

Modified PrefixSpan 法を用いた 頻出正規パターンの抽出をめざして

Mining of Sequential Patterns with Variable Wildcard Regions using Modified PrefixSpan Method

塔野 薫隆[†] 北上 始[‡]
田村 慶一[‡] 森 康真[‡]
黒木 進[‡]

Shigetaka TONO Hajime KITAKAMI
Keiichi TAMURA Yasuma MORI
Susumu KUROKI

著者らは、ある配列データベースから可変長ワイルドカード領域が含まれている頻出パターンを抽出するために、新 Modified PrefixSpan 法を提案する。この新しい方式は、既に提案した Modified PrefixSpan 法に、最大誤差数と呼ばれる入力パラメータを追加することによって開発されている。新方式の有効性を調べるために、PROSITE から Zinc Finger データセットを取り出し、可変長ワイルドカード領域が含まれる頻出パターンを抽出する能力の評価を行った。その結果、本提案方式が従来方式に比べて 9 倍の優れた抽出能力を持っている事を確認した。

In order to extract frequent patterns with a “variable wildcard region” from a sequence database, we propose the new Modified PrefixSpan method. The new method is enabled to develop by adding a new input parameter, called a maximum error count, to the former Modified PrefixSpan method. We verify this new method with a Zinc Finger dataset that is included in PROSITE. The results show that this new method has 9 times superior capacity for extraction of frequent patterns.

1. はじめに

配列データベースから頻出パターンを抽出する方法を用いると、多くの問題に答えることができる。この方法は、DNA やアミノ酸の配列データだけではなく、顧客の購買履歴、Web アクセス履歴、科学的な実験データ、病気の治療履歴、自然災害の履歴などを分析するのに有用であるといわれている。我々は、これまでに分子生物学の分野におけるモチーフ抽出の問題に着目してきた。現在、さまざまなモチーフが見つかっているが、各モチーフはあるアミノ酸配列データベース中に存在する特別な配列パターンであり、多様な蛋白質が持つ

機能の 1 つに関係し、生物の進化の過程で保存されてきたものであると考えられている。

本論文では、モチーフ発見を支援するために、アミノ酸の配列データベースからさまざまなワイルドカード領域をもつ頻出パターンの抽出方法を提案する。ワイルドカード領域には、固定長と可変長のワイルドカード領域の 2 種類があるが、著者らは PrefixSpan 法[1]を拡張し、両者を含む頻出パターンを抽出する。拡張された PrefixSpan 法の評価を行うために、Zinc Finger モチーフを含むデータセットを PROSITE[2]から取り出し、そのデータセットを配列データベースとして利用し、固定長ワイルドカード領域を抽出する処理と可変長ワイルドカード領域を抽出する処理の性能比較を行う。

2. 関連研究

PrefixSpan 法は、配列データベースから頻出パターンを抽出するアルゴリズムとして知られているが、抽出過程では配列の軸位置から最左端位置までの広い範囲を参照しなければならないため、多大な計算時間を要し、無意味な頻出パターンが数多く抽出されてしまうという問題がある。また、ワイルドカード x が明記された頻出パターンを抽出する機能を持たないため、正規表現で表されるモチーフを表現できないという問題がある。例えば、頻出パターンの抽出過程において、 $\langle AC \rangle$, $\langle Ax \rangle$ と $\langle Ax \rangle$ は、単に $\langle AC \rangle$ と解釈されてしまうので、モチーフ抽出問題には不向きである。

これらの問題を解決するために、PrefixSpan 法を改良した Modified PrefixSpan 法[3]では、頻出パターン中に複数存在する固定長ワイルドカード領域の最大長の設定を可能にし、頻出パターンの抽出過程において、 $\langle AC \rangle$, $\langle Ax \rangle$ と $\langle Ax \rangle$ を異なるパターンとして区別する機能を持たせた。

しかしながら、PROSITE に登録されているモチーフには、固定長ワイルドカード領域の他に可変長ワイルドカード領域をもつモチーフが数多く存在するにもかかわらず、可変長ワイルドカード領域を自動的に抽出する研究が行われていなかった。本論文では、Modified PrefixSpan 法を拡張し、可変長ワイルドカード領域を含む頻出パターンを抽出する方法について提案している。

3. 問題の形式

配列データベース $S = \{S_1, S_2, S_3, \dots, S_n\}$ の各配列 S_i は、あるアルファベット Σ 上で定義されるとする ($1 \leq i \leq n$)。例えば、蛋白質の場合は、以下の 20 文字のアルファベット Σ_{protein} 上で定義される配列である。

$$\Sigma_{\text{protein}} = \{A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y\}$$

Π を Σ 上で定義された曖昧文字の集合、 X をワイルドカード文字列の集合、 $\{*\}$ を任意長のワイルドカードを表す文字列を要素とする集合とする。PROSITE のパターンは、 $(\Sigma \cup \Pi \cup X \cup \{*\})$ 上の正規表現として書かれている[4][5]。我々は抽出するパターン表記として、PROSITE のパターン表記を用いる。

ここでは、 k 個のアルファベット文字要素をもつパターンを k -パターンと呼ぶ。 k -パターンが利用者によって与えられた最小支持数 min_sup 以上の配列にマッチするとき、そのパターンを k -頻出パターンと呼び、 pat^k で表現する。配列データベース S に含まれる k -パターン pat^k は、以下の正規表現とする。

$$pat^k = \langle A_1 - x(i_1, j_1) - A_2 - x(i_2, j_2) - \dots - x(i_{k-1}, j_{k-1}) - A_k \rangle$$

[†] 学生会員 広島市立大学大学院情報科学研究科
tono@db.its.hiroshima-cu.ac.jp

[‡] 正会員 広島市立大学情報科学部
{kitakami, ktamura, mori, kuroki}@its.hiroshima-cu.ac.jp

ここで、 $\langle \rangle$ 内のハイフン(-)は、隣同士がお互いに連続していることを表現するための記号であり、その記号を省略することがある。 A_i をアルファベット Σ 上の文字要素または単に要素と呼ぶ。また、文字要素 A_i が、1文字で表現されるとき単一要素、2文字以上で表現(たとえば[ILVF]など)されるとき曖昧要素と呼ぶ。 $x(i, j)$ の領域において、 $0 \leq i \leq j$ のとき、その領域を可変長ワイルドカード領域と呼ぶ。 $i=j$ のとき、それを固定長ワイルドカード領域と呼び、この領域を $x(i)$ で簡略表現することがある。以下では、固定長ワイルドカード領域だけを含むパターンを基本パターンと呼ぶ。また、 $x(0, \infty)$ は、任意長のワイルドカード文字列(*)と同じ意味である。

本論文では、配列データベース S から可変長ワイルドカード領域が含まれる頻出パターンの集合 P を抽出する方法に着目する。これによって抽出された各頻出パターンの形式は複数の基本パターンを1つにまとめられた形式を持つ。例えば抽出された頻出パターンの形式が $\langle F-x(2,5)-A \rangle$ であれば、これは、 $\langle F-x(2)-A \rangle, \langle F-x(3)-A \rangle, \langle F-x(4)-A \rangle, \langle F-x(5)-A \rangle$ の4つの基本パターンが含まれている事を表す。頻出パターン P を抽出するためには、最小支持数 min_sup , 最大誤差数 ϵ_{max} , 最大ワイルドカード数 max_wc を必要とする。支持数 cnt_r を持つ k -頻出パターン pat_r^k を $\langle pat_r^k \rangle : cnt_r$ で表現し、 m 個の k -頻出パターンが配列データベース S から抽出されたとする、それらの集まりとする集合 P_k は以下のように表現される。

$$P_k = \{ \langle pat_1^k \rangle : cnt_1, \langle pat_2^k \rangle : cnt_2, \dots, \langle pat_m^k \rangle : cnt_m \}$$

S から抽出される頻出パターンの最大要素数を q とすると、 P は $\{P_1, P_2, \dots, P_q\}$ で表現される。 $(k+1)$ -パターンの最右端文字は、 k -頻出パターンに対して、 Σ 上の1文字を追加することにより構成され得る。その1文字は、 S_i において、 k -頻出パターン pat_r^k の最右端文字よりも右側の位置に存在する可能性がある。その右側の位置を表現するために、以下の $PDB(pat_r^k)$ を導入する。

$PDB(pat_r^k) = \{ (i, j) \mid S_i \text{において、} pat_r^k \text{の最右端文字の右隣位置を} j \text{とする} \}$ 。ただし、 $1 \leq j \leq \| S_i \|$ である。以後、 $PDB(pat_r^k)$ を pat_r^k の射影データベース(Projected Database)と呼ぶ。

4. Modified PrefixSpan 法

PrefixSpan法を拡張した Modified PrefixSpan法は、固定長ワイルドカード領域が含まれる頻出パターンを抽出する方法である。

表1に示される配列データベース S から固定長ワイルドカード領域が含まれる頻出パターンの抽出方法を示そう。ただし、最小支持数 min_sup および最大ワイルドカード数 max_wc を各々3とする。また、表1は、 $S = \{S_1, S_2, S_3, S_4, S_5\}$ であり、 $S_1 = \text{FKYAKWL}, S_2 = \text{SFVKTA}, S_3 = \text{ALR}, S_4 = \text{MSKPL}, S_5 = \text{FSKFLMAW}$ であることを示している。配列データベース S から頻出パターンを抽出する処理は、 k -頻出パターンから $(k+1)$ -頻出パターンを作成することによって行われる($k \geq 1$)。この処理を $k=1$ と $k=2$ の場合について考えてみよう。

$k=1$ のとき、配列データベース中に存在するパターンは、 $\langle A \rangle : 4, \langle F \rangle : 3, \langle K \rangle : 4, \langle L \rangle : 4, \langle M \rangle : 2, \langle P \rangle : 1, \langle R \rangle : 1, \langle S \rangle : 3, \langle W \rangle : 2, \langle V \rangle : 1, \langle Y \rangle : 1$ である。これらにより、最小支持数を満たす1-頻出パターンを選び出すと、その集合は図1に表されているように、 $P_1 = \{ \langle A \rangle : 4, \langle F \rangle : 3, \langle K \rangle : 4, \langle L \rangle : 4, \langle S \rangle : 3 \}$ になる。例えば、1-頻出パターン $\langle A \rangle : 4, \langle F \rangle : 3, \langle K \rangle : 4$ の射影データベースは、各々、 $PDB(\langle A \rangle) = \{ (1, 5), (2, end), (3, 2), (5, 8) \}$ 、

表1 配列データベース

Table 1 Sequence Database

ID	Sequence
1	FKYAKWL
2	SFVKTA
3	ALR
4	MSKPL
5	FSKFLMAW

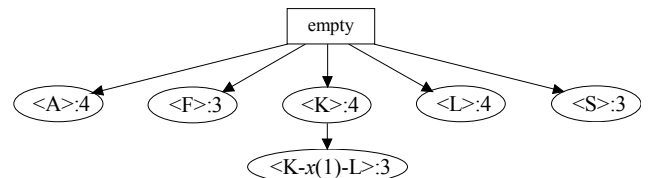


図1 従来方式で得られたマイニング木

Fig. 1 Mining Tree generated by the Modified PrefixSpan Method

$\langle F \rangle = \{ (1, 2), (2, 3), (5, 2), (5, 5) \}$, $PDB(\langle K \rangle) = \{ (1, 3), (1, 6), (2, 5), (4, 4), (5, 4) \}$ で表現される。 $PDB(\langle A \rangle)$ 内の $(1, 5)$ は、パターン $\langle A \rangle$ が1番目の配列 S_1 の4文字目にあり、 S_1 には $\langle A \rangle$ から始まる2-パターンの最右端文字は5文字目以降に存在する可能性があることを意味する。また、 $PDB(\langle A \rangle)$ 内の $(2, end)$ は $\langle A \rangle$ が2番目の配列 S_2 の最右端文字にあり、 S_2 には $\langle A \rangle$ から始まる2-パターンの最右端文字は存在しないことを意味する。

$k=2$ のとき、2-頻出パターンの集合 P_2 は P_1 中の各1-頻出パターンから構成することができる。例えば、 P_1 中の $\langle A \rangle$ の次にくる1文字は、 $\langle A \rangle$ の射影データベース $PDB(\langle A \rangle) = \{ (1, 5), (2, end), (3, 2), (5, 8) \}$ を用いて選択され得る。それらを集めた集合は $\{ S_1[5], S_3[2], S_5[8] \} = \{ \langle K \rangle, \langle L \rangle, \langle W \rangle \}$ になるが、その中で P_1 に含まれていない文字は頻出ではない。従って、2-頻出パターンの最右端文字は $\langle K \rangle, \langle L \rangle$ の2つの文字に絞られるが、残念なことに $\langle A-K \rangle : 1, \langle A-L \rangle : 1$ である。また、1~3文字のワイルドカードを含む文字列 $\langle A-x(1) \rangle, \langle A-x(2) \rangle, \langle A-x(3) \rangle$ のどれかを考えても、その次にくる1文字をさがし2-頻出パターンを構成することができない。以上から、 $\langle A \rangle$ から始まる2-頻出パターンは存在しない(図1)。

次に、 P_1 中の $\langle K \rangle$ から始まる2-頻出パターン(ワイルドカード数がゼロの場合)について考えてみよう。 $\langle A \rangle$ の場合と同様に、 $\langle K \rangle$ の次にくる1文字は $\langle K \rangle$ の射影データベース $PDB(\langle K \rangle)$ を用いて選択される。それらの中で P_1 に含まれる文字は $\{ S_2[5], S_4[4] \} = \{ \langle T \rangle, \langle F \rangle \}$ である。従って、 $\langle K-T \rangle : 1, \langle K-F \rangle : 1$ となるので、ワイルドカード数をゼロとする限り、2-頻出パターンは存在しない。次に、1文字のワイルドカードを含む2-頻出パターンについて考えてみよう。 $\langle K-x(1) \rangle$ の次にくる1文字は $PDB(\langle K-x(1) \rangle) = \{ (1, 3+1), (1, 6+1), (2, 5+1), (4, 4+1), (5, 4+1) \}$ を計算することにより選択される。それらの中で P_1 に含まれる文字は、 $\{ S_1[4], S_1[7], S_2[6], S_4[5], S_5[5] \} = \{ \langle A \rangle, \langle L \rangle, \langle A \rangle, \langle L \rangle, \langle L \rangle \}$ である。これにより、 $\langle K-x(1)-A \rangle : 2, \langle K-x(1)-L \rangle : 3$ を構成することができる。しかし、2文字および3文字のワイルドカードを含む文字列 $\langle K-x(2) \rangle, \langle K-x(3) \rangle$ のどちらかを考えても、もはや2-頻出パターンは存在しない。以上により、 $\langle K \rangle$ から始まる2-頻出パターンは、 $\langle K-x(1)-L \rangle : 3$ だけとなる(図1)。3-頻出パターンを構成するために必要な $\langle K-x(1)-L \rangle$ の射影データベース $PDB(\langle K-x(1)-L \rangle) = \{ (1, end), (4, end), (5, 6) \}$ となるので、この2-頻出パターンから3-頻出パターンを構成することができない。

5. 新 Modified PrefixSpan 法

新 Modified PrefixSpan 法の処理では、前章で説明した入力パラメータの他に、最大誤差数 ϵ_{max} を必要とする。最大誤差数がゼロのとき、従来の Modified PrefixSpan 法と同じである。また、 k -頻出パターン $\langle pat^k \rangle$ から可変長ワイルドカード領域が含まれる $(k+1)$ -パターンを抽出するときに、同じ最右端文字 α をもつ以下のパターンが数多く見つかる。

$\langle pat^k-x(i, i+\epsilon_1)-\alpha \rangle: cnt_{\epsilon_1}, \langle pat^k-x(i, i+\epsilon_2)-\alpha \rangle: cnt_{\epsilon_2}, \dots$
 $\epsilon_1 \leq \epsilon_2 (\leq \epsilon_{max})$ のとき、明らかに、 $cnt_{\epsilon_1} \leq cnt_{\epsilon_2}$ であり、 $\langle pat^k-x(i, i+\epsilon_1)-\alpha \rangle: cnt_{\epsilon_1}$ は、 $\langle pat^k-x(i, i+\epsilon_2)-\alpha \rangle: cnt_{\epsilon_2}$ に含まれる表現である。従って、このような同じ最右端文字のパターンに対しては、我々は、最も一般的な形式のパターンだけを候補パターンとして選び出している。これにより、冗長性のない $(k+1)$ -頻出パターンを抽出している。

以下、表 1 で示される配列データベースの例を用いて、可変長ワイルドカード領域が含まれる頻出パターンを抽出する処理について考えてみよう。ただし、最小支持数 min_sup 、最大ワイルドカード数 max_wc 、最大誤差数 ϵ_{max} を各々 3 とする。配列データベース S から頻出パターンを抽出する処理は、前章と同じであり、 k -頻出パターンから $(k+1)$ -頻出パターンを作成することによって行う ($k \geq 1$)。 $k=1$ に対応する 1-頻出パターンの構成方法は、図 2 に示されるように、従来の Modified PrefixSpan 法と同じである。 $k=2$ に対して、従来の方法では、 $\langle F \rangle$ から始まる 2-頻出パターンは抽出できなかったが、新方式ではそれが可能になっている。

以下では、2-頻出パターン $pat = \langle F-x(0,3)-K \rangle$ に着目し、この 2-頻出パターンから 3-頻出パターンを生成する方法について考えてみよう。図 3 に示されるように、ワイルドカード数 wc が 0 と 1 に対応する 3-頻出パターンが生成されるが、ワイルドカード数 wc が 2 と 3 に対応する 3-頻出パターンは生成されない。

例えば、2-頻出パターンを $pat = \langle F-x(0,3)-K \rangle$ とし、その 2-頻出パターンからワイルドカード数 wc が 1 の 3-頻出パターンを見つけてみよう。このために、誤差 ϵ を最小値 0 から最大値の 3 まで変化させる。即ち、 $0 \leq \epsilon \leq 3$ とし、3-頻出パターン $\langle pat-x(1, 1+\epsilon)-\alpha \rangle: cnt$ の最右端文字 α および支持数 cnt を計算する方法について考える。この 3-パターンは、4 つの基本パターン $\langle pat-x(1,1)-\alpha \rangle, \langle pat-x(2,2)-\alpha \rangle, \langle pat-x(3,3)-\alpha \rangle, \langle pat-x(4,4)-\alpha \rangle$ をまとめた表現である。各基本パターンの支持数は、誤差 ϵ を 0 から 3 まで順に変化させながら、誤差 ϵ ごとに、同じ α をもつ 3-パターンが異なるタプルに何回出現するかを正確に数え上げることで計算できる。この様子を図 3 の (a) と (b) に示す。

図 3 で、 $pat = \langle F-x(0,3)-K \rangle$ から始まる 3-パターンの最右端文字 α が $\langle A \rangle$ の場合について考えてみよう。図 3 の (a) で誤差

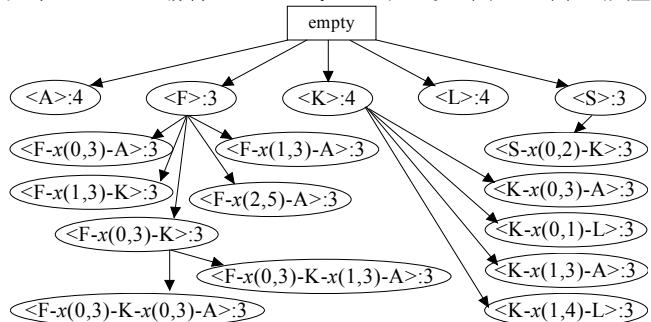


図 2 提案方式で得られたマイニング木

Fig. 2 Mining Tree generated by the new method

ID	ϵ			
	0	1	2	3
1	$pat-x-A$ $pat-x-L$	$pat-xx-K$	$pat-xxx-W$	$pat-xxxx-L$
2	$pat-x-A$	-	-	-
3	-	-	-	-
4	-	-	-	-
5	$pat-x-L$	$pat-xx-M$	$pat-xxx-A$	$pat-xxxx-W$

α	ϵ			
	0	1	2	3
$\langle A \rangle$	$\langle pat-x(1,1)-A \rangle:2$	-	$\langle pat-x(1,3)-A \rangle:3$	-
$\langle F \rangle$	-	-	-	-
$\langle T \rangle$	-	-	-	-
$\langle K \rangle$	-	$\langle pat-x(1,2)-K \rangle:1$	-	-
$\langle L \rangle$	$\langle pat-x(1,1)-L \rangle:2$	-	-	$\langle pat-x(1,4)-L \rangle:2$
$\langle M \rangle$	-	$\langle pat-x(1,2)-M \rangle:1$	-	-
$\langle W \rangle$	-	-	$\langle pat-x(1,3)-W \rangle:1$	$\langle pat-x(1,4)-W \rangle:2$

図 3 $\langle pat-x(1, 1+\epsilon)-\alpha \rangle: cnt$ 形式の 3-頻出パターンの計算 ($0 \leq \epsilon \leq \epsilon_{max}$)

Fig. 3 Extraction of 3-Frequent Patterns having $\langle pat-x(1, 1+\epsilon)-\alpha \rangle: cnt$ Expression ($0 \leq \epsilon \leq \epsilon_{max}$)

ϵ が 0 のとき、3-パターン $\langle pat-x-A \rangle$ は、配列データベース中に 2 箇所存在するが、それらは、お互いに異なるタプル (1 番目と 2 番目) に含まれているので、図 3 の (b) に示されるように $\langle pat-x(1,1)-A \rangle:2$ となる。図 3 の (a) で誤差 ϵ が 1 のとき、最右端文字が $\langle A \rangle$ であるパターン $\langle pat-xx-A \rangle$ はどのタプルにも存在しないが、誤差 ϵ が 2 のとき、 $\langle pat-xxx-A \rangle$ はそれらのタプルとは異なる 5 番目のタプルに含まれている。従って、図 3 の (b) に示されるように、ここまでの計算では、 $\langle pat-x(1,1)-A \rangle:2$ かつ $\langle pat-xxx-A \rangle:1$ であるので、 $\langle pat-x(1,3)-A \rangle$ の支持数は 3 になる。最後に、誤差 ϵ が 3 のときは、最右端文字が $\langle A \rangle$ で終わる 3-パターン $\langle pat-xxxx-A \rangle$ は存在しない。以上から、 $pat = \langle F-x(0,3)-K \rangle$ から始まる 3-パターンは、 $\langle pat-x(1,1)-A \rangle:2, \langle pat-x(1,3)-A \rangle:3$ の 2 つ存在することになるが、後者は前者を包含するので、前者を除外する。従って、ワイルドカード数を 1 に固定したとき、 $pat = \langle F-x(0,3)-K \rangle$ から始まる 3-頻出パターンは、 $\langle F-x(0,3)-K-x(1,3)-A \rangle:3$ となる。

ワイルドカード数を 1 に保持したままで、他の候補パターンについて考えると、図 3 の (b) に示されるように 6 個の 3-パターンが見つかる。そして冗長なパターンとして、 $\langle pat-x(1,1)-L \rangle:2$ および $\langle pat-x(1,3)-W \rangle:1$ が除外され、4 個の 3-パターンが残る。しかし、4 個のいずれも最小支持数を満たさないの、他の 3-頻出パターンは存在しない。

6. 計算結果と性能評価

PROSITE から Zinc Finger モチーフを含むデータセットを選択するために、アクセス番号として PS00028 を用いた。このデータセットは 744 件の配列が含まれている。Zinc Finger モチーフは $\langle C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H \rangle$ の形式で知られている。この形式は、72 種類の基本パターンを一つにまとめたものである。しかし、評価のために使われたデータセットには、その内の 34 種類の基本パターンだけが存在する。また、その 34 種類の基本パターンの内、配列データベース内で最も支持率の高いパターンの形式は $\langle C-x(2)-C-x(3)-F-x(8)-H-x(3)-H \rangle$ である。このパターンは、配列データベース内に 3299 箇所あり、支持率は 82% である。また、最も支持率の低いパターンは $\langle C-x(2)-C-x(3)-M-x(8)-H-x(4)-H \rangle$ である。このパターンの支持率は 0.13% である。

この配列データベースから頻出パターンを抽出するために、入力データとして max_wc を 8 にして、従来方式では ϵ_{max} を 0、提案方式では ϵ_{max} を 2 として入力し、表 2 でその計算結果が表されている。提案方式の計算では min_sup は 90% から 50% まで下げていって各支持率によって得られた基本パタ

表2 Zinc Finger の計算処理
Table 2 Zinc Finger Evaluation

比較項目	最小支持率	90%	80%	70%	60%	50%
	従来方式					
	正解の数(件)	0	1	1	1	1
	計算時間(sec)	2.23	3.46	6.7	15.25	49.89
	頻出パターン数(件)	81	277	1,108	3,998	19,104
提案方式						
	正解の数(件)	9	9	-	-	-
	計算時間(sec)	22.9625	128.616	-	-	-
	頻出パターン数(件)	1,042	9,669	-	-	-

一の種類数や頻出パターン数や計算時間を表している。また、提案方式では min_sup を 90% と 80% にして得られた基本パターンの種類数や頻出パターン数や計算時間を表している。

従来方式では、最小支持率が 80% のとき、Zinc Finger モチーフの一部分として、82% の支持率をもつ頻出パターン $\langle C-x(2)-C-x(3)-F-x(8)-H-x(3)-H \rangle$ が頻出パターンの集合に含まれていた。それは、1 件の正解に相当する。最小支持率を 50% まで下げたが、他の正解は抽出されなかった。一般に最小支持率を下げると、計算に必要なメモリー領域が増える。我々の計算機環境では、最小支持率を 50% 未満にすると、計算の途中でメモリー領域が不足したため、途中で計算を打ち切った。これに対して、提案方式では、最小支持率が 90% のとき、 $\langle C-x(2,4)-C-x(3,5)-F-x(8,10)-H-x(3,5)-H \rangle$:93% ほか 17 件の頻出パターンが正解を含んでいた。それらを基本パターンとして整理すると、9 件の正解(正解総数の 26%)に相当することがわかった。最小支持率を 80% 未満に設定すると、提案方式では計算の途中でメモリー領域が不足する。提案方式では最小支持率が 90% のときの計算時間は 23.0 秒であったが、従来方式でこれと同じ程度の計算時間に対応する最小支持率は 60%~50% である。このとき、従来方式では、正解が 1 件だけしか抽出されていないことがわかる。以上から、提案方式は従来の方式に比べて、正解の数が 9 倍多かったと言える。

7. まとめ

従来の *Modified PrefixSpan* 法は、最小支持数と最大ワイルドカード数の 2 つの設定パラメータに対して、固定長ワイルドカード領域を含む頻出パターンだけを抽出していたので、可変長ワイルドカード領域を含む頻出パターンの一部分になっている基本パターンを見落としてしまうという問題がある。本論文では、この問題を解決するために、*Modified PrefixSpan* 法に可変長ワイルドカード領域を含む頻出パターンを抽出する機能を持たせる方法を提案した。具体的には、従来方法に最大ワイルドカード数に対する最大誤差数のパラメータを新たに導入した。提案方法の有効性を確かめるために、PROSITE から Zinc Finger データセットを選び、計算結果を従来方法と比較した。そこでは、提案方法と従来方法の計算終了時間がお互いにほぼ同じになるような最小支持率を各々選択し、各々の方法で抽出された頻出パターン集合を調べた。各々の頻出パターン集合の中を調べた結果、提案方法が従来方法に比べて、正解としての基本パターンが 9 倍多く抽出されていることがわかった。提案方法に限らず従来方法では、一般に、最小支持数を下げれば、計算時間や計算に必要なメモリー容量が増える。今回、提案方法で抽出された基本パターンの中には、設定された最小支持数よりも小さな支持数が含まれている。極端な例になるが、Zinc Finger データセットの場合、抽出された可変長ワイルドカード領域を含む頻出パターンの中に、支持数が 1 の基本パターン

($\langle C-x(3)-C-x(3)-F-x(8)-H-x(5)-H \rangle$) が含まれていた。これを従来方法で最小支持数を 1 と設定して計算をすると、計算の途中で組み合わせ爆発が生じる。即ち、 n 本の配列の平均配列長を m とすると、時間計算量は $O(n^m)$ になってしまうので、従来方法を用いてこの基本パターンを見つけるのは現実的ではない。

今後の課題として、3 章で説明した「曖昧要素」を含む頻出パターンを抽出する方法の研究があげられる。

【謝辞】

本研究の一部は、広島市立大学特定研究費(一般研究費(コード番号: 3106))の支援により行われた。

【文献】

- [1] Jian Pei, Jiawei Han, Behzad Mortazavi-Asl, and Helen Pinto: PrefixSpan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth, Proc. of International Conference on Data Engineering (ICDE 2001), pp.215-224, IEEE Computer Society Press, 2001.
- [2] PROSITE: <http://kr.expasy.org/prosite>
- [3] Hajime Kitakami, Tomoki Kanbara, Yasuma Mori, Susumu Kuroki, and Yukiko Yamazaki: Modified PrefixSpan Method for Motif Discovery in Sequence Databases, Proc. of the 7th Pacific Rim International Conference on Artificial Intelligence (PRICAI2002), pp.482-491, Springer-Verlag, August 2002.
- [4] Inge Jonassen, John F. Collins, and Desmond G. Higgins: Finding Flexible Patterns in Unaligned Protein Sequences, Protein Science, pp.1587-1595, Cambridge University Press, 1995.
- [5] D. Gusfield: Algorithms on Strings, Trees, and Sequences, Cambridge University Press, 1997.

塔野 薫隆 Shigetaka TONO

広島市立大学大学院情報科学研究科研究生。2003 ポリビア国・カトリック大学サンタクルス校システム工学部卒業。日本データベース学会学生会員。

北上 始 Hajime KITAKAMI

広島市立大学情報科学部教授。1976 東北大学大学院工学研究科博士前期課程修了。博士(工学)。データベースシステムの研究・開発に従事。情報処理学会一般情報処理教育委員会委員および CE 研究会運営委員、日本データベース学会 BI 研究グループ運営委員など。

田村 慶一 Keiichi TAMURA

広島市立大学情報科学部助手。2000 九州大学大学院システム情報科学府修士課程修了。データベース並列処理の研究に従事。日本データベース学会、情報処理学会 各会員。

森 康真 Yasuma MORI

広島市立大学情報科学部助手。1994 北陸先端科学技術大学院大学情報科学研究科博士前期課程修了。データベースシステムの研究・開発に従事。日本データベース学会会員。

黒木 進 Susumu KUROKI

広島市立大学情報科学部助教授。1990 東京大学大学院工学系研究科修士課程修了。博士(工学)。空間データベースの研究に従事。日本データベース学会、情報処理学会、電子情報通信学会、ACM、IEEE 各会員。