

2004 年度卒業論文

アイテム間の距離を考慮した
Sequential Pattern Mining の提案

提出日：2005 年 2 月 2 日

指導：山名 早人 助教授

早稲田大学理工学部情報学科

学籍番号：1G01P043-6

小松 俊介

概要

近年の技術の進歩によって記憶装置の大容量化が進み、膨大な量のデータを人間が直接扱うことが困難になってきている。そこで、膨大なデータの中から有用な情報を取り出す技術としてデータマイニング技術が注目されている。データマイニングの中で重要とされる問題が、頻出パターンと呼ばれる大量のデータの中から頻繁に出現するアイテムの組合せを抽出する（頻出パターン抽出：**Frequent Pattern Mining**）問題である。**Frequent Pattern Mining**はユーザが与えた最小サポート値よりも高い頻度で出現するアイテムの組合せを抽出するが、アイテム同士の順序は考えられていない。しかし、アイテム間の順序が重要となる事例が世の中には数多くある。そこで、考え出されたのが **Sequential Pattern Mining** である。**Sequential Pattern Mining** ではアイテム間の順序を考慮してマイニングを行う。アイテム間の順序を考慮することによってより現実 に即したマイニングが可能となった。しかし、**Sequential Pattern Mining** ではアイテム間の距離（時間）については考慮していないため、抽出されたアイテムセット中の任意の2つのアイテム間に、どれだけの距離があるのかを区別することが出来ない。しかし、距離を区別することができれば、意味の違う行動をそれぞれ抽出することができる。

本論文では **Sequential Pattern Mining** におけるアイテム間の順序に加えて、アイテム間の距離（時間）を考慮してマイニングを行う方法を提案する。提案手法ではアイテム間の順序だけでなく、距離（時間）を考慮してマイニングを行うことにより、直後に起こるトランザクションとしばらくして起こるトランザクションとの区別を可能にすることができた。

目次

第1章	はじめに.....	4
第2章	Sequential Pattern Mining とは.....	5
2.1	データマイニングとは.....	5
2.1.1	KDD プロセス [7].....	5
2.1.2	相関ルール抽出問題.....	6
2.1.3	Frequent Pattern Mining から Sequential Pattern Mining へ.....	7
2.2	Sequential Pattern Mining	7
2.2.1	シーケンス.....	7
2.2.2	問題定義.....	8
2.2.3	マイニングの例.....	8
第3章	関連研究.....	10
3.1	Sequential Pattern Mining の歴史.....	10
3.2	GSP [3].....	10
3.2.1	Apriori の理論の適用.....	11
3.2.2	GSP アルゴリズムの流れ.....	12
3.2.3	GSP の欠点.....	16
3.3	SPADE [5].....	17
3.3.1	ID-List	17
3.3.2	Lattice	18
3.3.3	SPADE アルゴリズムの流れ.....	22
3.4	SPAM [6].....	24
3.4.1	ビットマップを用いたデータ構造.....	24
3.4.2	ビットマップ表現におけるシーケンスの結合.....	25
3.5	PrefixSpan [6].....	27
3.5.1	Prefix projection と Prefix データベース.....	27
3.5.2	PrefixSpan アルゴリズムの流れ.....	28
3.5.3	PrefixSpan の最適化.....	30
3.6	各手法の比較.....	32
3.7	Sequential Pattern Mining の拡張 [11].....	33
第4章	提案手法.....	35
4.1	Sequential Pattern Mining における順序と距離.....	35
4.2	アイテム間の距離（時間）を考慮した Sequential Pattern Mining	35
4.3	提案手法のアルゴリズムの流れ.....	37
4.4	Stream Mining の観点から見た提案手法.....	39

第 5 章	性能評価	41
5.1	実験環境	41
5.2	データセット	41
5.3	実験	43
5.3.1	実験方法	43
5.3.2	実験結果	43
5.4	考察	50
第 6 章	おわりに	56

第1章 はじめに

近年の技術の進歩により、記憶装置の容量が飛躍的に増大した。それによって、様々な場面で多くの情報を蓄積することが可能になった。しかし、蓄積される情報量が増大すればするほど、人間が直接データを扱うことが困難になる。情報を大量に集めることができても有効に活用できなければ無意味となってしまう。そこで、大量のデータの中から有用な情報を取り出す技術としてデータマイニング技術が注目されている。

データマイニングにおける重要な技術の一つとして頻出パターン抽出が挙げられる。頻出パターン抽出 (**Frequent Pattern Mining**) とは、データベースの中からユーザが与えた最小サポート値以上の頻度で出現するパターン (アイテムの組合せ) を見つける技術である。**Frequent Pattern Mining** では速度の向上のためのアルゴリズムや最小サポートにこだわらないアルゴリズムが数多く提案された。

しかし、現実の世界では、顧客の購買行動、自然災害 (地震など)、**Web** ページのクリックの流れ、株取引などアイテムの出現順序が重要となってくる場面が多く存在するが、**Frequent Pattern Mining** では順序を考慮せずにマイニングを行っている。そこで、より現実に応じたマイニングを行うために、**Sequential Pattern Mining** が 1995 年に **Agrawal** と **Srikant** によって提案された[1]。**Sequential Pattern Mining** ではアイテム間の順序を考慮してマイニングを行うことにより、上で述べたような事例において **Frequent Pattern Mining** より有用な情報を得ることが可能になった。**Sequential Pattern Mining** のアルゴリズムは現在までにいくつか提案されているが、有名なものは **GSP**[2]、**PrefixSpan**[5]、**SPADE**[3]、**SPAM**[4]が挙げられる。

しかし、**Sequential Pattern Mining** ではアイテム間の順序は考慮するが、アイテム間の距離 (時間) は考慮されていない。そのため、抽出されたアイテムセット中の任意の2つのアイテム間に、どれだけの時間間隔があるのかを区別することが出来ない。例えば、商品 **A** を購入したあと、1日後に商品 **B** を購入する人と1年後に商品 **B** を購入する人がいるとする。この場合、従来の **Sequential Pattern Mining** ではこの2人の行動は同じものとして同じパターンで扱われてしまう。しかし、2人の行動の持つ意味はそれぞれ違っている。そこで、本論文では、従来の **Sequential Pattern Mining** において、アイテム間の距離が考慮されていない問題を解決することを目的として、従来の **Sequential Pattern Mining** では別々のトランザクションとして扱われていたものを、ユーザが定義した一定の距離 (時間) 内に起こったトランザクションを同時に起こったものとみなしてマイニングを行う手法を提案する。

本論文では以下の形をとる。第2章で **Sequential Pattern Mining** について述べる。第3章で関連研究として、既存の **Sequential Pattern Mining** アルゴリズムについて述べる。第4章でアイテム間の距離 (時間) を考慮した **Sequential Pattern Mining** の提案をする。第5章で提案手法の評価と考察をし、最後に第6章でまとめを述べる。

第2章 Sequential Pattern Mining とは

本章では、**Sequential Pattern Mining** とは何かについて説明する。具体的には、まず知識抽出におけるデータマイニングについて述べる。次に、そして最後に、**Sequential Pattern Mining** の問題定義について述べる。

2.1 データマイニングとは

2.1.1 KDD プロセス[7]

大量のデータから有用な情報を取り出すプロセスである、**KDD(Knowledge Discovery and DataMining)**プロセスは、図 2.1 のような5段階で構成されている。以下に、図 2.1 における 5 段階のプロセスについて説明する。

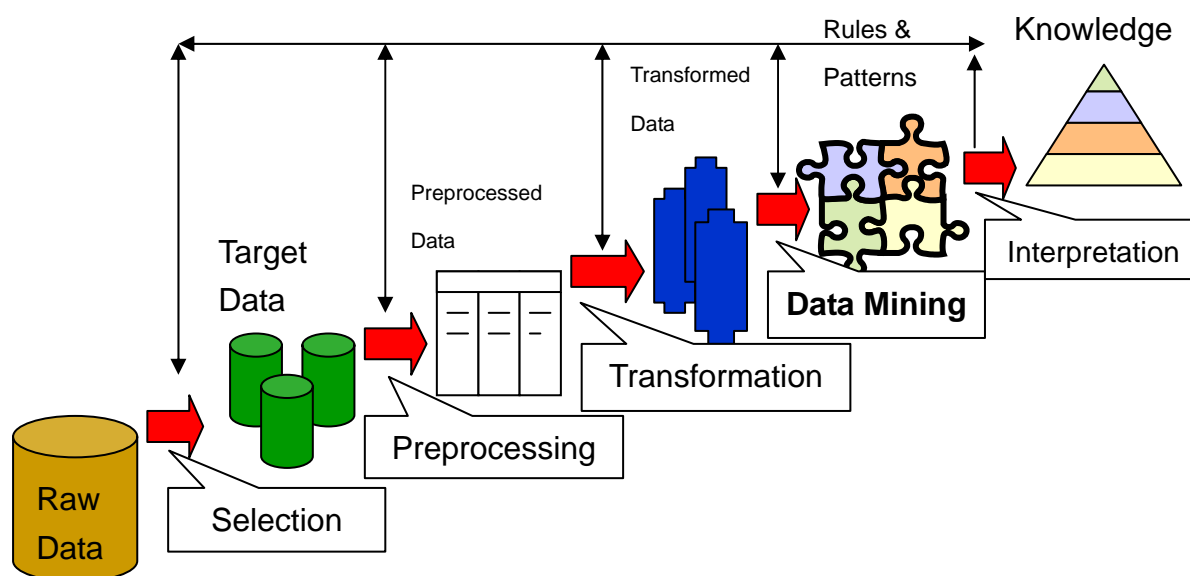


図 2.1 KDD プロセス

(1) Selection (データ選択)

大規模なデータベースの中から、どのデータに対して、知識発見をするかを選択することを **Selection** という。

(2) Preprocessing (前処理)

データの中の欠損値に対して、ダミーのデータを与えたり、欠損値が存在するレコードは消去するなど、マイニングの前にデータを完全な状態にしておくことを **Preprocessing** という。

(3) Transformation(データ変換)

マイニングを実行する前にデータベースフォーマットからマイニングしやすいフォーマット

トに変換することを **Transformation** という。

(4) **Mining** (データマイニング)

完全性の保たれているデータに対してルールやパターンを抽出することを **Mining** という。データマイニング手法には、相関ルール抽出、クラスタリング、決定木などがある。

(5) **Interpretation/Evaluation** (解釈・分析)

マイニングによって得られたルールやパターンを解釈、分析することを **Interpretation/Evaluation** という。

2.1.2 相関ルール抽出問題

KDD プロセスにおける **Mining** において、相関ルール抽出問題というのは重要な役割を担っている。相関ルール抽出問題とは何かを説明する前に、相関ルール抽出における定義、相関ルールの例を示す。

定義

$I=\{i_1, i_2, \dots, i_n\}$ をアイテム集合とする。データベース $D=\{t_1, t_2, \dots, t_n\}$ をトランザクションの集合とする。任意のトランザクション t はアイテムの集合で構成されている ($t_i \subseteq I (1 \leq i \leq n)$)。また、各トランザクションにはユニークな識別子 **TID(transaction id)** が付けられているとする。

相関ルールとは、 $X \subseteq I$ 、 $Y \subseteq I$ 、 $X \cap Y = \emptyset$ であるような任意のアイテム集合 X 、 Y を使って作られる $X \Rightarrow Y$ という表現のことである。

相関ルールは2つのパラメータ、確信度 $\text{conf}(X \Rightarrow Y)$ と、サポート $\text{sup}(X \Rightarrow Y)$ と呼ばれる2つのパラメータを持つ。確信度は、データベース D 中の X を含むトランザクションにおける、 Y を含むトランザクションの割合である。サポートはデータベース D における、 $X \cup Y$ を含む全トランザクションに対する割合である。

相関ルールの例

表 2.1 に、スーパーマーケットの購買データの例を示す。以下に、相関ルールについて、表 2.1 を元に説明していく。

表 2.1 スーパーマーケットの購買データ D

TID	アイテム
0001	A,C,D
0002	B,C,E
0003	A,B,C,E
0004	B,E

表 2.1 のデータベースにおいて、各行がトランザクションを表す。つまり、表 2.1 において、**TID(TransactionID)0001** 番の顧客は商品 **A,C,D** を購入し、**TID0002** 番の顧客は商品 **B,C,E**

を購入したことになる。 $X \Rightarrow Y$ という相関ルールの確信度が $c\%$ で、そのサポートが $s\%$ だとする。この場合、商品集合 X を購入したトランザクションのうち $c\%$ が商品の集合 Y も同時に購入していて、 $X \cup Y$ を購入していた顧客のトランザクション全体に対する割合は $s\%$ だったということである。

相関ルール抽出問題

相関ルール抽出問題では、データベース D 、確信度の最小値 min_conf 、そしてサポートの最小値 min_sup が与えられて、これらを満足する相関ルールを抽出する。相関ルール抽出は以下の2段階の処理に分けられる。

- ① ユーザが与えた min_sup を満たすアイテム集合（頻出アイテム集合という）を全て抽出する。
- ② 得られた頻出アイテム集合の中から min_conf を満たす相関ルールを得る。

②は①で求めた頻出アイテム集合を用いて相関ルールを得る処理であり、負荷が小さいので比較的早く算出することができる。しかし、①は大規模なデータベースに対し繰り返しスキャンを行い、それぞれのアイテム集合のサポートを計算するため、負荷が大きくなり算出に時間がかかる。つまり、相関ルール抽出アルゴリズムというのは、①の処理をいかに効率よく計算するかが主要な課題となっており、相関ルール抽出問題は頻出パターン抽出 (**Frequent Pattern Mining**) 問題に置き換えることができる。

2.1.3 Frequent Pattern Mining から Sequential Pattern Mining へ

Frequent Pattern Mining が対象とするデータにおいて、アイテム間の順序という概念は存在していない。例えば、表 2.1 の購入データにおいて、**TID0001** の顧客は **ACD** を購入しているが、**ACD** は同時に購入されている。つまり、**ACD** でも **CDA** でも構わないのである。同様に、抽出されるパターンもアイテム間の順序は考えられていない。

しかし、現実の世界ではアイテム間の順序が多くの場合で重要となってくる。そこで、**Frequent Pattern Mining** においてアイテム間の順序が考慮されていないという問題を解決するために、**Sequential Pattern Mining** が考案された。**Sequential Pattern Mining** は、顧客購買行動、医療、自然災害（地震など）、科学実験、株取引、電話をかけるパターン、Web ページ閲覧の流れ、**DNA** 配列、遺伝子構造などの分野に応用することができる。

2.2 Sequential Pattern Mining

2.2.1 シーケンス

Sequential Pattern Mining におけるパターンは、シーケンスである。ひとつのシーケンスは、“<”で開始し、“>”で終了する。シーケンスでは、“< >”の間にアイテムを時系列順に配置する。なお、同時期に発生したアイテムは、“()”でくくることによって表現される。シーケンスの例を以下に示す。

シーケンス：<(ef)(ab)(df)c b>

□を1つのトランザクションとして考える。()の中は同時に行われているので順序は問わない。しかし、□同士は時系列で並んでいるので入れ替えはできない。

サブシーケンス：<a (bc) d c>は<a (abc) (ac) d (cf)>のサブシーケンスである。

シーケンス **S** が **S'** のサブシーケンスとは、**S** のアイテムが **S'** 内において **S** と同じ順番で存在しているということを意味する。上の例では、<a (bc) d c>は<a (abc) (ac) d (cf)>内において、斜体字で示すように、a、(bc)、d、c の順番で存在しているので<a (bc) d c>は<a (abc) (ac) d (cf)>のサブシーケンスである。

2.2.2 問題定義

表2.2のようなデータベース **D** が与えられている。それぞれのトランザクションには顧客ID、トランザクションタイム、アイテム集合が含まれている。

シーケンスはアイテム集合の整列されたリストである。アイテム集合は空ではなく通常、完全に連続したものであると考える。アイテム集合 i は $(i_1 i_2 \dots i_m)$ と表される。ここで i_j ($1 \leq j \leq m$) は1つのアイテムである。さらにシーケンス s は $\langle s_1 s_2 \dots s_n \rangle$ と表される。ここで、 s_k ($1 \leq k \leq n$) は1つのアイテム集合である。

表2.2 元のデータベース

Transaction Time	Customer ID	Items Bought
June 10 '93	2	10,20
June 12 '93	5	90
June 15 '93	2	30
June 20 '93	2	40,60,70
June 25 '93	4	30
June 25 '93	3	30,50,70
June 25 '93	1	30
June 30 '93	1	90
June 30 '93	4	40,70
July 25 '93	4	90

2.2.3 マイニングの例

表2.3の顧客の購買データベースの例を用いて、**Sequential Pattern Mining** について説明する。この例において、**min_sup=2** とする。

2つのシーケンス<(30)(90)>と<(30)(4070)(90)>は、アイテム30を買い、その後にアイテム90を買っている。つまり<(30)(90)>は最小サポート値を満たしている。

同様に<(10 20)(30)(40 60 70)>と<(30)(40 70)(90)>は、アイテム30を買った後、アイ

テム40と70を同時に購入しているので、 $\langle (30)(40\ 70) \rangle$ は最小サポート値を満たしている。マイニングの結果を表2.4に示す。

表2.3 Customer-Sequence バージョンのデータベース

Customer ID	Customer Sequence
1	$\langle (30)(90) \rangle$
2	$\langle (10\ 20)(30)(40\ 60\ 70) \rangle$
3	$\langle (30\ 50\ 70) \rangle$
4	$\langle (30)(40\ 70)(90) \rangle$
5	$\langle (90) \rangle$

表2.4 マイニングの結果

Sequential Patterns with support >40%
$\langle (30)(90) \rangle$
$\langle (30)(40\ 70) \rangle$

第3章 関連研究

本章では、既存の **Sequential Pattern Mining** のアルゴリズムを示す。まず、**Sequential Pattern Mining** の歴史について説明する。次に、**Sequential Pattern Mining** の代表的なアルゴリズムである **GSP**[3]、**SPADE**[4]、**SPAM**[5]、**PrefixSpan**[6]について説明する。最後に4つのアルゴリズムをそれぞれ比較する。

3.1 **Sequential Pattern Mining** の歴史

Sequential Pattern Mining はアイテム間の順序を考慮したデータマイニングで、1995年に **Agrawal** と **Srikant** によって初めて提案された[2]。同時に、**Sequential Pattern Mining** のアルゴリズムとして **AprioriAll** という **Apriori** (**Apriori** については次節で説明する) 的なアルゴリズムが提案された。さらに 96 年に **Apriori** ベースのアルゴリズムである **GSP** (**Generalized Sequential Pattern mining**) [3]が **Agrawal** と **Srikant** によって提案された。**GSP** は **Apriori** の「あるシーケンス **S** のサブシーケンス **S'** が非頻出ならばそのシーケンスは非頻出である」という考えのもとに候補シーケンスを生成し、頻出シーケンスを抽出する手法である。しかし、**GSP** では候補シーケンスが膨大な数になってしまうという欠点がある。

GSP の候補シーケンスが膨大な数になるという欠点を解決するために、2000年に **Zaki** が **SPADE**(**Sequential Pattern Discovery using Equivalence classes**)[4]を提案した。**SPADE** は **Lattice** という概念を用い候補シーケンスをグループに分割する。また、**ID-List** というデータ構造を用いることによってサポート値をカウントする際のコストを削減することができる。

さらに、**SPADE** を改良したアルゴリズムとして、2001年に **Ayres** らが、**SPAM**(**Sequential Pattern Mining**)[5]を提案した。**SPAM** では、**ID-List** をビットマップを用いて表現することによりさらなる高速化を実現した。

一方で、候補シーケンスを生成せずにデータベースを射影することによって頻出シーケンスを抽出する手法が 2001年に **Pei** らによって提案された **PrefixSpan**[6]である。[8]によれば、**PrefixSpan** が 2004 年時点では最も高速なアルゴリズムである。

3.2 **GSP**[3]

GSP(**Generalized Sequential Pattern mining**)は 1996 年に **Agrawal**, **Srkant** によって提案されたアルゴリズムである。**GSP** は **Sequential Pattern Mining** において最初に提案された有名なアルゴリズムである。**GSP** は **Frequent Pattern Mining** のアルゴリズムである

Apriori[1]の考え方をベースにしている。本節では、まず **GSP** のベースとなっている概念である **Apriori** について説明する。次に、**GSP** における重要な動作である候補シーケンス生成について説明した後、**GSP** のアルゴリズムの流れを説明する。そして、最後に **GSP** の欠点を説明する。

3.2.1 Apriori の理論の適用

Apriori は 1994 年に Agrawal らによって提案された **Frequent Pattern Mining** のアルゴリズムである。Apriori は「長さ k の頻出でないパターンを含む長さ $k+1$ のパターンは頻出でない」という理論の元で、ボトムアップに頻出パターンを抽出するアルゴリズムである。

Sequential Pattern Mining においても Apriori と同様の理論が適用でき、「あるシーケンス S が頻出シーケンスならば S のスーパーシーケンスはすべて頻出ではない」。例えば、シーケンス $\langle hb \rangle$ が頻出でなければ、 $\langle hb \rangle$ のスーパーシーケンス（内部に $\langle hb \rangle$ を含むシーケンス）である、 $\langle hab \rangle$ と $\langle (ah)b \rangle$ も頻出でない。この理論を元にして、**GSP** はマイニングを行っている。

3.2.2 候補シーケンス生成

k 個のアイテムを持つ（長さが k である）シーケンスを **k-sequence** と呼ぶ。 L_k をすべての頻出 **k-sequence** の集合とし、 C_k を候補 **k-sequence** の集合とする。

候補シーケンス生成とは、与えられた頻出 **(k-1)-sequence** の集合 L_{k-1} が与えられたとして、すべての候補 **k-sequence** の集合のスーパーセットを生成することを指す。そこで、まず連続サブシーケンスの定義を述べる。

シーケンス $s = \langle s_1 s_2 \dots s_n \rangle$ と、サブシーケンス c が与えられているとする。以下のいずれかの条件が成り立つとき、 c は s の連続サブシーケンスである。

- 1、 s_1 または s_n からアイテムを削除することで c が s から導出される。
- 2、少なくとも 2 つのアイテムを持つ要素 s_i からアイテムを削除することで、 c から s が導出される。
- 3、 c が c' の連続サブシーケンスであり、 c' が s の連続サブシーケンスである。

例えば、シーケンス $s = \langle (1,2) (3,4) (5) (6) \rangle$ について考える。ここで、 $\langle (2) (3,4) (5) \rangle$ $\langle (1,2) (3) (5) (6) \rangle$ $\langle (3) (5) \rangle$ は s の連続サブシーケンスである。
しかし、 $\langle (1,2) (3,4)(6) \rangle$ $\langle (1) (5) (6) \rangle$ は s の連続サブシーケンスではない。

候補生成は **Join Phase** と **Prune Phase** の 2 つのフェーズに分かれている。

(1) Join Phase

L_{k-1} に含まれるシーケンス s_a と s_b を結合することによって候補シーケンスを生成する。ここで、 $s_a = \langle sa_1, sa_2, \dots, sa_m \rangle, s_b = \langle sb_1, sb_2, \dots, sb_n \rangle$ とする。 $sa_2 = sb_1, sa_3 = sb_2, \dots, sa_m = sb_{n-1}$ が成立する場合、候補シーケンスは s_a と s_b を結合し、 $\langle sa_1, sa_2, \dots, sa_m, sb_n \rangle$ となる。追加されるアイテムは、それが s_b において分離した要素

(最後のトランザクションが (ab) のときに b だけ取り出すと (b) になる。これを分離した要素という) は、分離した要素にする。それ以外のときは s_a の最後の要素の一部とする。 L_1 と L_1 の結合するときは、 s_b のアイテムをアイテム集合の一部として、また分離した要素としてとの両方で追加する必要がある。なぜなら $\langle (x) (y) \rangle$ と $\langle (x,y) \rangle$ は最初のアイテムを削除することで同一のシーケンス $\langle (y) \rangle$ を生成するためである。。

(2) Prune Phase

サポート値が最小サポートより小さい連続 $(k-1)$ -subsequence を持つ候補シーケンスを削除する。

以下で長さ 3 の頻出シーケンスが与えられたときに、長さ 4 の候補シーケンスを生成する例を説明する。

表 3.1 候補生成の例

Frequent 3-Sequences	Candidate	4-Sequence
	After join	After pruning
$\langle (1,2) (3) \rangle$	$\langle (1,2) (3,4) \rangle$	$\langle (1,2) (3,4) \rangle$
$\langle (1,2) (4) \rangle$	$\langle (1,2)(3)(5) \rangle$	
$\langle (1) (3,4) \rangle$		
$\langle (1,3) (5) \rangle$		
$\langle (2) (3,4) \rangle$		
$\langle (2) (3) (5) \rangle$		

Join Phase では、 $\langle (1,2) (3) \rangle$ と $\langle (2) (3,4) \rangle$ の結合によって $\langle (1,2) (3,4) \rangle$ が生成される。同様に $\langle (1,2) (3) \rangle$ と $\langle (2) (3) (5) \rangle$ の結合によって $\langle (1,2) (3) (5) \rangle$ が生成される。

次に、Prune Phase では $\langle (1,2) (3) (5) \rangle$ がサブシーケンスの一つである $\langle (1) (3) (5) \rangle$ が頻出でないので、削除される。よって長さ 4 の候補シーケンスは $\langle (1,2) (3,4) \rangle$ となる。

3.2.2 GSP アルゴリズムの流れ

GSP アルゴリズムの流れを以下に示す。

- 1、長さ k の頻出シーケンスから長さ $k+1$ の候補シーケンスを生成する
- 2、候補シーケンスのサポート値をカウントし、最小サポート値以上のシーケンスのみが長さ $k+1$ の頻出シーケンスとして抽出される。
- 3、候補シーケンスを生成できなくなる、または頻出シーケンスが抽出できなくなるまで 1, 2 の動作を繰り返す

表 3.2 のシーケンスデータベースが与えられたとき、上の各ステップでどのような動きになるのか説明する。なお、**min_sup=2** とする。

表 3.2 シーケンスデータベース

SID	Sequence
10	< (bd) c b (ac) >
20	< (bf) (ce) b (fg) >
30	< (ah) (bf) a b f >
40	< (be) (ce) d >
50	< a (bd) b c b (ade) >

表 3.3 長さ 1 の頻出シーケンス

Cand	Sup
< a >	3
< b >	5
< c >	4
< d >	3
< e >	3
< f >	2
< g >	1
< h >	1

表 3.4 長さ 2 の候補シーケンス

	<a>		<c>	<d>	<e>	<f>
<a>	<aa>	<ab>	<ac>	<ad>	<ae>	<af>
	<ba>	<bb>	<bc>	<bd>	<be>	<bf>
<c>	<ca>	<cb>	<cc>	<cd>	<ce>	<cf>
<d>	<da>	<db>	<dc>	<dd>	<de>	<df>
<e>	<ea>	<eb>	<ec>	<ed>	<ee>	<ef>
<f>	<fa>	<fb>	<fc>	<fd>	<fe>	<ff>
	<a>		<c>	<d>	<e>	<f>
<a>	<(aa)>	<(ab)>	<(ac)>	<(ad)>	<(ae)>	<(af)>
		<(bb)>	<(bc)>	<(bd)>	<(be)>	<(bf)>
<c>			<(cc)>	<(cd)>	<(ce)>	<(cf)>
<d>				<(dd)>	<(de)>	<(df)>
<e>					<(ee)>	<(ef)>
<f>						<(ff)>

表 3.5 長さ 2 の頻出シーケンス

	<a>		<c>	<d>	<e>	<f>
<a>	<aa>:2	<ab>:2	<ac>:2	<ad>:1	<ae>:1	<af>:1
	<ba>:2	<bb>:4	<bc>:4	<bd>:1	<be>:3	<bf>:2
<c>	<ca>:2	<cb>:2	<cc>:1	<cd>:2	<ce>:1	<cf>:1
<d>	<da>:2	<db>:2	<dc>:2	<dd>:1	<de>:1	<df>:0
<e>	<ea>:0	<eb>:1	<ec>:1	<ed>:1	<ee>:1	<ef>:1
<f>	<fa>:1	<fb>:2	<fc>:1	<fd>:0	<fe>:1	<ff>:2
	<a>		<c>	<d>	<e>	<f>
<a>	<(aa)>:0	<(ab)>:0	<(ac)>:1	<(ad)>:1	<(ae)>:1	<(af)>:0
		<(bb)>:0	<(bc)>:0	<(bd)>:2	<(be)>:1	<(bf)>:2
<c>			<(cc)>:0	<(cd)>:0	<(ce)>:2	<(cf)>:0
<d>				<(dd)>:0	<(de)>:1	<(df)>:0
<e>					<(ee)>:0	<(ef)>:0
<f>						<(ff)>:0

表 3.6 長さ 3 の候補シーケンス

candidate 3-sequence						
<aaa>	<aca>	<bbc>	<bfb>	<cda>	<dbc>	<ffb>
<aab>	<acb>	<bbe>	<bff>	<cdb>	<dca>	<fff>
<aac>	<baa>	<bbf>	<caa>	<daa>	<dcb>	<(bd)a>
<aba>	<bab>	<b(bf)>	<cab>	<dab>	<fbb>	<(bd)b>
<abb>	<bac>	<bca>	<cba>	<dac>	<fbf>	<(bd)c>
<abc>	<bba>	<bc b>	<cbb>	<dba>	<f(bd)>	<(bf)b>
<abe>	<bbb>	<b(ce)>	<c(bd)>	<dbb>	<f(bf)>	<(bf)f>

表 3.7 長さ 3 の頻出シーケンス

frequent 3-sequence						
<aaa>:0	<aca>:1	<bbc>:2	<bfb>:0	<cda>:0	<dbc>:2	<ffb>:0
<aab>:1	<acb>:1	<bbe>:1	<bff>:0	<cdb>:0	<dca>:2	<fff>:0
<aac>:0	<baa>:0	<bbf>:2	<caa>:0	<daa>:0	<dcb>:2	<(bd)a>:2
<aba>:2	<bab>:1	<b(bf)>:0	<cab>:0	<dab>:0	<fbb>:0	<(bd)b>:2
<abb>:2	<bac>:0	<bca>:2	<cba>:2	<dac>:0	<fbf>:2	<(bd)c>:2
<abc>:1	<bba>:2	<bc b>:2	<cbb>:0	<dba>:2	<f(bd)>:0	<(bf)b>:2
<abe>:1	<bbb>:1	<b(ce)>:2	<c(bd)>:0	<dbb>:1	<f(bf)>:0	<(bf)f>:2

表 3.8 長さ 4 の候補シーケンス

candidate 4-sequence		
<abba>	<dbca>	<(bd)bc>
<bbca>	<dcba>	<(bd)ca>
<bcba>	<(bd)ba>	<(bd)cb>

表 3.9 長さ 4 の頻出シーケンス

frequent 4-sequence		
<abba>:1	<dbca>:1	<(bd)bc>:2
<bbca>:1	<dcba>:2	<(bd)ca>:2
<bcba>:2	<(bd)ba>:2	<(bd)cb>:2

表 3.10 長さ 5 の候補シーケンス

candidate 5-sequence
<(bd)cba>

表 3.1 1 長さ 5 の頻出シーケンス

frequent 5-sequence
<(bd)cba>:2

まず表 3.2 のデータベースをスキャンし、長さ 1 の頻出シーケンスを抽出する。抽出した結果が表 3.3 である。次に抽出した長さ 1 の頻出シーケンスを元に、表 3.4 のような長さ 2 の候補シーケンスを生成する。そして、長さ 2 の候補シーケンスのサポート値をカウントし、長さ 2 の頻出シーケンスを抽出する（表 3.5）。同様に長さ 2 の頻出シーケンスを元に、長さ 3 の候補シーケンスを生成（表 3.6）し、頻出シーケンスを抽出する（表 3.7）。続いて長さ 3 の頻出シーケンスを元に、長さ 4 の候補シーケンスを生成し（表 3.8）、長さ 4 の頻出シーケンスを抽出する（表 3.9）。引き続き長さ 5 の候補シーケンスを生成し（表 3.10）、長さ 5 の頻出シーケンスが抽出される（表 3.11）。長さ 5 の頻出シーケンスから長さ 6 の候補シーケンスは生成できないので、ここでアルゴリズムは終了する。

3.2.3 GSP の欠点

GSP の欠点は以下の 2 つである。

・生成される候補シーケンスが莫大な数になる

GSP における候補シーケンスはシーケンスの中で起こり得るすべてのアイテムの組み合わせを含んでいる。そのため、候補シーケンスが莫大な数になってしまう可能性がある。例えば、長さ 1 の頻出アイテム(<a₁>, <a₂>...<a₁₀₀₀>)が 1000 個あるとき、長さ 2 の候補シーケンス(<a₁a₁>, <a₁a₂>...<a₁a₁₀₀₀>, <a₂a₁>...<a₁₀₀₀a₁₀₀₀>と<(a₁a₂)><(a₁a₃)>...<(a₉₉₉a₁₀₀₀)>)は $1000 \times 1000 + 1000 \times 999 \div 2 = 1499500$ 個生成されることになる。

また、抽出するシーケンスの長さが長くなればなるほど、必要となる候補シーケンスが指数的に増大する。例えば、min_sup=1（すべてのパターンが頻出のとき）が与えられたとき、長さ 100 のシーケンスをマイニングした場合、長さ 1 の候補シーケンスは 100 個、長さ 2 の候補シーケンスは $100 \times 100 + 100 \times 99 \div 2 = 14950$ 個、同様に長さ 3 の候補シーケンスは 161700 個.....、合計では $2^{100} - 1$ 個となり、記憶容量が多く必要になってしまう。

・データベースのスキャン回数が多くなる

GSP ではそれぞれの長さごとに候補シーケンスを生成するためにデータベースをスキャンしなければならない。そのため、データベースのスキャン回数が多くなってしまう。例えば、<(abc) (abc) (abc) (abc) (abc)>というパターンを見つけるためには 15 回もデータベースをスキャンしなければならないので、記憶容量が多く必要になってしまう。

3.3 SPADE[5]

SPADE(**S**equential **P**attern **D**iscovery using **E**quivalence classes)[5]は2000年にZakiに提案されたアルゴリズムである。SPADEは *lattice* (格子) という概念を用いて、候補シーケンスをアイテムごとにグループに分割し、それぞれのグループがメインメモリに完全に格納されることによって高速化を図っている。また、*ID-List* というデータフォーマットを用いることによってサポート値をカウントするコストを削減している。本節ではまず、*ID-List* と *Lattice* という概念について説明した後、具体的なSPADEアルゴリズムの流れを説明する。

3.3.1 *ID-List*

ID-List とは SPADE で用いられているデータフォーマットである。*ID-List* はシーケンスごとに分けられており、それぞれが **SID** と **EID** という2つの要素を持っている。**SID** はシーケンシャル **ID** であり、**EID** はトランザクションが行われた時間を格納するイベント **ID** である。表3.12のデータベースをもとに作成した長さ1の頻出アイテムの *ID-List* が表3.13である。例えば、アイテム **D** に関するトランザクションが行われたのは **SID** 1 で時間10と25と **SID** 4 で時間10であるということを表している。

表 3.1 2 元のデータベース

Sequential ID	Time	Items
1	10	C D
1	15	A B C
1	20	A B F
1	25	A C D F
2	15	A B F
2	20	E
3	10	A B F
4	10	D G H
4	20	B F
4	25	A G H

表 3.1 3 表 3.1 を元に作成した長さ 1 の頻出シーケンスの ID-List

<A>				<D>		<F>	
SID	EID	SID	EID	SID	EID	SID	EID
1	15	1	15	1	10	1	20
1	20	1	20	1	25	1	25
1	25	2	15	4	10	2	15
2	15	3	10			3	10
3	10	4	20			4	20
4	25						

3.3.2 Lattice

SPADE では **Lattice** という概念に基づいてシーケンスを列挙している。**Lattice** とは先ほど述べたように格子という意味で、候補シーケンスをアイテムごとにグループに分割し、それぞれのグループをメインメモリに完全に格納することにより高速化を実現している。表 3.1 4 のデータベースを例に **Lattice** の説明を行う。

表 3.1 4 シーケンスデータベース(min_sup=2)

Sequence ID	Sequence
1	<(C D) (A B C) (A B F) (A C D F)>
2	<(A B F) E>
3	<(A B F)>
4	<(D G H) (B F) (A G H)>

ここで、最小サポートを満たす長さ 1 のアイテム集合を $F_1 = \{A, B, D, F\}$ とする。この 4 つのアイテムを元に最小サポートに関係なく、長さが 2、3、4…のシーケンスを作り、それを **Lattice** で表現したものが図 3.1 である。どのようにして長さ 2、3、4…のシーケンスを生成するかは後ほど説明する。図 3.1 の格子の底は“{}”になっているが、格子は無限なので上限は存在しない。

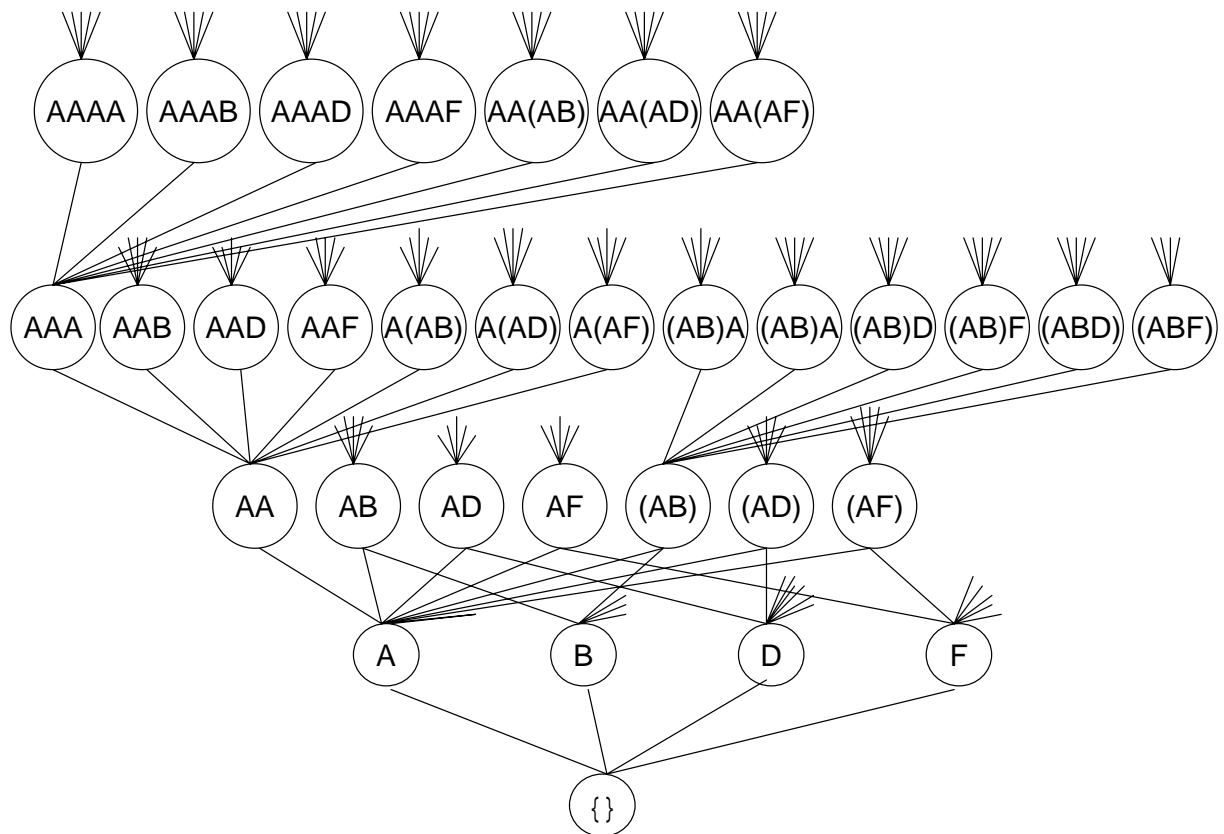


図 3.1 長さ 1 の頻出アイテムをもとにシーケンスを格子状に表現 ([5]より引用)

また、表 3.1 4 の例において最小サポート考慮して頻出シーケンスのみで作成した **Lattice** が図 3.2 である。この図を元に、**Lattice** を用いた候補シーケンスのグループ分けをした探索について説明する。

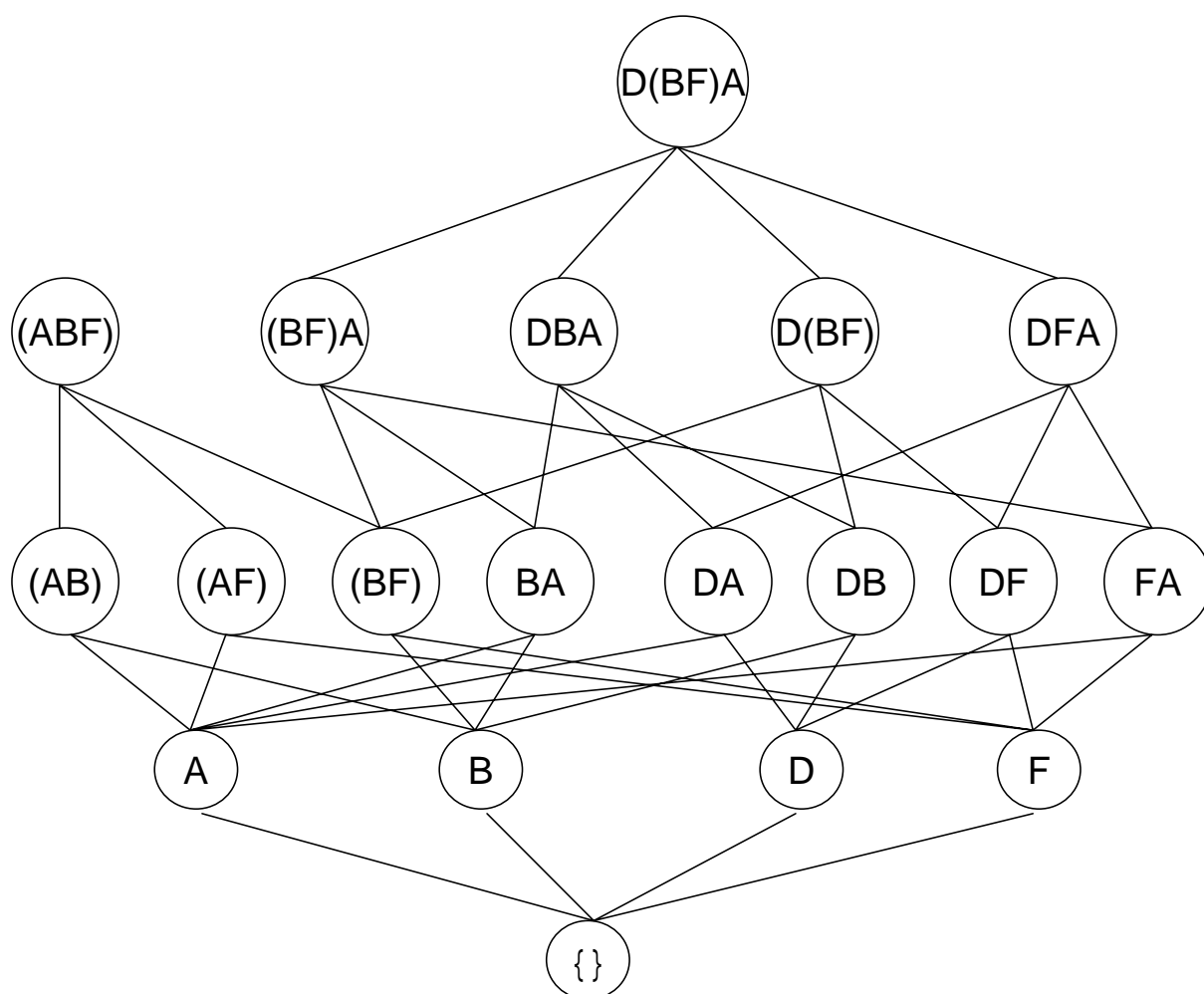


図 3.2 極大の頻出シーケンス<(ABF)>と<D (BF) A>によって作成された格子([5]より引用)

SPADE では深さ優先探索でマイニングを行う。図 3.3 では探索の順序を 4 つに分割している。<A>から始まる頻出シーケンス (**class[A]**) の抽出、から始まる頻出シーケンス (**class[B]**) の抽出、<C>から始まる頻出シーケンス (**class[C]**) の抽出、<D>から始まる頻出シーケンス (**class[D]**) の抽出の順でマイニングを行う。同様に、図 3.4 は<D>から始まるシーケンスの抽出の順序をさらに<DA>で始まるもの (**class[DA]**)、<DB>で始まるもの (**class[DB]**)、<DF>で始まるもの (**class[DF]**) の 3 つに分割している。そして、図 3.5 では、<D>から始まる候補シーケンスの流れについてを表している。<D>とその他の長さ 1 の頻出アイテム<A>、、<F>とを結合することによって長さ 2 の候補シーケンス<DA>、<DB>、<DF>が得られる。同様に探索していくと、<DBA>、<D(BF)>、<D(BF)A>、<DFA>の順に抽出される。

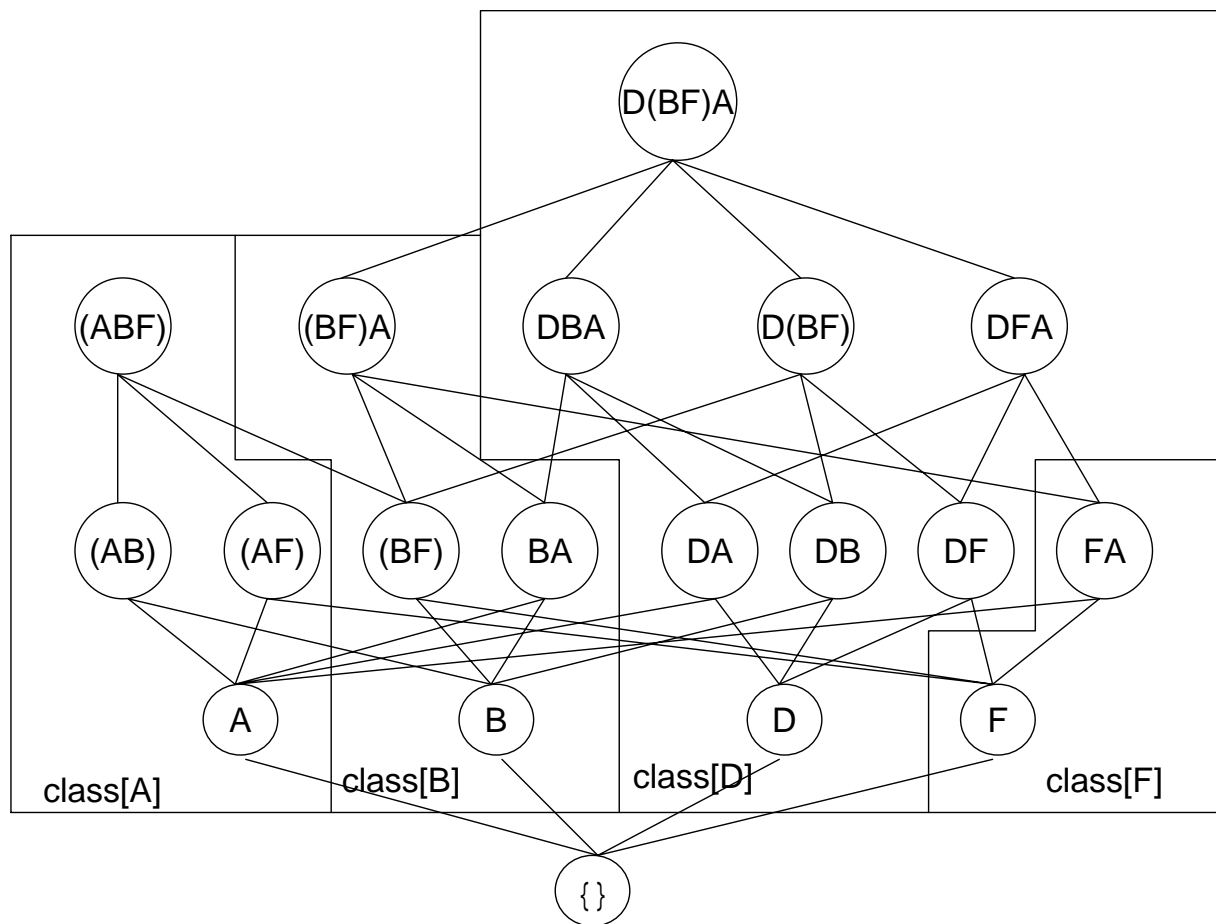


図 3.3 Lattice によるグループ分け

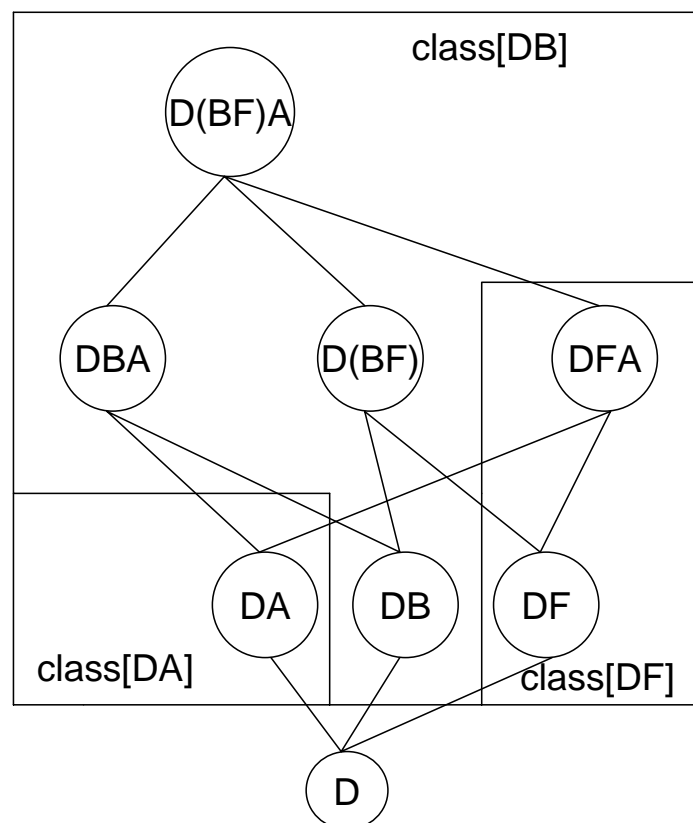


図 3.4 クラス $[D]$ におけるグループ分け

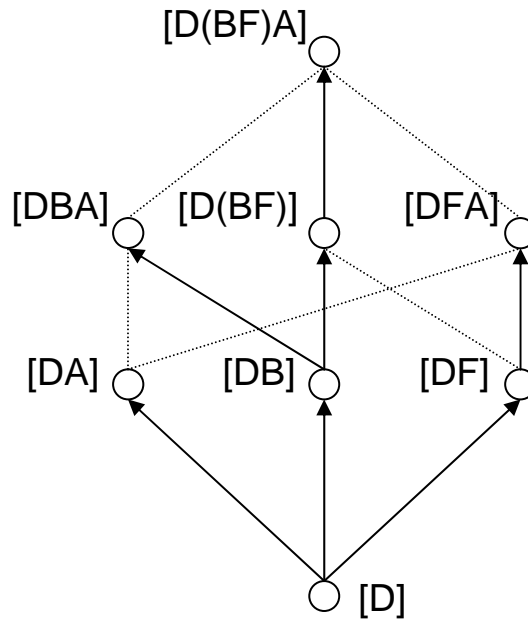


図 3.5 <D>から始まる頻出シーケンスの抽出における探索の流れ

3.3.3 SPADE アルゴリズムの流れ

SPADE では表 3.1 3 のような ID-List のデータベースを用いる。SPADE の主要な流れは以下のとおりである。

- 1、長さ 1 の頻出アイテムを抽出する
- 2、長さ 2 の頻出シーケンスを抽出する
- 3、深さ優先探索で長さ k ($k \geq 2$) の頻出シーケンス同士を結合し、長さ $k+1$ の頻出シーケンスを抽出する

ID-List の結合

図 3.6 のシーケンス<PA>と<PF>の結合の例について考える。<PA>と<PF>との結合からは、<PAF>、<PFA>、<P(AF)>の 3 つのシーケンスを生成することができる。まずは、<PAF>を生成する場合の SID が 1 のトランザクションについて見てみると、PA では時間 (EID) 20, 30, 40 に起こっていて、PF は時間(EID) 70, 80 に起こっている。ここで、<PAF>とは P、A、F の順でトランザクションが起こっているので、PF より前の時間(EID)に、PA が起こっていれば、結合が可能となる。よって、<PAF>の結合は、図 3.6 のようになる。また、<P (AF)>を生成する場合は、<PA>と<PF>が同じ時間(EID)のものを探さことになる。よって、<PA>、<PF>に共通に存在している SID が 8 で時間 (EID)が 30, 50, 80 のトランザクションが<P (AF)>の ID-List となる。

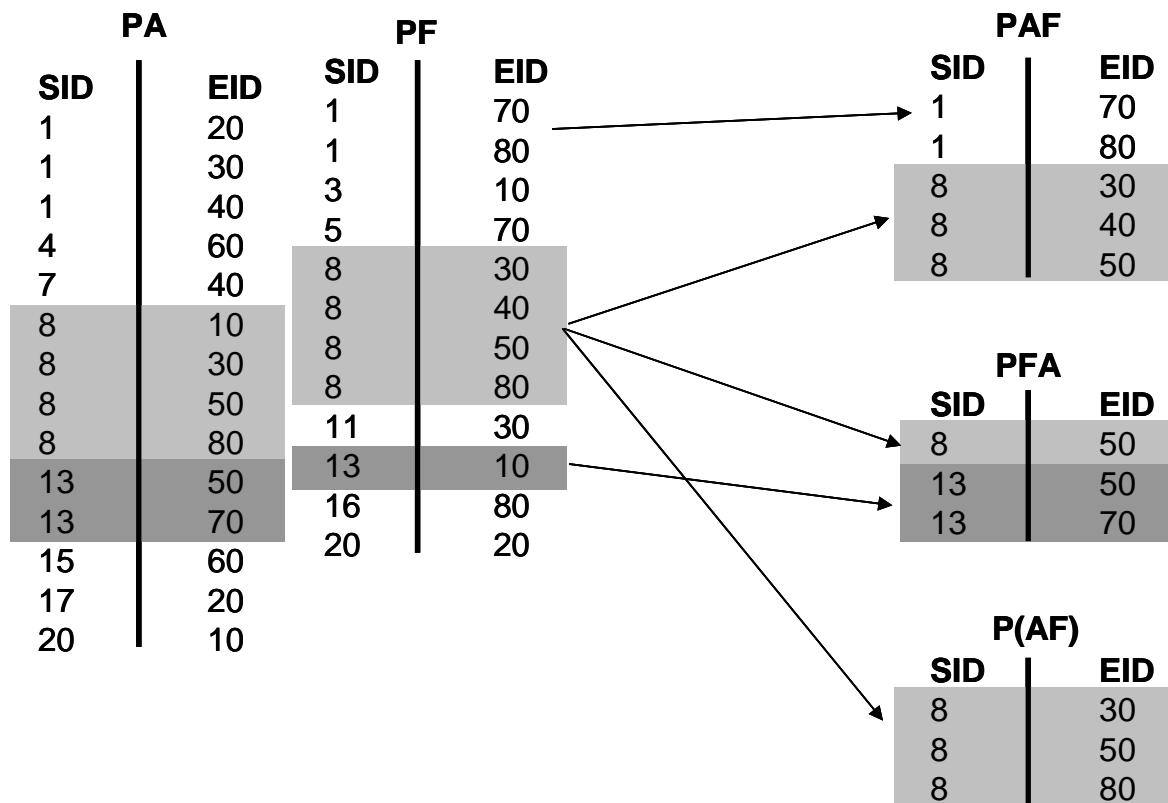


図 3.6 ID-List の結合

表 3.5 のマイニングにおける、ID-List の結合を図 3.7 に示す。

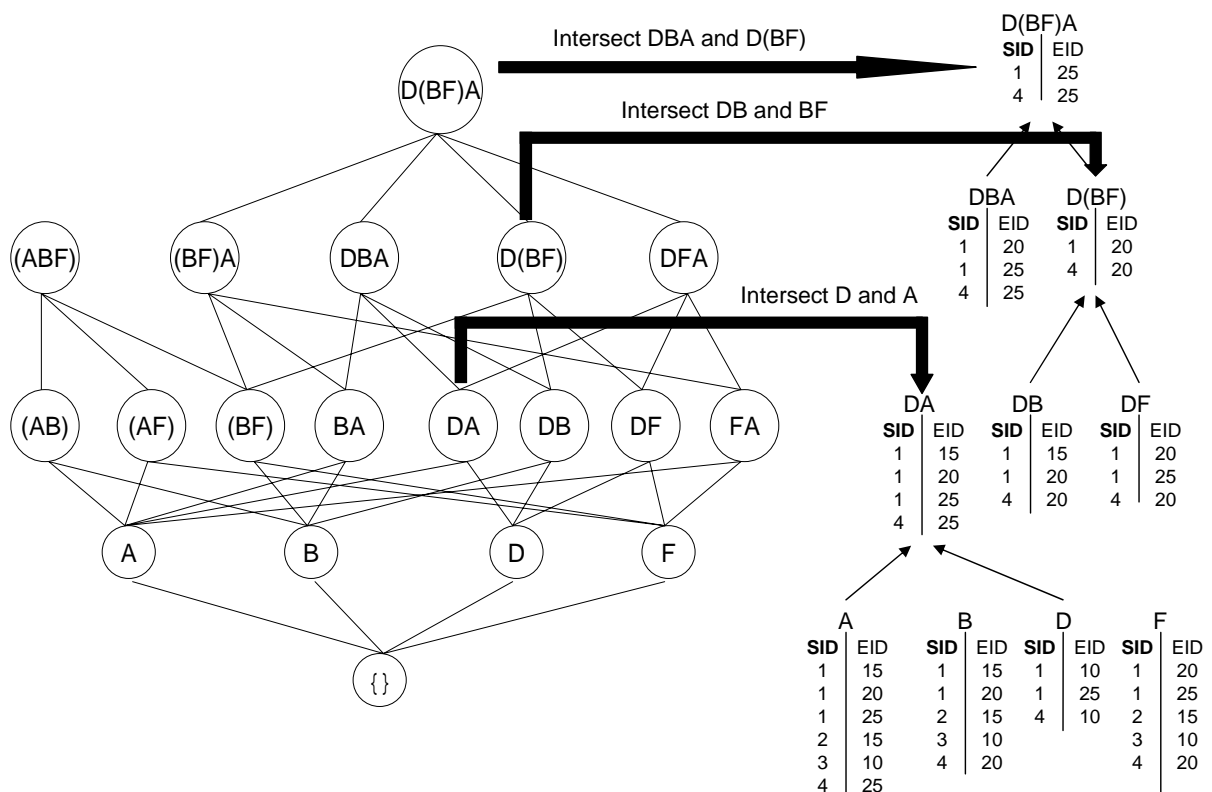


図 3.7 表 3.5 のマイニングの流れにおける ID-List の結合

＜A＞と＜D＞を結合することによって＜DA＞（図中では DA）を生成している。同様に、＜DB＞（図中では DB）と＜DF＞（図中では DF）の結合から＜D (BF)＞（図中では D(BF)）、＜DBA＞（図中では DBA）と＜D (BF)＞（図中では D(BF)）の結合から＜D (BF) A＞（図中では D(BF)A）が生成される。

3.4 SPAM[6]

SPADE では、頻出シーケンスを結合するところが主要なコストになっている。そのため、何度も結合を繰り返すと、候補シーケンスが膨大な数になってしまい、マイニングの時間が長くなってしまう。以上の問題を解決するために Ayres らによって 2001 年に提案された手法が SPAM(Sequential Pattern Mining)[6]である。

SPAM は SPADE と同様に Lattice の考え方をを用いている。また、アルゴリズムの流れも SPADE と基本的には同じで、長さ $k - 1$ の頻出シーケンスを深さ優先探索で結合して、長さ k の頻出パターンを抽出していく。SPADE との大きな違いは ID-List を **vartical** なビットマップで表現するところにある。以下で、ビットマップを用いたデータ構造について説明する。

3.4.1 ビットマップを用いたデータ構造

表 3.1 5 のデータベースを例に説明する。表 3.1 5 のデータベースをビットマップで表現したものが表 3.1 6 である。

表 3.1 5 元のデータベース

CID	EID	Item
1	1	a,b,d
1	3	b,c,d
1	6	b,c,d
2	2	b
2	4	a,b,c
3	5	a,b
3	7	b,c,d

表 3.1 6 表 3.1 5 のビットマップ表現

CID	EID	a	b	c	d
1	1	1	1	0	1
1	3	0	1	1	1
1	6	0	1	1	1
-	-	0	0	0	0
2	2	0	1	0	0
2	4	1	1	1	0
-	-	0	0	0	0
-	-	0	0	0	0
3	5	1	1	0	0
3	7	0	1	1	1
-	-	0	0	0	0
-	-	0	0	0	0

ビットマップ表現では、トランザクション内にアイテムが存在すれば1で、存在しなければ0で表現している。例えば、表 3. の **CID** が 1、**EID** が 3 のトランザクションでは、**b,c,d** が行われているので、**a** のところは 0、**b,c,d** のところは 1 になっている。

3.4.2 ビットマップ表現におけるシーケンスの結合

SPAM におけるシーケンスの結合方法は 2 種類存在する。一つ目は別々の時間に起こったトランザクションの結合である。例えば、 $\langle a \rangle$ と $\langle b \rangle$ (**a** と **b** が別の時間に起こっている) から $\langle ab \rangle$ を生成するといった結合を指す。二つ目は同じ時間に起こったトランザクションの結合である。例えば、 $\langle ab \rangle$ と $\langle d \rangle$ (**b** と **d** が同じ時間に起こっている) から $\langle a(bd) \rangle$ を生成するといった結合を指す。この 2 種類のシーケンスの結合方法についてそれぞれ以下で説明する。

別々の時間に起こったトランザクション同士の結合

表 3.1 6 のビットマップの $\langle a \rangle$ と $\langle b \rangle$ から $\langle ab \rangle$ を生成する例を考える。結合の流れを表したものが図 3.8 である。

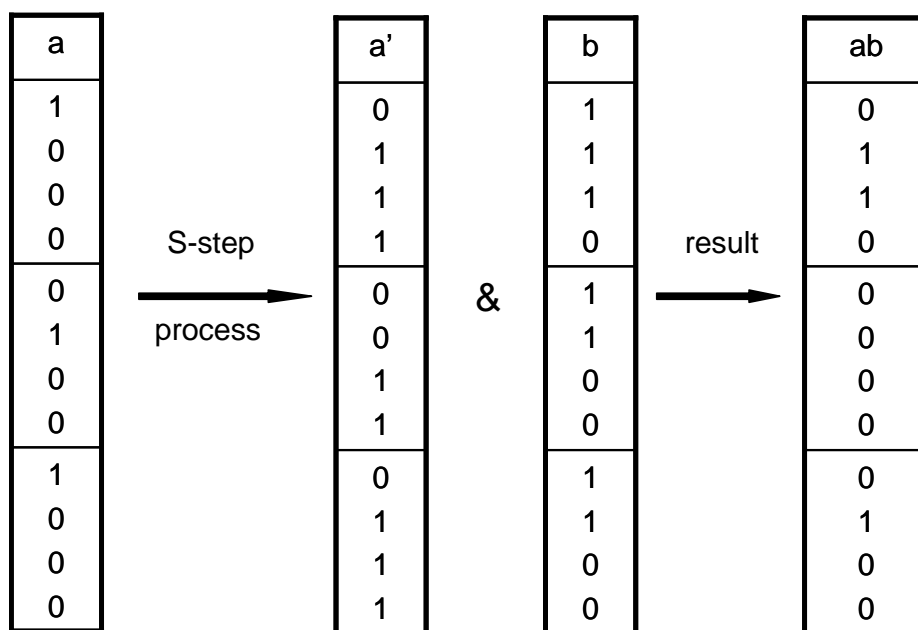


図 3.8 $\{a\}$ と $\{b\}$ を結合して $\{a\}, \{b\}$ を生成する

結合の流れを説明する。まず、 $\langle a \rangle$ について、それぞれのセクションではじめて 1 が出てくるまで 0 にし、それより後を 1 に置き換える（図 3.8 における **S-step process** で、 a から a' への変換）。例えば、 a の一番上のセクションを見てみると、上から 1, 0, 0, 0 となっている。1 は最初に出現しているので一番上だけ 0 にし、後はすべて 1 にする。同様に a の 2 つ目のセクションを見てみると、0, 1, 0, 0 となっており、2 番目に 1 が出現している。よって変換後は、0, 0, 1, 1 となる。そして変換したビットマップ（図中では a と b との論理積が $\langle ab \rangle$ を表している（図 3.8 中では ab が結合した結果を表している）。図 3.8 では **CID** が 1 と 3 のシーケンスで $\langle ab \rangle$ が存在していることを表している。つまりサポート値が 2 であることが分かる。

同時に起こったトランザクションの結合

図 3.8 で結合したシーケンス $\langle ab \rangle$ と $\langle d \rangle$ を結合して $\langle a(bd) \rangle$ を生成する例を考える。結合の流れを図 3.9 に示す。

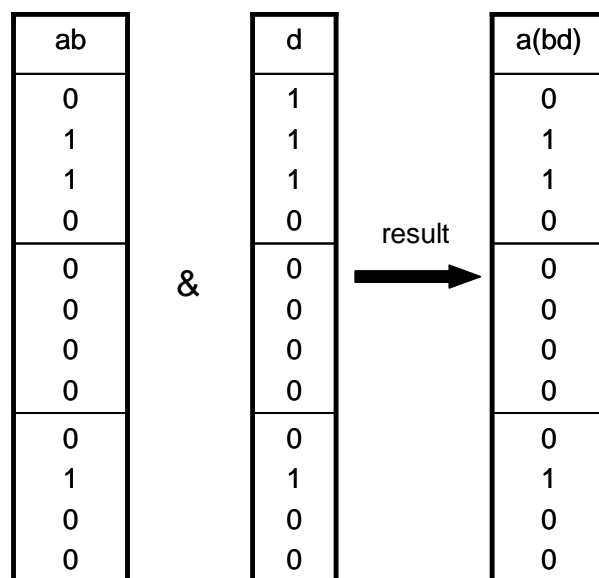


図 3.9 <ab>と<d>を結合して<a(bd)>を生成する

同時に起こったトランザクションの結合はそれぞれのビットマップの論理積をとることで生成することができる。図 3.9 の例では、<a (bd)>（図中では a(bd)）のサポート値が 2 であることが分かる。

3.5 PrefixSpan[6]

PrefixSpan は 2001 年に **Pei** らによって提案された射影ベースのアルゴリズムである。射影ベースのアルゴリズムとは、候補シーケンスを生成せずに、データベースを射影することによって頻出シーケンスを抽出するアルゴリズムである。すなわち、**PrefixSpan** は候補シーケンスを生成せずに、**Prefix projection** という特殊な射影方法と射影によって生成される **Prefix** データベースを用いることで、マイニングの高速化を実現している。本節では、まず **Prefix projection** と **Prefix** データベースについての説明をした後、具体的なアルゴリズムの流れを説明する。また、**PrefixSpan** を最適化したものとして **bi-level projection** と、**Pseudo-Projection** の 2 つを紹介する。

3.5.1 Prefix projection と Prefix データベース

PrefixSpan では **Prefix Projection** と呼ばれる射影を行っている。**Prefix Projection** とは、射影元のシーケンスから射影対象のシーケンスより後ろに存在するアイテムからなるシーケンスのみを抽出する射影である。例えばシーケンス<a (abc) (ac) d(cf)>の射影について考えてみる。ここで、<a>、<aa>、<ab>について射影を行った結果が表 3.17 である。

表 3.1 7 Prefix projection

Prefix	射影後のシーケンス
<a>	<(abc) (ac) d (cf)>
<aa>	<(_bc) (ac) d (cf)>
<ab>	<(_c) (ac) d (cf)>

<a>について射影を行ったときは先頭の a を省いたシーケンスが射影後のシーケンスとなる。同様に、<aa>についても最初の ab が除かれている。<ab>についての射影では、シーケンス内で最初に出てくる ab 以降が射影後のシーケンスとなる。そして、与えられたデータベースに対し、射影を行った結果のデータベースを **Prefix** データベースという。

3.5.2 PrefixSpan アルゴリズムの流れ

PrefixSpan アルゴリズムの流れを以下に示す。

- 1、長さ 1 の頻出シーケンスを抽出する
- 2、深さ優先探索で射影を行い、マイニングを行う。

表 3.1 8 のデータベースが与えられたときの **PrefixSpan** の動きを以下に示す。なお、**min_sup=2** とする。

表 3.1 8 シーケンスデータベース

SID	Sequence
10	<a (abc) (ac) d (cf)>
20	<(ad) c (bc) (ae)>
30	<(ef) (ab) (df) c b>
40	<e g (af) c b c>

1、長さ 1 の頻出シーケンスを抽出する

データベースをスキャンし、長さ 1 の頻出アイテムを抽出する。表 3.1 8 の例では長さ 1 の頻出シーケンスは<a>:4、:4、<c>:4、<d>:3、<e>:3、<f>:3 の 6 つとなる。ここで、各シーケンスの後ろの数字はサポート値である。

2、深さ優先探索で射影を行い、マイニングを行う（図 3.9 参照）

例ではまず、<a>について射影を行い、射影されたデータベースを構築し、a を先頭とする長さ 2 の頻出シーケンスを抽出する。ここでは<aa>:2、<ab>:4、<(ab)>:2、<ac>:4、<ad>:2、<af>:2 の 6 つが抽出される。**PrefixSpan** は深さ優先探索でマイニングを行うので、次は<aa>について射影を行う。射影によって構築されたデータベースが表 3.1 9 である。ここでそれぞれのサポート値をカウントすると、<aaa>:1、<a(ab)>:1、<a(ac)>:1、<aac>:1、<aad>:1、<aae>:1、<aaf>:1 となり最小サポート以上のシーケンスは存在しないので、<aa>より深く射影を行うことは不可能である。よって次は、<ab>について

射影を行う。以降も同様に深さ優先で射影したデータベースにおいて、最小サポート以上のアイテムが存在しなくなるまでマイニングを行っていく。具体的には表 3.20 のようになる。表 3. の **Frequent Patterns** の $\langle a \rangle$ 、 $\langle aa \rangle$ 、 \dots 、 $\langle af \rangle$ 、 $\langle b \rangle$ 、 \dots 、 $\langle ec \rangle$ の順に射影を行う。また、表 3.18 のマイニングの全体像を図 3.10 に示す。

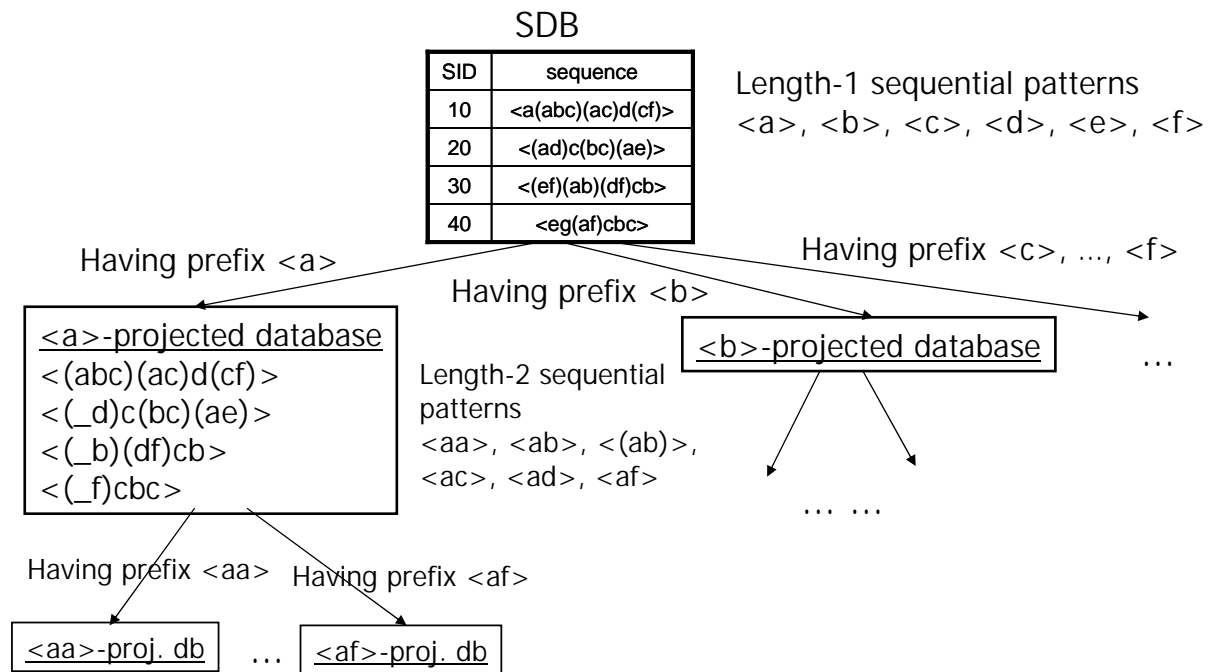


図 3.10 PrefixSpan の流れ

表 3.19 $\langle aa \rangle$ について射影した結果

$\langle aa \rangle$ -projected database
$\langle (_bc)(ac)d(cf) \rangle$
$\langle (_e) \rangle$

表 3.20 射影されたデータベースとシーケンシャルパターン

Prefix	Projected(postfix) database	Frequent Patterns
<a>	<(abc)(ac)d(cf)>、<(_d)c(bc)(ae)>、 <(_b)(df)cb>、<(_f)cbc>	<a>、<aa>、<ab>、<a(bc)>、 <a(bc)a>、<aba>、<abc>、<(ab)>、 <(ab)c>、<(ab)d>、<(ab)dc>、 <(ab)f>、<ac>、<aca>、<acb>、 <acc>、<ad>、<adc>、<af>
	<(_bc)(ac)d(cf)>、<(_c)(ae)>、 <(df)cb>、<c>	、<ba>、<bc>、<(bc)>、 <(bc)a>、<bd>、<bdc>、<bf>
<c>	<(ac)d(cf)>、<(bc)(ae)>、 、<bc>	<c>、<ca>、<cb>、<cc>
<d>	<(cf)>、<c(bc)(ae)>、<(_f)(cb)>	<d>、<db>、<dc>、<dcb>
<e>	<(_f)(ab)(df)cb>、<(af)cbc>	<e>、<ea>、<eab>、<eac>、 <each>、<eb>、<ebc>、<ec>、 <ecb>、<ef>、<efb>、<efc>、 <efcb>
<f>	<(ab)(df)cb>、<cbc>	<f>、<fb>、<fbc>、<fc>、<fcb>

3.5.3 PrefixSpan の最適化

PrefixSpan における主要なコストは、射影されたデータベースを構築することにある。つまり、射影されたデータベースのサイズを減らすことができれば、PrefixSpan の速度を向上させることができる [4]。[4]では、PrefixSpan の最適化として **bi-level projection** と **Pseudo-Projection** の 2 つの手法が提案されている。

bi-level projection

bi-level projection は長さが奇数の頻出シーケンスを抽出するときは通常の PrefixSpan と同じように射影を行い、射影データベースを構築するが、長さが偶数の頻出シーケンスを抽出する際にはデータベースをスキャンし、三角行列形式のデータベースを構築する。このデータベースは通常の PrefixSpan の射影により構築されるデータベースよりもサイズを削減することができる。

再び表 3. の例を用いて説明する。また、**min_sup=2** とする。まずは、通常の PrefixSpan と同様に、データベースをスキャンし、長さ 1 の頻出アイテムを抽出する。抽出されるのは <a>、、<c>、<d>、<e>、<f> の 6 つである。

次のステップでは、通常の PrefixSpan では射影を行い、新たなデータベースを構築するが、**bi-level projection** では、表 3.1 4 のような **6×6** の三角行列を構築する。ここでは構築された行列を **M** とする。

表 3.1 4 長さ 2 における行列 **M**

a	2					
b	(4,2,2)	1				
c	(4,2,1)	(3,3,2)	3			
d	(2,1,1)	(2,2,0)	(1,3,0)	0		
e	(1,2,1)	(1,2,0)	(1,2,0)	(1,1,0)	0	
f	(2,1,1)	(2,2,0)	(1,2,1)	(1,1,1)	(2,0,1)	1
	a	b	c	d	e	f

それぞれの行列の要素は長さ 1 の頻出アイテムを組み合わせて生成した、長さ 2 のサポート値を表している。また、行列には値を 3 つ含んでいるもの($M[a,c]$ (**a** と **c** の要素)=(4,2,1))と、1 つだけ含んでいるものの 2 種類が存在する。値が 1 つだけのものは対角線上に存在し、例えば、 $M[a,a]=2$ は $\langle aa \rangle$ のサポート値が 2 であることを表している。また、値を 3 つ含んでいるものは対角線以外の場所に存在する。例えば、 $M[a,c]=(4,2,1)$ における $\langle ac \rangle$ のサポート値が 4、 $\langle ca \rangle$ のサポート値が 2、 $\langle (ac) \rangle$ のサポート値が 1 であることを表している。 $M[c,a]$ の表す値は $M[a,c]$ の値と対称になっているので、表現する必要が無い。よって、データベースを削減することが可能になる。

長さ 2 の頻出シーケンスについては通常の **PrefixSpan** と同様に射影を行い、データベースを構築する。例えば表 3.1 4 における頻出シーケンスの一つである $\langle ab \rangle$ について射影を行い、構築されたデータベースが表 3.1 5 である。

表 3.1 5 $\langle ab \rangle$ について射影した結果

projected database
$\langle (_c)(ac)(cf) \rangle$
$\langle (_c)a \rangle$
$\langle c \rangle$

ここで、表 3.1 5 のデータベースに対しスキャンを行い、長さ 1 の頻出シーケンスを抽出する。ここで抽出される長さ 1 の頻出シーケンスは $\langle a \rangle$ 、 $\langle c \rangle$ 、 $\langle (_c) \rangle$ の 3 つである。これをもとに表 3.1 6 のような $\langle ab \rangle$ について射影した行列を生成する。表 3.1 6 中の Φ という記号は組み合わせを生成できないという意味である。

表 3.1 6 $\langle ab \rangle$ について射影した行列

a	0		
c	(1,0,1)	1	
(_c)	(Φ ,2, Φ)	(Φ ,1, Φ)	Φ
	a	c	(_c)

ここで、最小サポート値を満たす組み合わせは $\langle _c \rangle a$ だけなので、これ以上射影を行う必要は無い。

以降も長さが奇数の頻出シーケンスを抽出するときは通常の **PrefixSpan** と同じように射影を行い、長さが偶数の頻出シーケンスを抽出するときは三角行列形式のデータベースを構築する。そして、通常の **PrefixSpan** と同様に深さ優先探索でマイニングを行っていく。この結果、表 3. のデータベースのマイニングにおいて、通常の **PrefixSpan** では合計で 53 回射影されたデータベースを構築しなければならなかったが、**bi-projection** では、射影されたデータベースの構築は 22 回ですむ。

Pseudo Projection

Prefix Projection とは、シーケンス $s_1 = \langle a(abc)(ac)d(cf) \rangle$ が与えられた場合、 $\langle a \rangle$ について射影した結果、 $\langle (abc)(ac)d(cf) \rangle$ 、 $\langle ab \rangle$ について射影した結果、 $\langle _bc \rangle (ac)d(cf)$ が得ることができる射影である。**Pseudo Projection** では射影されたデータベース（シーケンス）を冗長なものとし、別の表現を用いることによって高速化を図っている。具体的には、それぞれの射影を、射影元のシーケンスへのポインタと射影されたシーケンスのオフセットという 2 つの情報で表現する。ポインタとは射影する対象、オフセットは射影されたシーケンスが射影元のシーケンスの中で何番目以降が射影されているかを表している。 s_1 の例を用いて説明すると、 s_1 における $\langle a \rangle$ についての射影はポインタが s_1 、オフセットが 2 と表される。 $\langle a \rangle$ について射影したシーケンスは $\langle (abc)(ac)d(cf) \rangle$ なので、 s_1 の中で 2 番目以降のアイテムで表されている。よってオフセットは 2 となる。同様に、 s_1 の $\langle ab \rangle$ についての射影はポインタが s_1 、オフセットが 4 と表すことができる。

3.6 各手法の比較

本章で紹介した 4 つの手法について、候補シーケンス削減、データベース分割、シーケンスの縮小という 3 つの戦略から比較してみる。

(1) 候補シーケンス削減

候補シーケンス削減とは、頻出シーケンスに為り得ない候補シーケンスをできるだけ早い段階で削除する戦略で、プロセスのコストとサポート値カウントの際のオーバーヘッドを減らすことができる。候補シーケンス削減の戦略は **GSP**[3]、**SPADE**[5]、**SPAM**[6]、**PrefixSpan**[4]のすべてのアルゴリズムで用いられている。

(2) データベース分割

データベース分割は、データベースを複数のグループに分割する戦略で、それぞれのグループをメインメモリにうまく割り当てることができるようになる。**SPADE**[5]と **SPAM**[6]では、候補シーケンスをもとにしてデータベースを分割している。**PrefixSpan**[4]では、射影によって得られた頻出アイテムをもとにデータベースを分割している。**GSP**[3]ではデータベー

ス分割の戦略は用いられていない。

(3) シーケンスの縮小

シーケンスを減らす戦略とは、できるだけ多くのシーケンスを減らし、プロセスのコストを現象させる戦略である。この戦略は **PrefixSpan[4]**のみで用いられており、射影によってシーケンスの縮小を実現している。例えば、シーケンス<(f) (ag) (bfh) (bf)>に対し、<a>について射影を行うと<(g) (bfh) (bf)>となり、シーケンスを縮小することができる。

以上の比較をまとめたものが表 3.17 である。

表 3.17 アルゴリズムと戦略

	候補シーケ ンス削減	データベー ス分割	シーケンス の縮小	アルゴリズムの特徴
GSP[3](1996)	✓			Apriori をベースに候補 シーケンスを生成
SPADE[4](2000)	✓	✓		Lattice と ID-List を用 い高速化を実現
SPAM[5](2001)	✓	✓		ID-List をビットマップ で表現
PrefixSpan[6](2001)	✓	✓	✓	Prefix Projection によ ってシーケンスデータ ベースを縮小

3.7 Sequential Pattern Mining の拡張[11]

近年、**Sequential Pattern Mining** の拡張手法として、**Stream Mining** という研究が盛んになっている([12][13])。 **Stream Mining** とは、無限に入ってくるデータ(**data stream**)をマイニングの対象とした手法で、**Web** ページ閲覧の流れの解析、エネルギー消費測定、ネットワーク上のデータの解析、株取引といったデータが継続して入ってくる分野に応用が可能である。 **Stream Mining** の概要を以下で説明する。

上述したように、**Stream Mining** が対象とするデータは、無限に流入してくるデータである。しかし、これに対して計算機の資源は有限であり、**Stream Mining** では有限な計算機資源を用いてどのようにマイニングを行なうかということが課題となっている。

この課題に対する解決方法として以下の2つの方法が挙げられる。

① 近似的な解を求める

図 3.11 のようにデータ要約という小さなデータ構造を主記憶上に保持し、近似的な解を求める方法である。

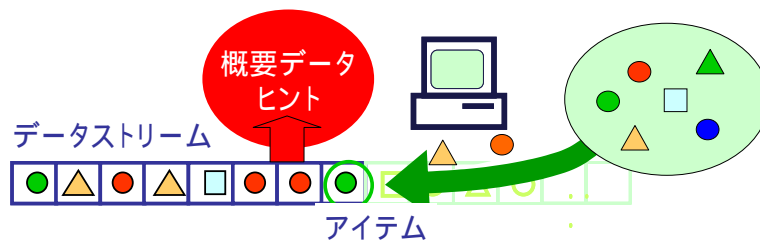


図 3.1 1 Data Stream から近似解を得る

② ある時点にスポットを当ててマイニングを行なう

図 3.1 2 のように対象をすべての過去のデータではなくある一定の時間内のデータに絞ってマイニングを行う方法である。

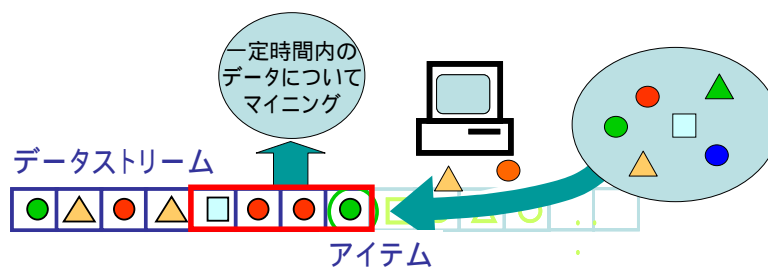


図 3.1 2 一定時間内のデータに絞ってマイニング

第4章 提案手法

本章では提案手法であるアイテム間の距離（時間）を考慮した **Sequential Pattern Mining** について述べる。まず、**Sequential Pattern Mining** における順序と時間について述べた後、提案手法についての具体的な説明をする。

4.1 **Sequential Pattern Mining** における順序と距離

Sequential Pattern Mining の最大の特徴はアイテム間の順序を考慮するところにある。現実の世界にはアイテム間の順序が重要となってくる事例が数多く存在する。例えば、最初にパソコンを購入し、その後 **CD-ROM** を購入した人と、最初に **CD-ROM** を購入し、その後パソコンを購入した人がいるとする。**Frequent Pattern Mining** では、両者は同じ組合せとして扱われている。しかし、最初に **CD-ROM** を購入し、その後パソコンを購入するという行動は考え難い。そこで、アイテム間の順序を考慮する点が特徴である **Sequential Pattern Mining** では、両者を別の組合せとして、区別してマイニングを行なう。

しかし、従来の **Sequential Pattern Mining** では、アイテム間の距離が考慮されないという問題がある。距離（時間）を考慮しないでマイニングを行なうと、短期的な行動と長期的な行動を区別することができない。シーケンス $\langle a\ b \rangle$ では **a** と **b** の間の距離は分からないのである。上の例を用いると、最初にパソコンを購入し、次に **CD-ROM** を購入するという行動の中でも1週間後に **CD-ROM** を購入する人と、1年後に **CD-ROM** を購入する人とはそれぞれ行動の持つ意味は異なっているのだが、**Sequential Pattern Mining** では両者の行動を同じものとして扱っている。

4.2 アイテム間の距離（時間）を考慮した **Sequential Pattern Mining**

4.1 で説明したように、**Sequential Pattern Mining** では、アイテム間の距離（時間）によってトランザクションが区別されない。そのため、短期的な行動と長期的な行動という持つ意味の異なる2つの行動を区別することができない。図4.1にレンタルビデオ店の顧客の行動を例に示す。上の人物はシリーズ物のビデオを一度に8本まとめて借りている。それに対し、下の人物は、2本だけ借りて、見終わったらまた2本借りている。このように、同じシリーズ物をレンタルする場合でも、人によって借り方が異なってくる。従来の **Sequential Pattern Mining** では、図4.1の2人の行動は全く別のものとして扱われてしまう。

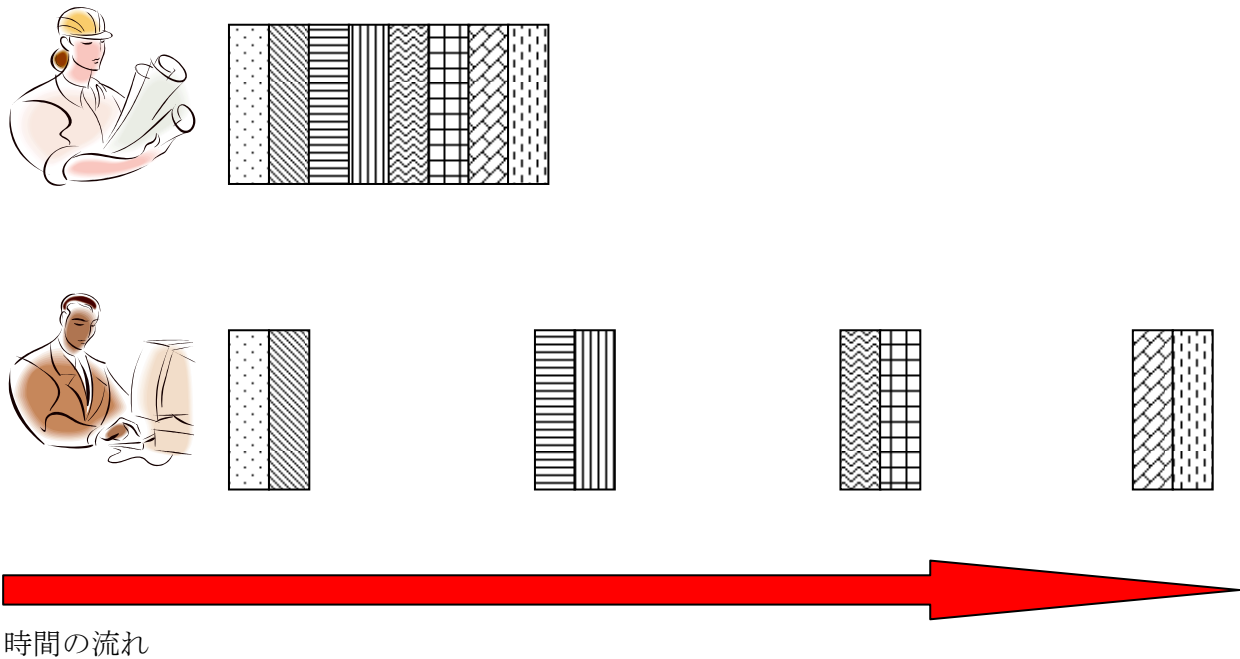


図 4.1 シリーズ物の借り方の違い

以上のような問題を解決するために、本論文では、アイテム間の距離（時間）を考慮した **Sequential Pattern Mining** を提案する。

提案手法は、ユーザが定義した制限距離（時間）以下のトランザクションは、同時に起こったものとみなしてマイニングを行う。つまり、従来の **Sequential Pattern Mining** では別々のものとして扱われていたトランザクションの中で、ユーザにとっては同時に行われているとみなしてもよい距離のトランザクションを結合し、マイニングを行う。以下で、制限距離について定義を行なう。

シーケンス $\langle AB \rangle$ が与えられたとする。また、 AB 間の距離を a ($a \geq 0$) とする。ここで、制限距離 x が与えられた場合について考える。 $a \leq x$ の時は、図 4.2 のように AB が同時に起こったとみなして $\langle (AB) \rangle$ とする。

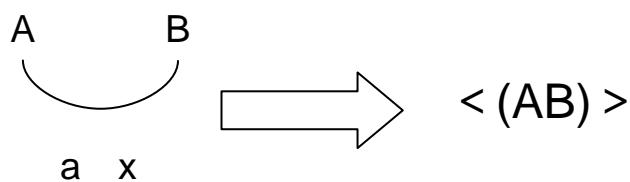


図 4.2 アイテム間の距離が制限距離以下の場合

しかし、 $a > x$ の時は、図 4.3 のように $\langle AB \rangle$ のままとする。

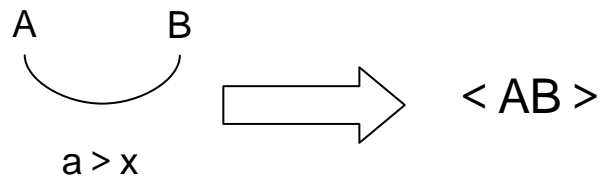


図 4.3 アイテム間の距離が制限距離より大きい場合

また、制限距離以下のトランザクションが連続した場合は、連続した制限距離以下のトランザクションをすべて一つの“()”でまとめる。ここでシーケンス $\langle ABCD \rangle$ が与えられている例について考える。 AB 間の距離を a 、 BC 間の距離を b 、 CD 間の距離を c とする。制限距離が x が与えられたときに、 $a \leq x, b \leq x, c > x$ ならば、シーケンスは図 4.4 のように $\langle (ABC)D \rangle$ となる。

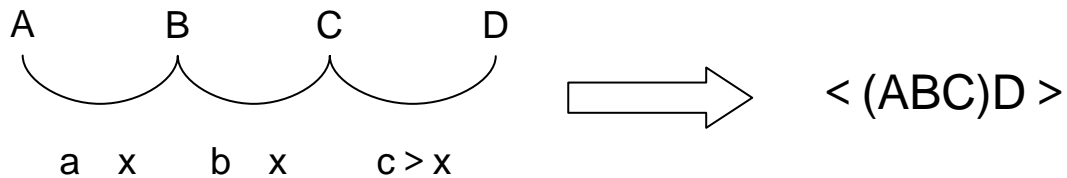


図 4.4 制限距離以下のトランザクションが連続した例

4.3 提案手法のアルゴリズムの流れ

提案手法は、データ変換とマイニングの2つのフェーズに分けることができる。

データ変換のフェーズでは、制限距離（時間） n ($n \geq 0$) を与え、マイニングの対象となるデータベースから、距離が n 以下のトランザクションは同時に起こったものとしてシーケンスデータベースを作成する。

マイニングのフェーズでは、データ変換のフェーズで得られたシーケンスデータベースを元に実際にマイニングを行なう。提案手法では **Pseudo-Projection** 方式の **PrefixSpan**[9]を用いてマイニングを行なう。**PrefixSpan** の詳細な説明は3.5章において行なっている。

以下に提案手法の流れを例を用いて示す。

データ変換

マイニング対象のデータベースは表 4.1 のようなものである。

表 4.1 マイニングに使用するデータベースの例

ユーザ ID	トランザクションタイム	トランザクション
1 0	4 月 3 0 日	A
1 0	4 月 3 0 日	B
1 0	7 月 2 4 日	D
2 0	5 月 1 3 日	B
2 0	9 月 1 7 日	C
3 0	4 月 2 日	A
3 0	6 月 2 7 日	B
3 0	6 月 2 8 日	C
3 0	7 月 1 日	D

表 4.1 のデータベースをもとにシーケンスを作成する。提案手法を用いずに従来どおりにシーケンスを作成したものが表 4.2 である。

表 4.2 通常通り作成したシーケンスデータベース

SequenceID	Sequence
1 0	< (A B) D >
2 0	< B C >
3 0	< A B C D >

提案手法では、表 4.1 のような基のデータベースにおいて、ユーザが定義した距離内のトランザクションを同時に起こったものとして変換を行なう。表 4.1 のデータベースを制限距離を 1 0 日にしてシーケンスを作成したシーケンスデータベースが表 4.3 である。

表 4.3 制限距離をもとに作成したシーケンスデータベース

SequenceID	Sequence
1 0	< (A B) D >
2 0	< B C >
3 0	< A (B C) D >

表4.2と表4.3のシーケンスデータベースを比較してみると、シーケンス ID 30 のシーケンスにおいて、従来手法では $\langle A B C D \rangle$ になっているが、提案手法では、 $\langle A (B C) D \rangle$ となっている。表4.1のデータベースを見てみると、トランザクション B は 6 月 27 日に行われており、トランザクション C は 6 月 28 日に行われている。制限距離（時間）は 10（日）なので、提案手法では、トランザクション B と C は同時に行われたと見なし、一つのトランザクション (B C) として扱う。

マイニング

提案手法によって、変換されたシーケンスデータベースについて、実際にマイニングを行う。上で述べたように、提案手法では **Pseudo Projection** 方式の **PrefixSpan[9]** を用いてマイニングを行う。

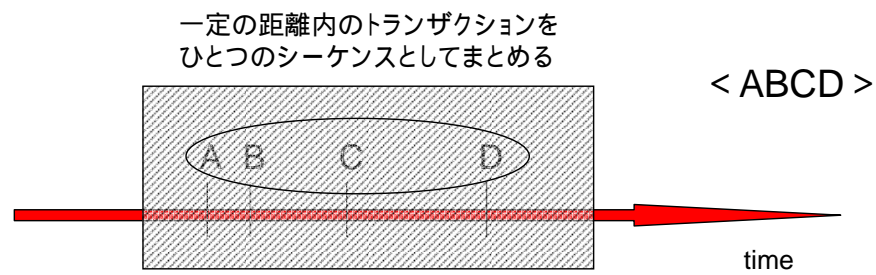
4.4 Stream Mining の観点から見た提案手法

関連研究でも紹介した **Sequential Pattern Mining** の拡張手法である **Stream Mining** の観点から見た提案手法はどのような位置づけになるのかについて説明する。

提案手法は制限距離をユーザが定義し、制限距離以下のトランザクション同士を一つのトランザクションにまとめる手法であるが、制限距離を設けて時間で区切るという点については **Stream Mining** の方法の一つ（関連研究における②の方法）と共通している。では、提案手法と **Stream Mining** の違いはどのような点にあるのだろうか。

Stream Mining ではある一定の距離（時間）内に起こったトランザクションを一つのシーケンスにまとめてマイニングを行なっている。それに対し、提案手法では、制限距離内にあるトランザクションを一つのトランザクションにまとめている。図4.5の例では A、B、C、D の4つのトランザクションが起こっているが、**Stream Mining** では制限時間内に入っているすべてのトランザクションをひとつのシーケンスとしてまとめるので、シーケンス $\langle ABCD \rangle$ になる。しかし、提案手法ではトランザクション間の距離に着目してその制限距離以下のトランザクションを一つのトランザクションにまとめるので、シーケンスは $\langle (AB)CD \rangle$ となる。

Stream Mining



提案手法

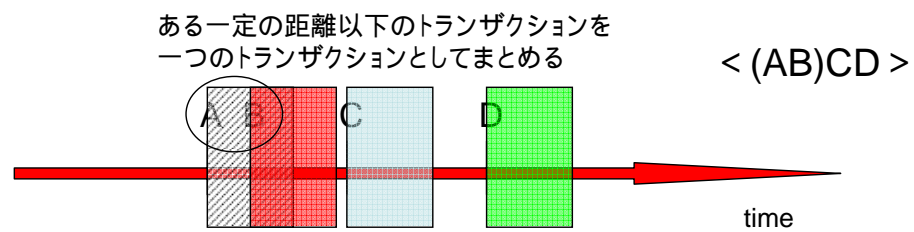


図 4.5 制限距離の区切りによるトランザクションのまとめ方の違い

つまり、**Stream Mining** の観点から提案手法をみた場合、制限距離に対するトランザクションのまとめ方という点において **Stream Mining** とは異なっている。

第5章 性能評価

本章では、第4章で述べた提案手法を実装し、抽出されたシーケンスについて評価を行う。

5.1 実験環境

実験を行った計算機の環境を以下に示す。

表 5.1 計算機の環境

CPU	Celeron M 1.4GHz
メモリ	752MB
OS	Windows XP Home Edition
コンパイラ	Borland C++ 5.5.1 for Win32

5.2 データセット

今回、実験に用いたデータセットとして、MUSASHI (Mining Utilities and System Architecture for Scalable processing of Historical data) [10] の顧客購買履歴データ生成スクリプトによって生成された人工データを使用した。MUSASHI では、店、日付、時間、レシート、顧客、商品、大分類、中分類、小分類、細分類、メーカー、ブランド、仕入単価、単価、数量、金額、仕入金額、粗利金額を含むデータを生成することができるが、今回はこの中から、日付、顧客、商品の3つのデータを切り取って使用した。ここで、使用した3つのデータについて生成する際に調整可能なパラメータについて表5.2に示す。

表 5.2 調整可能なパラメータ

項目	計算ルール	使用したパラメータ
店別顧客人数	手入力にて店コードと来店人数を指定する	店舗数:1 顧客数:2000 人
顧客の一来店あたりの購入数量	正規分布に基づく乱数にて決定	nrnd(nrnd(5 個,2 個),3) 但し、0 以下は省く
顧客の来店間隔	正規分布に基づく乱数にて決定	nrnd(30 日,10 日)
顧客の来店開始日	ある日付を中心に正規分布するよう に乱数で決める	07/01+nrnd(0,200 日)
顧客の来店終了日	ある日付を中心に正規分布するよう に乱数で決める	07/01+nrnd(0,200 日)
購入商品	商品によって売れる頻度を乱数で 変える	nrnd(商品番号の中点, 商品 数/3)

生成されたデータセットには表 5.3 のような顧客 ID、購入日、購入アイテムが含まれている。

表 5.3 生成されたデータセットの例

顧客 ID	購入日(月/日)	購入アイテム
1	1/26	234
1	1/26	73
2	2/9	547
3	2/1	170
3	2/9	365

上述の **MUSASHI** によって生成されたデータセットをもとにして、顧客 ID ごとにシーケンスを作成する。作成するシーケンスの形式をシーケンス<(2 5 7) (1 2) (3 9) 4>を例に説明する。

2 5 7 -1 19 1 2 -1 6 3 9 -1 30 4 -1 -2

－1 はトランザクションの終了、－2 はシーケンスの終了を表す。また、－1 の直後の数字（上の数列では 1 9 と 6 と 3 0）はトランザクション間の距離を表している。例えば、上の数列において、トランザクション**(2 5 7)**と**(1 2)**の距離は 1 9 である。また、各トランザクション内のアイテムは昇順にソートされている。

このような形式の数列からトランザクション間の距離を除いた形に変換する。上の例の場合

合、以下のような形になる。

2 5 7 -1 1 2 -1 3 9 -1 4 -1 -2

ここで、提案手法を用いた場合、ユーザが定義した制限距離以下のトランザクションは結合される。上の例において、制限距離を10とした場合、以下のような形になる。

2 5 7 -1 1 2 3 9 -1 4 -1 -2

距離が10以下のトランザクション間は**(1 2)**と**(3 9)**なので、この2つが結合され、**<(2 5 7) (1 2 3 9) 4>**となる。

なお、実際にマイニングを行う際には、作成したシーケンスデータベースはバイナリ形式になっている。

今回の実験では、シーケンス数 **2000**、トランザクションが行なわれる期間が **1** 年、アイテム数 **577** となっている。

5.3 実験

5.3.1 実験方法

MUSASHI によって生成したデータセットに対して提案手法を用い、シーケンスを変換してマイニングを行なった場合と、従来手法を用いてマイニングを行なった場合について比較する。ここでいう従来手法とは、与えられたデータに対し、前処理を何もせずにマイニングを行なう手法を指す。また、それぞれの手法に対し、最小サポート値を変えて、マイニングを行なった結果についても比較を行なう。

5.3.2 実験結果

提案手法を用い、シーケンスの変換を行なった場合と、従来手法（提案手法を用いない場合）について、それぞれマイニングを行ない、比較を行なった。最小サポート値は **100%**～**82%**まで変化させ、それぞれの抽出時間、抽出シーケンス数、**1** シーケンスあたりの抽出時間について比較したものをまとめたのが表5.4である。また、抽出時間、頻出シーケンス数、**1** シーケンスあたりの抽出時間のグラフをそれぞれ図5.1、図5.2、図5.3に示す。

表 5.4 従来手法と提案手法の比較

min_sup	従来手法			提案手法		
	時間 (sec)	頻出シー ケンス数	1 シーケンスあ たりの抽出時間 (Sec)	時間 (sec)	頻出シー ケンス数	1 シーケンスあ たりの抽出時間 (Sec)
100%	0.01	3	0.003333	0.01	2	0.005
99%	0.02	9	0.002222	0.04	20	0.002
98%	0.02	10	0.002	0.04	21	0.001905
97%	0.02	10	0.002	0.04	21	0.001905
96%	0.03	13	0.002308	0.05	24	0.002083
95%	0.03	13	0.002308	0.05	31	0.001613
94%	0.13	71	0.001831	0.07	46	0.001522
93%	0.13	71	0.001831	0.07	47	0.001489
92%	0.17	96	0.001771	0.08	47	0.001702
91%	0.18	96	0.001875	0.08	47	0.001702
90%	0.18	96	0.001875	0.12	83	0.001446
89%	0.18	97	0.001856	0.12	83	0.001446
88%	1.832	1195	0.001533	0.821	658	0.001248
87%	3.575	2225	0.001607	1.001	786	0.001274
86%	8.292	5243	0.001582	1.502	1093	0.001374
85%	28.811	25589	0.001126	2.463	2035	0.00121
84%	38.715	28118	0.001377	4.266	3350	0.001273
83%	207.528	164854	0.001259	7.270	6787	0.001071
82%	901.607	1113153	0.00081	22.472	23913	0.00094

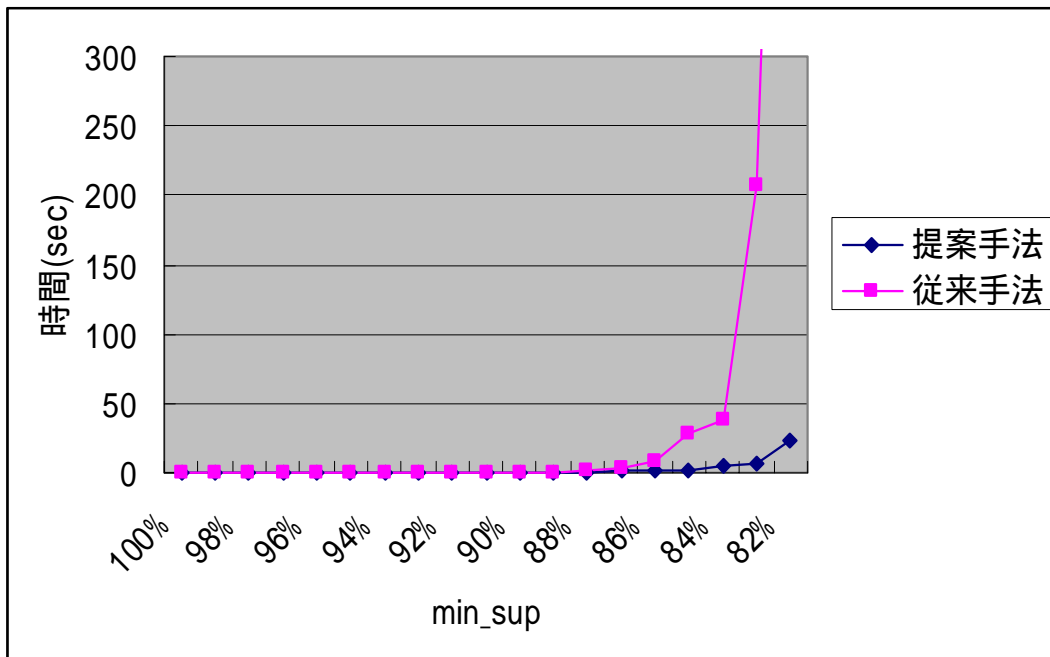


図 5.1 抽出時間の比較のグラフ

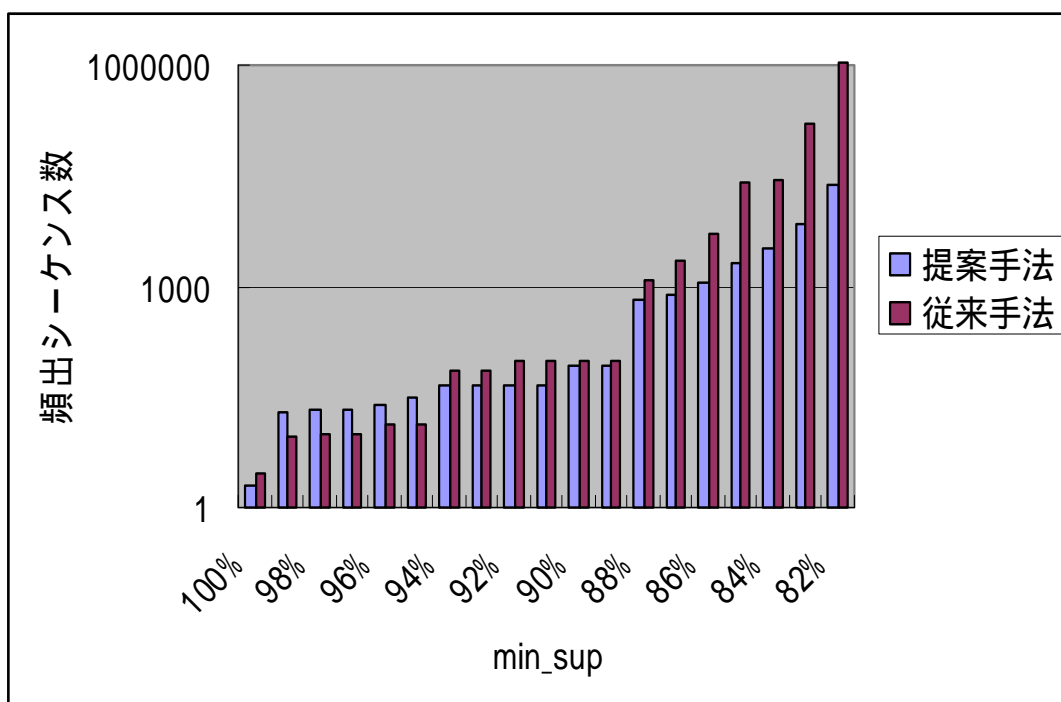


図 5.2 抽出される頻出シーケンスの数の比較

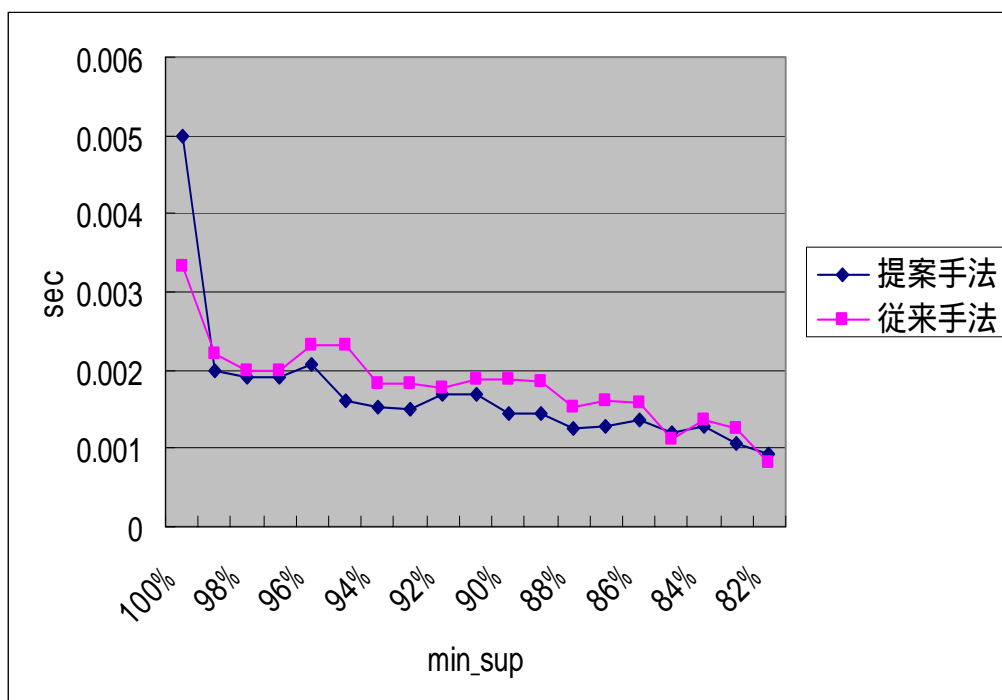


図 5.3 1 シーケンスあたりの抽出時間

次に、提案手法と従来手法では、抽出されるシーケンスの長さに変化があるのかを比較してみる。まず、提案手法を用いた場合に各最小サポート値ごとに抽出されるシーケンスを長さ別の個数を表 5.5 に示す。なお、表の縦軸は最小サポート値、横軸はシーケンスの長さである。

表 5.5 提案手法を用いた場合に抽出されるシーケンスの長さごとの個数

	1	2	3	4	5	6	7	8	9	10	11
100%	2	1									
99%	8	7	4	1							
98%	9	7	4	1							
97%	9	7	4	1							
96%	9	9	5	1							
95%	9	12	8	2							
94%	14	15	11	5	1						
93%	14	16	11	5	1						
92%	14	16	11	5	1						
91%	14	16	11	5	1						
90%	14	22	24	16	6	1					
89%	14	22	24	16	6	1					
88%	25	132	193	166	95	37	9	1			
87%	28	175	240	192	103	38	9	1			
86%	32	254	343	268	138	47	10	1			
85%	38	367	617	562	317	110	22	2			
84%	39	421	883	966	660	288	79	13	1		
83%	44	538	1402	1879	1583	893	344	89	14	1	
82%	50	750	2932	5379	6045	4556	2381	861	207	30	2

同様に従来手法の場合における、抽出される頻出シーケンスの長さ別の個数を表 5.6 に示す。

表 5.6 従来手法を用いたときに抽出されるシーケンスの長さごとの個数

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
100%	2	1															
99%	8	1															
98%	9	1															
97%	9	1															
96%	9	3	1														
95%	9	3	1														
94%	14	38	17	2													
93%	14	38	17	2													
92%	14	39	29	12	2												
91%	14	39	29	12	2												
90%	14	39	29	12	2												
89%	14	40	29	12	2												
88%	25	170	332	344	218	85	19	2									
87%	28	227	584	653	448	208	64	1									
86%	32	314	1033	1359	1209	790	367	115	22	2							
85%	38	465	2137	4330	5915	5747	4010	2018	723	177	27	2					
84%	39	498	2388	5000	6656	6257	4247	2090	736	178	27	2					
83%	44	682	4513	14378	27109	35398	34215	25316	14444	6257	1992	441	61	4			
82%	50	960	8259	36177	92303	161485	211702	216652	176240	114739	59724	24646	7922	1922	333	37	2

次に、各最小サポート値における抽出される頻出シーケンスの平均長を比較したグラフを図5.4に示す。

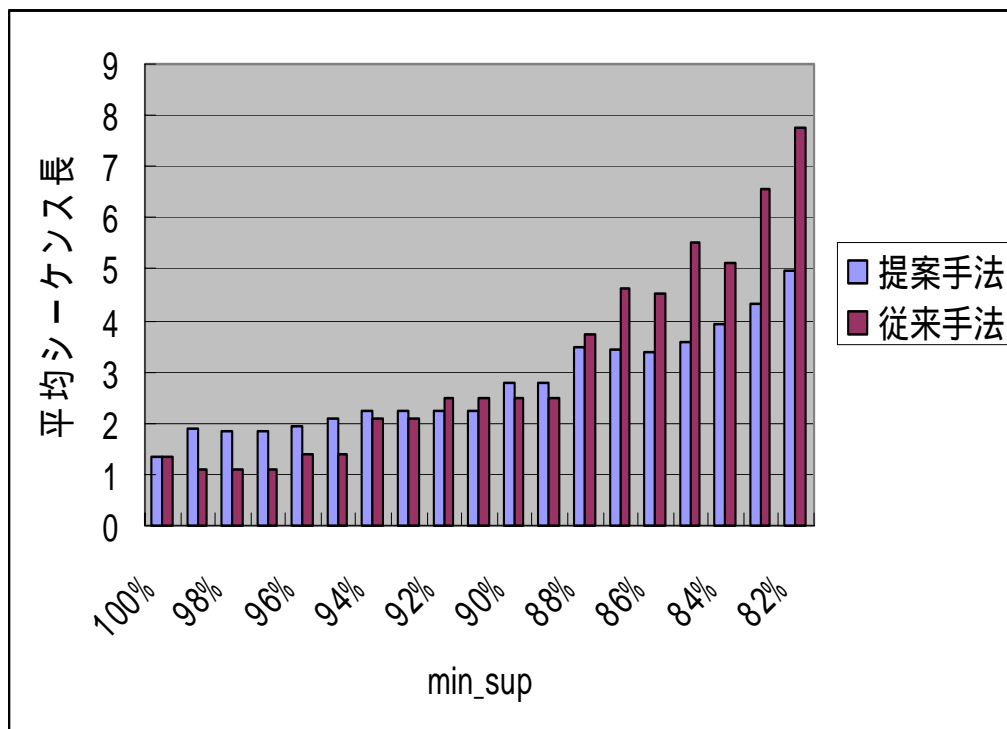


図5.4 各最小サポート値における頻出シーケンスの平均長

5.4 考察

提案手法のシーケンスの結合を行なうことによって起こる変化は以下の 2 つである。

① 従来手法では非頻出であったシーケンスが結合によって頻出になる

2 つのシーケンス $\langle 3\ 4 \rangle$ と $\langle (3\ 4) \rangle$ があつたとする。従来手法では当然 2 つのシーケンスは別々だが、前者のシーケンスにおいて、3 と 4 が結合されると、2 つとも $\langle (3\ 4) \rangle$ となる。

② 従来手法では頻出であったシーケンスが結合によって非頻出になる

従来手法ではシーケンス $\langle 1\ 2 \rangle$ は頻出であつたとする。しかし、提案手法を用いると、一部のシーケンス ID では $\langle (1\ 2) \rangle$ として扱われてしまうため、 $\langle 1\ 2 \rangle$ が頻出でなくなる。

表 5.4 と図 5.2 から分かるように、最小サポート値が高いときには、提案手法を用いた場合のほうが、従来手法に比べて抽出される頻出シーケンスが多くなっている。しかし、最小サポート値が下がるにつれて、提案手法を用いた場合のほうが、従来手法に比べ、抽出される頻出シーケンス数が少なくなっている。これは、①の現象は頻出でないものが結合によって頻出になる現象であり、①の現象によって頻出になるシーケンスは最小サポート値の高い時に起こりやすいといえる。また、②の現象は頻出であつたものが非頻出になる現象であるために、高いサポート値では非頻出になるために高いサポート値では提案手法のほうが抽出される頻出シーケンスが多くなっている。しかし、最小サポート値が下がると、②で非頻出であつたものも頻出になるため、従来手法で抽出される頻出シーケンス数が増大すると考えられる。

また、抽出される頻出シーケンス数が従来手法より少なくなったということは、従来手法では、同じ行動として扱われていたものを、提案手法を用いることにより、区別することができたということである。区別されるということはそれぞれのサポート値が下がるということであり、区別された結果の頻出シーケンスというものは従来手法で抽出されるシーケンスに比べ質の高いものであるといえることができる。

さらに、表 5.4 と図 5.1 が示すように、抽出される頻出シーケンスが少なくなった分、抽出時間も削減することができた。しかし、提案手法を用いた場合と用いない場合では、データの構造自体には変わりはない。そのため、表 5.4 と図 5.3 が示すように、抽出されるシーケンス 1 つあたりの抽出時間においては大きな差は見られなかった。

次に抽出された頻出シーケンスについて着目してみる。抽出される頻出シーケンスの長さは、表 5.5、表 5.6、図 5.4 からサポート値が高い時は、提案手法の方が従来手法に比べて長くなっていることが分かる。これは、先に述べたように、提案手法では最小サポート値が高いときは、結合によって頻出になるシーケンスが多い。結合されるということは、少なくとも長さ 2 以上であるので、頻出シーケンスの長さは長くなる。

また、頻出シーケンスのトランザクション数についてみると、提案手法はトランザク

ション同士を結合する手法であるので、抽出される頻出シーケンスは従来手法に比べて、” () “の数が少なくなることが予想される。56 ページから 59 ページに最小サポート値が 90 %の時に、提案手法によってマイニングし、抽出された頻出シーケンスと、従来手法によってマイニングを行い、抽出されたシーケンスを載せてある。これらを比較したときに、提案手法を用いた場合の方が、従来手法に比べて、頻出シーケンスの” () “の数が少なくなっていることが確認できた。例えば、提案手法を用いた場合、**(353 424 454 549)**といったように、” () “の少ない頻出シーケンスが多く見られる。サポート値 90%の場合、頻出シーケンスの平均の” () “の数は **1.00** であった。しかし、従来手法の場合、**(243 424) (383 491)**のように” () “の数が多い頻出シーケンスが多く見られるということである。最小サポート値 90%の時の頻出シーケンスの平均の” () “の数は **1.21** であった。最小サポート値を変えても、同様に提案手法を用いた場合の方が、” () “の少ない頻出シーケンスを多いことが確認できた。これは、当初の予想通り、提案手法は結合を行なっている手法であり、上述の 2 つの現象①、②においても、頻出になる場合は” () “がついている。そのため、提案手法を用いると、” () “の少ない頻出シーケンスが多く抽出されるということがいえる。

提案手法を用いてマイニング(**min_sup=90%**)を行い抽出された頻出シーケンスを以下に示す。(頻出シーケンス:サポート値)

(218) : 0.944839	(253 353 549) : 0.992900
(243) : 1.000000	(253 353 424 454) : 0.905516
(253) : 0.992900	(253 353 424 494) : 0.908247
(315) : 1.000000	(253 353 424 549) : 0.908247
(335) : 0.944839	(253 353 424 454 494) : 0.905516
(353) : 0.992900	(253 353 424 454 549) : 0.905516
(383) : 0.944839	(253 353 424 454 494 549) : 0.905516
(419) : 0.944839	(253 353 424 494 549) : 0.908247
(424) : 0.992900	(253 353 454 494) : 0.949208
(454) : 0.990169	(253 353 454 549) : 0.990169
(491) : 0.981431	(253 353 454 494 549) : 0.949208
(494) : 0.992900	(253 353 494 549) : 0.951939
(526) : 0.944839	(253 424 454) : 0.905516
(549) : 0.992900	(253 424 494) : 0.908247
(218 335) : 0.944839	(253 424 549) : 0.908247
(243 315) : 1.000000	(253 424 454 494) : 0.905516
(243 424) : 0.966685	(253 424 454 549) : 0.905516
(243 494) : 0.906062	(253 424 454 494 549) : 0.905516
(243 315 424) : 0.966685	(253 424 494 549) : 0.908247
(243 315 494) : 0.906062	(253 454 494) : 0.949208
(243 315 424 494) : 0.906062	(253 454 549) : 0.990169
(243 424 494) : 0.906062	(253 454 494 549) : 0.949208
(253 353) : 0.992900	(253 494 549) : 0.951939
(253 424) : 0.908247	(315 424) : 0.966685
(253 454) : 0.990169	(315 494) : 0.906062
(253 494) : 0.951939	(315 424 494) : 0.906062
(253 549) : 0.992900	(353 424) : 0.908247
(253 353 424) : 0.908247	(353 454) : 0.990169
(253 353 454) : 0.990169	(353 494) : 0.951939
(253 353 494) : 0.951939	(353 549) : 0.992900

(353 424 454) : 0.905516

(353 424 494) : 0.908247

(353 424 549) : 0.908247

(353 424 454 494) : 0.905516

(353 424 454 549) : 0.905516

(353 424 454 494 549) : 0.905516

(353 424 494 549) : 0.908247

(353 454 494) : 0.949208

(353 454 549) : 0.990169

(353 454 494 549) : 0.949208

(353 494 549) : 0.951939

(383 491) : 0.944839

(424 454) : 0.905516

(424 494) : 0.932277

(424 549) : 0.908247

(424 454 494) : 0.905516

(424 454 549) : 0.905516

(424 454 494 549) : 0.905516

(424 494 549) : 0.908247

(454 494) : 0.949208

(454 549) : 0.990169

(454 494 549) : 0.949208

(494 549) : 0.951939

提案手法を用いずに（従来手法で）マイニング(**min_sup=90%**)を行なった際に抽出された頻出シーケンスを以下に示す。（頻出シーケンス:サポート値）

(218) : 0.944839	(243 315 424) : 0.966667
(243) : 1.000000	(243 315) (491) : 0.944809
(253) : 0.992900	(243 315) (526) : 0.944262
(315) : 1.000000	(243 315) (218 335) : 0.944809
(335) : 0.944839	(243 315) (383 491) : 0.944809
(353) : 0.992900	(243 315 424) (218) : 0.921311
(383) : 0.944839	(243 315 424) (335) : 0.921311
(419) : 0.944839	(243 315 424) (383) : 0.921311
(424) : 0.992900	(243 315 424) (419) : 0.921311
(454) : 0.990169	(243 315 424) (491) : 0.921311
(491) : 0.981431	(243 315 424) (526) : 0.920765
(494) : 0.992900	(243 315 424) (218 335) : 0.921311
(526) : 0.944839	(243 315 424) (383 491) : 0.921311
(549) : 0.992900	(243) (383 491) : 0.944809
(218 335) : 0.944839	(243 424) (218) : 0.921311
(243) (218) : 0.944809	(243 424) (335) : 0.921311
(243 315) : 1.000000	(243 424) (383) : 0.921311
(243) (335) : 0.944809	(243 424) (419) : 0.921311
(243) (383) : 0.944809	(243 424) (491) : 0.921311
(243) (419) : 0.944809	(243 424) (526) : 0.920765
(243 424) : 0.966667	(243 424) (218 335) : 0.921311
(243) (491) : 0.944809	(243 424) (383 491) : 0.921311
(243) (526) : 0.944262	(253) (218) : 0.944809
(243) (218 335) : 0.944809	(253) (335) : 0.944809
(243 315) (218) : 0.944809	(253) (419) : 0.944809
(243 315) (335) : 0.944809	(253) (526) : 0.944262
(243 315) (383) : 0.944809	(253) (218 335) : 0.944809
(243 315) (419) : 0.944809	(315) (218) : 0.944809

(315) (335) : 0.944809
(315) (383) : 0.944809
(315) (419) : 0.944809
(315 424) : 0.966667
(315) (491) : 0.944809
(315) (526) : 0.944262
(315) (218 335) : 0.944809
(315) (383 491) : 0.944809
(315 424) (218) : 0.921311
(315 424) (335) : 0.921311
(315 424) (383) : 0.921311
(315 424) (419) : 0.921311
(315 424) (491) : 0.921311
(315 424) (526) : 0.920765
(315 424) (218 335) : 0.921311
(315 424) (383 491) : 0.921311
(353) (526) : 0.944262
(383 491) : 0.944809
(424) (218) : 0.944809
(424) (335) : 0.944809

(424) (383) : 0.944809
(424) (419) : 0.944809
(424) (491) : 0.944809
(424) (526) : 0.944262
(424) (218 335) : 0.944809
(424) (383 491) : 0.944809
(454) (526) : 0.920765
(494) (218) : 0.944809
(494) (335) : 0.944809
(494) (383) : 0.944809
(494) (419) : 0.944809
(494) (491) : 0.944809
(494) (526) : 0.944262
(494) (218 335) : 0.944809
(494) (383 491) : 0.944809
(549) (218) : 0.944809
(549) (335) : 0.944809
(549) (419) : 0.944809
(549) (526) : 0.944262
(549) (218 335) : 0.944809

第6章 おわりに

本論文では、ユーザが定義した制限距離（時間）以下のトランザクションは同時に起こったとみなして、トランザクションの結合を行ない、得られたデータベースに対してマイニングを行なう手法を提案した。その結果、従来手法と比較した時、頻出シーケンスの長さについては、最小サポート値 **90%** で頻出シーケンスの平均長が、提案手法では **2.771**、従来手法では **2.469** となった。また、トランザクション数については、最小サポート値 **90%** で、提案手法では平均トランザクション数が **1.00**、従来手法では **1.21** となり、質の異なる頻出シーケンスを抽出することができたといえる。つまり、短期的な行動と長期的な行動を区別してマイニングを行なうことができたといえる。また、抽出時間については、**1** シーケンスあたりの抽出時間には変化はないが、全体の抽出時間については最小サポート値 **90%** で提案手法では **0.12** 秒、従来手法では **0.18** 秒となり、削減することができた。

今後、提案手法の拡張を考えた場合、長期的行動に着目して新たな処理を行なう方法が考えられる。具体的には、トランザクション間の距離（時間）が、ユーザが定義した制限距離（制限時間）より離れている場合は、相関関係がないと判断し、サポート値にカウントしないようにする。これにより、短期的な行動と長期的な行動を、さらに厳密に区別することが可能になる。そして、ユーザにより詳細な行動の情報を与えることができると考えられる。

さらに、提案手法ならびに拡張手法を実際に実世界の例に応用するということも考えられる。

謝辞

本研究を行なうにあたり、数々のご指導を頂いた山名早人助教授に深く感謝致します。また、様々な面でお世話になった斉田直幸先輩、平手勇宇先輩、並びにユビキタス班の方々に厚く御礼申し上げます。

参考文献

- [1] J.Pei, “Frequent Pattern Mining”, <http://www.cs.buffalo.edu/faculty/jianpei/>
- [2] R.Agrawal and R.Srikant. “Mining Sequential Patterns”. In *Proc. of the 11th international conference on Data Engineering*, Mar 1995.
- [3] R.Srikant and R.Agrawal ,” Mining Sequential Patterns:Generalizations and Performance Improvements”In *Proc. of 5th Int. Conf. Extending Database Technology*, 1996
- [4] J.Pei, J.Han, B.Mortazavi-Asl, H.Pinto, Q.Chen, U.Dayal and M.Hsu, “PrefixSpan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth” In. *Proc. of 2001 Int. Conf. Data Engineering (ICDE'01)*, pp. 215-224 , April 2001
- [5] M.J.Zaki, “SPADE: An efficient algorithm for mining sequences”,In *Proc. of Machine Learning*,42(1/2):31-60, 2001
- [6] Ayres J, Gehrke J, Yu T and Flannick J, “Sequential PAttern Mining using Bitmap Represetation”, In *Proc. of ICDE* ,2001
- [7] 福田剛志 森本康彦 徳山豪 “データマイニング” 共立出版,2001
- [8] D.Chiu, Y.Wu and A.Chen,”An Efficient Algorithm for Mining Frequent Sequences by a New Strategy without Support Counting”In *Proc. of ICDE*, 2004
- [9] J. Han, <http://www-sal.cs.uiuc.edu/~hanj/>
- [10] “Musashi 人工データ”, <http://musashi.sourceforge.jp/artificialData/index.html>
- [11] 有村博紀 , “ストリームデータマイニングに関する研究動向”, <http://www.i.kyushu-u.ac.jp/~arim/jtalks/arimura-dews04-mini-survey.ppt>
- [12] T.Asai, K.Abe, S.Kawasoe, H.Arimura, H.Sakamoto and S.Arikawa. “Online Algorithms for Mining Semistructured Data Stream” In *Proc. of DOI Technical Report*, June 2002
- [13] Q.Zheng, K.Xu and S.Ma “When to Update the Sequential Patterns of Stream Data” In *Proc. of 7th Pacific-Asia Conference on Knowledge Discovery and Data Mining*,2003