

半構造データに対する効率良い近接パターン照合アルゴリズム

和佐 州洋[†] 金田 悠作[†] 宇野 毅明^{††} 有村 博紀[†]

[†] 北海道大学大学院情報科学研究科 〒060-0814 札幌市北区北14条西9丁目

^{††} 国立情報学研究所 〒101-8430 東京都千代田区一橋2-12

E-mail: †{wasa,y-kaneta,arim}@ist.hokudai.ac.jp, ††uno@nii.ac.jp

あらまし 本稿では、近接パターン照合問題のひとつとして、木に対するグラフモチーフ問題に着目する。そのための重要な副問題である k -部分木列挙問題を考察し、その効率良い列挙アルゴリズムを与える。本論文では、問題の定義とわれわれのアルゴリズムを述べ、アルゴリズムの正当性と時間計算量について考察する。

キーワード グラフマイニング, 近接パターン照合, 双方向ビットカウンタ, 部分木列挙, 文字列アルゴリズム

1. はじめに

文字列や、木、グラフのような非定型構造データを半構造データと呼ぶ。近接パターン照合問題(proximity pattern match problem) [3], [4], [6], [7] は、与えられた半構造データ上で、制約を満たしたアイテムの出現を発見する問題である。これは、検索とマイニングの中間に位置する問題であり、大規模文書検索や生物情報学で研究されている近接文字列照合問題 [6], [7] やグラフモチーフ問題 [3], [4] は近接パターン照合問題の例である。これらの問題では、パターンに対する最もゆるい制約として、連結性制約が用いられる。

グラフモチーフ問題(graph motif problem) [4] は、近接パターン照合問題のひとつである。これは、頂点にラベルを持つグラフ $G = (V(G), E(G))$ (テキスト) とラベルの多重集合 M (モチーフ) を入力として受け取り、 G の連結部分グラフ S で、そのラベルの多重集合が M に一致するようなものを見つける問題である。

本稿では、木に対するグラフモチーフ問題について着目し、その重要な副問題である k -部分木列挙問題(k -subtree enumeration problem) の効率良いアルゴリズムについて考察する。このとき、 G は根付き木であり、 S は G の(連結)部分木である。主結果として、頂点数 n の木と正整数 $k \geq 1$ が与えられた時、解1個あたりならし定数時間で頂点数 k の部分木を重複なく全て列挙する効率良いアルゴリズムを与える。このような部分木を用いたDB処理として、Saito と、Morishita [8] などがある。

2. 準備

2.1 木の定義

根付き木 $T = (V(T), E(T))$ とは、ノード集合 V と辺集合 $E \in V^2$ から成り立つ。 V の要素数は $n = |V|$ で表す。また、根付き木のサイズを $|T| = |V| = n$ と定義する。各辺 $(u, v) \in E$

に対して、 u を v の親ノードといい、 $\text{pnode}(v) = u$ と書く。また、 v を u の子ノードという。 u の子ノードの数を u の次数という。根ノード $\text{root}(T)$ 以外は1つの親ノードを持つ。子ノードのないノードを葉という。根ノードと葉以外のノードを中間ノードという。同じ親ノードを持つノードを兄弟という。 T の葉の集合を $L(T)$ とする。ノード v に対して、 $v[i]$ を v の i 番目の子供とする。

兄弟間に大小関係のある根付き木を順序木という。大きいノードほど左側にあるとする。以下、順序木を木という。ノード u, v に対して、 u から v への長さ $k \geq 0$ のパスとは、ノードの列 $\pi = (v_0 = u, v_1, \dots, v_k = v)$ で、各 $i = 1, \dots, k$ に対して、 $(v_{i-1}, v_i) \in E$ であるものである。木 T で、任意の2ノードの間にパスが存在する時、 T は連結(connected)であるという。ノード v の先祖とは、ノード v と根ノードを結ぶパス上の v を除くノード集合をいう。

ここで、左優先深さ優先探索順(以下、DFS順という)を考える。木 $T = (V(T), E(T))$ の各ノード v のDFS順の番号を $\text{dfs}(v) \in \{1, \dots, |V|\}$ とおき、 $\text{dfs}(v)$ を v の番号という。以下では、文脈から明らかならば、ノードと番号を同一視する。よって、頂点数 n の順序木 T の頂点集合を $V(T) = \{1, \dots, n\}$ とみなす。

2.2 部分木の定義

木 $T = (V(T), E(T))$ を固定する。 $V(T)$ の部分集合 S が誘導する T の連結な部分木を $T(S) = (S, E(S))$ とする。ここで、 $E(S) = \{(u, v) \in E(T) \mid \{u, v\} \subseteq S\}$ とする。以後、文脈から明らかならば、部分集合 S を部分木ということがある。部分木 S のサイズをノード数 $|S|$ と定義する。サイズ k の部分木を k -部分木という。 $T(v)$ は、ノード v を根ノードとする T の部分木を表す。ノード u の子孫 $\text{dec}(u)$ とは、 $\text{dec}(u) = \{v \in T(u) \mid u \neq v\}$ で表されるノード集合をいう。

2.3 列挙アルゴリズム

列挙とは、与えられた条件を満たす解を漏れ無く重複なく重

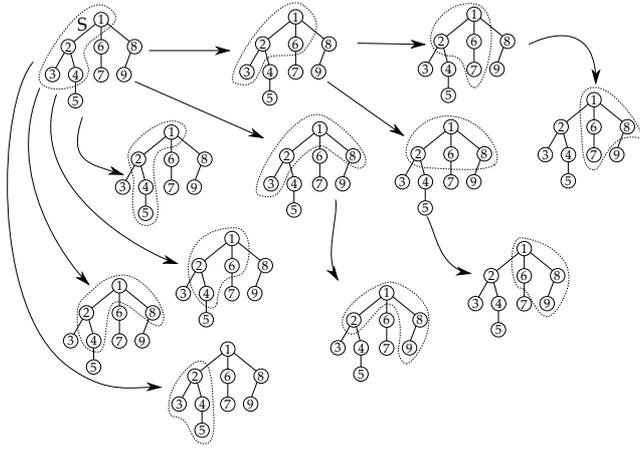


図 1 探索木の例.

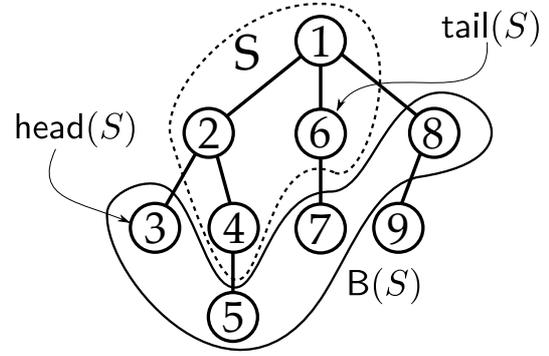


図 2 番号 1, 2, 4, 6 をノードとする部分木 S に対する $\text{tail}(S)$, $B(S)$, $\text{head}(S)$ の例である. $\text{tail}(S)$ は番号 6 を持つノード, $B(S)$ は番号 3, 5, 7, 8 を持つノードの集合, $\text{head}(S)$ は番号 3 を持つノードである.

複することである [1], [9]. 列挙アルゴリズムの計算量は, 解一個あたりを出力するのに必要な計算時間 (遅延) で評価する. 本論文では, Ferreira と, Grossi, Rizzi ら [5] が研究したグラフ上の k -部分木列挙問題を, 入力グラフ G が木の場合に制限して次の問題を考察する. 以下では, k を任意の正整数とする.

問題 1 (k -部分木列挙問題). 入力として与えられた木 $T = (V(T), E(T))$ においてサイズ k の部分木を重複なくすべて列挙する.

計算モデルとしては, 通常の命令を持つ RAM [2] を用いる.

3. 部分木同士の親子関係

われわれの列挙アルゴリズムは逆探索手法 [1], [9] に基づく. これは, 与えられた解空間上に探索木 (家系木) を定義し (図 1), 根から出発して, バックトラックしながら, 全ての解を列挙する手法である. $T = (V(T), E(T))$ をサイズ n の任意の根付き木とする. 以下では, 与えられた根付き木のノードに対して, DFS 順の番号が振られていると仮定する. ノードとその番号を同一視する. よって, ノード u がノード v より大きいなどという. 初めに, 出発点となる初期木を定義する. その準備として, 連続木を定義する.

定義 1 (連続木). 木 T の任意のノード $r \in V(T)$ が $|T(r)| \geq k$ を満たすとき, サイズ k の連続木は $\text{con}_k(r) = \{r, \dots, r+k-1\}$ である.

ここで, ノードと番号を同一視していることに注意されたい.

定義 2 (初期木). 木 T の連続木 $\text{con}_k(1)$ を特に初期木といい, $\perp(k) = \text{con}_k(1) = \{1, \dots, k\}$ と書く.

補題 1. 木 T の任意のノード r が $|T(r)| \geq k$ を満たすとき, $S = \text{con}_k(r)$ はサイズ k の連結した部分木である.

証明. 仮定から, S に含まれるノード全ては $T(r)$ に含まれる. ノードの番号付けが DFS 順であることを考慮すると, これらから S は明らかに連結な部分木である. ■

次に, 与えられた部分木 $S \subseteq V(T)$ の「親」となる部分木 $P(S)$

を定義する. T 上のサイズ $k \leq |T|$ の部分木 S に対して $\text{tail}(S)$ と, $\text{head}(S)$, $B(S)$ を以下のように定義する (図 2).

定義 3 (境界). S に含まれるノードの子ノードで, かつ S に含まれないノードの集合を S の境界といい, $B(S) = \{v \in V(T) \mid u \in S, (u, v) \in E(T), v \notin S\}$ とする.

定義 4 (先頭). $B(S)$ の中で最も小さいノードを先頭といい, $\text{head}(S) = \underset{v \in B(S)}{\text{argmin}} v$ とする.

定義 5 (末尾). S の中で最も大きいノードを末尾といい, $\text{tail}(S) = \underset{v \in S}{\text{argmax}} v$ とする.

補題 2. 木 T のサイズ k の部分木 S における $\text{tail}(S)$ は葉である. ただし, $k \leq |T|$ とする.

証明. $v = \text{tail}(S)$ が S において中間ノードであるとする. 子ノード u を持つ. このとき, $u > v$ であるので, $\text{tail}(S)$ は S 中で最大という定義に反する. よって, $\text{tail}(S)$ は葉である. ■

定義 6 (親木). $n \geq k$ を満たす木 T のサイズ k の部分木 S の親木 $P(S)$ は, 以下で定義される. ただし, $S \neq \perp(k)$ とする.

$$P(S) = \begin{cases} (S \setminus \{\text{tail}(S)\}) \cup \{\text{head}(S)\}, & \text{if } S \text{ が連続木でない} \\ (S \setminus \{\text{tail}(S)\}) \cup \{\text{pnode}(\text{root}(S))\}, & \text{if } S \text{ が連続木である} \end{cases}$$

定義 7 (子木). 部分木 R に対して, ある手続きを行い出力された部分木を S とする. このとき, $P(S) = R$ となるような S を R の子木という. また, このような R と S に対して, R と S に親子関係がある, という.

補題 3. 木 T の任意のノード r が $|T(r)| \geq k$ を満たすとする. このとき, 根ノードが r でサイズ k の部分木 S を考える. (1) S がサイズ k の連続木ではないならば, その時に限り, (2) $\text{head}(S) < \text{tail}(S)$ である.

証明. (1) \Rightarrow (2): 木 T の連続木でない部分木 $S = \{x_1, x_2, \dots, x_k\}$ を考える. ただし, $x_1 < x_2 < \dots < x_k$ とす

る。このとき、 $1 \leq i < k$ に対して、 $x_i \in S$ かつ $y = x_i + 1 \notin S$ かつ $y < x_{i+1}$ となる i が存在する。ここで、DFS の番号付けから、 y は x_i の先祖の子、または x_i の子となる。また、 $x_k = \text{tail}(S)$ である。よって、このような y の中で最小のノードが定義 4 より $\text{head}(S)$ となるので、 $\text{head}(S) < \text{tail}(S)$ となる。

(2) \Rightarrow (1): S が連続木であるとする。このとき、 $\text{tail}(S) = r + k - 1$ である。 $S = \{r, \dots, \text{tail}(S)\}$ と連続であり、かつ $\text{head}(S) \notin S$ であるので、 $r < \text{head}(S)$ から $\text{tail}(S) < \text{head}(S)$ となる。対偶より、 $\text{head}(S) \neq \text{tail}(S)$ から $\text{head}(S) < \text{tail}(S)$ であるとき、 S は連続木ではない。 ■

また、補題 3 より、自明に系 1 が導かれる。

系 1. 木 T で任意の連続木でない部分木 S に対して、 $(\text{tail}(S), \text{head}(S)) \notin E(T)$ である。

補題 4. T のサイズ k の任意の部分木 S に対して、もし、 S が初期木でないならば $P(S)$ は正しく定義される。

証明. $P(S)$ は S が連続木であるときと、そうでないときに異なる。よって、場合分けをして証明する。

S が連続木でないとき: S において、 $\text{head}(S)$ は S のノードの子ノードであり $\text{tail}(S)$ は補題 2 より葉ノードである。また、 S は連続木ではないので、系 1 より $\text{head}(S)$ が $\text{tail}(S)$ の子ノードとなることはない。よって、 $\text{tail}(S)$ を除き $\text{head}(S)$ を加えた部分木は連結であるので、 $P(S)$ は正しく定義される。

S が連続木であるとき: S において、 $\text{tail}(S)$ は補題 2 より葉ノードである。また、 S は初期木ではないので、 $\text{pnode}(\text{root}(S))$ は必ず存在し、また、 $\text{root}(S)$ の親ノードである。よって、 $\text{tail}(S)$ を除き、 $\text{pnode}(\text{root}(S))$ を加えた部分木は連結であるので、 $P(S)$ は正しく定義される。 ■

4. 列挙アルゴリズム

前章では、子木から親木を求める方法について考察した。本章では、親木から子木を求める方法について、子木が連続木であるときとそうでないときに分けて考察する。そのあとで、木 T の k -部分木を定数遅延時間で列挙するアルゴリズムを示し、その正当性と時間計算量について考察する。

4.1 連続木でない子木の生成

木 T の部分木 R から子木 S を求めるには、子木 S からその親木である R を求める逆の手順を踏めば良い。つまり、適当なノード $u \in L(R)$ を削除し、 R に含まれるノードの子ノードでかつ $\text{tail}(R)$ よりも大きく R に含まれないノード v を追加すれば良い。しかし、ノード u, v の候補は一つとは限らない。そこで、 u の候補となるノードの集合を $\text{can}(R)$ 、 v の候補となるノードの集合を $\text{LB}(R)$ とし、以下のように定義する (図 3)。

定義 8. 部分木 R の葉の中で、 $\text{head}(R)$ よりも小さいノードの集合を $\text{can}(R) = \{u \in L(R) \mid u < \text{head}(R)\}$ とする。

定義 9. 部分木 R のボーダーの中で、 $\text{tail}(R)$ よりも大きいノードの集合を $\text{LB}(R) = \{v \in B(R) \mid \text{tail}(R) < v\}$ とする。

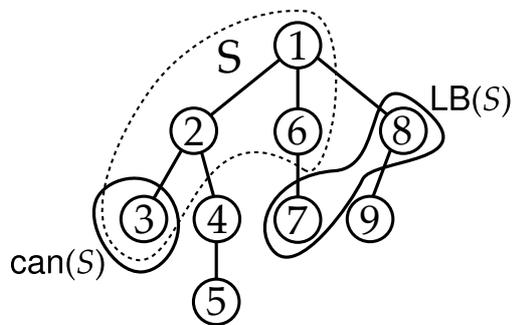


図 3 部分木 S に対する $\text{can}(S)$ と $\text{LB}(S)$ の例. $\text{can}(S)$ は番号 3 を持つノードの集合、 $\text{LB}(S)$ は番号 7, 8 を持つノードの集合である。

定義 10. 木 T の部分木 R に対して、 $u \in \text{can}(R)$ 、 $v \in \text{LB}(R)$ 、 $(u, v) \notin E(T)$ であるような u, v に対して、

$$C_1(R) = (R \setminus \{u\}) \cup \{v\}$$

を R の子木とする。

補題 5. 木 T の部分木を R とする。ここで、 $\text{can}(R) \neq \emptyset$ かつ $\text{LB}(R) \neq \emptyset$ のとき、 $C_1(R) = (R \setminus \{u\}) \cup \{v\}$ は木 T の連結な部分木である。ただし、 $u \in \text{can}(R)$ 、 $v \in \text{LB}(R)$ 、 $(u, v) \notin E(T)$ とする。

証明. $(u, v) \notin E(T)$ であり $\text{can}(R) \subseteq L(R)$ 、 $\text{LB}(R) \subseteq B(R)$ であるので、 $C_1(R)$ は連結な部分木である。 ■

補題 6. 木 T の任意の部分木を R とし、 T の連続木でない部分木を S とする。ここで、(1) $u \in \text{can}(R)$ かつ $v \in \text{LB}(R)$ かつ $(u, v) \notin E(T)$ であるような u, v に対して、木 T の部分木 S が $S = C_1(R) = (R \setminus \{u\}) \cup \{v\}$ が R の子木であるならば、その時に限り、(2) $\text{head}(S) \in \text{can}(R)$ かつ $\text{tail}(S) \in \text{LB}(R)$ である。

証明. (1) \Rightarrow (2): 補題 5 より、 S は連結な部分木である。対偶を用いて示す。 $\text{head}(S) \notin \text{can}(R)$ または $\text{tail}(S) \notin \text{LB}(R)$ であるとする。ここで、 $P(S) = (S \setminus \{\text{tail}(S)\}) \cup \{\text{head}(S)\}$ である。しかし、 $R = (R \setminus \{v\}) \cup \{u\}$ であるが、 $u \in \text{can}(R)$ かつ $v \in \text{LB}(R)$ であるので、 $P(S)$ と R は一致しない。よって、 S は R の子木でない。よって、(1) \Rightarrow (2) が示された。

(2) \Rightarrow (1): $\text{head}(S) \in \text{can}(R)$ かつ $\text{tail}(S) \in \text{LB}(R)$ と仮定する。ここで、 $u = \text{head}(S)$ かつ $v = \text{tail}(S)$ とおく。このとき、 S は連続木ではないので、系 1 から $(v, u) \notin E(T)$ である。よって、 S は連結な部分木である。ここで、

$$\begin{aligned} P(S) &= (S \setminus \{\text{tail}(S)\}) \cup \{\text{head}(S)\} \\ &= ((R \setminus \{\text{head}(S)\}) \cup \{\text{tail}(S)\}) \setminus \{\text{tail}(S)\} \cup \{\text{head}(S)\} \\ &= R \end{aligned}$$

であるので、 S と R は親子関係である。よって、(2) \Rightarrow (1) が示された。 ■

補題 7. T 上の連続木でないサイズ k の S について、 $(u, v) \notin E(T)$ である。ただし、 $u \in \text{can}(S)$ 、 $v \in \text{LB}(S)$ とする。

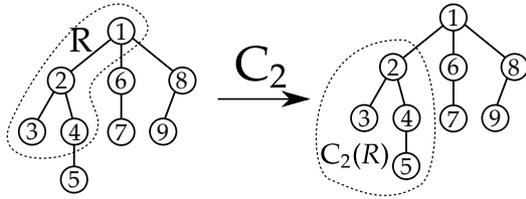


図 4 擬連続木 R から得られる連続木 $C_2(R)$ の例.

証明. S は連続木ではないので, 補題 3 より $\text{head}(S) < \text{tail}(S)$ である. よって, $\text{can}(S)$ の定義から, $\text{tail}(S)$ と $\text{tail}(S)$ の先祖は $\text{can}(S)$ に含まれない. 一方, $v \in \text{LB}(S)$ は定義より $v > \text{tail}(S)$ であるので, v の親ノードは $\text{tail}(S)$ または, $\text{tail}(S)$ の先祖である. これらから, 任意の v は任意の u の子ノードでない. ■

補題 8. T 上のサイズ k の連続木 S について, $(u, v) \in E(T)$ であるのは, $u = \text{tail}(S)$ のときである. ただし, $u \in \text{can}(S)$, $v \in \text{LB}(S)$ とする.

証明. 任意の $v \in \text{LB}(S)$ に対して, v の親ノード y が S 中で葉であるときとそうでないときの場合分けし, y と u が一致する場合について考察する. $y \notin L(S)$ のとき, $\text{can}(S)$ の定義から $u \in L(S)$ であるので, 明らかに $(u, v) \notin E(T)$ である. $y \in L(S)$ のとき, $\text{LB}(S)$ の定義から y は $\text{tail}(S)$ と等しくなくては行けない. よって, $(u, v) \in E(T)$ であるのは, $y = u$, つまり $u = \text{tail}(S)$ のときである. よって, これらから, v が u の子ノードであるのは, $u = \text{tail}(S)$ のときである. ■

4.2 連続木である子木の生成

次に, 初期木でない連続木 S をある木の子木として求める方法について考察する. その準備として, 擬連続木を導入する.

定義 11 (擬連続木). 木 T の部分木 R に対して, 部分木 R の根ノード $\text{root}(R)$ が $|T(v)| \geq k$ を満たすノード v を子ノードとしてただひとつに持ち, かつ, $\text{tail}(R) - v = k - 2$ であるような R を擬連続木という.

定義 12. 木 T のサイズ k の擬連続木であるような部分木 R に対して,

$$C_2(R) = (R \setminus \{\text{root}(R)\}) \cup \{\min \text{LB}(R)\}$$

を R の子木とする.

擬連続木 R に対して, $C_2(R)$ とすることで, 連続木が得られることを以下で示していく 4.

補題 9. 木 T の任意の部分木 S に対して, $\text{LB}(S) \neq \emptyset$ であるとき, $\text{tail}(S) + 1 = \min \text{LB}(S)$ である.

証明. $l = \text{tail}(S) + 1$ とする. l は, DFS 番号付けと $\text{tail}(S)$ の定義から $\text{tail}(S)$ の子ノード, もしくは $\text{tail}(S)$ の先祖の子ノードである. また, $\text{tail}(S)$ が S で最大のノードであるので, $l \notin S$ である. よって, $l \in \text{LB}(S)$ である. ここで, ノードの番号付けが整数であること考えると, $\text{tail}(S)$ と l の間にノードは存在しないので, 定義 9 から $l = \min \text{LB}(S)$ である. ■

Algorithm 1 木 T のサイズ k の部分木 S を列挙

```

1: procedure MAIN( $T, k$ )
2:   Input:  $T$ : 木,  $k$ : サイズ;
3:    $T$  に対して各ノードに対し, DFS 順に番号をつける;
4:    $S \leftarrow$  初期木;
5:    $\text{can}(S)$ ,  $\text{LB}(S)$  の計算;
6:   ENUMSUBTREECONSTANT( $S, T, k$ );

1: procedure ENUMSUBTREECONSTANT( $R, U, k$ )
2:   Input:  $R$ : 部分木,  $U$ : 木,  $k$ : サイズ;
3:   print  $R$ ;
4:   if  $R$  が擬連続木である then
5:      $R' \leftarrow C_2(R)$ 
6:      $\text{can}(R')$ ,  $\text{LB}(R')$  の再計算;
7:     ENUMSUBTREECONSTANT( $R', U, k$ );
8:   for each  $u \in \text{can}(R)$  do
9:     for each  $\{v \in \text{LB}(R) \mid (\text{tail}(R), v) \notin E(R)\}$  do
10:       $S \leftarrow (R \setminus \{u\}) \cup \{v\}$ ;
11:       $\text{can}(S)$ ,  $\text{LB}(S)$  の再計算;
12:      ENUMSUBTREECONSTANT( $S, U, k$ );

```

補題 10. 木 T の擬連続木 S に対して, $|S| \geq k$ のとき, $C_2(S)$ は連続木である.

証明. 補題 9 より, $\text{tail}(S) + 1 = \min \text{LB}(S)$ であるので, $C_2(S)$ は定義より連続木である. ■

補題 11. 木 T のサイズ k の初期木でない連続木 S には $C_2(R) = S$ となる擬連続木 R がただひとつ存在する.

証明. $C_2(R) = S$ となるような擬連続木 R の根ノードを α とする. 定義 11 より, $S \cap R = \{\alpha, \alpha + 1, \dots, \alpha + k - 2\}$ である. ここで, $A = (A \cap R) \cup \{\min \text{LB}(A)\}$ である. また, R は初期木でない連続木であるので, α には T 上で親ノードがただひとつ存在する. よって, これらから R に対応する A はただひとつ存在する. ■

4.3 アルゴリズム

これらから, 入力として与えられた木 $T = (V(T), E(T))$ の k -部分木 $S = (S, E(S))$ をならし定数時間で列挙できるアルゴリズムを Algorithm 1 に挙げる.

Algorithm 1 は与えられた木の初期木の子木を再帰的に出力する. その子木の中に擬連続木が存在すれば, その擬連続木に対応する連続木を求め, その連続木の子木を再帰的に出力することで, 全ての k -部分木を列挙するアルゴリズムである.

4.4 正当性

本節では, Algorithm1 の正当性について考察する. MAIN は木 T とサイズ k を受け取り, T 上のサイズ k の部分木をすべて列挙するメインルーチンである. 3 行目で, T に DFS 番号付けをし, 6 行目から, 初期木で ENUMSUBTREECONSTANT を呼び出している. 初期木は家系木の根となる.

ENUMSUBTREECONSTANT の 4-7 行目で連続木を擬連続木か

ら呼び出し、呼び出した連続木に対して `ENUMSUBTREECONSTANT` を再帰的に呼び出す。また、8-12 行目より R の子木を再帰的に出力している。正当性は、4-7 行目の処理で、 T に含まれる連続木を再帰的にすべて呼び出すことができることが補題 11 からわかる。また、8-12 行目の処理で、処理全体を通して家系木を構築し、 T の部分木をすべて列挙できることが補題 7,8 からわかる。

よって、これらから Algorithm1 は MAIN を実行することで、木 T とサイズ k を入力とした時、サイズ k の T 中に含まれる部分木をすべて列挙する。

4.5 時間計算量

本節では、Algorithm1 の時間計算量について考察する。ここで、部分木 S の葉集合の保存には双方向連結リスト dl_S を、 $LB(S)$ の保存には双方向連結リストを db_S を用いる。 dl_S を求めることで、 $head(S)$ を用いると、 $can(S)$ が定義より得られる。また、 $dl_S(db_S)$ の i 番目の要素は $dl_S[i](db_S[i])$ と書く。ただし、添字 i は 1 から始まるものとする。

補題 12. 木 T の初期木に対して $can(\perp(k))$ は $O(|T|)$ 時間で得られる。

証明. DFS 順に葉を $dl_{\perp(k)}$ に追加すると、 $O(|T|)$ 時間で初期木の葉の昇順のリスト $dl_{\perp(k)}$ を構築できる。 ■

補題 13. 木 T の初期木に対して $LB(\perp(k))$ は $O(|T|)$ 時間で得られる。

証明. DFS 順で $tail(\perp(k))$ に到達したあと、行き掛け順に $LB(\perp(k))$ の定義を満たすノードをリストに追加していけば、 $O(|T|)$ 時間で、 $LB(\perp(k))$ の昇順のリスト $db_{\perp(k)}$ が得られる。 ■

補題 12, 13 から、初期化として MAIN の 3-5 行目の処理には $O(|T|)$ 時間かかる。

`ENUMSUBTREECONSTANT` の 4 行目の判定には根ノードの子の数が 1 つであること、またその子ノード f と $tail(R)$ まで連番であること、つまり $tail(R) - f = k - 2$ であることを確認し、 $B(R)$ に適切なノードが入っているか確認すればよいので、 $O(1)$ 時間である。また、5 行目の初期木の生成には定数回の操作で済むので $O(1)$ 時間である。また、8-12 行目の処理で 12 行目が呼び出される間の遅延を $O(1)$ 時間にするには、 u が $tail(R)$ であるときに、 u の子であるようなノード v をすべてスキップすればよい。このスキップは u の子ノードで最も大きいノードへのアクセスが $O(1)$ 時間で出来れば、 db_R が昇順で並んでいることを考慮すると、12 行目の呼び出しの遅延は $O(1)$ 時間である。6 行目の再計算、11 行目の再計算は以下の補題に関して考察する必要がある。

補題 14. 木 T 上の擬連続木 S に対して、 $S' = C_2(S)$ を求めた後、 $can(S')$ の再計算にかかる時間は $O(1)$ 時間である。

証明. $v = \min LB(S)$ とする。 S' は連続木であるので、補題 3 より、 $head(S') > tail(S')$ となる。また、 $dl_{S'}$ は dl_S の末尾

に v が追加し、もし、 $pnode(v)$ が S で葉であった時は、 dl_S から $pnode(v)$ を除いたリストである。また、 $head(S')$ は $LB(S')$ の先頭であるので、 $LB(S')$ の操作が $O(1)$ 時間で再計算できれば、 $can(S')$ の再計算は $O(1)$ 時間で実行可能である。 ■

補題 15. 木 T の擬連続木 S に対して、 $S' = C_2(S)$ を求めた後、 $LB(S')$ の再計算にかかる時間は $O(1)$ 時間である。

証明. $v = \min LB(S)$ とする。 $LB(S')$ は $T(\text{root}(S'))$ の子孫ノードであるので、 $db_{S'}$ は db_S から子孫でないノードを削除する。この操作は、 db_S が昇順であることを考慮すると、 $O(1)$ 時間で可能である。また、 S' に新しく追加するノード v を db_S から除いてから、 v の子ノードすべてを db_S の先頭に加えれば維持できる。よって、再計算にかかる時間は $O(1)$ 時間である。 ■

補題 16. 木 T の部分木 S に対して、 $S' = C_1(S)$ を求めた後、 $can(S')$ の再計算にかかる時間は $O(1)$ 時間である。

証明. まず S の葉集合の管理に関して考察する。ここで、 u は S の子木を求めるときに削除した葉ノードであるとする。 $dl_{S'}$ は dl_S から求める場合は、以下の 2 つに場合分けをする。 $pnode(u)$ が u のみを S 中で子ノードに持つ時は、 dl_S 中の u を $pnode(u)$ に置き換える。もし、 $pnode(u)$ が S 中で 2 つ以上の子ノードを持つ時は、 u を dl_S 中から削除する。よって、 $dl_{S'}$ は dl_S から求める時、 $O(1)$ 時間で可能である。よって、 u が $head(S')$ と等しくなることを考慮すると、 $can(S')$ の再計算には $dl_{S'}$ と u を持つておけば良いので、再計算にかかる時間は $O(1)$ 時間である。 ■

補題 17. 木 T の部分木 S に対して、 $S' = C_1(S)$ を求めた後、 $LB(S')$ の再計算にかかる時間は $O(1)$ 時間である。

証明. $C_1(S)$ で S に追加したノード v が $db_S[j] = v$ であるとする。このとき、 v が木 T の中で葉かどうかで場合分けして考察する。 v が T の中で葉の時、 $LB(S)$ から $LB(S')$ に再計算される際に追加されるノードはない。 v は定義から $tail(S') = v$ となるので、 $LB(S)$ から $LB(S')$ に再計算される際、 v より小さいノードが $LB(S)$ から削除される。よって、更新された $db_{S'}$ は $db_{S'}[1] = db_S[j+1], db_{S'}[2] = db_S[j+2], \dots$ となる。 v が T の中で中間ノードの時、 $LB(S)$ から $LB(S')$ に再計算される際に追加されるノードは v の T の子ノード全てである。 v は定義から $tail(S') = v$ となるので、 $LB(S)$ から $LB(S')$ に再計算される際、 v より小さいノードが $LB(S)$ から削除される。 v の次数を d とすると、更新された $db_{S'}$ は $db_{S'}[1] = v[1], db_{S'}[2] = v[2], \dots, db_{S'}[d] = v[d], db_{S'}[d+1] = db_S[j+1], \dots$ となる。これらの操作は、リストのポインタの張替えで行えるので、 $O(1)$ 時間で実行できる。よって、再計算にかかる時間は $O(1)$ 時間である。 ■

よって、`ENUMSUBTREECONSTANT` の 6 行目の再計算は補題 14, 15, 11 行目の再計算は補題 16, 17 から $O(1)$ 時間であることがわかる。今までの議論から以下の定理が得られる。

定理 1. *Algorithm1* は前処理に $O(|T|)$ 時間, 解 1 個あたりならし定数時間で k -部分木をすべて列挙する.

5. ま と め

本稿では, 木に対するグラフモチーフ問題を着目し, そのための重要な手続きである根付き木に対する部分木列挙問題を考察した. さらに与えられた木のサイズ k の部分木を定数時間遅延で列挙するアルゴリズムを与えた.

文 献

- [1] D. Avis and K. Fukuda. Reverse search for enumeration. *Discrete Applied Mathematics*, 65:21–46, 1993.
- [2] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. The MIT Press, 2 edition, 2001.
- [3] R. Dondi, G. Fertin, and S. Vialette. Finding approximate and constrained motifs in graphs. In *the 22nd Annual Symposium on Combinatorial Pattern Matching (CPM'11)*, volume 6661 of *Lecture Notes in Computer Science*, pages 388–401. Springer, 2011.
- [4] M. Fellows, G. Fertin, D. Hermelin, and S. Vialette. Sharp tractability borderlines for finding connected motifs in vertex-colored graphs. In *Proceedings of 34th International Colloquium on Automata, Languages and Programming (ICALP'07)*, volume 4596 of *Lecture Notes in Computer Science*, pages 340–351. Springer, 2007.
- [5] R. Ferreira, R. Grossi, and R. Rizzi. Output-sensitive listing of bounded-size trees in undirected graphs. In *Proceedings of the 19th European conference on Algorithms (ESA'11)*, volume 6942 of *Lecture Notes in Computer Science*, pages 275–286. Springer, 2011.
- [6] S.-R. Kim, I. Lee, and K. Park. A fast algorithm for the generalized k -keyword proximity problem given keyword offsets. *Inf. Process. Lett.*, 91(3):115–120, 2004.
- [7] K. Sadakane and H. Imai. Fast algorithms for k -word proximity search. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, E84-A(9):2311–2318, 2001.
- [8] T. L. Saito and S. Morishita. Relational-style xml query. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, SIGMOD '08, pages 303–314, New York, NY, USA, 2008. ACM.
- [9] A. Shioura, A. Tamura, and T. Uno. An optimal algorithm for scanning all spanning trees of undirected graphs. *SIAM J. Comput.*, 26(3):678–692, 1997.