

構造圧縮された木構造データからの頻出部分木枚挙アルゴリズム

Algorithm for enumerating all frequent subtrees from structurally compressed tree-structured data

堀部 智也 *1
Tomoya Horibe糸川 裕子 *2
Yuko Itokawa内田 智之 *1
Tomoyuki Uchida鈴木 祐介 *1
Yusuke Suzuki宮原 哲浩 *1
Tetsuhiro Miyahara

*1 広島市立大学大学院情報科学研究科

Graduate School of Information Sciences, Hiroshima City University

*2 広島国際大学心理学部

Faculty of Psychology, Hiroshima International University

Structurally compressing large tree-structured data without loss of information leads to reduction in running time needed to extract structural features. We first introduce a compression tree, which represents an edge-labeled ordered tree structurally compressed by a Lempel-Ziv compression scheme. Then, given a set D of compression trees, we propose a time- and memory-efficient algorithm for enumerating all frequent subgraphs, which have tree structures, in the set U_D without decompressing compression trees in D , where U_D is the set of edge-labeled ordered trees obtained by decompressing compression trees in D . Finally, we discuss the implementation of the proposed algorithm on a computer, explain the experimental results obtained by applying it to a set of synthetic edge-labeled ordered trees, and provide discussion on its evaluation.

1. はじめに

Web 文書, L^AT_EX ソースファイル, 自然言語などの木構造データは順序木で表現されることが多い。また, それらの木構造データの数日々増大している。それら巨大化した木構造データから共通する構造的特徴を抽出するには計算時間が非常にかかることが多い。そのため同じ特徴を繰り返し枚挙することのない効率的なマイニング手法が提案されている [Asai et al. 2004, Zaki 2005, Itokawa et al. 2010]。文字列上の圧縮手法である LZ 圧縮手法 [Ziv and Lempel 1977] は, 繰返し出現する部分文字列を最初に出現する位置へのリファレンスで置き換えることでサイズの可逆圧縮を行う手法である。Itokawa ら [Itokawa et al. 2010] は, この LZ 圧縮手法を木構造データの構造的繰返しの圧縮に拡張した可逆圧縮手法を提案している。木構造データ D を可逆圧縮したデータ $Compress(D)$ から, 解凍することなく D に共通する構造的特徴を抽出するグラフマイニング手法の開発を本研究の目的とする。これにより抽出時間の短縮がもたらされ, より大きな木構造データの解析が可能となると考えられる。

この目的を達成するために, 論文 [Horibe et al. 2017] において, 木構造データを表現する辺ラベル付き順序木 T の構造的繰返しを可逆的に圧縮して得られる辺ラベル付き順序木である圧縮木を定義し, 順序木に対する簡潔データ表現である DFUDS [Jansson et al. 2012, Raman et al. 2013] を圧縮木に拡張した簡潔データ表現を与えた。さらに, T に対する圧縮木の簡潔データ表現に基づく T の圧縮コーディングを提案した。図 1 に, 辺ラベル付き順序木 T に対する圧縮木 t とその簡潔データ表現及び T の圧縮コーディングの例を示す。論文 [Horibe et al. 2017] では, 辺ラベル付き順序木 T の圧縮コーディング $Code(T)$ と最小出現回数 k が与えられたとき, $Code(T)$ を解凍することなく k 回以上出現するパスを全て枚挙するアルゴリズム $EnumFreqPaths$ を提案した。本稿では, 論文 [Horibe et al. 2017] の結果を拡張し, 辺ラベル付き順序木の圧縮コーディングの集合 $Compress(D)$ と実数 σ ($0 < \sigma \leq 1$)

が与えられたとき, $Compress(D)$ 中の圧縮コーディングを解凍することなく, $Compress(D)$ が表す辺ラベル付き順序木集合 D に含まれる $[|Compress(D)| \times \sigma]$ 個以上の辺ラベル付き順序木に共通して出現する順序木構造を持つ部分グラフを全て枚挙する効率的なグラフマイニングアルゴリズムを与える。最後に, 提案アルゴリズムを計算機上に実装し, 人工的なデータを用いた評価実験を行ったので, その評価実験結果の報告と考察を行う。

2. 準備

本節では, 論文 [Horibe et al. 2017] で与えた, 構造圧縮された木構造データのモデルである圧縮木について述べ, 簡潔データ表現の 1 つである DFUDS (Depth-First Unary Degree Sequence, [Jansson et al. 2012, Raman et al. 2013]) を用いた圧縮木の簡潔データ表現を紹介する。

2.1 圧縮木

Λ をアルファベットとする。辺ラベル付き順序木とは, 任意の内部ノードが順序付けられた子を持ち, 任意の辺が Λ 内の要素でラベル付けられている根付き木をいう。辺ラベル付き順序木 T のノード集合と辺集合はそれぞれ $V(T)$ と $E(T)$ で表される。ノード u がノード v の親であるラベル a をもつ $E(T)$ 内の辺を 3 つ組 $(u, a, v) \in V(T) \times \Lambda \times V(T)$ で表す。 T の内部ノード u と u の連続した子 u_1, u_2, \dots, u_k からなるラベル ℓ をもつリスト $p = (u, \ell, u_1, u_2, \dots, u_k)$ をポートリストという。このときの u を p の親ポート, u_1, u_2, \dots, u_k を子ポートという。 T の 2 つのポートリスト $p = (u, \ell_p, u_1, u_2, \dots, u_k)$ と $p' = (u', \ell_{p'}, u'_1, u'_2, \dots, u'_{k'})$ が互いに素であるとは, 以下の 2 つの条件を満たすときをいう。(1) $\{u_1, u_2, \dots, u_k\} \cap \{u'_1, u'_2, \dots, u'_{k'}\} = \emptyset$. (2) もし u と u' が同じノードであれば, u_k は u'_1 より兄弟であるか, もしくは u_1 は $u'_{k'}$ の弟妹である。 T のノード v に対し, v を根とする T の部分木を $T[v]$ で表す。さらに, $V(T)$ の部分集合 U に対して, U と両端点が U に含まれる辺からなる $E(T)$ の部分集合 E_U からなる T の誘導部分グラフを $T[U] = (U, E_U)$ と表す。 T のノード u とその子孫 v に対して, v から u へのパスを $P_{v,u}$ で表す。以下の 3 条件を満たす T のノードリス

連絡先: 内田智之, 広島市立大学大学院情報科学研究科, 〒731-3194 広島市安佐南区大塚東 3 丁目 4-1, uchida@hiroshima-cu.ac.jp

トとリファレンスの全種類のリストのコーディングとする。また、 C_T を T の圧縮木 t の簡潔データ表現のコーディングとする。このとき、 T の圧縮コーディング $Code(T)$ とは、区切り記号 "!" で連結されたコーディング列 $EL_T \circ "!" \circ D_T \circ "!" \circ C_T$ をいう。なお、 \circ は文字列の連結を行う 2 項演算を表す。以後、 T の圧縮コーディングを $Code(T) = \langle EL_T, D_T, C_T \rangle$ で表す。

EL_T, D_T, C_T は T の深さ優先走査に基づき作成されるため、容易に辺ラベル、リファレンスを返すハッシュ関数を作ることができる。図 1 に辺ラベル付き順序木 T の圧縮コーディングを与える。なお、 C_T のコーディングにおいて、 \square , "#", "a", "b", "c" はそれぞれ 0, 1, 2, 3 に、"\$", "P" はそれぞれ "-1", "-2" にコード化されている。

3. 頻出部分順序木枚挙アルゴリズム

辺ラベル付き順序木の集合 D と実数 σ ($0 < \sigma \leq 1$) に対して、順序木 s が D における σ 頻出であるとは、 s を部分グラフとして持つ D 中の辺ラベル付き順序木の数が $\lceil |D| \times \sigma \rceil$ 以上であるときをいう。辺ラベル付き順序木に対する圧縮コーディングの集合が与えられたとき、順序木構造を持つ σ 頻出な部分グラフを全て枚挙する問題を以下に定義する。

圧縮コーディング集合からの頻出部分順序木枚挙問題

インスタンス：圧縮コーディング集合 U と

実数 σ ($0 < \sigma \leq 1$)

問題： U の各圧縮木を解凍することなく D_U における σ 頻出な順序木構造を有する部分グラフを全て枚挙せよ。
なお、 D_U は U の各圧縮木を解凍して得られる辺ラベル付き順序木の集合である。

この問題を解く *EnuFreqSubTrees* アルゴリズムを Algorithm 1 に示す。*EnuFreqSubTrees* アルゴリズムは、最右拡張 [Asai et al. 2004, Zaki 2005, Itokawa et al. 2010] により順序木を枚挙する手法を用い、辺数に関しインクリメントしながら頻出な部分順序木を枚挙するアルゴリズムである。*EnuFreqSubTrees* アルゴリズムでは、 j 番目のコーディング上の位置 i に順序木 p の最右葉に付随する辺が出現することを保持するために 4 つ組 (j, i, p, OP_i^j) で表現される出現位置管理構造を用いる。ここで、 OP_i^j は位置 i がリファレンス木内のとき、そのリファレンス木へのリファレンスをもつ辺のリストである。*EnuFreqSubTrees* アルゴリズムの関数 *MAKECANDFREQSUBTREES* は、辺数 ln の σ 頻出な順序木の出現位置管理構造の集合 Z_{ln} と辺数 1 の σ 頻出な順序木の出現位置管理構造の集合 Z_1 が入力として与えられたとき、辺数 $ln + 1$ の順序木の出現位置管理構造の集合 Z_{ln+1} を返す関数である。*MAKECANDFREQSUBTREES* では、 Z_{ln} 内の各出現位置管理構造 (j, i, p, OP_i^j) に対して、 j 番目の圧縮コーディングの位置 i から始まる p の最右パス上の各ノードから子 u を伸ばしてできるそれぞれの順序木について、次の 4 つの場合に分けて処理を行う。

(Case 1) p から u を伸ばすことで i が含まれる辞書内のリファレンス木から u が出てしまう場合、

(Case 2) 辞書内のリファレンス木内に位置 i が含まれかつ p' がそのリファレンス木内に止まる場合、

(Case 3) 位置 i を含むリファレンス木が辞書内にないが、 u が辞書内のリファレンス木に入ってしまう場合、

(Case 4) 位置 i を含むリファレンス木が辞書内になく u を伸ばしても、どのリファレンス木にも入ることがなく、ポートリストが貼り付けられている辺にもならない場合。

Algorithm 1 *EnuFreqSubTrees*

Require: 圧縮コーディング集合 U と実数 σ ($1 < \sigma \leq 1$)

Ensure: 辺ラベル付き順序木集合 $D = \{T \mid \text{辺ラベル付き順序木 } T \text{ の圧縮コーディングが } U \text{ に含まれる}\}$ における全 σ 頻出順序木集合 F

```

1:  $P_1$  を  $D$  における 1 辺のみからなる  $\sigma$  頻出順序木の集合;
2:  $Z_1$  を  $P_1$  中の辺ラベルの出現位置管理構造の集合;
3:  $F = P_1$  and  $ln = 1$ 
4: while  $P_{ln} \neq \emptyset$  do
5:    $P_{ln+1} = \emptyset$ 
6:    $Z_{ln+1} = \text{makeCandFreqSubTrees}(Z_{ln}, Z_1)$ 
7:   for all  $z = (j, i, p, OP_i^j) \in Z_{ln+1}$  do
8:     if  $p$  が  $D$  において  $\sigma$  頻出 then
9:        $P_{ln+1} = P_{ln+1} \cup \{p\}$ 
10:       $Z_{ln+1} = Z_{ln+1} \cup \{z\}$ 
11:     end if
12:   end for
13:    $F = F \cup P_{ln+1}$ 
14:    $ln++$ 
15: end while
16: return  $F$ 

```

全ての場合において、 OP_i^j で管理されている各位置 k に対して、位置 i が辞書のリファレンス木内に含まれている、または u を伸ばすことでリファレンス木に入ったり、リファレンス木を参照するポートリストを表す辺に対応する場合には、同様の場合分け処理が必要である。Case 2 や Case 3 の場合には、解凍して得られる辺ラベル付き順序木において最右拡張した順序木が出現するか否かを確認する時間より余分に時間がかかる。しかし、Case 1 や Case 4 の場合で、は位置 i の処理と OP_i^j 内の位置 k の処理が同一となり、位置 k の処理を省くことができる。これにより計算時間が短縮できる。また、*EnuFreqSubTrees* を簡潔データ構造を用いて実装することにより、SDSL (Succinct Data Structure Library) [SDSL] を使用することができ、*EnuFreqSubTrees* の高速化と省メモリ化を図ることができる。

4. 実験と考察

本節では、前節で提案したアルゴリズム *EnuFreqSubTrees* を計算機上へ実装し、人工データを用いた評価実験を行ったので、その実験結果を報告し考察する。

EnuFreqSubTrees を、主記憶メモリ 32GB、OS が MacOS Sierra で、4GHz の Intel Core i7 プロセッサを持つ計算機上に C++ で実装した。辺ラベル付き順序木の圧縮コーディング集合に対するデータ構造として簡潔データ構造を採用し、各種関数は SDSL [SDSL] を使用した。また、圧縮木に現れる辞書のサイズを 5、辺ラベル種を 2 種類、解凍時のノード数を約 1000、圧縮率 (= ノード数減少率) を C に設定した 5 個の圧縮木を根で連結させた圧縮木 $D_C^5(1000)$ を、 C が 20 から 80 まで 20 刻みでそれぞれランダムに作成した。さらに、圧縮率を 50 とし、ノード数を約 N に設定した圧縮木 $D_{50}^5(N)$ を、 N が 1000 から 3000 まで 500 刻みでそれぞれランダムに作成した。加えて、上記の圧縮木 D に対応する非圧縮木もそれぞれ用意した。まず、4 つの圧縮木 $D_C^5(1000)$ ($C = 20, 40, 60, 80$) とそれぞれ対応する非圧縮木に *EnuFreqSubTrees* を入力して与えた時の実行時間と圧縮木の圧縮率との関係性を評価した。そ

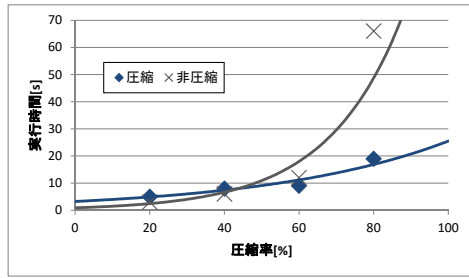


図 2: 実行時間と圧縮率の関係

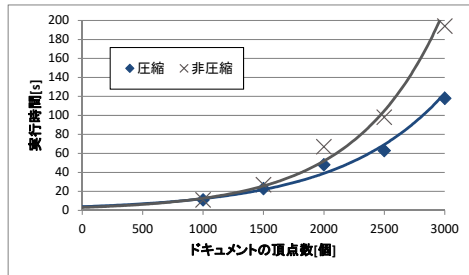


図 3: 実行時間とノード数の関係

の実験結果として、*EnuFreqSubTrees* の実行時間と圧縮率の関係を示したグラフを図 2 に示す。図 2 より、低圧縮率の圧縮木を与えた場合には、圧縮する前の辺ラベル付き順序木 (非圧縮木) を与えた場合より低速であるが、圧縮率が高くなるにつれて圧縮木を与えた方が非圧縮木の場合より高速なことがわかる。これは、低圧縮の場合、非圧縮木と比較してノード数の減少が少なく、圧縮木中のリファレンスに対する操作にかかる時間の方が実行時間に大きな影響を与えたものと考えられる。高圧縮率によるノード数の減少が実行時間に影響し、圧縮木に対する *EnuFreqSubTrees* アルゴリズムの実行時間が非圧縮木の場合より高速となった。さらに、6 つの圧縮木 $D_{50}^5(N)$ ($N = 1000, 1500, \dots, 3000$) とそれに対応する非圧縮木を用いて、ノード数と実行時間の関係を示したグラフを図 3 に示す。図 3 より、ノード数の増加に従って実行時間の差が開いていくことがわかる。圧縮率を固定してノード数を増やすと、辞書にあたるリファレンス木が大きくなる。このため、圧縮木の場合の頻出部分グラフの出現位置数と非圧縮木の場合の出現位置数の差が大きく開くことになり、圧縮木の場合の方が、非圧縮木の場合より高速になっていくと考えられる。これらの実験により、大きな繰返し構造を持つ巨大な辺ラベル付き順序木に対して *EnuFreqSubTrees* はより有効であることが実証された。

5. おわりに

本稿では、LZ 圧縮法を辺ラベル付き順序木の構造圧縮に拡張し、辺ラベル付き順序木 T を可逆的圧縮した圧縮木 t を紹介し、簡潔データ表現 DFUDS を拡張した圧縮木の簡潔データ表現を与えた。また、圧縮木 t の簡潔データ表現を用いた T の圧縮コーディングを与えた。さらに、辺ラベル付き順序木の圧縮コーディングの集合 $Compress(D)$ から解凍することなく $Compress(D)$ に対応する辺ラベル付き順序木集合 D に頻出する順序木構造を有する部分グラフを枚挙するアルゴリズム *EnuFreqSubTrees* を与えた。

今後の課題としては、実データを用いた評価実験を行うこと

や構造的変数を持つ順序木パターン p [Suzuki et al. 2015] と圧縮木 t が与えられたとき、 t を解凍することなく t が表す辺ラベル付き順序木と p がマッチするかを判定するパターンマッチング手法への拡張が考えられる。

参考文献

- [Asai et al. 2004] Asai, T., Abe, K., Kawasoe, K., Sakamoto, H., Arimura, H. and Arikawa, S.: Efficient substructure discovery from large semi-structured data. *IEICE TRANSACTIONS on Information and Systems* E87-D(12): pp.2754-2763 (2004).
- [Horibe et al. 2017] Horibe, T., Itokawa, Y., Uchida, T., Suzuki, Y. and Miyahara, T.: Algorithm for enumerating all frequent paths from structurally compressed tree-structured data. *2017 IAENG International Conference on Computer Science* [2017 年 3 月 15-17 日, 香港, 発表予定] (2017).
- [Itokawa et al. 2010] Itokawa, Y., Katoh, K., Uchida, T. and Shoudai, T.: Algorithm using expanded LZ compression scheme for compressing tree structured data. *Lecture Notes in Electrical Engineering*, Springer, pp.333-346 (2010).
- [Itokawa et al. 2010] Itokawa, Y., Miyoshi, J., Wada, M. and Uchida, T.: Succinct Representation of TTSP Graphs and its Application to the Path Search Problem. *Proc. 6th IASTED International Conference on Advances in Computer Science and Engineering (ACSE 2010)*, pp.33-40 (2010)
- [Jansson et al. 2012] Jansson, J., Sadakane, K. and Sung, W.-K.: Ultra-succinct representation of ordered trees with applications. *Journal of Computer and System Sciences* 78, pp.619-631 (2012)
- [Raman et al. 2013] Raman, R. and Rao, S.S.: Succinct representations of ordinal trees. *Space-efficient data structures, streams, and algorithms, Proc. LNCS* Vol. 8066, pp.319-332 (2013)
- [SDSL] SDSL: Succinct data structure library. <http://simongog.github.io/sdsl/>
- [Suzuki et al. 2015] Suzuki, Y., Shoudai, T., Uchida, T. and Miyahara, T.: An efficient pattern matching algorithm for ordered term tree patterns. *IEICE TRANSACTIONS on Fundamentals of Electronics, Communications and Computer Sciences*, E98-A(06):1197-1211 (2015)
- [Zaki 2005] Zaki, M. J.: Efficiently mining frequent trees in a forest: algorithms and applications. *IEEE Transactions on Knowledge and Data Engineering*, 17(8):1021-1035 (2005)
- [Ziv and Lempel 1977] Ziv, J. and Lempel, A.: A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, IT-23(3):337-343 (1977)