての集合を  $A^*$  で表す。 $A$  上の言語とは  $A^*$  の部分集合のことを指す。すなわち  $L \subseteq A^*$  となる  $L$  を言語と呼ぶ。任意の2つの言語  $L, M$  の接続を以下のように語の接続を自然に拡張して定義する：

$$L \cdot M \triangleq \{uv \in A^* \mid u \in L, v \in M\}. \quad (2)$$

言語  $L$  に対し、 $L^n$  で言語  $L$  の  $n$  回の接続

$$L^0 \triangleq \{\varepsilon\} \quad L^n \triangleq L \cdot L^{n-1} \quad (3)$$

を表す。言語  $L$  に対してその Kleene 閉包  $L^*$  を

$$L^* = \bigcup_{n=0}^{\infty} L^n \quad (4)$$

で定義する。有限要素の言語を全て含み、和集合、接続および Kleene 閉包に閉じた最小の言語族を正則言語族と呼ぶ。

## 2.1 原始語

**定義 2.1** (原始語).  $A$  上の語  $w \in A^*$  が原始語であるとは、 $w$  より短い語の繰り返して表せられないこと：

$$\text{任意の } u \in A^* \text{ について } w = u^n \Rightarrow n = 1 \quad (5)$$

ことを言う。 $A$  上の全ての原始語の集合を  $\mathbb{Q}_A$  で表す。

例えば語  $abbab$  は原始的であるが、 $abbabb = (abb)^2$  は原始的ではない。また、語  $w$  の長さが素数  $p$  の場合は  $a^p$  という形の語を除いて原始的となる：

$$\text{任意の素数 } p \text{ について } A^p \cap \mathbb{Q}_A = A^p \setminus \{a^p \mid a \in A\}. \quad (6)$$

さらに、任意の非空語  $w \in A^* (|w| \geq 1)$  には原始語  $u$  と自然数  $n \geq 1$  のペアが一意に存在し  $w = u^n$  が成り立つ。

## 2.2 文脈自由言語

**定義 2.2** (文脈自由文法). 文脈自由文法とは、文字集合  $A$  と

- 変数集合と呼ばれる  $A$  と素な有限集合  $V$
- 導出規則と呼ばれる関係  $\xrightarrow{R} \subseteq V \times (V \cup A)^*$
- 初期変数と呼ばれる  $S \in V$

から成る4つ組  $G = (V, A, \xrightarrow{R}, S)$  である。

**定義 2.3** (文脈自由言語). 文脈自由文法  $G = (V, A, \xrightarrow{R}, S)$  が表現する言語  $\llbracket G \rrbracket$  を初期変数  $S$  から導出規則

$\xrightarrow{R}$  の反射推移閉包  $\xrightarrow{R^*}$  によって導出される  $A$  上の言語として定義する.

$$\llbracket G \rrbracket \triangleq \{w \in A^* \mid S \xrightarrow{R^*} w\}. \quad (7)$$

言語  $L$  が文脈自由とは,  $L = \llbracket G \rrbracket$  となる文脈自由文法  $G$  が存在することを言う.

文脈自由言語の典型的な例としては,  $A_0 = \{(\,,)\}$  上の次の文法  $G_D$  が挙げられるだろう.

$$G_D = (\{S\}, A = \{(\,,)\}, \xrightarrow{R} = \{(S, \varepsilon), (S, (S)S)\}, S). \quad (8)$$

文法  $G_D$  によって生成される語として

$$S \xrightarrow{R} (S)S \xrightarrow{R} (\varepsilon)S \xrightarrow{R} ()(S)S \xrightarrow{R} ()(\varepsilon)S \xrightarrow{R} ()()S \xrightarrow{R} ()()(\varepsilon)S \xrightarrow{R} ()()() \varepsilon = ()()() \quad (9)$$

などがある. 直感的には  $\llbracket D \rrbracket$  は括弧の対応がきちんと取れている語の集合である. なお, 数式に用いる括弧  $()$  と混同しないように太字の括弧  $\mathbf{()}$  を用いている.  $\llbracket D \rrbracket$  は Dyck 言語とも呼ばれ, 正則言語ではないが文脈自由言語である例として代表的な言語となっている.

### 2.3 Dömösi-Horváth-Ito 予想

原始語と文脈自由言語の定義がそろったので, 本稿で解説する Dömösi-Horváth-Ito 予想を以下に述べよう.  
**予想 (Dömösi-Horváth-Ito).**  $A$  が 2 つ以上の文字を含む ( $\#A \geq 2$ ) とき, 全ての原始語の集合  $\mathbb{Q}_A$  は文脈自由ではない.

この予想のように「ある言語がある言語族に属さない」という形の命題を否定命題と呼ぶ. 形式言語理論には否定命題の一般的な証明技法がそれぞれの言語族についていくつか種類がある. 続く 3 章では母関数を用いた証明技法 (Chomsky-Schützenberger の定理) を, 最後の 4 章では繰り返し構造を用いた証明技法 (ポンピング補題) と測度を用いた証明技法を紹介する.

予想の初出である Dömösi, Horváth, Ito らの 1991 年の論文 (4) から, 原始語と文脈自由言語に関する研究が続けられているが現時点で未解決である.

**補足 2.1.**  $A$  が単一の文字  $a$  のみ含む場合 ( $A = \{a\}$ ) は, 定義より

$$\mathbb{Q}_{\{a\}} = \{a\} \quad (10)$$

となってしまうため, これは明らかに文脈自由言語 (文法  $G = (\{S\}, \{a\}, \{(S, a)\}, S)$  の言語) となる.

$w$  が原始語であるかどうかの判定は, 素朴に  $w$  の全ての真の接頭辞 ( $w = uv$  ( $|u|, |v| \geq 1$ ) となる  $u$ ) に対してその繰り返し  $w$  と一致するかどうかを判定すれば良い. そのため  $w$  が  $\mathbb{Q}_A$  に属するかどうかの判定は線形領域で計算可能であるため,  $\mathbb{Q}_A$  が文脈依存言語であることは自明である.

## 3 原始語と無曖昧文脈自由言語

原始語の集合が文脈自由言語かどうかについてはまだ未解決であるが, 文脈自由言語よりも表現力の弱い言語クラスである無曖昧文脈自由言語ではないことが知られている (文献 (5) の 8.3 章を参照せよ). 本節では, 既存の証明とは少し違った形で (不完全ではあるが) 原始語の集合が無曖昧文脈自由でないことを示すためのアイデアを解説したいと思う. まずは無曖昧文脈自由言語の定義 (およびそれに必要な諸定義) から述べよう.

### 3.1 無曖昧文脈自由言語と Chomsky-Schützenberger の定理

**定義 3.1** (最左導出). 文脈自由文法  $G = (V, A, \xrightarrow{R}, S)$  導出規則  $\xrightarrow{R}$  は次の最左導出規則  $\xrightarrow[\text{left}]{R} \subseteq (V \cup A)^* \times (V \cup A)^*$  に自然に拡張することができる.

$$\xrightarrow[\text{left}]{R} \triangleq \{(uXw, uvw) \mid u \in A^*, v, w \in (V \cup A)^* X \in V \text{ such that } X \xrightarrow{R} v\}. \quad (11)$$

$u \xrightarrow[\text{left}]{R}^* v$  となる場合に  $G$  は  $u$  から  $v$  を最左導出すると言う.

**定義 3.2** (文脈自由文法と曖昧性). 文脈自由文法  $G = (V, A, R, \xrightarrow{R})$  が無曖昧であるとは, 任意の  $w \in \llbracket G \rrbracket$  について  $S \xrightarrow[\text{left}]{R}^* w$  となる最左導出が唯一存在することである. すなわち任意の  $w \in \llbracket G \rrbracket$  に対して

$$S \xrightarrow[\text{left}]{R} w_1 \xrightarrow[\text{left}]{R} w_2 \xrightarrow[\text{left}]{R} \cdots \xrightarrow[\text{left}]{R} w_n \xrightarrow[\text{left}]{R} w \quad (12)$$

となる  $n \in \mathbb{N}$  と  $w_1, \dots, w_n \in (V \cup A)^*$  が唯一存在することを言う.

文脈自由言語  $L$  が無曖昧であるとは,  $L = \llbracket G \rrbracket$  となる無曖昧な文脈自由文法  $G$  が存在することを言う. 文脈自由言語  $L$  が本質的に曖昧であるとは,  $L = \llbracket G \rrbracket$  となる無曖昧な文脈自由文法  $G$  が存在しないことを言う.

例えば前節で紹介した Dyck 言語を生成する文法

$$G_D = (\{S\}, \{(\cdot)\}, \{(S, \varepsilon), ((S)S)\}, S) \quad (13)$$

は無曖昧である (実際には機能法などを用いて示す) が, 同じく Dyck 言語を生成する次の文法

$$G_{D'} = (\{S\}, \{(\cdot)\}, \xrightarrow{R} = \{(S, \varepsilon), (S, SS), (S, (S))\}, S) \quad (14)$$

は曖昧である. なぜなら,  $G_{D'}$  においては例えば  $()$  という語に対して

$$S \xrightarrow[\text{left}]{R} SS \xrightarrow[\text{left}]{R} (S)S \xrightarrow[\text{left}]{R} (\varepsilon)S \xrightarrow[\text{left}]{R} (\varepsilon)\varepsilon = (), \quad (15)$$

$$S \xrightarrow[\text{left}]{R} SS \xrightarrow[\text{left}]{R} \varepsilon S \xrightarrow[\text{left}]{R} \varepsilon(S) \xrightarrow[\text{left}]{R} \varepsilon(\varepsilon) = () \quad (16)$$

という異なる 2 つの最左導出列が存在するためである.

「ある言語が無曖昧文脈自由言語でない」という否定命題を示すための強力な道具として, 言語の数え上げ関数に対する次の古典的な Chomsky-Schützenberger の定理がある.

**定義 3.3** (数え上げ関数と母関数). 言語  $L$  について, その数え上げ関数  $\Gamma_L : \mathbb{N} \rightarrow \mathbb{N}$  とは

$$\Gamma_L(n) \triangleq \#\{w \in L \mid |w| = n\} = \#(L \cap A^n) \quad (17)$$

で定義される関数である. 言語  $L$  について

$$\sum_{n \geq 0} \Gamma_L(n) z^n \quad (18)$$

で定義される 1 変数の有理数係数級数を  $L$  の母関数と呼ぶ.

**定理 3.1** (Chomsky-Schützenberger(3)). 無曖昧文脈自由言語の母関数は代数関数.



$\overline{\mathbb{Q}'_A}$  の数え上げ関数は

$$\Gamma_{\overline{\mathbb{Q}'_A}}(n) = 0 \iff n = 1 \text{ または } n \text{ は素数} \quad (30)$$

という性質を満たすことがわかる．つまり  $\overline{\mathbb{Q}'_A}$  の数え上げ関数の零点の集合は (1 を含むことを除いて) 素数全体と一致するのである．

級数  $F(z) = \sum_{n=0}^{\infty} c_n z^n$  において、係数  $c_n$  が 0 となる  $n$  を  $F$  の消失点と呼び、消失点全体の集合を  $\mathcal{Z}(F)$  と置くことにしよう．ある種の関数のクラスにおいては、その Taylor 展開  $F(z)$  の消失点の集合  $\mathcal{Z}(F)$  が等差数列の有限和となることが知られている．

**定義 3.4** (周期的集合)．自然数  $c, d \geq 0$  について  $\{cn + d \mid n \geq 0\}$  の形でかける自然数の集合を等差数列と呼ぶ．自然数の等差数列の有限和を周期的集合と呼ぶ．

**補足 3.1**．上の定義では、公差  $c = 0$  の場合も許しているため 1 点集合  $\{d\}$  は等差数列となり、よって任意の有限集合は周期的集合となる．

**定理 3.2** (Skolem-Mahler-Lech (文献 (1) を参照せよ))．有理関数  $f(z)$  の Taylor 展開の消失点の集合  $\mathcal{Z}(f)$  は周期的集合．

上述した Chomsky-Schützenberger の定理から、無曖昧文脈自由言語の母関数は代数関数となる．無曖昧文脈自由言語は正則言語との和において閉じているため、もしも  $\mathbb{Q}_A$  が無曖昧文脈自由言語であれば  $\mathbb{Q}'_A = \mathbb{Q}_A \cup \{a^n \mid a \in A, n \geq 2\}$  も  $(\{a^n \mid a \in A, n \geq 2\})$  が正則言語であり、無曖昧文脈自由言語は正則言語との非交和について閉じているため) 無曖昧文脈自由言語となり、代数関数は差において閉じているため  $\overline{\mathbb{Q}'_A}$  の母関数も代数関数となるはずである．一方  $\overline{\mathbb{Q}'_A}$  の消失点の集合は素数全体であるため周期的集合にはならない (素数は無限個あり、さらに素数の集合は無限長の等差数列を含まないため)．そのため、上述した Skolem-Mahler-Lech の定理の言明 (有理関数の消失点は周期的集合) を代数関数にまで一般化することができれば、 $\overline{\mathbb{Q}'_A}$  の母関数の消失点に関する上記の考察から  $\mathbb{Q}_A$  の本質的な曖昧性を示すことができる．Skolem-Mahler-Lech の定理の一般化については、代数関数を含む広い関数クラス (ホロノミック関数) について Bell ら (2) が技術的な条件付きで一般化について成功している．完全な一般化については今後の課題である．

#### 4 関連研究と課題

前節では原子語の集合の無曖昧性についての解説を行った．しかし、肝心の予想である「原子語の集合が文脈自由言語ではない」は未解決である．文脈自由言語の母関数は一般に超越関数になり得るため (文献 (7) を参照せよ)、無曖昧文脈自由言語における Chomsky-Schützenberger の定理のような「否定命題に対する強力な道具」が欠けていることが、予想の解決を困難にしている原因の 1 つである．

また、予想の解決を困難にしている他の原因としてに原子語の集合が「非常に大きい」ことが挙げられる．ここで言う「非常に大きい」の意味は後で説明するが、直感的には「とても多くの語を含む」と考えてもらいたい．文脈自由言語の否定命題に用いられる道具としては次の有名なポンピング補題がある．

**補題 4.1** (ポンピング補題)．任意の文脈自由言語  $L$  に対して、ある自然数  $p \geq 1$  が存在し、 $|u| \geq p$  となる任意の  $u \in A^*$  は  $u = vxyz$  と次の条件を満たす  $v, w, x, y, z \in A^*$  に分解できる：

$$(1) |w| + |y| \geq 1 \quad (2) |wxy| \leq p \quad (3) \text{ 任意の自然数 } i \geq 0 \text{ について } vw^i xy^i z \in L \quad (31)$$

ポンピング補題の性質上、否定命題の対象となる言語にはあまり語が含まれていないことが好ましい．なぜなら言語  $L$  が非常に多くの語を含む (= 大きい) 場合、ポンピング補題の言明で保証されるべき条件 ( $xy^i z \in L$  for all  $i \geq 0$ ) が成立しやすくなるためである．実際、原始語の集合はポンピング補題で要求される性質を満たすため、ポンピング補題を用いて予想を示すことはできないのである．ポンピング補題にはさまざま

な拡張があるが、 $\mathbb{Q}_A$  はことごとくそれらの補題をかいくぐるのである (詳しくは文献 (5) の 4 章を参照せよ).

さて、ここで本節冒頭で述べた「 $\mathbb{Q}_A$  は非常に大きい」の意味を説明することにしよう. 言語  $L$  の “大きさ” を測る尺度として「ランダムに選んだ語が  $L$  に属する確率」を表す

$$\mu(L) = \lim_{n \rightarrow \infty} \frac{\#(L \cap A^n)}{\#(A^n)} \quad (32)$$

を用いることはしばしばある ( $\mu(L)$  を「言語  $L$  の測度」と呼ぶ). 原始語の集合  $\mathbb{Q}_A$  については性質 (28) が成り立つため、

$$\limsup_{n \rightarrow \infty} \frac{\#(\mathbb{Q}_A \cap A^n)}{\#(A^n)} = 1 \quad (33)$$

であることが簡単に示せるが、実際にはより詳細な解析を行うことで  $\mu(\mathbb{Q}_A) = 1$  を示すことができる. そういった意味で  $\mathbb{Q}_A$  は「非常に大きい」のである. その他にも、( $a^n$  ( $n \neq 2$ ) という形の 1 つの文字の繰り返しを除く) 全ての語は 2 つの原始語の接続に分解できる、つまり

$$\mathbb{Q}_A^2 = A^* \setminus \{a^n \mid n \neq 2, a \in A\} \quad (34)$$

が成り立つということもわかって (文献 (6)) おり、そういった意味でも  $\mathbb{Q}_A$  は大きな集合と言える.

ポンピング補題のような否定命題に対する既存の道具では、予想を解決することはこれらの理由で困難だと思われる. 今後の課題として、 $\mathbb{Q}_A$  のような大きな言語においても有用な、否定命題に対する道具を開発する必要があると著者は考えている. 著者はこれまで正則言語における研究を行っており、測度 1 の言語 (= 非常に大きい言語) に対する次の形の否定命題の道具を開発した:

**定理 4.1** (文献 (7) を参照せよ).  $A$  上の正則言語  $L$  に対して、次の条件は同値:

1.  $L$  は測度 1 ( $\mu(L) = 1$ )
2. ある語  $w \in A^*$  が存在して  $A^*wA^* \subseteq L$

例えばこの定理によって Dyck 言語や回文の集合、さらには「素数の 1 進数表記の集合  $\{a^p \mid p \text{ は素数}\}$ 」(及びこれらの補集合の言語) が正則言語ではないことが簡単に示せる (文献 (8)).

文脈自由言語は正則言語にくらべはるかに難しい構造を持っているため、正則言語の理論や道具を文脈自由言語の世界に輸入することは容易ではないが、先に述べた「大きい言語に対する文脈自由言語の否定命題のための道具」を開発するためには測度 1 の文脈自由言語に対する定理 4.1 の拡張について研究を進めていくことは重要であり、著者の課題である.

#### 参考文献

- (1) Bell, J. (2005): A generalised Skolem-Mahler-Lech theorem for affine varieties, *Journal of the London Mathematical Society*, Volume 73, Issue 2, pp. 367–379.
- (2) Bell, J., Burris, N. S., and Yeats, K. (2012): On the set of zero coefficients of a function satisfying a linear differential equation, *Mathematical Proceedings of the Cambridge Philosophical Society*, Volume 153, Issue 2, pp. 235–247.
- (3) Chomsky, N. and Shützenberger, M. (1963): The Algebraic Theory of Context-Free Languages, *Studies in Logic and the Foundations of Mathematics*, Volume 35, pp. 118–161.
- (4) Dömösi, P., Horváth, S., and Ito, M. (1991): On the connection between formal languages and primitive words, *Proc. First Session on Scientific Communication (University of Oradea, Romania)*, pp. 59–67.
- (5) Dömösi, P., Horváth, S., and Ito, M. (2014): Context-Free Languages and Primitive Words, World

Scientific Publishing Company Pte Limited.

- (6) Reis, C. M. and Shyr, H. (1978): Some Properties of Disjunctive Languages on a Free Monoid, *Information and Control*, Volume 37, Issue 3, pp. 334–344.
- (7) 新屋良磨 (2017): オートマトン理論再考, コンピュータソフトウェア, 34 巻, 3 号, pp. 3–35.
- (8) 新屋良磨 (2017): 言語の測度に基づく非正規性の証明技法, コンピュータソフトウェア, 34 巻, 1 号, pp. 119–124.