

基本パターンの出現を保存した融合に基づく 関係型パターンマイニング手法

Multi-Relational Pattern Mining Based-on Combination of Properties
with Preserving Their Structure in Examples

中野 裕介*¹
Yusuke Nakano

犬塚 信博*¹
Nobuhiro Inuzuka

*¹名古屋工業大学 大学院 工学研究科 情報工学専攻

Department of Computer Science and Engineering, Graduate School of Engineering, Nagoya Institute of Technology

We propose a new algorithm for the problem of multi-relational pattern mining through the problem established in WARMR. In order to overcome the combinatorial problem of large pattern space, another algorithm MAPIX restricts patterns into combination of basic patterns, called properties. A property is defined as a set of literals appeared in examples and is an extended form of the attribute-value form. MAPIX enumerates patterns made from conjunction of the properties. Although the range of patterns is clear and MAPIX enumerates them efficiently, a large part of patterns are out of the range. Advantage of MAPIX is to make patterns from pattern fragments occurred in examples. Many patterns which are not appeared in examples are not tested. The proposing algorithm keeps this advantages and extends the way of combination of properties. The algorithm adopts a way of combination it combines properties as they appeared in examples.

1. はじめに

大量のデータから隠された知識や新しい規則を発見するプロセスをデータマイニングとよび、その中でも複数の関係表に跨るようなパターンを扱うものを関係型データマイニング (Multi-Relational Data Mining: MRDM) とよぶ。従来のデータマイニングの多くは1つの関係表を用いたマイニングであるが、一般的なデータは複数の関係表に渡るデータベースで表されるため MRDM を行う必要がある。

MRDM は帰納論理プログラミング (Inductive Logic Programming: ILP) の枠組みで行われてきた。ILP は述語論理を用いた記述により、豊かな表現力と高い可読性を持っており、MRDM の有力な手法と考えられている。

MRDM の代表的な手法として WARMR[2] がある。これは、単純なパターンから複雑なパターンへとトップダウンに頻出パターンを探索するアルゴリズムである。しかし、生成される候補が膨大にあり計算時間が大きい。

これに対し、元山らはボトムアップに基本パターンを生成し、その組み合わせによりパターンを探索する MAPIX (Mining Algorithm by Property Item eXtraction)[4] を考案した。これは、引数のモードやタイプ等で制限を行いながら、事例に関する事実を取り出し、それらをアイテムとして興味深いルールを導出するものである。これにより探索の高速化に成功したが、出力されるパターンに制限を受けてしまった。

そこで、MAPIX を拡張させた EQUIPIX[3] は基本パターンを構造に基づいて組み合わせるオペレーションにより探索の範囲を広げた。しかし、そこで用いられているオペレータは全ての組み合わせを見ておらず改善の余地がある。

本研究では EQUIPIX を発展させたアルゴリズムを提案する。提案手法では MAPIX を2回用いる。1回目は従来と同様に頻出基本パターンの組み合わせを枚挙する。2回目の MAPIX を行う前に1回目の出力結果を用いて構造的な組み合わせを行い、その出力結果を用いて再度 MAPIX を実行する。

2. 準備

ILP は述語論理の一般的な定義に基づく。MRDM では、それぞれの関係表を述語論理の形式で表現しパターンをマイニングする。例として家族関係データベース R_{fam} について考える (図1)。 R_{fam} は4つの関係表をもち、それぞれを次のように述語論理の形式で表現する。 $parent(x, y)$ は x が y の親であることを、 $male(x)/female(x)$ は x が男性/女性であることを、 $grandfather(x)$ は x が祖父であることを意味する。01 から 24 は人物を表すインデックスである。以下では省略して gf, p, m, f を用いる。まず MAPIX におけるパターンを定義する。ここで $target$ は特徴を取り出したい事例の述語で構成されたりテラルであり目標事例 (目標述語) とよぶ。これは WARMR における key の概念に相当する。

定義1 (パターン) パターンは次のような節である。

$$P = target \leftarrow char_1 \wedge \dots \wedge char_n$$

$e = head(P)\rho$ となる最汎単一化代入 ρ を使って、 $R = body(P)\theta$ となる代入 θ が必ず存在するとき、事例 e がデータベース R 上でパターン P を満たすという。

定義2 (頻出パターン) パターン P の支持度とは、 P を満たす事例の割合であり $supp(P)$ と表記する。 $supp(P) \geq sup_{min}$ のときパターン P は頻出であるという。ただし、 sup_{min} はユーザにより与えられる最低支持度である。

例えば、次の式は R_{fam} におけるパターンである。

$$P = gf(A) \leftarrow m(A) \wedge p(A, B) \wedge f(B).$$

このパターンを事例 $e = gf(01)$ が満たしているかどうか調べるとき、 $e = head(P)\rho$ となるような代入 $\rho = \{A/01\}$ を考える。 $body(P)\rho = m(01) \wedge p(01, B) \wedge f(B)$ の中の全ての述語が R_{fam} に現れるような代入 $\theta = \{B/01\}$ が存在するので、事例 $e = gf(01)$ はパターン P を満たす。同様に事例 $gf(07)$, $gf(12)$, $gf(20)$ も P を満たすが、 $gf(19)$ はこれを満たさない。つまり P の支持度は $supp(P) = 4/5 = 80\%$ となり、最低支持度が 80% 以下のとき P は頻出パターンとなる。

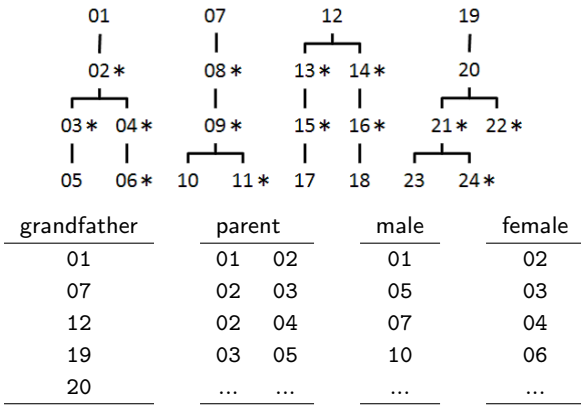


図 1: 4つの関係系 grandfather, parent, male, female からなる家族関係のデータベース R_{fam} . * がついているものは女性を, それ以外は男性を表す.

MRDM では, 効率よく探索を行うために論理的に意味の同じものを出力しないようにすることがある. そこで, 同値なパターンについて述べるため包摂関係を定義する.

定義 3 (包摂関係, 同値) 節 C, D について $C\theta \subseteq D$ をみたす代入 θ が存在するとき, 節 C は D を包摂するといいい $C \succeq D$ と表す. また, $C \succeq D$ かつ $D \succeq C$ が成り立つとき, 節 C と D は包摂に関して同値であるといいい $C \sim D$ と表す.

また, ILP におけるアルゴリズムの多くはパターンを制限するために述語の引数にモードを利用している. これは各引数が入力引数であるか出力引数であるかを表す情報で, 前者を入力モード, 後者を出力モードという. R_{fam} に現れる述語にはそれぞれ $p(+, -), m(+), f(+)$ のようにモードが与えられている. ただし $+/-$ は入力/出力を意味する.

さて, 述語のモードについて注目すると, 述語は 2 つのクラスに分けることができる. 例えば $p(+, -)$ のように少なくとも 1 つの $(-)$ -引数を持つ述語を経路述語とよび, これは入力の項から出力の項へ導く機能をもつ. $m(+)$ と $f(+)$ のように $(-)$ -引数を持たない述語を判定述語とよび, その項を持つ特徴を表す機能をもつ. また, 経路/判定述語で構成されたリテラルを経路/判定リテラルとよぶ.

以上で説明した概念を用いて MAPIX における基本パターンである性質と, それを変数化した性質アイテムを定義する.

定義 4 (性質) データベース R について事例 e に関する性質 pr は, 以下の条件を満たす R 中のリテラルの極小集合である.

1. pr は, 必ずただ一つの判定リテラルを含む.
2. pr は, 全てのリテラル $l (l \in pr)$ の入力引数に現れる項が, そのリテラルより前のリテラルの引数に現れるように順序づけを与えることができる.

定義 5 (変数化) ある基礎節 α について, 以下の条件を満たすような節 β は α の変数化である.

1. β は基礎項を持つようなリテラルを含まない.
2. 以下の条件を満たす代入 $\theta = \{v_1/t_1, \dots, v_n/t_n\}$ が存在する.
 - (a) $\alpha = \beta\theta$
 - (b) t_1, \dots, t_n は, すべて異なる項である.

定義 6 (性質アイテム) 基礎リテラルの集合 $L = \{l_1, \dots, l_m\}$ と事例 e について, $e \leftarrow l_1 \wedge \dots \wedge l_m$ を変数化したものを $\text{var}(e \leftarrow L)$ と表記する. L が e の性質であるとき, $\text{var}(e \leftarrow L)$ を性質アイテムと呼ぶ.

表 1: $e = \text{gf}(01)$ の性質と性質アイテム

性質 (properties)
$pr0 = \{m(01)\}$
$pr1 = \{p(01,02), f(02)\}$
$pr2 = \{p(01,02), p(02,03), f(03)\} (\{p(01,02), p(02,04), f(04)\})$
$pr3 = \{p(01,02), p(02,03), p(03,05), m05\}$
$pr4 = \{p(01,02), p(02,04), p(04,06), f06\}$
性質アイテム (property items)
$it0 = \text{gf}(A) \leftarrow m(A).$
$it1 = \text{gf}(A) \leftarrow p(A, B) \wedge f(B).$
$it2 = \text{gf}(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C).$
$it3 = \text{gf}(A) \leftarrow p(A, B) \wedge p(B, C) \wedge (C, D) \wedge m(D).$
$it4 = \text{gf}(A) \leftarrow p(A, B) \wedge p(B, C) \wedge (C, D) \wedge f(D).$

R_{fam} において, $pr = \{p(01,02), p(02,03), f(03)\}$ は事例 $e = \text{gf}(01)$ の性質である. これから生成される性質アイテムは $it = \text{var}(e \leftarrow pr) = \text{gf}(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C)$ となる. 表 1 に事例 $\text{gf}(01)$ から抽出される性質と性質アイテムを全て示す. $pr2$ のように同値な性質が存在する場合は同値な性質アイテムを生成しないようにしてパターンの重複を避ける.

以下に MAPIX アルゴリズムの概要を示す:

1. いくつかの事例をサンプリングする
2. サンプリングした事例から性質を全て抽出する
3. 性質を変数化し性質アイテムを生成する
4. アプリアリアルゴリズム [1] と同様の手法でアイテムの頻出な組み合わせを枚挙する

事例 $\text{gf}(01)$ から生成される性質アイテムについて, 最低支持度を 60% とするとこれらの性質アイテムは全て頻出である. しかし, 事例 $\text{gf}(20)$ から生成される性質アイテム $\text{gf}(A) \leftarrow p(A, B) \wedge p(B, C) \wedge m(C)$ は頻出ではない.

性質アイテムを組み合わせるとき, MAPIX では連言で性質アイテムを結合している. 例えば「孫娘をもつ」という意味を表す $it2$ と「ひ孫娘をもつ」という意味を表す $it4$ を組み合わせたパターンは以下ようになる.

$$\langle it2, it4 \rangle = \text{gf}(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge p(A, D) \wedge p(D, E) \wedge (E, F) \wedge f(F).$$

これは「孫娘とひ孫娘をもつ」という意味である. このように MAPIX では, それぞれの性質アイテムについて目標述語に含まれる変数以外を他の性質アイテムの変数と共通しないように変数化し性質アイテム同士は独立であるようにしている. つまり, $it2$ と $it4$ を組み合わせても実際に事例に現れる「娘と孫娘をもつ子供をもつ」や「娘を持つ孫娘をもつ」という意味を表すパターンは生成されない.

3. 提案手法

MAPIX では性質アイテム同士を独立に組み合わせていたため生成できないパターンがあった. そこで, 性質アイテムの構造的な組み合わせを用いて MAPIX で生成されないパターンを生成する手法を提案する. まず, 性質アイテムの事例への出現の仕方を保存するために性質アイテムの影を定義する.

定義 7 データベース R と性質アイテム it について, 以下のような集合 $\text{shadow}(it, R)$ を性質アイテム it の影という.

$$\text{shadow}(it, R) = \{ 'e \leftarrow L' \in T \times 2^R \mid$$

L は事例 e に関する性質であり $\text{var}(e \leftarrow L) \sim it \}$ ただし, T は目標述語の事例集合であり, 2^X は集合 X のべき集合である.

次に構造的な組み合わせの方法である性質の融合と、それにより生成される分子アイテムを定義する。

定義 8 (性質の融合) データベース R について、事例 e に関する変数化をしていない性質の集合を P とする。 $\{p_1, \dots, p_n\} \subseteq P$ について、 $L_P = p_1 \cup \dots \cup p_n$ とし、リテラル集合 $L_P = \{\ell_1, \dots, \ell_k\}$ から節 $e \leftarrow \ell_1 \wedge \dots \wedge \ell_k$ を生成し変数化することを、性質 p_1, \dots, p_n の融合という。

定義 9 性質アイテムセット $\langle it_{i_1}, \dots, it_{i_n} \rangle$ に対し、影の組み合わせ、

$$\begin{aligned} & \langle e \leftarrow pr_{i_1}, \dots, e \leftarrow pr_{i_n} \rangle \\ & \in \text{shadow}(it_{i_1}, R) \times \dots \times \text{shadow}(it_{i_n}, R) \\ & \text{s.t. } \bigcap_{j=1, \dots, n} (\text{terms}(pr_{i_j}) - \text{terms}(e)) \neq \emptyset \end{aligned}$$

が存在するとき、この性質アイテムセットは融合可能であるという。ただし、 $\text{terms}(p)$ は p の中の全ての項の集合である。

定義 10 (分子アイテム) 性質アイテムセット $is = \langle it_{i_1}, \dots, it_{i_n} \rangle$ について融合可能な影の組み合わせ $\langle e \leftarrow pr_{i_1}, \dots, e \leftarrow pr_{i_n} \rangle$ が存在するとき、

$$\text{var}(e \leftarrow \bigcup_{j=1, \dots, n} pr_{i_j})$$

を is から生成された分子アイテムとよぶ。

例えば、性質アイテム it_2 と it_4 の影は次のようになる。

$$\begin{aligned} \text{shadow}(it_2, R_{\text{fam}}) = \{ & \\ & gf(01) \leftarrow \{p(01, 02), p(02, 03), f(03)\}, \\ & gf(01) \leftarrow \{p(01, 02), p(02, 04), f(04)\}, \\ & gf(07) \leftarrow \{p(07, 08), p(08, 09), f(09)\}, \\ & gf(12) \leftarrow \{p(12, 13), p(13, 15), f(15)\}, \\ & gf(12) \leftarrow \{p(12, 14), p(14, 16), f(16)\}, \\ & gf(19) \leftarrow \{p(19, 20), p(20, 21), f(21)\}, \\ & gf(19) \leftarrow \{p(19, 20), p(20, 22), f(22)\}, \\ & gf(20) \leftarrow \{p(20, 21), p(21, 24), f(24)\} \} \end{aligned}$$

$$\begin{aligned} \text{shadow}(it_4, R_{\text{fam}}) = \{ & \\ & gf(01) \leftarrow \{p(01, 02), p(02, 04), p(04, 06), f(06)\}, \\ & gf(07) \leftarrow \{p(07, 08), p(08, 09), p(09, 11), f(11)\}, \\ & gf(19) \leftarrow \{p(19, 20), p(20, 21), p(21, 24), f(24)\} \} \end{aligned}$$

性質アイテムセット $\langle it_2, it_4 \rangle$ の融合を考えると、融合可能な影の組み合わせの有無をみる。 $\langle it_2, it_4 \rangle$ は以下の 5 つの融合可能な影の組み合わせを持つ。

$$\begin{aligned} & \langle gf(01) \leftarrow \{p(01,02), p(02,03), f(03)\}, \\ & gf(01) \leftarrow \{p(01,02), p(02,04), p(04,06), f(06)\} \rangle \\ & \langle gf(01) \leftarrow \{p(01,02), p(02,04), f(04)\}, \\ & gf(01) \leftarrow \{p(01,02), p(02,04), p(04,06), f(06)\} \rangle \\ & \langle gf(07) \leftarrow \{p(07,08), p(08,09), f(09)\}, \\ & gf(07) \leftarrow \{p(07,08), p(08,09), p(09,11), f(11)\} \rangle \\ & \langle gf(19) \leftarrow \{p(19,20), p(20,21), f(21)\}, \\ & gf(19) \leftarrow \{p(19,20), p(20,21), p(21,24), f(24)\} \rangle \\ & \langle gf(19) \leftarrow \{p(19,20), p(20,22), f(22)\}, \\ & gf(19) \leftarrow \{p(19,20), p(20,21), p(21,24), f(24)\} \rangle \end{aligned}$$

これらの組み合わせから融合を行い同値なものを生成しないようにすると、次の 2 つの分子アイテムが生成される。

$$\begin{aligned} it_{2-4} &= gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge p(C, D) \wedge f(D). \\ it_{2-4'} &= gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge \\ & p(B, D) \wedge p(D, E) \wedge f(E). \end{aligned}$$

it_{2-4} は「娘をもつ孫娘をもつ」という意味を表し、 $it_{2-4'}$ は「娘と孫娘をもつ子供をもつ」という意味を表す。これらは MAPIX では生成できなかったパターンである。 R_{fam} における頻出な性質アイテムと分子アイテムの全てを表 2 に示す。

表 2: 最低支持度が 60% のときの、 R_{fam} における頻出な性質アイテムと分子アイテム。

it0	=	$gf(A) \leftarrow m(A).$
it1	=	$gf(A) \leftarrow p(A, B) \wedge f(B).$
it2	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C).$
it3	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge p(C, D) \wedge m(D).$
it4	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge p(C, D), f(D).$
it1-2	=	$gf(A) \leftarrow p(A, B) \wedge f(B) \wedge p(B, C) \wedge f(C).$
it1-3	=	$gf(A) \leftarrow p(A, B) \wedge f(B) \wedge$ $p(B, C) \wedge p(C, D) \wedge m(D).$
it2-3	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge$ $p(C, D) \wedge m(D).$
it2-3'	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge$ $p(B, D) \wedge p(D, E) \wedge m(E).$
it2-4	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge$ $p(C, D) \wedge f(D).$
it2-4'	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge$ $p(B, D) \wedge p(D, E) \wedge f(E).$
it3-4	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge p(C, D) \wedge f(D) \wedge$ $p(B, E) \wedge p(E, F), m(F).$
it1-2-3	=	$gf(A) \leftarrow p(A, B) \wedge f(B) \wedge$ $p(B, C) \wedge f(C) \wedge p(C, D) \wedge m(D).$
it1-2-3'	=	$gf(A) \leftarrow p(A, B) \wedge f(B) \wedge p(B, C) \wedge f(C) \wedge$ $p(B, D) \wedge p(D, E) \wedge m(E).$
it2-3-4	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge f(C) \wedge$ $p(C, D) \wedge f(D) \wedge p(B, E) \wedge p(E, F) \wedge m(F).$
it2-3-4'	=	$gf(A) \leftarrow p(A, B) \wedge p(B, C) \wedge p(C, D) \wedge f(D) \wedge$ $p(B, E) \wedge f(E) \wedge p(E, F) \wedge m(F).$

また、性質アイテムと分子アイテムを独立に組み合わせるとあらたなパターンを得ることができる。例えば、 it_1 と it_{2-4} を組み合わせたアイテムセットは次のようになる。

$$\begin{aligned} \langle it_1, it_{2-4} \rangle &= gf(A) \leftarrow p(A, B) \wedge f(B) \wedge \\ & p(A, C) \wedge p(C, D) \wedge f(D) \wedge p(D, E) \wedge f(E). \end{aligned}$$

これは「娘と娘を持つ孫娘をもつ」という意味を表す。

以上の考えを用いて提案手法の概要を次に示す(詳細は表 3):

1. 目標事例をサンプリングする
2. サンプリングした事例から性質を抽出し性質アイテムを生成する
3. アプリオリ同様のアルゴリズムで性質アイテムの頻出な組み合わせを全て枚挙する
4. 性質アイテムの頻出な組み合わせから分子アイテムを生成する
5. 再度アプリオリと同様のアルゴリズムで性質アイテムと分子アイテムの頻出な組み合わせを全て枚挙する

ステップ 4 では頻出な性質アイテムセットから分子アイテムを生成する。これは性質アイテムセット $\langle it_1, \dots, it_N \rangle$ と分子アイテム $it_{1-\dots-N}$ の支持度が $\text{supp}(\langle it_1, \dots, it_N \rangle) \leq \text{supp}(it_{1-\dots-N})$ の関係にあるためである。この式は性質アイテムセットよりそれから生成される分子アイテムの方が支持度が低くなるということを表している。つまり頻出でない性質アイテムセットから生成される分子アイテムもまた頻出ではないので生成する必要がない。

4. 実験とまとめ

2 つのデータを用いて実験を行った。1 つの実験は R_{fam} を用いたもので提案手法により従来よりも多種のパターンが生

表 3: 提案手法のアルゴリズム.

input R : 関係 DB; T : 目標事例; sup_{\min} : 最低支持度 (%);
output Freq : 頻出アイテムセットの集合;
 1. サンプルングする目標事例を選択 $T' \subseteq T$;
 2. $\text{Items} := \emptyset$; $\mathcal{P} := \emptyset$; $\text{Freq} := \emptyset$;
 3. **For each** $e \in T'$ **do** $\mathcal{P} := \mathcal{P} \cup \{e \leftarrow \text{pr} \mid \text{pr} \text{ は } e \text{ の性質}\}$;
 4. **For each** ' $e \leftarrow \text{pr}$ ' $\in \mathcal{P}$ **do**
 5. **If** $\exists I \in \text{Items}, I \sim \text{var}(e \leftarrow \text{pr})$ **then**
 6. $S[I] := S[I] \cup \{e \leftarrow \text{pr}\}$; % $S[I] = \text{shadow}(I, R)$
 7. **else** $I' = \text{var}(e \leftarrow \text{pr})$;
 8. $S[I'] := \{e \leftarrow \text{pr}\}$;
 9. $\text{Items} := \text{Items} \cup \{I'\}$;
 10. $k := 1$;
 11. $\mathcal{F}_1^1 := \{\langle I \rangle \mid I \in \text{Items} \text{ かつ } \text{supp}(I) \geq \text{sup}_{\min}\}$;
 12. $\text{Freq} := \mathcal{F}_1^1$;
 13. **While** $\mathcal{F}_k^1 \neq \emptyset$ **do**
 14. $\mathcal{C}_{k+1} := \text{CANDIDATE}(\mathcal{F}_k^1, \mathcal{F}_k^1)$;
 15. $\mathcal{F}_{k+1}^1 := \{IS \in \mathcal{C}_{k+1} \mid \text{supp}(IS) \geq \text{sup}_{\min}\}$;
 16. $\text{Freq} := \text{Freq} \cup \mathcal{F}_{k+1}^1$;
 17. $k := k + 1$;
 18. $\text{Comb} := \text{CANDICOMB}(\text{Freq})$;
 19. $k := 1$;
 20. $\mathcal{F}_1^2 := \{\langle I \rangle \mid I \in \text{Comb} \text{ かつ } \text{supp}(I) \geq \text{sup}_{\min}\}$;
 21. $\text{Freq} := \text{Freq} \cup \mathcal{F}_1^2$;
 22. **While** $\mathcal{F}_k^2 \neq \emptyset$ **do**
 23. $\mathcal{C}_{k+1} := \text{CANDIDATE}(\mathcal{F}_k^2, \mathcal{F}_k^2)$;
 24. $\mathcal{F}_{k+1}^2 := \{IS \in \mathcal{C}_{k+1} \mid \text{supp}(IS) \geq \text{sup}_{\min}\}$;
 25. $\text{Freq} := \text{Freq} \cup \mathcal{F}_{k+1}^2$;
 26. $k := k + 1$;
 27. **Return** Freq ;

$\text{CANDIDATE}(\mathcal{F}_k^1, \mathcal{F}_k^2)$:

input $\mathcal{F}_k^1, \mathcal{F}_k^2$: あるレベルの頻出アイテムセットの集合;
output \mathcal{C}_{k+1} : 次のレベルの候補アイテムセットの集合;
 1. $\mathcal{C}_{k+1} := \emptyset$;
 2. **For each** $\langle \langle I_1, \dots, I_k \rangle, \langle I'_1, \dots, I'_k \rangle \rangle \in \mathcal{F}_k^1 \times \{\mathcal{F}_k^1 \cup \mathcal{F}_k^2\}$,
 ただし $I_1 = I'_1, \dots, I_{k-1} = I'_{k-1}$ かつ $I_k < I'_k$ **do**
 3. $\mathcal{C}_{k+1} := \mathcal{C}_{k+1} \cup \{\langle I_1, \dots, I_{k-1}, I_k, I'_k \rangle\}$;
 4. **For each** $IS \in \mathcal{C}_{k+1}$ **do**
 5. **If** $k = 1$ かつ $(I \leq I'$ または $I' \leq I)$,
 ただし $IS = \langle I, I' \rangle$ **then**
 6. \mathcal{C}_{k+1} から IS を取り除く;
 7. **For each** $I \in IS$ **do if** $IS - \{I\} \notin \mathcal{F}_k^1 \cup \mathcal{F}_k^2$ **then**
 8. \mathcal{C}_{k+1} から IS を取り除く;
 9. **Return** \mathcal{C}_{k+1} ;

$\text{CANDICOMB}(\text{Freq})$:

input Freq : 頻出性質アイテムセットの集合;
output Comb : Freq から生成される分子アイテムの集合;
 1. $\text{Comb} := \emptyset$;
 2. **For each** $\langle I_1, \dots, I_n \rangle \in \text{Freq}$ **for** $n \geq 2$ **do**
 3. **For each** $\langle e_1 \leftarrow \text{pr}_1, \dots, e_n \leftarrow \text{pr}_n \rangle$
 $\in S[I_1] \times \dots \times S[I_n]$ **do**
 4. **If** $e_1 = \dots = e_n$ かつ
 $\bigcap_{j=1, \dots, n} (\text{terms}(\text{pr}_j) - \text{terms}(e_j)) \neq \emptyset$ **then**
 5. $\text{Comb} := \text{Comb} \cup \{\text{var}(e_1 \leftarrow \bigcup_{j=1, \dots, n} \text{pr}_j)\}$;
 6. **For each** $I \in \text{Comb}$ **do**
 7. **If** $\exists I' \in \text{Comb}, I' \neq I$ かつ $I' \sim I$ **then**
 8. Comb から I を取り除く;
 9. **Return** Comb ;

表 4: R_{fam} を用いた実験結果.

	最低支持度	20%	40%	60%	80%
MAPIX	パターン数	55	31	23	11
[4]	時間 (sec.)	0.04	0.01	0.01	0.01
EQUIVPIX	パターン数	441	153	51	15
[3]	時間 (sec.)	6.23	0.45	0.06	0.03
提案手法	パターン数	4601	1063	109	17
	時間 (sec.)	9.55	0.54	0.08	0.01

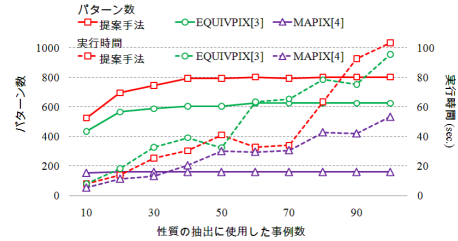


図 2: Bongard を用いた実験結果.

成されることを確認した. 表 4 は最低支持度をを変えながら, それぞれのアルゴリズムにより生成されるパターン数と実行時間を示している. これらのパターンの中には同値であるという意味での重複がないことを確認している. 提案手法により多種のパターンが生成されていることを確認できる.

もう一つの実験は図形に関するデータ Bongard を用いたものである. 図 2 ではサンプルングする事例数を変えながら, 生成されるパターン数と実行時間を示している. これらはそれぞれ 10 回実行したときの平均である. 支持度を計算するには全ての事例を使うが, 性質を抽出するときはサンプルングされたいくつかの事例を用いる. 提案手法では 80 事例をサンプルングすることで全事例 392 をサンプルングして得られる全パターン (802 個) を生成した. またこれらのパターンに重複がないことを確認している.

本研究では MRDM における MAPIX アルゴリズムを発展させた手法を提案した. 提案手法により従来手法よりも多種のパターンを出力することが可能になった. しかし, サンプルング事例のサイズが大きくなると処理時間が爆発的に増えてしまった. 実際, Bongard で全事例 392 をサンプルングして実行すると処理時間が 6556 秒かかった. これはサンプルング事例数が増加するに伴い性質アイテムの影の組み合わせが膨大になるためだと思われる. 効率的な分子アイテムの生成方法が今後の課題として挙げられる.

参考文献

- [1] R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules. VLDB, pp. 487–499, 1994.
- [2] L. Dehaspe and L. De Raedt. Mining association rules with multiple relations. ILP97, pp.125–132, 1997.
- [3] N. Inuzuka, J. Motoyama, S. Urazawa and T. Nakano. Relational pattern mining based on equivalent classes of properties extracted from samples. PAKDD2008, pp. 582–591, 2008.
- [4] J. Motoyama, S. Urazawa, T. Nakano and N. Inuzuka. A mining algorithm using property items extracted from sampled examples, ILP2006, pp.335–350, 2007.