



# HOKKAIDO UNIVERSITY

Title	分散データベースからの 頻出飽和アイテム集合のプライバシー保護発見
Author(s)	山本, 章博
Description	ERATO 세미나2010 : No.12. 2010年8月6日
Relation	2010年度科学技術振興機構ERATO湊離散構造処理系プロジェクト講究録. p. 79-90.
Issue Date	2011-06
Doc URL	<a href="https://hdl.handle.net/2115/48473">https://hdl.handle.net/2115/48473</a>
Type	conference presentation
File Information	12_all.pdf



ERATO セミナ 2010 - No. 12

# 分散データベースからの $k$ 頻出飽和アイテム集合の プライバシー保護発見

山本章博

京都大学 大学院情報学研究科

2010/8/6

## 概要

水平分割され分散配置されたトランザクションデータベースに、全体として頻出する飽和アイテム集合を暗号化と組み合わせて発見する手続きを提案する。頻出飽和アイテム集合は、頻出アイテム集合の圧縮形とみなせることから、分散配置されたデータベースからのプライバシー保護発見に有用と考えられる。本講では基本アイデアを実験結果とともに報告する。

## 分散データベースからの 頻出飽和アイテム集合の プライバシー保護発見

06 Aug. 2010 ERATO セミナ

京都大学情報学研究科

久野慎弥\* 土井晃一郎\*\* 山本章博

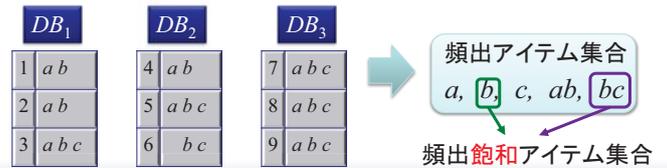
\*現在 野村総合研究所

\*\*現在 東京大学新領域創成科学研究科

### 概要

問題: 大規模データベースから興味深いパターンを発見すること

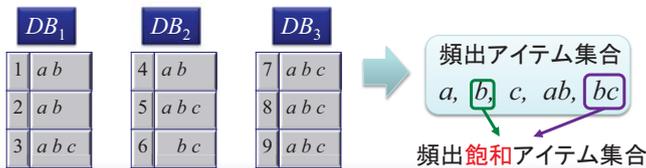
- “データベース”...分散配置された**トランザクションDB**
  - ・トランザクションDB: **バスケット分析** (買い物かごや取引レコードで**アイテム**購入の相関ルールを発見する問題) で用いられるモデル
- “興味深いパターン”...出現**頻度**の高いアイテム集合 (頻出アイテム集合) の中で代表的なもの (**頻出飽和アイテム集合**)



### 概要

問題: 大規模データベースから興味深いパターンを発見すること

- “データベース”...分散配置された**トランザクションDB**
  - ・トランザクションDB: **バスケット分析** (買い物かごや取引レコードで**アイテム**購入の相関ルールを発見する問題) で用いられるモデル
- 本研究では、分散配置されたトランザクションDBから**頻出飽和アイテム集合**を発見することを目標とする。



### 概要

問題: 大規模データベースから興味深いパターンを発見すること

- ・分散配置されたデータベースから、プライバシー保護のため**暗号化**を用いて**頻出飽和アイテム集合**をすべて発見する手続き**SFC-Combine**を構築
  - 通信を考慮しない手続き**FC-Combine**を構築し、これを基にする
- ・従来研究の手続きとの性能比較を行い、従来より**通信コスト**を抑えられていることを確認



### 発表の概要

- ・ 本研究の背景
- ・ 準備
  - 頻出アイテム集合
  - 頻出飽和アイテム集合
  - 問題例
- ・ マイニングの手続き
  - FC-Combine
  - SFC-Combine
- ・ 性能比較
- ・ 今後の課題

### 研究の背景

- ・ ネットワーク上に分散されたデータベースから興味深いパターンを発見する手法に近年注目が集まっているが、幾つかの問題点が指摘されている。
  - 分散データベースに蓄えられているデータ量の増大ペースに対して、ネットワーク帯域のコストパフォーマンス改善が追いついていない。[SCS2009 Konishi]
- **通信コストを削減して全体として蓄えられたデータを分析する発見する手法の必要性**
- データベースから個人情報漏洩することの危険性が大きい、個々のデータベースを死蔵させることの損失もまた大きい。[Kantarcioglu et al.][Verykios et al.]
- **プライバシーを保護してデータマイニングを行う手法の必要性**

## 頻出アイテム集合発見問題

- 頻出アイテム集合発見問題・・・データマイニング研究の発展の契機となった問題

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

### 頻出アイテム集合発見問題

与えられたトランザクションデータベースDBと最小支持度minsupに対して、出現頻度がminsup以上であるアイテム集合を全て発見せよ。

$$\text{アイテム集合} X \text{の出現頻度} = \frac{X \text{が出現するトランザクションの数}}{\text{DBのトランザクションの数}}$$

## 頻出アイテム集合発見問題

- 頻出アイテム集合発見問題・・・データマイニング研究の発展の契機となった問題

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

頻出アイテム集合の例:

a (sup: 60%)      b (sup: 80%)  
 ac (sup: 60%)    bc (sup: 60%)  
 be (sup: 80%)    ce (sup: 60%)  
 bce (sup: 60%)  
 ...  
 ( $\{i_1, i_2, \dots, i_k\}$  を  $i_1i_2\dots i_k$  と表す)

## 頻出アイテム集合発見問題

- 頻出アイテム集合発見問題・・・データマイニング研究の発展の契機となった問題

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

頻出アイテム集合の例:

a (sup: 60%)      b (sup: 80%)  
 アイテムの種類に対して指数オーダー数の解が存在する  
 →出力を絞り込むことが必要  
 [Uno et al. 05]

( $\{i_1, i_2, \dots, i_k\}$  を  $i_1i_2\dots i_k$  と表す)

## 頻出飽和アイテム集合

アイテム集合Cが飽和アイテム集合である



Cを含むトランザクション全てに共通して含まれるアイテム集合がC自身である。

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

頻出飽和アイテム集合: c, ac, be, bce  
 頻出であるが飽和でない例: a, bc, ...

## 頻出飽和アイテム集合

アイテム集合Cが飽和アイテム集合である



Cを含むトランザクション全てに共通して含まれるアイテム集合がC自身である。

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

頻出飽和アイテム集合: c, ac, be, bce  
 頻出であるが飽和でない例: a, bc, ...

## 頻出飽和アイテム集合

アイテム集合Cが飽和アイテム集合である



Cを含むトランザクション全てに共通して含まれるアイテム集合がC自身である。

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

頻出飽和アイテム集合: c, ac, be, bce  
 頻出であるが飽和でない例: a, bc, ...

## 頻出飽和アイテム集合

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

アイテム集合Cが**飽和アイテム集合**である



Cを含むトランザクション全てに共通して含まれるアイテム集合がC自身である。

頻出飽和アイテム集合: c, ac, be, bce

頻出であるが飽和でない例: a, bc, ...

## 頻出飽和アイテム集合

DB	
1	a c d
2	b c e
3	a b c e
4	b e
5	a b c e

minsup = 50%

アイテム集合Cが**飽和アイテム集合**である



Cを含むトランザクション全てに共通して含まれるアイテム集合がC自身である。

頻出飽和アイテム集合: c, ac, be, bce

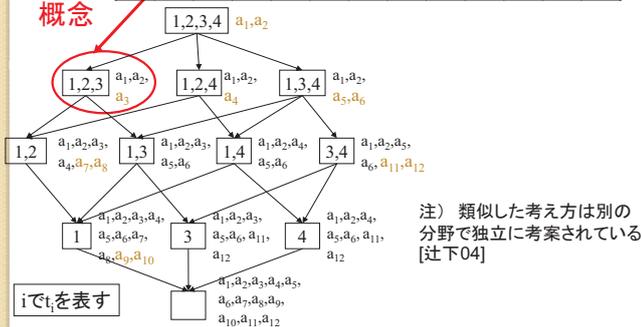
頻出であるが飽和でない例: a, bc, ...

**補題:** すべての頻出アイテム集合は**頻出飽和アイテム集合**から求めることができる。

## Formal Concept Analysis [Wille 84]

	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>	a <sub>4</sub>	a <sub>5</sub>	a <sub>6</sub>	a <sub>7</sub>	a <sub>8</sub>	a <sub>9</sub>	a <sub>10</sub>	a <sub>11</sub>	a <sub>12</sub>
t <sub>1</sub>	●	●	●	●	●	●	●	●	●	●		
t <sub>2</sub>	●	●	●	●	●	●	●	●				
t <sub>3</sub>	●	●	●	●	●	●	●	●			●	●
t <sub>4</sub>	●	●	●	●	●	●	●	●			●	●

概念



## 用いるモデル

- アイテム全体の集合をIとする。
- アイテムの集合とトランザクションIDの組をトランザクションとよび、トランザクションの集合を**トランザクションDB**とよぶ。
- 分散DBのモデルとして、 $DB_i \cap DB_j = \emptyset$  なるM個のトランザクションDB,  $DB_1, DB_2, \dots, DB_N$  を用いる。
- $DB_1, DB_2, \dots, DB_N$  を**部分DB**とよび、 $DB_i$ を管理するサーバをサイトiとよぶ。
- $DB_{1\dots N} = DB_1 \cup \dots \cup DB_N$ を**全体DB**とよぶ。

## 問題例

	DB <sub>1</sub>	DB <sub>2</sub>	DB <sub>3</sub>
1	a b	4 a	7 a b c
2	a b	5 a	8 a b c
3	a c	6 b c	9 a b c

minsup=50%

全体DBの頻出飽和アイテム集合を全て発見せよ。

## 飽和アイテム集合に注目する理由

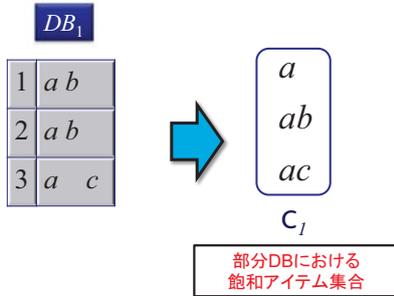
- アイテムの種類に対して指数オーダーとなる**頻出アイテム集合の数を抑えられる**。 [Uno et al. 05]
- 集合族間の二項演算子  $\hat{\otimes}$  (後述) で、部分DBの飽和アイテム集合から、**全体DBの飽和アイテム集合**を求めることができる。 [Lucchese et al. 06]
- 全体DBの頻出飽和アイテム集合は、**少なくとも一つの部分DBで頻出な飽和アイテム集合の部分集合**である。



サイト間の通信路における**通信コスト**を抑えた発見手続きを構築できることが期待される。

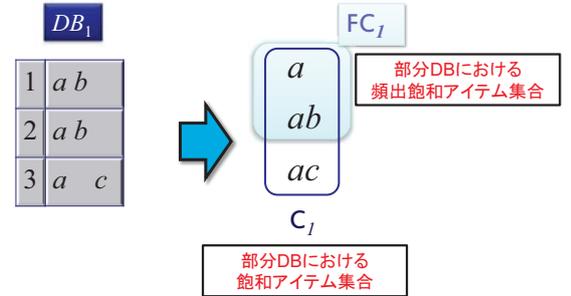
### 用いる記号

minsup=50%



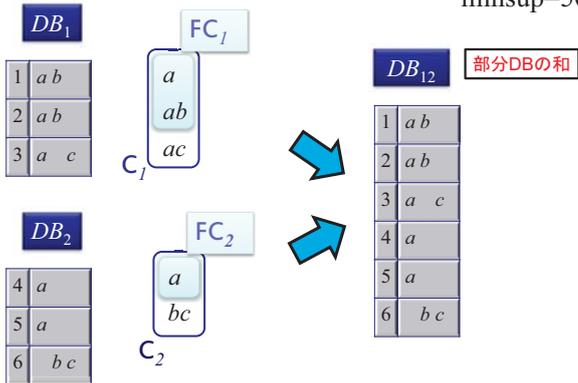
### 用いる記号

minsup=50%



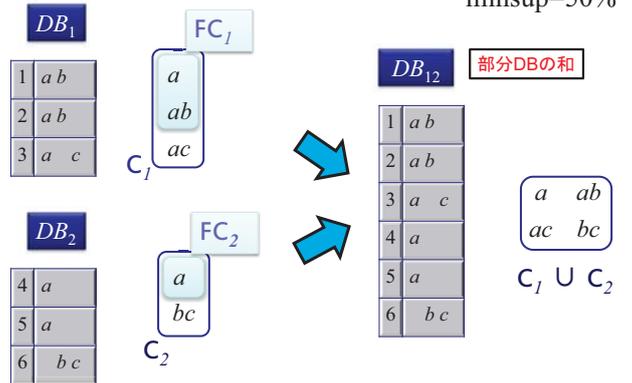
### 部分DBの和

minsup=50%



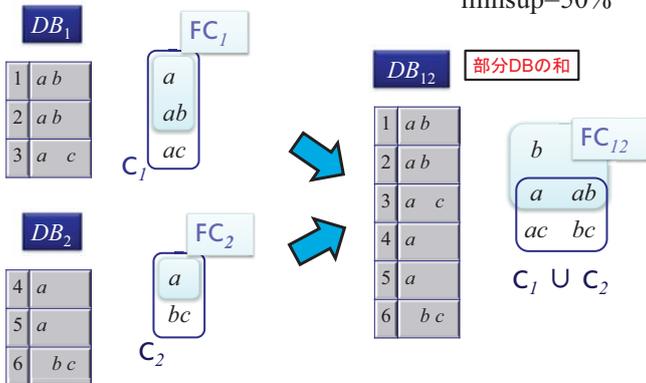
### 部分DBの和

minsup=50%



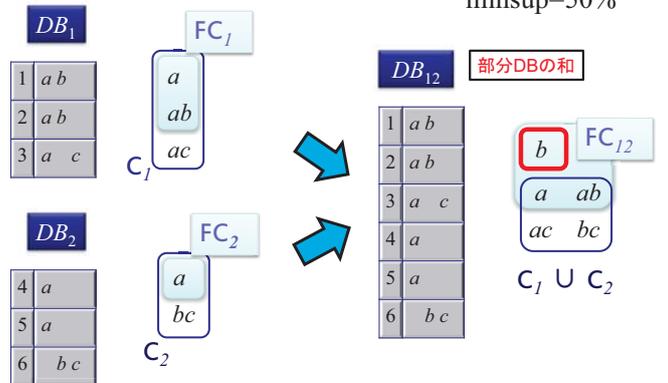
### 部分DBの和

minsup=50%



### 部分DBの和

minsup=50%



### 部分DBの和における飽和アイテム集合

Iをアイテム全体の集合とする.

$S_1, S_2 \subseteq 2^I$ に対して

$$S_1 \otimes S_2 = \{X \cap Y \mid X \in S_1, Y \in S_2\}$$

$$S_1 \hat{\otimes} S_2 = S_1 \cup S_2 \cup (S_1 \otimes S_2) \text{ と定義する.}$$

補題:  $S_1, S_2, S_3 \subseteq 2^I$ に対して

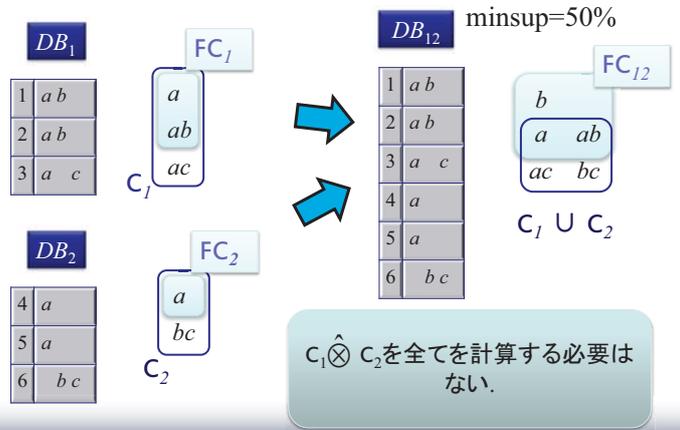
$$(S_1 \hat{\otimes} S_2) \hat{\otimes} S_3 = S_1 \hat{\otimes} (S_2 \hat{\otimes} S_3) \text{ が成り立つ.}$$

$$S_1 \hat{\otimes} S_2 \hat{\otimes} S_3 = (S_1 \hat{\otimes} S_2) \hat{\otimes} S_3 \text{ とおく.}$$

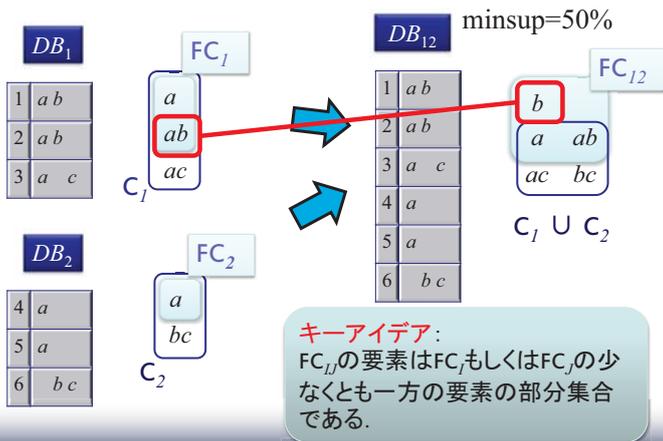
定理[Lucchese et al.]:  $DB_1, \dots, DB_N$ に対して

$$C_{1..N} = C_1 \hat{\otimes} C_2 \hat{\otimes} \dots \hat{\otimes} C_N \text{ が成り立つ.}$$

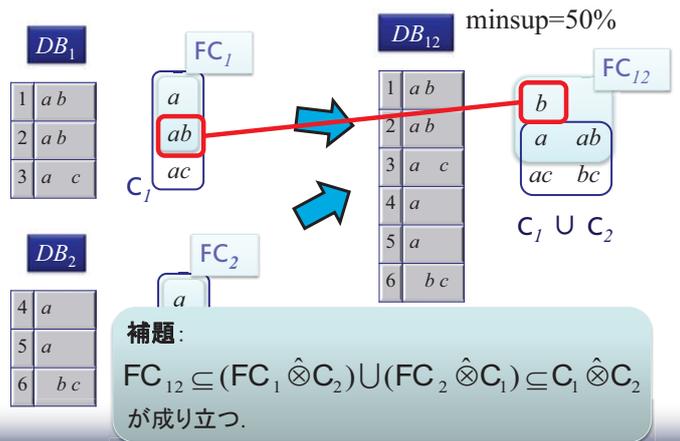
### 手続きの方針



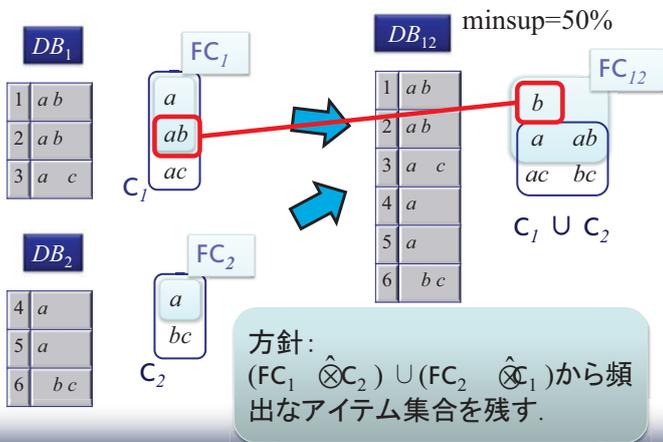
### 手続きの方針



### 手続きの方針



### 手続きの方針



### 手続き FC-Combine<sub>0</sub> の概要

**FC-Combine<sub>0</sub>** (FC-Combineの部分手続き):

$I, J \subseteq \{1, \dots, N\}$ とする.

- 入力  
 $(DB_I \text{から}) FC_I \cup_{k \in I} C_k$  ( $DB_J \text{から}) FC_J, \cup_{k \in J} C_k$
- 出力  
 $FC_{IJ}, (\cup_{k \in I} C_k) \cup (\cup_{k \in J} C_k)$
- 手続き:
  1.  $S_I := \text{ClosedsetSum}(FC_I, \cup_{k \in J} C_k)$ ;
  2.  $S_J := \text{ClosedsetSum}(FC_J, \cup_{k \in I} C_k)$ ;
  3.  $FC_{IJ} := \text{freq}(S_I \cup S_J)$ ;
  4. return  $FC_{IJ}, (\cup_{k \in I} C_k) \cup (\cup_{k \in J} C_k)$

## 手続き FC-Combine<sub>0</sub> の概要

**FC-Combine<sub>0</sub>** (FC-Combineの部分手続き):

$I, J \subseteq \{1, \dots, N\}$

• 入力

$(D_1, \dots, D_N)$

• 出力

$(FC_I, FC_J)$

• 手続き

1.  $S_I := \text{ClosedsetSum}(FC_{I_1}, \dots, FC_{I_m})$

2.  $S_J := \text{ClosedsetSum}(FC_{J_1}, \dots, FC_{J_n})$

3.  $FC_{IJ} := \text{freq}(S_I \cup S_J)$

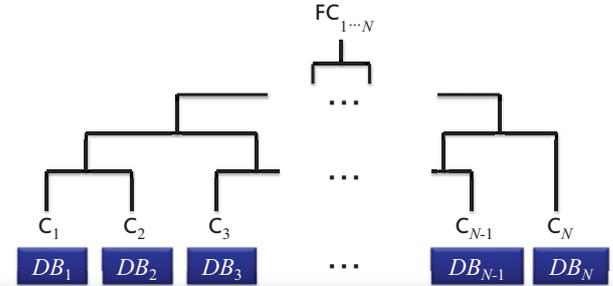
4. return  $FC_{IJ}, (\bigcup_{k \in I} UC_k) \cup (\bigcup_{k \in J} UC_k)$

### FC-Combine:

$FC_1, \dots, FC_N, C_1, \dots, C_N$ から  
トーナメント方式で  
FC-Combine<sub>0</sub>を行っていき、  
FC<sub>1...N</sub>を出力する手続き

## トーナメント方式

トーナメント方式で手続きを行っていき、FC<sub>1...N</sub>を出力する。  
(演算子の性質から、どのような順序で行ってもよい)



## 手続き SFC-Combine

- **SFC-Combine:** 手続きFC-Combineを暗号化を用いて行う。
- **Pohlig-Hellman encryption scheme**
  - ある大きい素数 $p$ をお互いのサイトで共有し、それぞれのサイトは $e \times d = 1 \pmod{p-1}$ となるような $e, d$ の組をランダムに選ぶ。
  - 平文 $M$ の暗号化:  $E(M) = M^e \pmod{p}$  (秘密鍵)
  - 暗号文 $C$ の複合化:  $D(C) = C^d \pmod{p}$
  - 補題:  $E_1(E_2(M)) = E_2(E_1(M))$  が成り立つ。

## 手続き SFC-Combine

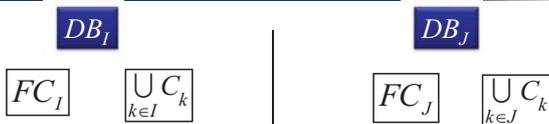
- **SFC-Combine:** 手続きFC-Combineを暗号化を用いて行う。

### SFC-Combine:

入力 $FC_I, FC_J, \bigcup_{k \in I} UC_k, \bigcup_{k \in J} UC_k$ 内のそれぞれの  
アイテムを互いの暗号鍵で暗号化し、  
手続きFC-Combineを行う手続き。

- 補題:  $E_I(E_J(M)) = E_J(E_I(M))$  が成り立つ。

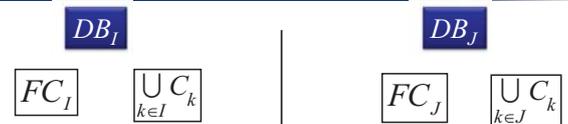
## 手続き SFC-Combine<sub>0</sub>



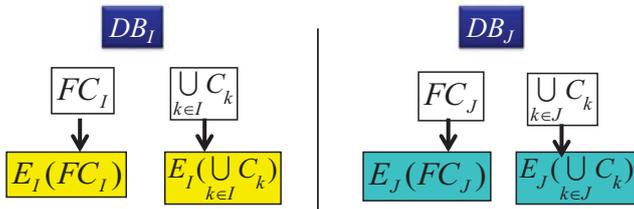
### SFC-Combine<sub>0</sub> :

暗号化を用いてこれらのアイテム集合の集合を  
隠匿しながら、FC-Combine<sub>0</sub>を行う手続き。

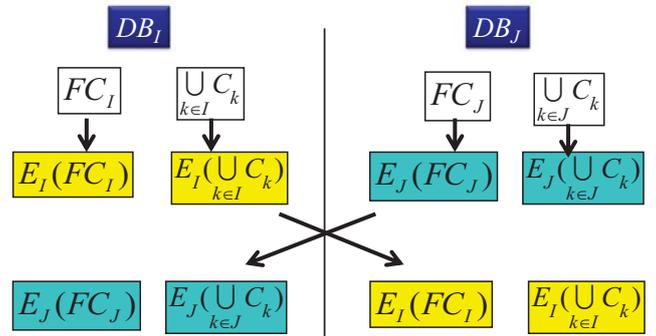
## 手続き SFC-Combine<sub>0</sub>



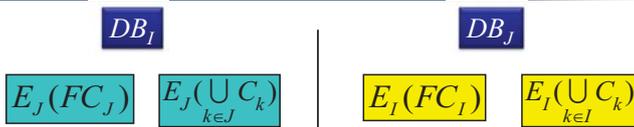
手続きSFC-Combine<sub>0</sub>



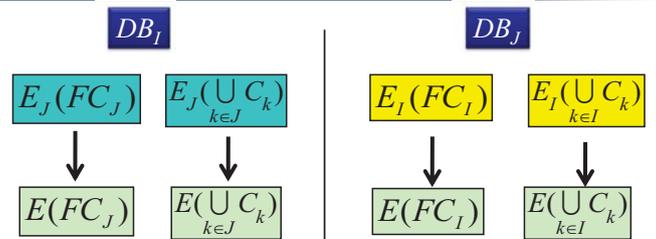
手続きSFC-Combine<sub>0</sub>



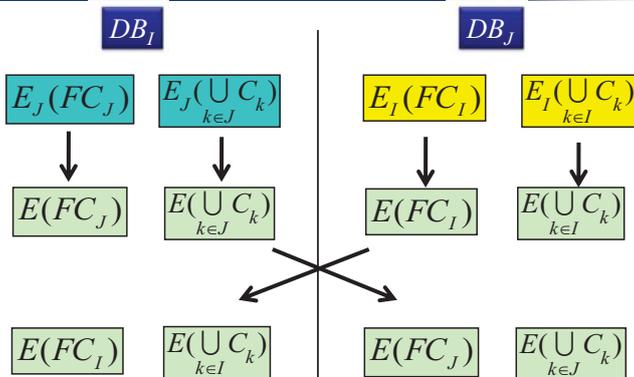
手続きSFC-Combine<sub>0</sub>



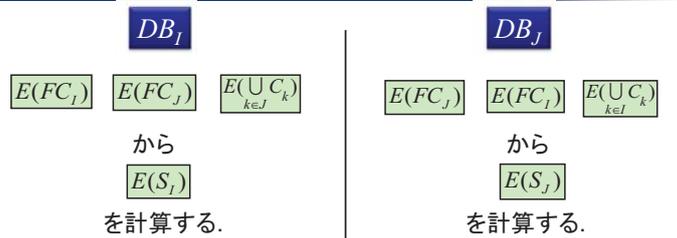
手続きSFC-Combine<sub>0</sub>



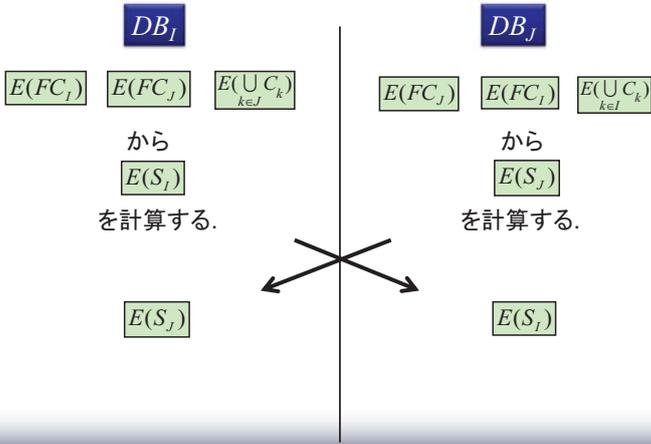
手続きSFC-Combine<sub>0</sub>



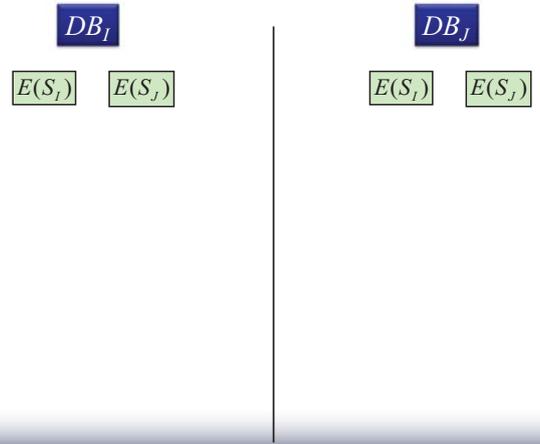
手続きSFC-Combine<sub>0</sub>



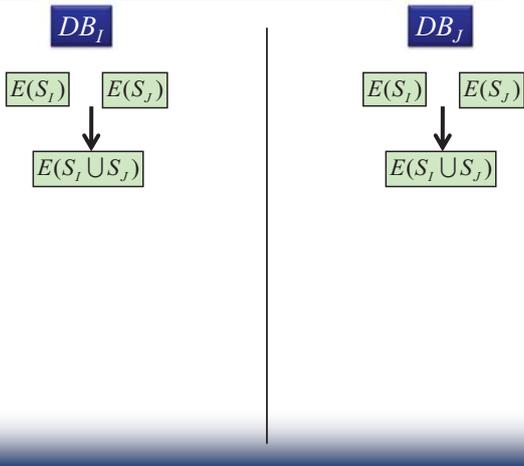
手続きSFC-Combine<sub>0</sub>



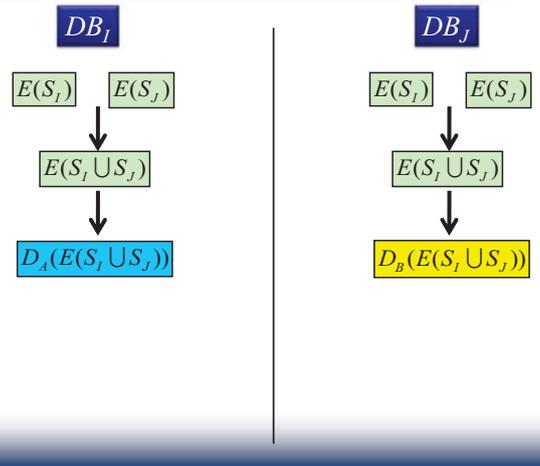
手続きSFC-Combine<sub>0</sub>



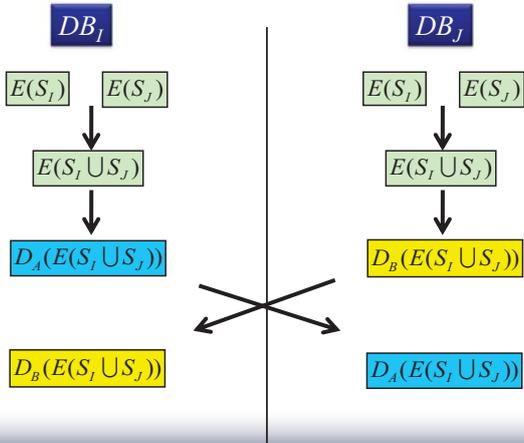
手続きSFC-Combine<sub>0</sub>



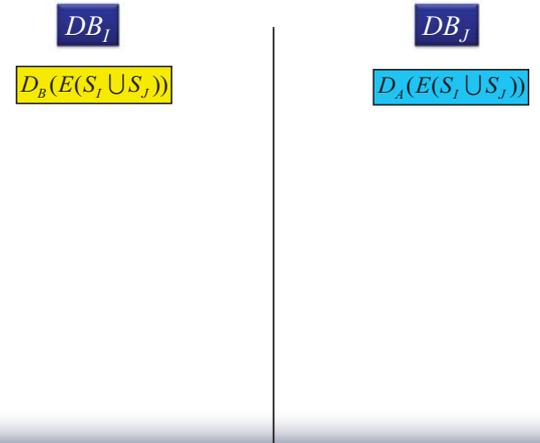
手続きSFC-Combine<sub>0</sub>



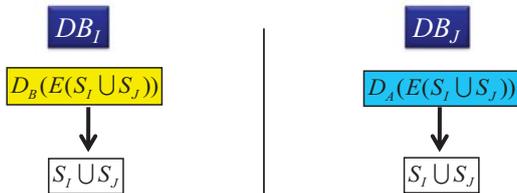
手続きSFC-Combine<sub>0</sub>



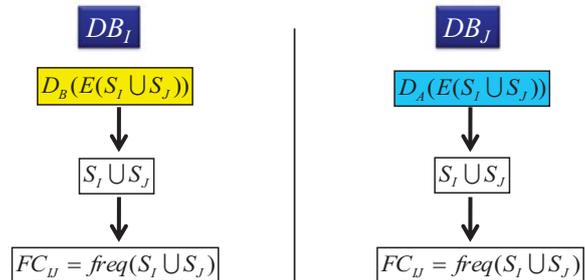
手続きSFC-Combine<sub>0</sub>



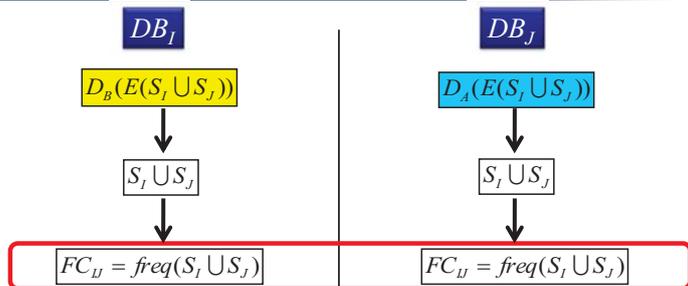
### 手続きSFC-Combine<sub>0</sub>



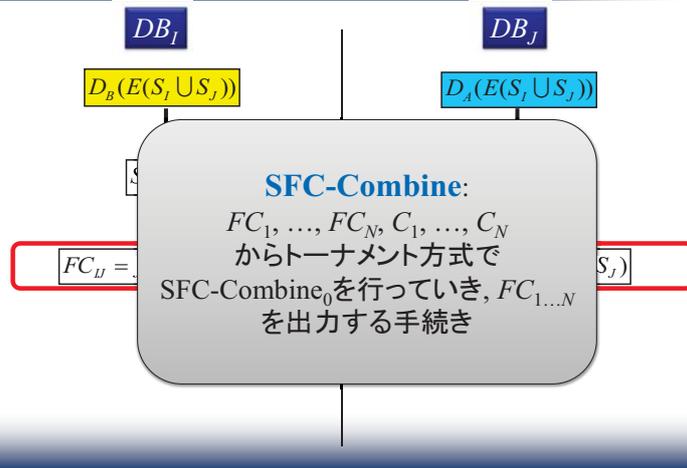
### 手続きSFC-Combine<sub>0</sub>



### 手続きSFC-Combine<sub>0</sub>



### 手続きSFC-Combine<sub>0</sub>



### 安全性

それぞれのDBがプロトコルで決められた振る舞いから逸脱しないことを仮定する。  
出力以外に漏れる情報(演算時):

SFC-Combine	(1) $C_1 \otimes C_2$ の演算時, 任意の $X \in C_1, Y \in C_2$ に対する $ X \cap Y $ の大きさ (2) $C_1 \cup C_2$ の演算時, $C_1 \cap C_2$ の大きさ
SFDM	

### 安全性の比較

	SFC-Combine (本研究)	SFDM (従来研究)	マルチパーティ プロトコル
他サイトと共通しない アイテム	○(隠匿できる)	○	○
他サイトと共通する アイテム		○	○
他サイトと共通しない アイテム集合	○	○	○
他サイトと共通する アイテム集合			○

SFC-Combineがアイテムごとに暗号化を行っているのに対して、SFDMはアイテム集合ごとに暗号化を行っているため

## 性能評価

- 以下の3つの手法に対して、通信コストおよび安全性について比較
  - 本提案 **SFC-Combine**
  - 分散DBから暗号化を用いて**頻出アイテム集合**を求める手続き **SFDM** [Kantarcioglu et al. 03]
  - マルチパーティ・プロトコル**を用いて頻出飽和アイテム集合を求める手続き(本研究で一部を構築)
    - マルチパーティ・プロトコル[Goldreich 87]: 各参加者  $1, 2, \dots, N$  が、秘密情報  $x_1, x_2, \dots, x_N$  を隠したまま、関数  $f(x_1, \dots, x_N)$  を計算するプロトコル

## 通信コスト(分析)

- 《分析1》全ての通信で発生するメッセージの総数
- 《分析2》同時発生する通信のメッセージを同一とみなした場合のメッセージの総数

	分析1	分析2
<b>SFC-Combine</b>	$10N-10$	$5 \times \lceil \log_2 N \rceil$
<b>SFDM</b>	$N$ (偶数の場合: $11N-12$ )	$\lceil \log_2 N \rceil$

1回の手続きが送られて、手続きが...  
1回の手続きで5回メッセージが送られており、トーナメントの高さは  $\lceil \log_2 N \rceil$  である。

## 通信コストの比較

	<b>SFC-Combine</b> (本研究)	<b>SFDM</b> (従来研究)	<b>マルチパーティ・プロトコル</b>
<b>通信コスト (bit)</b>	$O(N \cdot  C_{12\dots N} )$	$O(N^2 \cdot (\sum_{k=1}^{m+1}  Cand_{(k)} ))$	少なくとも $O(N^2 \cdot w^N)$

**m**: 全体DBにおける頻出アイテム集合の大きさの最大値

**Cand<sub>(k)</sub>**: Aprioriアルゴリズム [Agrawal et al. 94]を用いてつくられる、大きさkの候補アイテム集合全体。

**w**:  $|C_1|, |C_2|, \dots, |C_N|$ のいずれより大きな数

## 通信コストの比較

	<b>SFC-Combine</b> (本研究)	<b>SFDM</b> (従来研究)	<b>マルチパーティ・プロトコル</b>
<b>通信コスト (bit)</b>	$O(N \cdot  C_{12\dots N} )$	$O(N^2 \cdot (\sum_{k=1}^{m+1}  Cand_{(k)} ))$	少なくとも $O(N^2 \cdot w^N)$

$|Cand_{(k)}|$ は、ユーザが与える閾値の**最小支持度が小さいとき**、爆発的に増加する性質があることが知られている [Uno et al. 03].  
→ 最小支持度が小さいとき、本研究は従来手法と比べて通信コストを抑えられることが期待される。

## 通信コストの比較(実装)

それぞれの手続きにおいて通信するアイテム集合の総数

mushroom.dat N = 4	minsup = 1%	2%	5%	10%	20%	50%
	<b>SFC-Combine</b>	679,238	493,334	321,320	243,030	207,962
<b>SFDM</b>	721,325,038	177,450,024	27,312,860	4,603,279	432,581	1,240

Frequent Itemset Mining Implementations Repository : <http://fimi.cs.helsinki.fi/data/>  
Intel Core 2 Duo 2.33GHz, 2GB RAM, Windows Vista Business SP1

**minsupが小さい**とき、本提案は通信コストが小さい  
→ 多くの興味深いパターンを小さい通信コストで発見できる

## 通信コストの比較(実装)

それぞれの手続きにおいて通信するアイテム集合の総数

mushroom.dat N = 4	minsup = 1%	2%	5%	10%	20%	50%
	<b>SFC-Combine</b>	679,238	493,334	321,320	243,030	207,962
<b>SFDM</b>	721,325,038	177,450,024	27,312,860	4,603,279	432,581	1,240

Frequent Itemset Mining Implementations Repository : <http://fimi.cs.helsinki.fi/data/>  
Intel Core 2 Duo 2.33GHz, 2GB RAM, Windows Vista Business SP1

**minsupが大きい**とき、本提案は通信コストが大きい  
→ 暗号化のため、手続きの最初と最後に大きなコストがかかる

## 安全性の比較

	SFC-Combine (本研究)	SFDM (従来研究)	マルチパーティ・ プロトコル
手続きの進行	2者間の手続きによる トーナメント方式	全体並列方式	全体並列方式

2者間の手続きでは、手続き中に漏れてしまう情報の他にも、Outputからも情報が漏れてしまう可能性がある[Kantarcioglu et al. 04].

→例えば、ある  $X \in FC_{i,j}$  に対し、ある  $i \in X$  がサイト  $I$  で頻出でない場合、 $i$  がサイト  $J$  で頻出であることを、サイト  $I$  は知ることができる。

## まとめ・今後の課題

- まとめ
  - 分散データベースから暗号化手法を用いて頻出飽和アイテム集合を全て発見する手続きSFC-Combineを構築した
  - 従来手法より通信コストを抑えられていることを確認した
- 今後の課題
  - 複数のアイテムをまとめた上での暗号化
    - 集合演算  $\otimes$  に対応できるような暗号化が必要
  - 全体並列方式により進行する手続きの構築