

## 2部クリークを用いた closed item set の効率的な列挙

宇野 毅明<sup>†</sup> 有村 博紀<sup>††</sup> 浅井 達哉<sup>††</sup>

<sup>†</sup> 国立情報学研究所, 〒 101-8430 東京都千代田区一ツ橋 2-1-2

<sup>††</sup> 九州大学システム情報科学研究院, 〒 812-0053 福岡県福岡市東区箱崎 6-10-1

E-mail: <sup>†</sup>uno@nii.jp, <sup>††</sup>{arim,t-asai}@i.kyushu-u.ac.jp

あらまし 顧客の売上データのような、アイテムの部分集合族により定められるデータに対して、その族のある一定数以上の要素に含まれるアイテム集合を頻出集合とよぶ。頻出集合はデータマイニングの分野への応用を持つ。近年、極大な頻出集合と、同じような頻出集合を1つにまとめて扱う closed item set が注目されている。それぞれを列挙することにより、頻出集合の中から意味のある部分を効率よく抽出できるからである。本稿では、極大2部クリークを列挙することにより、closed item set を高速に列挙する手法と、それを用いて極大頻出集合を列挙する手法を提案する。さらに、ベンチマーク問題を用いた計算実験により、既存の手法よりも高速であることを示す。

キーワード 数え上げ, 発生, 計算量, アルゴリズム, 疎グラフ, 頻出集合発見, データマイニング

## Enumerating Closed Item Sets via Maximal Bipartite Cliques

Takeaki UNO<sup>†</sup>, Hiroki ARIMURA<sup>††</sup>, and Tatsuya ASAI<sup>††</sup>

<sup>†</sup> National Institute of Informatics, 2-1-2, Hitotsubashi, Chiyoda-ku, Tokyo, JAPAN, 101-8430

<sup>††</sup> Information Science and Electrical Engineering, Kyushu University, 6-10-1, Hakozaki, Higashi-ku, Fukuoka-shi, Fukuoka, JAPAN, 812-0053

E-mail: <sup>†</sup>uno@nii.jp, <sup>††</sup>{arim,t-asai}@i.kyushu-u.ac.jp

**Abstract** For an item set and its subset family, a frequent set is a subset of the item set included in at least a specified number of elements of the subset family. Frequent sets have some applications in data mining. Recently, maximal frequent set and closed item sets are remarked where a closed item set is a group of similar frequent sets, since we can obtain necessary frequent sets by enumerating closed item sets. In this paper, we propose an efficient method for enumerating closed item sets by using enumeration of maximal bipartite cliques, and enumeration algorithm for maximal frequent sets via closed item sets. We evaluate their performances by computational experiments using benchmark problems, and show that our algorithms are faster than existing algorithms.

**Key words** listing, generation, frequent set mining, algorithm, complexity, sparse graph, data mining

### 1. はじめに

近年、データマイニングの分野で、データの中から頻出する部分構造を見つけ出す問題が脚光を浴びている。頻出する部分構造は、アソシエーションルール発見などに応用があるほか、データの特徴を捉えるためにも有効だからである。この種の問題の中で最も古くから研究されており、かつ最も重要な問題が、頻出集合列挙問題である [1], [2]。この問題は、与えられたアイテム集合と、その部分集合（トランザクションとよぶ）の集合からなるデータに対して、アイテムセットの部分集合で、ある一定数以上のトランザクションに含まれるものを列挙する問題

である。しかし、頻出集合の数は巨大であることが多く、通常は頻出集合の極大元のみを列挙するという手法が用いられてきた。

しかし最近、極大元のみでは元データの性質を十分に表現できていないのではないかという疑問が提示され始めている。そこで注目されてきているものが、closed item set とよばれるものである。詳しい定義は後述するが、簡単に説明すると、含まれるトランザクション集合が等しい頻出集合を同一視したものである。closed item set は頻出集合と極大頻出集合の間にあるモデルであり、頻出集合に比べて、それほど情報を失っていないと考えられ、頻出集合ほど大量には存在しないと考えられる。

closed item set の列挙アルゴリズムはいくつか提案されているが [7], [10], これらのアルゴリズムの closed item set 1 つあたりの計算時間は入力の多項式で押さえられていない．本研究では, 宇野 [9] により提案された極大 2 部クリーク列挙アルゴリズムをもとにして, closed item set と極大頻出集合を列挙するアルゴリズムを提案する．closed item set / 極大頻出集合 1 つあたりの計算時間は  $O(\text{トランザクションのサイズの最大}^3)$  である．宇野のアルゴリズムは, 疎なデータに対して高速であることが実験的に確認されており, 本研究のアルゴリズムも現実の疎な頻出集合列挙問題に対して高速であると考えられる．さらに, 提案するアルゴリズムを計算機に実装し, KDD-cup 2000 [6] で用いられたベンチマーク問題を解いた．計算結果より, これらのデータに対する本手法の計算時間が既存アルゴリズムよりも大幅に短いことを示す．

## 2. 記法と問題の導入

グラフ  $G = (V, A)$  の頂点集合  $V$  が  $V_1$  と  $V_2$  に分割できて, 任意の  $V_1$  の頂点の組, および任意の  $V_2$  の頂点の組の間に枝がないとき,  $G$  は 2 部グラフであるという．本稿では, 全ての  $V_1$  の頂点と隣接する  $V_2$  の頂点は存在しないとする． $G$  の最大次数を  $\Delta$  で表記する．

$G$  の頂点部分集合  $H \subseteq V_1, K \subseteq V_2$  に対して,  $H$  の任意の頂点と  $K$  の任意の頂点の間に枝があるとき,  $H$  と  $K$  を合わせた頂点集合を 2 部クリークとよぶ． $K = \emptyset, H = V_1$  である場合, あるいはその逆である場合も, 2 部クリークとよぶ．ある 2 部クリークが他の 2 部クリークに含まれないとき, その 2 部クリークを極大 2 部クリークとよぶ．

$E$  をアイテムの集合,  $\mathcal{F}$  を  $E$  の部分集合族とし,  $\mathcal{F}$  の要素をトランザクションとよぶ． $E$  の部分集合  $K$  に対して,  $X(K) = \{F | K \subseteq F, F \in \mathcal{F}\}$  とする．与えられた定数  $\alpha$  に対して,  $|X(K)| \geq \alpha$  を満たす  $K$  は頻出集合とよばれる． $K$  が他の頻出集合の部分集合で無いなら,  $K$  は極大頻出集合とよばれる．あるトランザクションの集合  $B \subseteq \mathcal{F}$  に対して,  $\{K | X(K) = B\}$  を closed item set とよぶ．任意の closed item set に対して, その極大元は 1 つしか存在しない．証明を以下に示そう．

証明:

任意の  $K, K', K \neq K', K, K' \subseteq E$  に対して  $K$  と  $K'$  が同じ closed item set に属するならば,  $X(K) = X(K')$  であり, 任意の  $F \in X(K)$  は  $K$  と  $K'$  を含む．よって,  $X'' = K' \cup K$  も,  $X(K'') = X(K)$  を満たし, closed item set に含まれるので,  $K$  と  $K'$  の少なくとも片方は極大元でない．

頻出集合は, データマイニングにおけるアソシエーションルール発見問題に応用がある．しかし, 頻出集合の数は通常巨大であるので, 極大頻出集合のみを列挙して使用するというモデルが主流であり, そのための列挙アルゴリズムの研究も盛んに行われている．しかし, 極大頻出集合の列挙はさほど効率良く行えないこと (極大頻出集合 1 つあたり多項式時間のアルゴリズムはまだ提案されていない), および極大頻出集合は頻出

集合に比べて数が少ないため, 頻出集合の持つ情報を完全には表現できていないと考えられることから, 近年は closed item set に注目が集まっている．closed item set は「同じトランザクション集合に含まれるアイテム集合」の集合であるので, ある意味で同じ特徴を持つアイテム集合の集合と考えられる．また, closed item set の極大元の集合は極大頻出集合の集合を含むため, closed item set は極大頻出集合の情報を完全に含む．つまり, closed item set は頻出集合と極大頻出集合の間にあるモデルであり, 頻出集合に比べて, それほど情報を失っていないと考えられるからである．

近年, closed item set の極大元の中で頻出であるものを列挙する問題が研究されている．アルゴリズムも, Pei, J. Han, R. Mao による closet [7], M. J. Zaki, C. Hsiao による CHARM [10] などが提案されている．しかし, これらは出力線形時間ではない．つまり, 出力する closed item set 1 つあたりの計算時間が入力の多項式時間で押さえられていない．近年, 牧野ら [4] などは, closed item set の列挙問題と 2 部グラフの極大 2 部クリークの列挙問題が等価であることを示した．極大 2 部クリークは, 築山らのアルゴリズム [8] などにより, 極大 2 部クリーク 1 つあたり  $O(\text{枝数} \times \text{頂点数})$  の時間で列挙できる．また, グラフが疎である場合には, 宇野のアルゴリズム [9] を用いて, 1 つあたり  $O(\text{最大次数}^3)$  時間で列挙できる．[9] のアルゴリズムを実装した計算機実験では, ランダムグラフに対する実用上の計算時間は  $O(\text{平均次数}^2)$  程度であることが示されている．

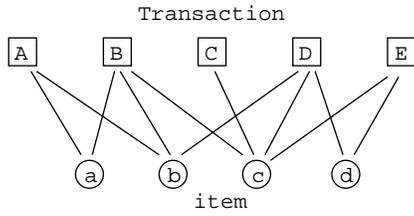
本研究では, 頻出な closed item set の極大元のみを出力するように宇野のアルゴリズムを改良し, さらに実装上のアルゴリズム的な工夫を加え, 高速な列挙アルゴリズムを構築した．通常, closed item set 列挙問題の入力は疎なデータである．つまり, 各トランザクションは,  $E$  の要素のごく一部しか含まないことが多い．このような疎な問題は, 最大次数・平均次数の小さな 2 部グラフの 2 部クリーク列挙問題に変換される．極大 2 部クリーク列挙アルゴリズムを用いることによって, 理論的にも, 実用上も, 高速な closed item set の極大元の列挙ができると考えられる．

また, 極大頻出集合は closed item set の極大元であるため, closed item set を列挙し, その中で極大なものだけを出力することにより, 極大頻出集合を列挙することができる．極大頻出集合の列挙が困難とされる, 頻出集合数が巨大な問題に対しても, closed item set の数はそれほど大きくない場合が多く, 効率の良い列挙が行えると考えられる．

## 3. 極大 2 部クリーク列挙問題への変換

与えられたアイテム集合  $E$  とトランザクションの集合  $\mathcal{F}$  に対して, closed item set の極大元は, 以下のように構築した 2 部グラフの極大 2 部クリークと 1 対 1 対応する．

- ・頂点集合  $V_1 = \mathcal{F}$  とする．つまり  $V_1$  の各頂点は各トランザクションと対応する
- ・頂点集合  $V_2 = E$  とする． $V_2$  の各頂点は各アイテムと対応する



Transactions	Maximal bipartite cliques with two or more transactions
A = {a, b}	{a, b, A, B}
B = {a, b, c}	{a, b, A, B}    {b, A, B, D}
C = {c}	{b, c, B, D}    {c, B, C, D, E}
D = {b, c, d}	{c, d, D, E}
E = {c, d}	

図 1 頻出集合問題を 2 部グラフで表現: トランザクションを 2 つ以上含む極大 2 部クリークは  $\{a, b, A, B\}$ ,  $\{b, A, B, D\}$ ,  $\{b, c, B, D\}$ ,  $\{c, B, C, D, E\}$ ,  $\{c, d, D, E\}$  の 5 個で, それぞれの  $V_2$  に含まれる頂点 (例えば  $a, b$ ) が closed item set の極大元に対応

・アイテム  $e \in E$  がトランザクション  $F$  に含まれるとき, またそのときに限り,  $e \in V_2$  と  $F \in V_1$  を枝で結ぶ.

図 1 に一例を示した. このグラフの極大 2 部クリークが closed item set と対応することを以下で示そう. 集合  $K$  に対して,  $X(K)$  は  $V_1$  の頂点集合に対応する. グラフの構築方法から,  $K$  の任意の頂点と,  $X(K)$  の任意の頂点の間には枝が張られている. そのため,  $K$  と  $X(K)$  はグラフの 2 部クリークになる.  $X$  がある closed item set の極大元であるならば, 任意のアイテム  $e \notin K$  に対して,  $X(K \cup \{e\}) \neq X(K)$  である. よって, あるトランザクション  $F \in X(K)$  が存在して,  $e$  と  $X(K)$  の間には枝が存在しない. また,  $X(K)$  の定義より, 任意のトランザクション  $F \notin X(K)$  に対して, あるアイテム  $e \in K$  が存在して,  $F$  と  $e$  の間には枝が存在しない. これは,  $K$  と  $X(K)$  を合わせた頂点集合が極大 2 部クリークであることを意味する.

逆に, 頂点集合  $H \subseteq V_1$  と  $K \subseteq V_2$  が極大 2 部クリークであれば,  $H = X(K)$  が成り立つ. さらに,  $X(K) = H$  となるような  $V_1$  の部分集合の中で,  $K$  が極大であることから,  $K$  は closed item set の極大元となる. 以上より, 構築されたグラフの極大 2 部クリークと, closed item set の極大元とが 1 対 1 対応することが示された.

#### 4. 宇野アルゴリズム

次に, [9] で提案されている極大 2 部クリーク列挙アルゴリズムを簡単に解説しよう. そのために, いくつかの記法を定義する.

2 部グラフ  $G = (V_1 \cup V_2, A)$  に対して,  $V_1 = \{v_1, \dots, v_m\}$ ,  $V_2 = \{u_1, \dots, u_n\}$  とする. 頂点集合  $H \subseteq V_1$  に対して,  $X(H)$  を  $H$  の全ての頂点に隣接する  $V_2$  の頂点の集合とする.  $K \subseteq V_2$  に対しても同様に  $X(K)$  を定める.  $H \subseteq V_1$  と  $K \subseteq V_2$  を合わせた頂点集合が極大 2 部クリークであれば, またそのときに限り,  $X(H) = K$  かつ  $X(K) = H$  が成り立つ. そこで, 以後極大 2 部クリークを, その  $V_2$  に含まれる頂点のみで定める. 頂点  $v$  に対して,  $X(\{v\})$  は  $v$  に隣接する頂点の

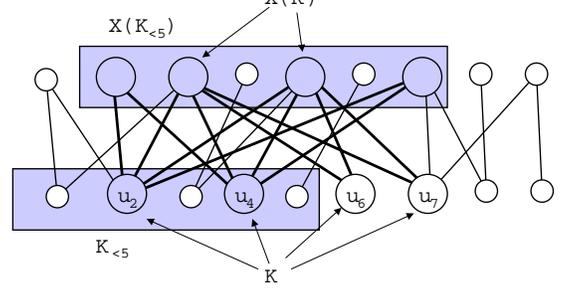


図 2 極大 2 部クリーク  $K$  と  $K_{\leq 5}, X(K), X(K_{\leq 5})$

集合になる. 頂点集合  $K$  に対して,  $K_{\leq i}$  を添え字が  $i$  以下の  $K$  の頂点の集合とする.

極大 2 部クリーク  $K \subseteq V_2$  に対して, その親を  $X(X(K_{\leq i}))$  で定める. ただし  $i$  は  $X(K_{\leq i-1}) \neq X(K)$  となる最大の  $i$  とし, この  $i$  を  $K$  の親添え字とよぶ. 任意の極大 2 部クリーク  $K \neq \emptyset$  に対して, その親が唯一的に定まる. また,  $K_{\leq i} \subset K$  であることから,  $X(K_{\leq i-1}) \neq X(K)$  であるなら,  $X(K) \subset X(K_{\leq i})$  であり, よって

$$X(X(K_{\leq i})) \subset X(X(K)) = K$$

となる. ゆえに, この親子関係は非巡回的であり, 極大 2 部クリーク  $\emptyset$  を根とする木構造を導出する. この木構造を列挙木とよぶ. 列挙木を根から出発して深さ優先探索することにより, 全ての極大 2 部クリークを列挙できる.

列挙木を深さ優先探索するには, 与えられた極大 2 部クリークの子供を列挙するアルゴリズムを構築し, 根の子供を列挙した後に, 再帰呼び出しによってそれら子供の子孫を列挙すれば良い. 計算を簡略化するため, 極大 2 部クリーク  $V_2$  は列挙する対象から除外する.

極大 2 部クリーク  $K$  に対して,  $K[i]$  を

$$K[i] = X(X(K_{\leq i} \cup \{u_i\}))$$

で定める. 極大 2 部クリーク  $H$  が  $K$  の子供であり,  $H$  の親添え字が  $i$  であるとする. このとき, 親の定義から  $K = X(X(H_{\leq i-1}))$  であり,  $u_i \in H, u_i \notin K$  であることから,

$$\begin{aligned} X(K) \cap X(\{u_i\}) &= X(K_{\leq i-1} \cup \{u_i\}) \\ &= X(H_{\leq i}) \end{aligned}$$

が成り立つ. よって,  $H = K[i]$  が成り立つ. 逆に,  $K$  の親添え字よりも大きな添え字  $i$  に対して  $K[i]_{\leq i-1} = K_{\leq i-1}$  が成り立つならば,  $X(K[i]) = X(K[i]_{\leq j})$  が任意の  $j \geq i$  について成り立つので,  $K[i]$  は  $K$  の子供である. よって,  $K$  の子供は,  $K[i]$  を各  $i$  について生成し,  $K[i]_{\leq i-1} = K_{\leq i-1}$  が成り立つかどうかを判定することにより, 全てを見つけることができる. この判定は全ての  $i$  について行う必要は無い. 以下の条件を満たすものだけで十分である.

- ・  $i$  は  $K$  の親添え字よりも大きい
- ・  $u_i$  は  $X(K)$  のどれかの頂点と隣接する

極大 2 部クリーク  $K$  が  $K = \emptyset$  でなければ, この 2 つの条件を満たす  $u_i$  は高々  $\Delta^2$  個しか存在しない ( $\Delta$  はグラフの最大次数). 各  $u_i$  について  $K[i]_{\leq i-1} = K_{\leq i-1}$  の判定には  $O(|X(K[i])|\Delta)$  時間かかる.  $|X(K[i])|$  は  $u_i$  に隣接する  $X(K)$  の頂点の数である. よって, 各  $u_i$  について  $|X(K[i])|$  の合計を取ると, それは  $|X(K)|\Delta$  を超えない. これにより, 判定にかかる時間の合計は  $O(\Delta^3)$  となる.

$K = \emptyset$  の場合,  $K$  の子供は高々  $n$  個である. 各添え字  $i$  に対して  $K[i]$  が子供であるかどうかの判定は,  $O(\Delta X(\{u_i\}))$  時間かかる. 全ての  $i$  について,  $|X(\{u_i\})|$  の合計を取ると, グラフの枝数と等しくなる. よって, 計算時間は  $O(|A|\Delta)$  となる.

子供を見つけるアルゴリズムを使えば, 列挙木を深さ優先探索できる. 深さ優先探索は, 各頂点からその子供に移動し, 全ての子供の探索が終了したら親に戻る, というアルゴリズムである. このアルゴリズムで必要となる作業は, 各頂点の子供を順に見つけることだけである. よって, 列挙木の根, つまり極大 2 部クリーク  $\emptyset$  から出発し, 現在訪れている頂点の子供を列挙し, 各子供に対してその子孫を再帰呼び出しで列挙することにより, 列挙木の全ての頂点を探索することができる. 列挙木の各頂点は極大 2 部クリークに対応するので, 列挙木の深さ優先探索により, 全ての極大 2 部クリークを列挙することができる.

探索中, 各反復では必ず 1 つ極大 2 部クリークが出力される. 先の議論により, 各反復の計算時間は, 列挙木の根以外では  $O(\Delta^3)$  である. よって, このアルゴリズムの極大 2 部クリーク 1 つあたりの計算量は,  $O(\Delta^3)$  となる. この他に, 根の反復でのみ  $O(|A|\Delta)$  時間がかかる. メモリの使用量は  $O(|V| + |A|)$  である.

さらに, [9] での計算実験では, ランダムに生成したグラフでの極大 2 部クリーク 1 つあたりの計算時間は, およそ  $O(\Delta^2)$  であることが示されている.

## 5. 本研究での改良と極大頻出集合の列挙

本研究では, [9] のアルゴリズムにいくつかの改良を加え, より高速に closed item set の高速な列挙を行う. 改良点は以下の 5 点である.

- (1) 頂点の添え字を, 次数の小さいものから昇順でつける
- (2) グラフの隣接行列の, 次数が大きな頂点に対応する行をメモリに持つ
- (3) 各反復で, 各  $X(K[i])$  を  $O(|X(K)|\Delta)$  時間で求める
- (4)  $K[i]_{\leq i-1}$  の頂点を添え字順で求め, 逐次  $K_{\leq i-1}$  と比較することにより,  $K[i]_{\leq i-1} = K_{\leq i-1}$  の判定を行う
- (5)  $|X(K)|$  が閾値より小さくなら, その子孫は列挙しない

(1) により, 極大 2 部クリーク  $K$  の親添え字よりも添え字が大きく, かつ  $X(K)$  に隣接するような頂点が少なくなり, 各反復で子供の判定をする回数が増える.

(2) により, 次数が大きな頂点に対しては, 他の頂点と隣接するかどうかの判定が定数時間でできるようになる. 次数の小

さな頂点に対しては, 隣接行列を用いずとも, 隣接の判定は短時間でできる.

(3) 上記の通り, 各  $u_i$  について  $|X(K[i])|$  の合計を取ったものは,  $|X(K)|\Delta$  を超えない. これを利用して, 以下のように全ての  $X(K[i])$  を  $O(|X(K)|\Delta)$  時間で作成する.

最初, 各頂点  $u_i$  について  $S_i = \emptyset$  とする. 次に,  $X(K)$  の各頂点  $v$  について,  $v$  が隣接する頂点  $u_j$  に対して  $S_j = S_j \cup \{v\}$  とする. この作業の終了後, 各  $u_i$  について,  $S_i = X(K[i])$  が成り立つ. この方法により, 判定が必要な  $X(K[i])$  全てを生成する時間は  $O(\Delta^2)$  となる.

(4) により,  $K[i]_{\leq i-1} = K_{\leq i-1}$  が成り立たない場合には, 早期に判定の手続きを終了することができ, 結果として  $K[i]$  の構築時間を短縮することができる.

(5) 極大 2 部クリーク  $K$  とその子供  $K'$  に対して,  $X(K') \subset X(K)$  が成り立つ. よって,  $K$  が頻出でなければ, その子孫も全て頻出でない. (5) の改良により, 頻出でない closed item set の極大元の列挙を省略することができる.

これらの改良に加え, 本研究では極大頻出集合を列挙する方法も開発した. 任意の極大頻出集合は, ある closed item set の極大元になっている. そこで, closed item set の中で極大頻出集合になっているもののみを出力することで, 極大頻出集合列挙アルゴリズムが得られる. closed item set の極大元が極大頻出集合であるかどうかの判定は以下のように行う.

ある closed item set の極大元  $K$  が極大頻出集合でないならば, ある closed item set の極大元  $K'$  が存在して,  $K \subset K'$  が成り立つ. このとき, ある  $i$  に対して,

$$K \subset K[i] \subset K', \text{ および}$$

$$X(K') \subset X(K[i]) \subset X(K)$$

が成り立つ. よって, 全ての  $K[i]$  が頻出集合でなければ, またそのときに限り,  $K$  は極大頻出集合である. この判定は, 全ての  $X(K[i])$  を生成することにより行えるので,  $O(\Delta^2)$  時間で行える.

## 6. 計算実験

本研究では, これらの工夫を加えたアルゴリズムを実装し, 計算機実験を行った. 使用機器は PentiumIII 500MHz と 256MB のメモリを搭載した PC, プログラム言語は C, OS は Linux である. 実験を行った問題は, KDD-cup2000 [6] で用いられた, 商品売り上げデータ (BMS-POS), Web ネットワークから抽出したデータ (BMS-WebView1, 2), 人工的に作られたデータ (IBM-artificial) の 4 つである. これらは, 広く世界的に使用されている問題である. 既存のアルゴリズムとの比較は, Z. Zheng, R. Kohavi, L. Mason [11] による比較実験の結果を用いた. この論文では以下の 4 つの頻出集合 / closed item set 列挙アルゴリズムの比較を行っている.

- R. Agrawal, H. Mannila, R. Srikant, H. Toivonen, A. I. Verkamo による apriori [2] (頻出集合)
- J. Han, J. Pei, Y. Yin による FP-growth [5] (頻出集合)

- J. Pei, J. Han, R. Mao による closet [7] (closed item set)
- M. J. Zaki, C. Hsiao による CHARM [10] (closed item set)

[11] の実験で用いられた機器は、CPU クロックが 550MHz の Duron と 1GB のメモリを搭載した PC であり、本研究で用いた機器より多少性能の良いマシンである。

本研究の計算実験では、上記 4 つの問題に対するこれら 4 つのアルゴリズムの計算時間と、本研究で提案する closed item set の列挙アルゴリズム (Ours)、および極大頻出集合列挙アルゴリズム (Ours(max)) の計算時間を比較した。それぞれの問題で、閾値を 0.1, 0.08, 0.06, 0.04, 0.02, 0.01 に設定した場合について、計算を行った。結果は、以下の表にまとめた。表中、「-」が書かれている欄は [11] に結果が記載されていない項目である。本研究で提案するアルゴリズム以外の 4 つのアルゴリズムの計算時間は [11] の論文に掲載されているグラフから定規を当てて読み出した数値であるので、1 割程度の誤差があることを注意しておく。また、IBM-Artificial は問題生成プログラムを用いて生成したものであり、使用する乱数の違いから、同じ問題は作成できない。そのため [11] で解かれた問題と本研究で解いた問題は異なる問題である。しかし、これらのグラフの頂点数・枝数は等しく、また同じ方法を使って作成されたことから、closed item set 数もほぼ同じであると考えられる。

また、本研究のアルゴリズムでは、閾値を用いずに全ての closed item set を列挙することができた。この結果も合わせて報告する。なお、表中の CIS は closed item set の略である。

**問題: BMS-Web-View1:** アイテム数約 500, トランザクション数約 6 万, トランザクションの平均の大きさ約 2.5

閾値 (%)	0.1	0.08	0.06	0.04	0.02	0.01
Apriori	1.1	3.6	113	-	-	-
FP-growth	1.2	1.8	51	-	-	-
Closet	33	74	-	-	-	-
Charm	2.2	2.7	7.9	133	422	-
Ours	1.0	1.3	3.9	10	16	38
Ours(max)	1.7	3.1	17	39	63	125
CIS 数	3974	9391	64762	155651	422692	1240700
頻出集合数	3992	10287	461522	-	-	-
極大頻出集合数	2067	4028	15179	12956	84833	129754

closed item set の総数: 1427216 個, 計算時間 46 秒

**問題: BMS-Web-View2:** アイテム数約 3300, トランザクション数約 8 万, トランザクションの平均の大きさ約 5

閾値 (%)	0.1	0.08	0.06	0.04	0.02	0.01
Apriori	13.1	15	29.6	58.2	444	-
Fp-growth	7.03	10	17.2	29.6	131	763
Closet	1500	2250	3890	6840	25800	-
Charm	5.82	6.66	7.63	13.8	27.2	76
Ours	2.9	3.4	4.0	5.0	9.2	16
Ours(max)	6.9	8.8	11	15	28	43
CIS 数	22976	37099	60352	116540	343818	754924
頻出集合数	24143	42762	84334	180386	1599210	9897303
極大頻出集合数	3901	5230	7841	16298	43837	118022

closed item set の総数: 1691051 個, 計算時間 38 秒

**問題: BMS-POS:** アイテム数約 1650, トランザクション数約 51 万, トランザクションの平均の大きさ約 6

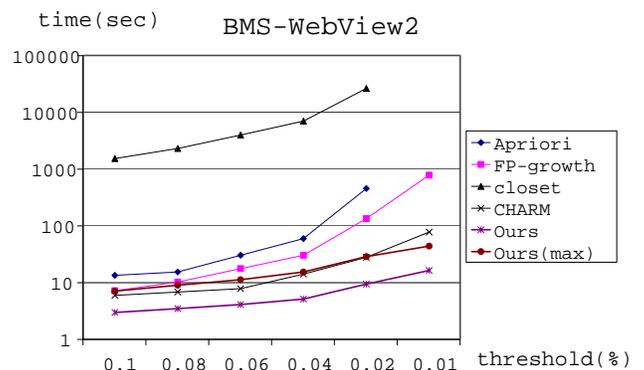
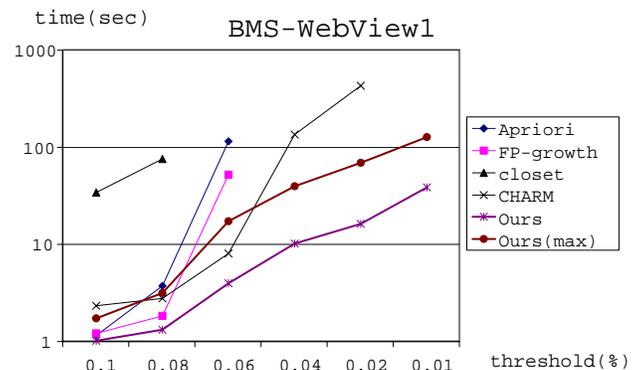
閾値 (%)	0.1	0.08	0.06	0.04	0.02	0.01
Apriori	251	341	541	1000	2371	10000
Fp-growth	196	293	398	671	1778	6494
Closet	-	-	-	-	-	-
Charm	100	117	158	215	541	3162
Ours	49	57	74	109	228	551
Ours(max)	261	343	486	837	2171	6147
CIS 数	121879	200030	378217	840544	1742055	21885050
頻出集合数	121956	200595	382663	984531	5301939	33399782
極大頻出集合数	30564	48015	86175	201306	891763	4280416

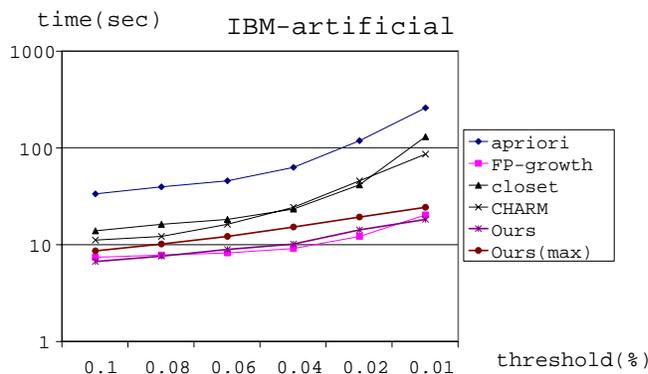
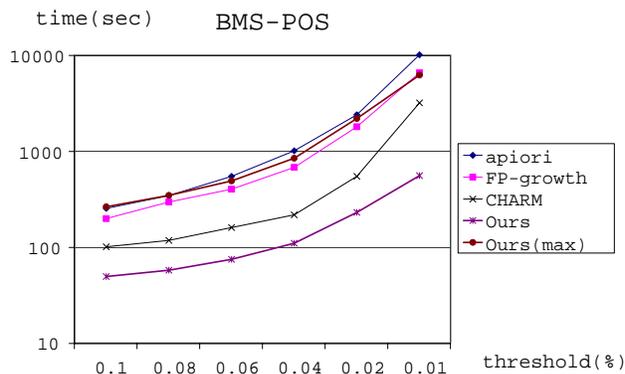
closed item set の総数: 1787673748 個, 計算時間 31259 秒

**問題: IBM-Artificial:** アイテム数 1000, トランザクション数 10 万, トランザクションの平均の大きさ 10

閾値 (%)	0.1	0.08	0.06	0.04	0.02	0.01
Apriori	33	39	45	62	117	256
Fp-growth	7.3	7.7	8.1	9.0	12	20
Closet	13	16	18	23	41	130
Charm	11	13	16	24	45	85
Ours	6.6	7.5	8.8	10	14	18
Ours(max)	8.5	10	12	15	19	24
CIS 数	13773	22943	38436	67536	131341	229029
頻出集合数	15010	28059	46646	84669	187679	335183
極大頻出集合数	7853	11311	16848	25937	50232	114114

closed item set の総数: 3160745 個, 計算時間 49 秒





ほとんど全ての問題において、本研究で提案した極大2部クリーク列挙アルゴリズム (Ours) が短時間で問題を解いた。特に closed item set を列挙する closet, CHARM より、どの問題でも高速であり、計算時間が増大するほどその差は開いている。極大頻出集合列挙アルゴリズム (Ours(極大)) は、極大頻出集合であるかどうかの判定に多大な時間がかかっているため、極大2部クリーク列挙アルゴリズムより時間がかかっている。しかし、いくつかの問題では、他のアルゴリズムよりも短時間で計算が終了している。

今回実験に用いたベンチマーク問題では、closed item set の数と頻出集合の数に2倍以内の差しかない場合が多い。このような問題では、closed item set であるかどうかの判定を行っている分、closed item set 列挙アルゴリズムのほうが、頻出集合列挙アルゴリズムより計算時間が余計にかかると思われる。しかし、本研究のアルゴリズムは、極大2部クリークの子供の生成方法を工夫するなどして1反復の計算時間を短縮したため、これらの問題に対しても頻出集合列挙アルゴリズムと大差ない時間で計算が終了している。BMS-WebView で閾値を0.04%以下に設定した問題など、closed item set 数が頻出集合数よりも5倍以上大きくなる問題では、明らかに極大2部クリーク列挙アルゴリズムが高速である。

## 7. まとめ

本稿では、極大2部クリークを列挙するアルゴリズムに改良を加え、closed item set および極大頻出集合を高速に列挙する方法を提案した。また、計算実験により、ベンチマーク問題で既存のアルゴリズムより高速であることを示した。

## 8. 謝 辞

この研究は、国立情報学研究所共同研究費の補助を受けた。

### 文 献

- [1] R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules in Large Databases," *In Proceedings of VLDB '94*, pp. 487-499, 1994.
- [2] R. Agrawal, H. Mannila, R. Srikant, H. Toivonen and A. I. Verkamo, "Fast Discovery of Association Rules," *In Advances in Knowledge Discovery and Data Mining*, MIT Press, pp. 307-328, 1996.
- [3] D. Avis and K. Fukuda, "Reverse Search for Enumeration," *Discrete Applied Mathematics*, Vol. 65, pp. 21-46, 1996.
- [4] E. Boros, V. Gurvich, L. Khachiyan, and K. Makino, "On the Complexity of Generating Maximal Frequent and Minimal Infrequent Sets," *STACS 2002*, pp. 133-141, 2002.
- [5] J. Han, J. Pei, Y. Yin, "Mining Frequent Patterns without Candidate Generation," *SIGMOD Conference 2000*, pp. 1-12, 2000
- [6] R. Kohavi, C. E. Brodley, B. Frasca, L. Mason and Z. Zheng, "KDD-Cup 2000 Organizers' Report: Peeling the Onion," *SIGKDD Explorations*, 2(2), pp. 86-98, 2000.
- [7] J. Pei, J. Han, R. Mao, "CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets," *ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery 2000*, pp. 21-30, 2000.
- [8] S. Tsukiyama, M. Ide, H. Ariyoshi and I. Shirakawa, "A New Algorithm for Generating All the Maximum Independent Sets," *SIAM Journal on Computing*, Vol. 6, pp. 505-517, 1977.
- [9] 宇野 毅明, "大規模グラフに対する高速クリーク列挙アルゴリズム," 電気通信学会コンピュータ研究会, 京都大学, 2003.
- [10] M. J. Zaki, C. Hsiao, "CHARM: An Efficient Algorithm for Closed Itemset Mining," *2nd SIAM International Conference on Data Mining (SDM'02)*, pp. 457-473, 2002.
- [11] Z. Zheng, R. Kohavi and L. Mason, "Real World Performance of Association Rule Algorithms," *KDD 2001*, pp. 401-406, 2000.