



列挙学校：第3コマ パターンマイニングにおける 列挙アルゴリズム

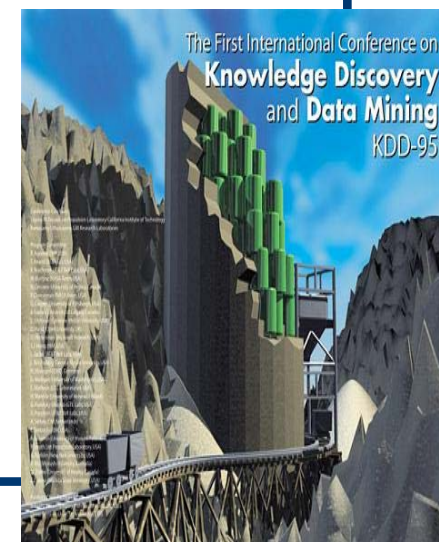
有村 博紀

北海道大学大学院 情報科学研究科
コンピュータサイエンス専攻

<http://www-ikn.ist.hokudai.ac.jp/~arim>

3コマ目の目的: データマイニングを題材に

- 列挙はどんなことに使えそうか？
- 応用分野では, 列挙はどのように使われているか？
- データマイニングと機械学習における列挙に関連した話題を提供



パート1: データマイニングと頻出集合発見

- データマイニング
- 頻出集合マイニング

パート2: 系列とグラフのマイニング

- 系列とグラフのマイニング概観
- 列挙によるパターンのマイニング
- 列挙とマイニングに関連した話題

列挙学校：第3コマ

パート1: データマイニングと頻出集合発見

有村 博紀

北海道大学大学院 情報科学研究科 コンピュータサイエンス専攻

email: {arim,kida}@ist.hokudai.ac.jp

<http://www-ikn.ist.hokudai.ac.jp/~arim>

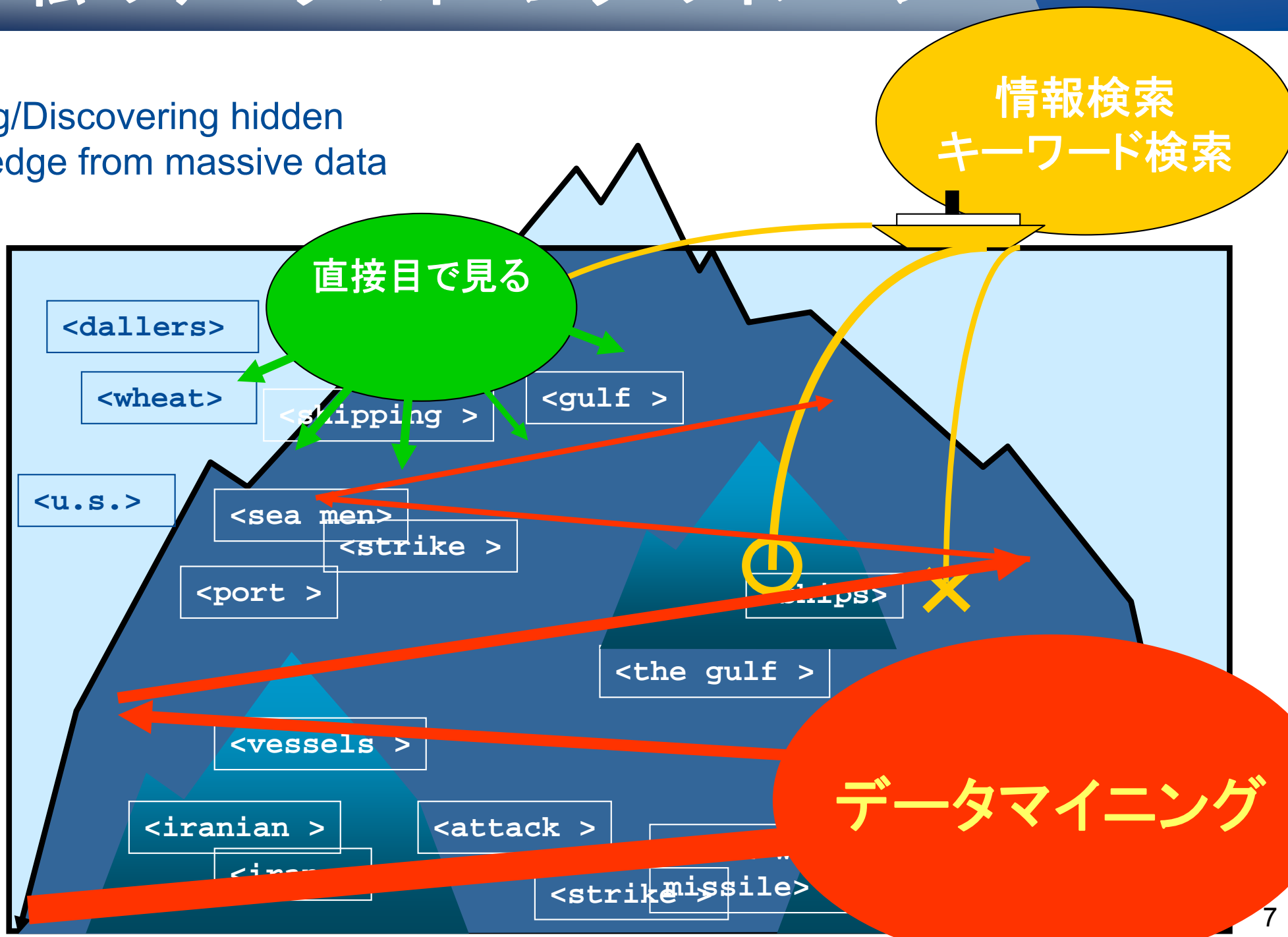
データマイニング

- 大量のデータから人間にとって**有用なパターン**や**規則**を半自動的にとりだす方法の研究
- 1990年代半ばから研究が盛んになった
 - Apriori algorithm [Agrawal, Srikant, VLDB1994]
- 潜在的には古くからある研究の集大成.
 - ただし, 大量データに対する効率的計算に重点
- 機械学習・数理統計学・データベース技術の境界分野

データマイニングのプロセスの全体

- 1.対象領域の理解
- 2.データ集合の前処理
- 3.パターンの発見(狭義のデータマイニング)
- 4.得られたパターンの解析
- 5.解析結果の利用

Finding/Discovering hidden knowledge from massive data



4テラバイト = 4,000,000,000,000バイト

- 400字詰めA4原稿用紙で309km²
= 山手線の内側の約3倍の面積
- 1秒間に10文字入力すると12,683年
- CD 6153枚 = 音楽として聴くと316日間
- 記憶装置からの読み込みに 15.2時間
- ちょっとした計算(N²時間)
2.2 GFLOPS × 32PE のスーパーコンピュータで
800万年以上かかる。
- 高速なアルゴリズムが必要！！！！

パターン発見

- トランザクションデータから共通して出現する規則性を発見する
- 頻出パターン発見 [Agrawal et al. '94]
- 最適化マイニング [森下 '96, '98, '00]

予測学習・自動分類

- 不完全なデータから、未知の規則を学習する
- SVM [Vapnik '96],
- Boosting [Shapire & Kearns '96]
- C4.5 [Quinlan '96]

構造マイニング

- 非定型構造データから特徴的な部分構造を規則性を発見する
- グラフマイニング [Washio & Motoda '00], [Zaki '02], [Uno, Asai, Arimura, '02, '03]

クラスタリング

- データを類似したものどうしグルーピングする。
- 大規模・不完全なデータからの高速クラスタリング
- K-means, CLARANS, DBSCAN

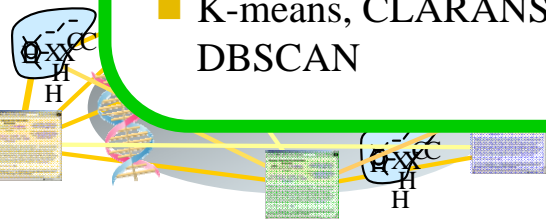
確率モデリング

- 高次元大規模データから不確実な現象を予測・モデル化する
- ベイジアンネットワーク [Pearl '90s]
- HMM [Asai], MCMC, ベイズ推定・MDL・AIC

新しいタイプのデータマイニング

- テキストマイニング
自然言語テキスト
情報抽出
意味マイニング
- ストリームマイニング
センサー監視

有用
規則・
ン・
知識
イ





Frequent Itemset Mining Definitions




Frequent Itemset Mining

- Finding **all "frequent" sets of elements** (items) appearing **no more than σ times** in a given transaction data base.
- Introduced by Agrawal and Srikant [VLDB'94]
- One of the most popular data mining problem
- Basis for more complicated / sophisticated data mining problems

- Association Rule Mining [Agrawal 1993/1994]
 - Finding combination of “items” frequently appearing in a given database

- トランザクションデータ
- バスケットデータ
- 二値データベース

- レコード/タプル
- バスケット



ID	Chips	Mustard	Sausage	Softdrink	Beer
001	1	0	0	0	1
002	1	1	1	1	1
003	1	0	1	0	0
004	0	0	1	0	1
005	0	1	1	1	1
006	1	1	1	0	1
007	1	0	1	1	1
008	1	1	1	0	0
009	1	0	0	1	0
010	0	1	1	0	1

- 
- カラム/属性
 - アイテム

トランザクション／レコードの意味

「レコード003の顧客は、ポテトチップとソーセージを一緒に買った」

- Frequent Itemset $X = \{ \text{Mustard, Sausage, Beer} \}$
 - with support/frequency 40%

ID	Chips	Mustard	Sausage	Softdrink	Beer
001	1	0	0	0	1
002	1	1	1	1	1
003	1	0	1	0	0
004	0	0	1	0	1
005	0	1	1	1	1
006	1	1	1	0	1
007	1	0	1	1	1
008	1	1	1	0	0
009	1	0	0	1	0
010	0	1	1	0	1

←
•アイテム集合 X

• X の出現リスト

$\text{Occ}(X) = \{002, 005, 006, 010\}$

• X の頻度(サポート)

$\text{freq}(X) = |\text{Occ}(X)|$
 $= 4/10 = 40\%$

アイテム集合 X の意味

$X =$ 「マスタードとソーセージ, ビールを一緒に買う人」が全体の40%いた

Frequent Itemset Mining Problem

- Given: A database **DB** over a set Σ of items, and a number $\sigma \geq 0$ called “minsup” (minimum support)
- Problem: Find all itemsets $X \subseteq \Sigma$ appearing in no less than σ records of DB

- A set $\Sigma = \{ 1, \dots, n \}$ of items (elements)
- Transaction database
 - A set $\mathbf{T} = \{ t_1, \dots, t_m \}$ of subsets of Σ
 - Each subset $t \subseteq \Sigma$ is called a **tuple** (record)

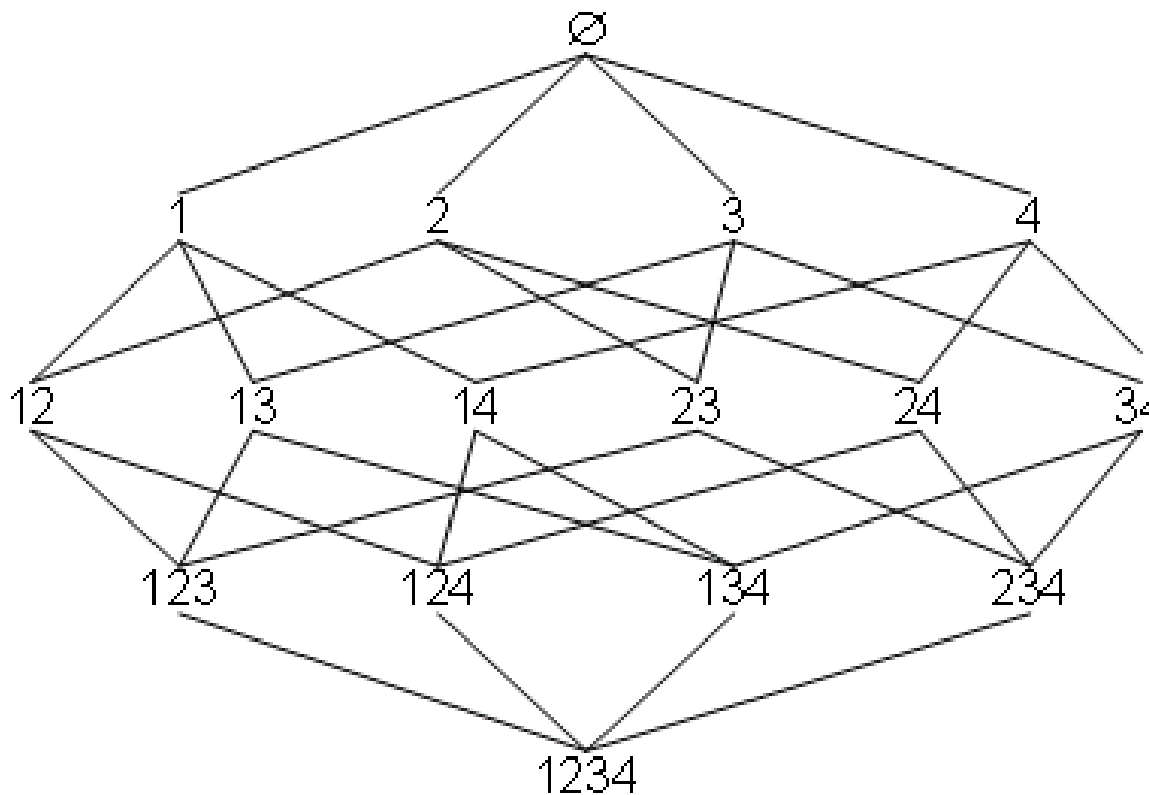
Alphabet of items

$$I = \{1, 2, 3, 4\}$$

Transaction database

id	tuples
1	1, 3
2	2, 4
3	1, 2, 3, 4
4	1, 2, 4

- Item set: any subset $X \subseteq \Sigma = \{1, \dots, n\}$
- (Item) set lattice $L = (2^\Sigma, \subseteq)$
 - The power set $2^\Sigma = \{X : X \subseteq \Sigma\}$
 - The subset relation \subseteq over 2^Σ



Example:
The set lattice
for $\Sigma = \{1, 2, 3, 4\}$



- An itemset X **appears** in a tuple t : $X \subseteq t$
- The **occurrence** of X in a database T :
 $Occ(X, T) = \{ t \in T : X \subseteq t \}$
- The **frequency** of X : $Fr(X, T) = | Occ(X, T) |$
- Minimum support (minsup): an integer $0 \leq \sigma \leq |T|$
- X is **σ -frequent (frequent)** in T if $Fr(X, T) \geq \sigma$.

Alphabet of items

$$I = \{A, B, C, D\}$$

Transaction database

id	tuples
1	1, 3
2	2, 4
3	1, 2, 3, 4
4	1, 2, 4

Occurrences and frequencies of itemsets

$$Occ(\mathbf{3}, T) = \{1, 3\}$$

$$Fr(\mathbf{3}, T) = 2$$

$$Occ(\mathbf{24}, T) = \{2, 3, 4\},$$

$$Fr(\mathbf{24}, T) = 3$$

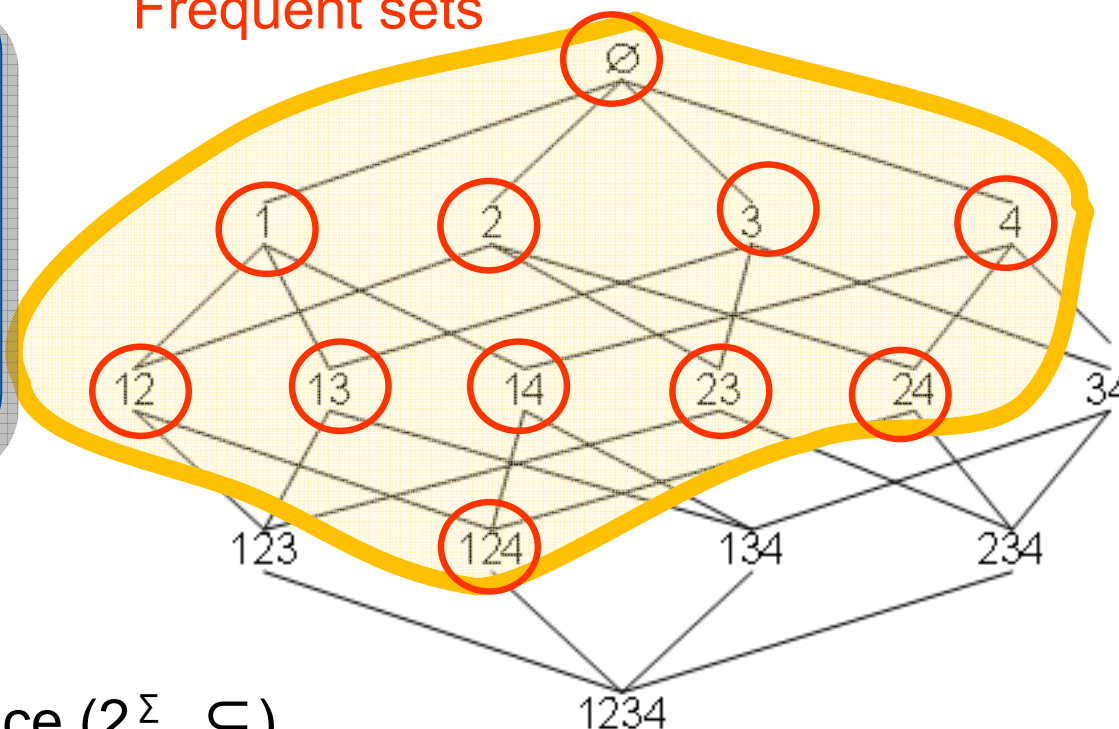
- The **occurrence** of X in a database T :
 $Occ(X, T) = \{ t \in T : X \subseteq t \}$
- X is **σ -frequent (frequent)** in T if $Fr(X, T) = | Occ(X, T) | \geq \sigma$.

minsup $\sigma = 2$

Frequent sets

\emptyset ,
1, 2, 3, 4,
12, 13, 14,
23, 24, 124

Frequent sets



	1	2	3	4	5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

database

The itemset lattice ($2^{\Sigma}, \subseteq$)



Frequent Itemset Mining Problem

- Given: A transaction database T and a non-negative integer $0 \leq \sigma \leq |T|$
- Task: Enumerate **all "frequent" itemsets X** in T that have frequency at least σ ($\text{Fr}(X) \geq \sigma$)
- \mathbf{F} : the class of all σ -frequent itemsets
- The number $|\mathbf{F}|$ of solutions is **exponential** in the **number n** of items.
- a typical **enumeration problem**.

Frequent Itemset Mining Problem

- Given: A database **DB** over a set Σ of items, and a number $\sigma \geq 0$ called “minsup” (minimum support)
- Problem: Find all itemsets $X \subseteq \Sigma$ appearing in no more than σ records of DB

Applications to more sophisticated data mining problems

- **Association rule** [Agrawal, Srikant '94]
 - {Mustard, Sausage, Beer => PotatoChips } **with frequency 40%**
- **Optimized classification rule** [Sese & Morishita PODS'90]
 - **If gene0001 & gene0012 then diabetes with classification error 8.5%**
- **Itemset boosting/SVM** [Saigo, Uno, Tsuda Bioinformatics **23(18)** 2007]
 - Learning a linear classifier over itemsets as composite features
- **Weighted substructure mining** [Nowozin, Tsuda, Uno, Kudo, Bakir, CVPR'07]
 - Application to image processing

ブースティング (Boosting) [Freund, Shapire 1996]

- 多数の機械学習アルゴリズムを統合して高精度予測
- オンライン予測の理論と深い関連

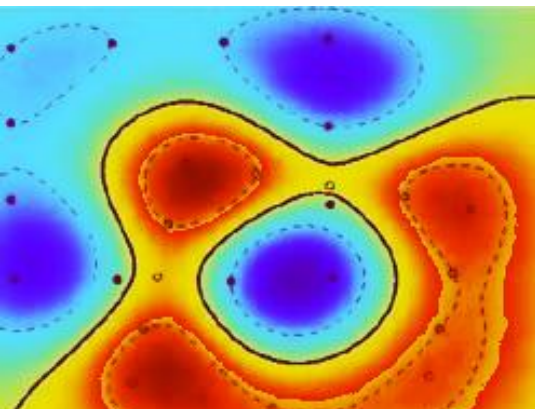
SVM (Support Vector Machines) [Vapnik 1996]

- マージン最大化による現在の state-of-the-art methods
- カーネル法を用いた高次元空間と多様なデータへの拡張

これらの学習アルゴリズムと、重みつきアイテム集合マイニング
を組み合わせる

キーワード: Boosting, SVM, オンライン予測, ニューラルネット, 計算学習理論
国際会議: NIPS, ICML, COLT, ALT

- V. Vapnik, Statistical Learning Theory, Wiley, 1998. (textbook)
- N. Cristianini and J. Shawe-Taylor, An introduction to support vector machines and other kernel-based learning methods, Cambridge, 2000. (textbook)
- Y. Freund and R. E. Schapire, A decisiontheoretic generalization of on-line learning and an application to boosting, JCSS, 55, 119-139, 1997. (AdaBoost)
- 金森, 畑埜, 渡辺, ブースティング: 学習アルゴリズムの設計技法, 森北出版 (text book)





How to model efficient data mining algorithms?

- **Light-weight**
- **High-throughput**

Modeling data mining as enumeration

- Idea: Measure the computation time per solution

■ Output-polynomial (OUT-POLY)

- Total time = $\text{poly}(N, M)$

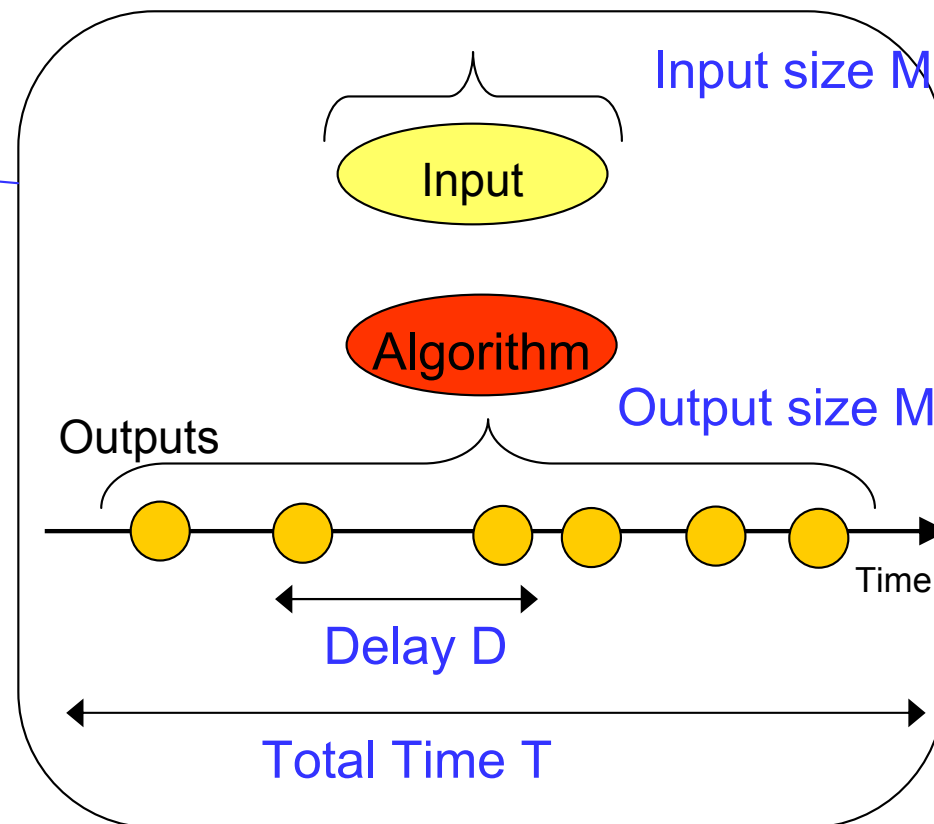
■ polynomial-time enumeration, or amortized polynomial-delay (POLY-ENUM)

- Amortized delay is $\text{poly}(\text{Input})$, or
- Total time = $M \cdot \text{poly}(N)$

■ polynomial-delay (POLY-DELAY)

- Maximum of delay is $\text{poly}(\text{Input})$

+ polynomial-space (POLY-SPACE)



Modeling data mining as enumeration

- Idea: Measure the computation time per solution

Ultimate Goal:

To design
polynomial delay and
polynomial space algorithm
 for a given data mining problem

Output-polynomial (POLY-OUTPUT)

- Total time is $O(\text{Output} \cdot \text{poly}(\text{Input}))$

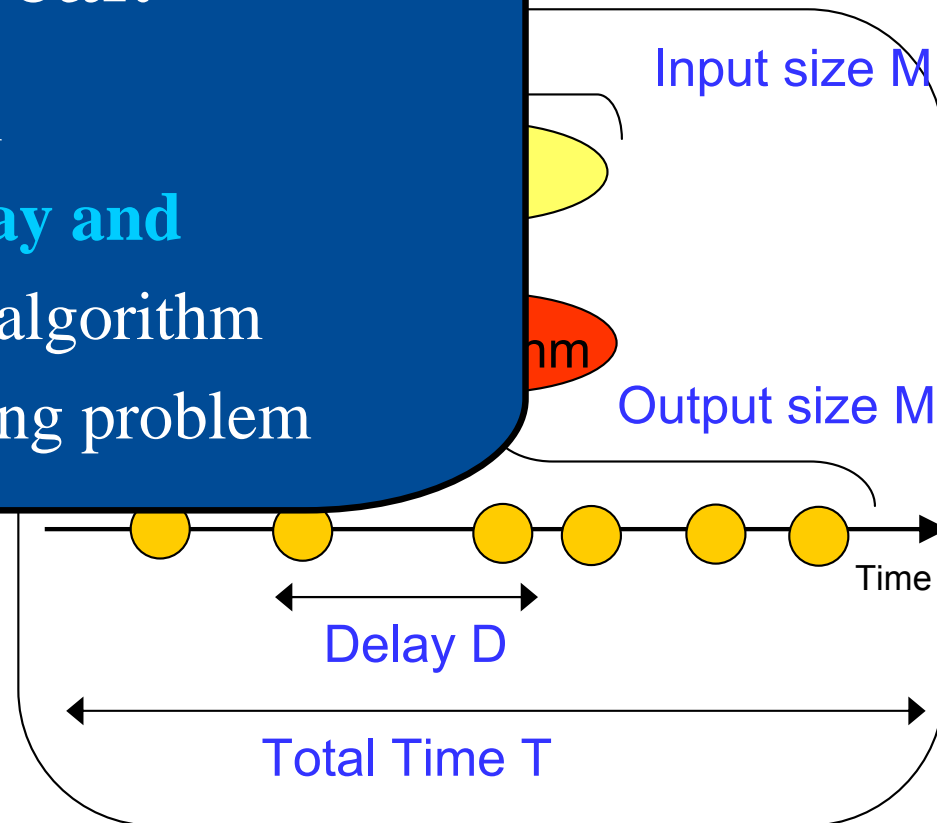
polynomial-delay (POLY-DELAY)

- Amortized delay is $O(\text{poly}(\text{Input}))$
- Total time is $O(\text{Output} \cdot \text{poly}(\text{Input}))$

polynomial-delay (POLY-DELAY)

- Maximum of delay is $\text{poly}(\text{Input})$

+ **polynomial-space (POLY-SPACE)**





Frequent Itemset Mining Algorithms

Problems

- 1. How to enumerate all frequent subsets without duplicates?

→ 列举の問題！

- 2. How to compute the frequencies quickly for each subset?

- **Lemma:** For every minsup σ , if Y is σ -frequent then any subset X of Y is also σ -frequent in DB \mathcal{T} .
- **Corollary:** Every non-empty σ -frequent set Y is obtained from some σ -frequent set X by adding some new element.
- The class F_σ of all frequent sets is monotone subclass of $(2^\Sigma, \subseteq)$

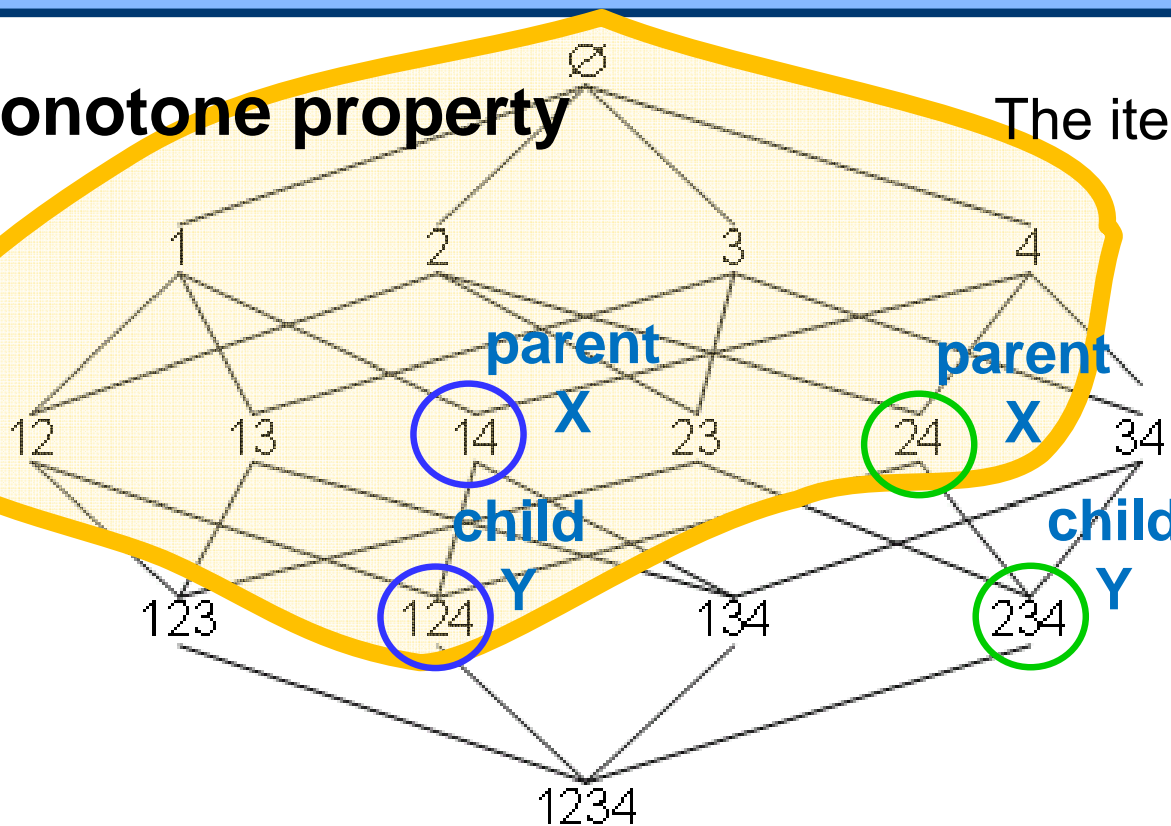
Anti-monotone property

The itemset lattice $(2^\Sigma, \subseteq)$

minsup $\sigma = 2$

Frequent sets

\emptyset ,
1, 2, 3, 4,
12, 13, 14,
23, 24, 124

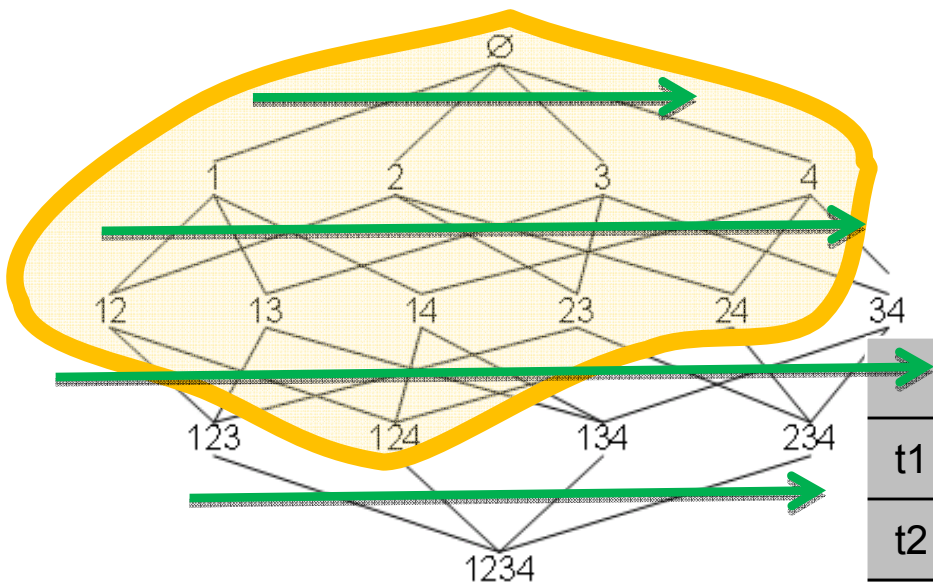


	1	2	3	4	5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

database

Apriori algorithm [1994]

- Breadth-first search (BFS)
- Horizontal data layout



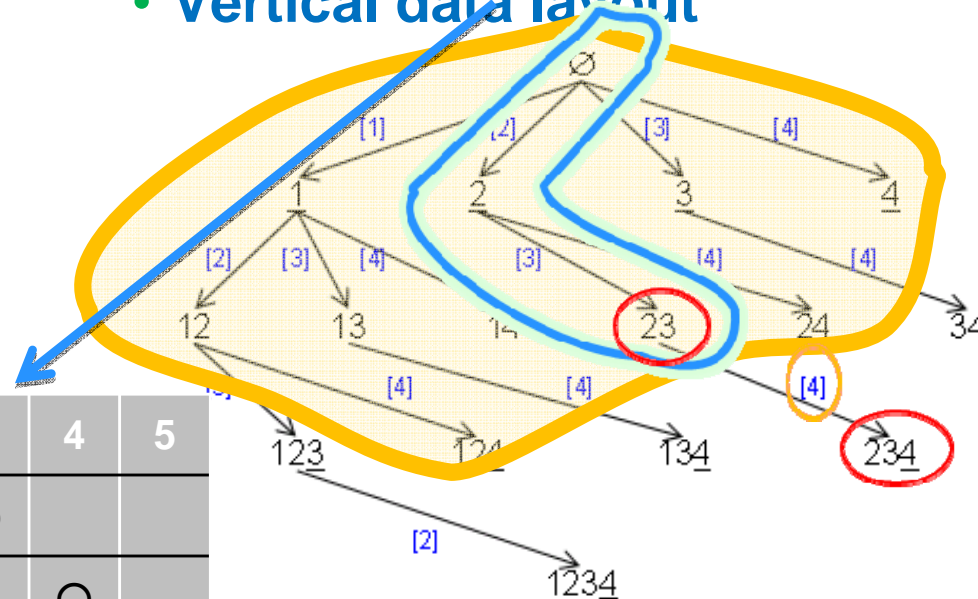
- 1st generation
- External memory algorithm

	1	2	3	4	5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

database

Backtrack algorithm [1997-1998]

- Depth-first search (DFS)
- Vertical data layout

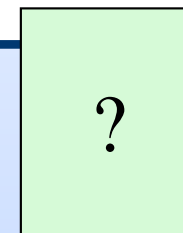


- 2nd generation
- In-core algorithm
- Space efficient

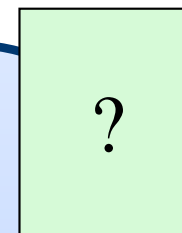


APRIORI Algorithm

- Most popular, the 1st generation frequent itemset miner
- [Agrawal & Srikant, VLDB'94; Mannila, Toivonen, Verkamo, KDD'94]



Dr. Agrawal
IBM Almaden Lab.



Prof. Mannila
Helsinki Inst. Tech
Univ. Helsinki

Candidate Generation: BFS (Breadth-first search)

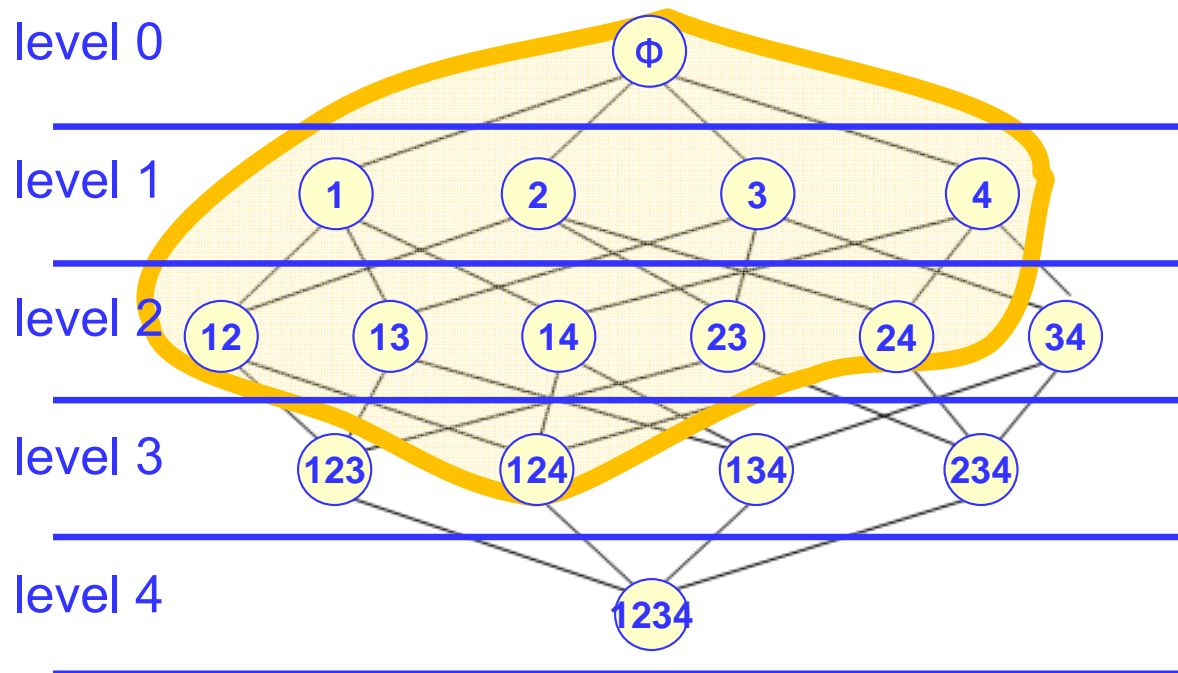
- Starting from the smallest element, search the itemset lattice **from smaller to larger**
- Generate itemsets **level by level** (level-wise algorithm)

Frequency Counting: Horizontal Data Layout

- Sequentially scanning **a large database** placed on **a hard disk** from left to right to compute the frequencies
- All the candidate itemsets is stored in a dictionary (a **hash tree**) in main memory.

Apriori algorithm [Agrawal, Srikant, 1994] (Also called Level-wise [Mannila et al., 1994])

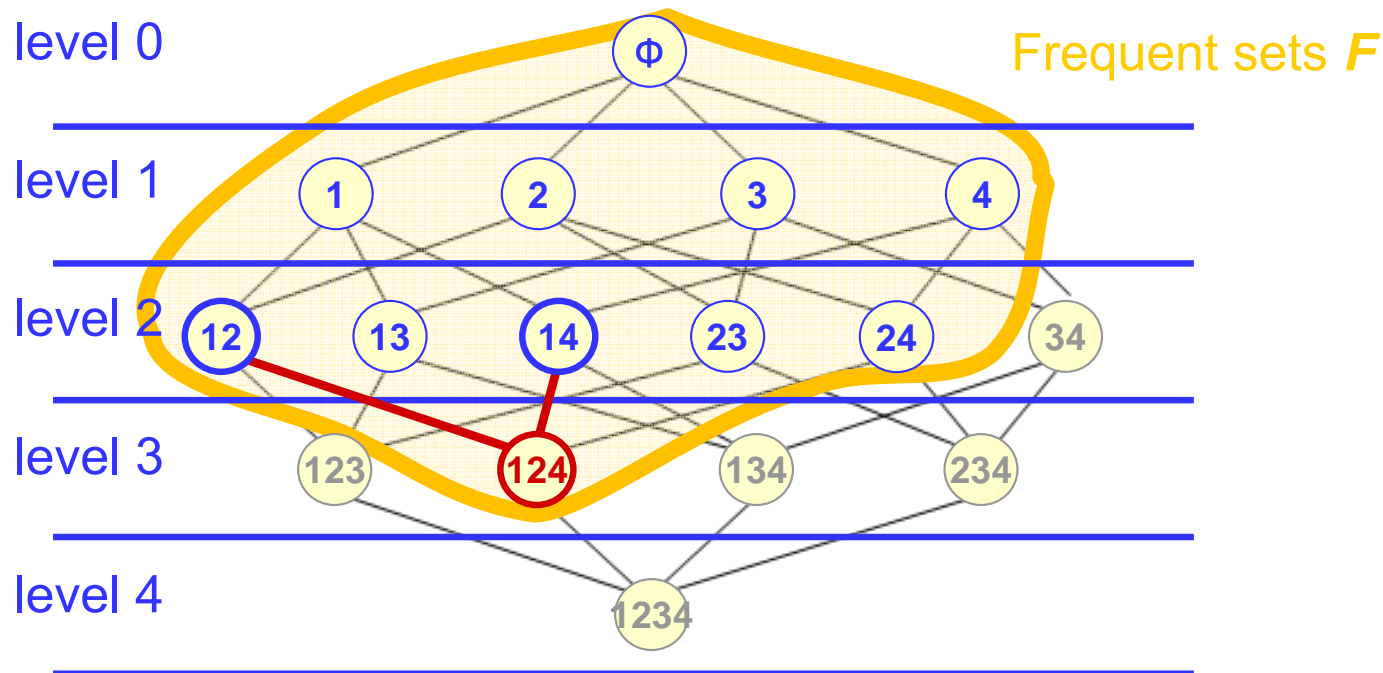
- Slice the item set-lattice 2^{Σ} into levels $L_0, L_1, \dots, L_n()$
- Compute $F_k = \{ X \in L_k : X \text{ is a frequent set of size } k \}$ for $k = 0, 1, 2, \dots$ in a bottom-up manner
- Generation of the k -th candidate set C_{k+1} from F_k :
 - For each $a_1 \dots a_k$ and $b_1 \dots b_k \in F_i$ with $a_1 = b_1 \dots a_{k-1} = b_{k-1}$ and $a_k < b_k$, add $a_1 \dots a_k b_k$ into C_i



	1	2	3	4	5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

Apriori algorithm [Agrawal, Srikant, 1994] (Also called Level-wise [Mannila et al., 1994])

- Slice the item set-lattice 2^{Σ} into levels $L_0, L_1, \dots, L_n()$
- Compute $F_k = \{ X \in L_k : X \text{ is a frequent set of size } k \}$ for $k = 0, 1, 2, \dots$ in a bottom-up manner
- Generation of the k -th candidate set C_{k+1} from F_k :
 - For each $a_1 \dots a_k$ and $b_1 \dots b_k \in F_i$ with $a_1 = b_1 \dots a_{k-1} = b_{k-1}$ and $a_k < b_k$, add $a_1 \dots a_k b_k$ into C_i

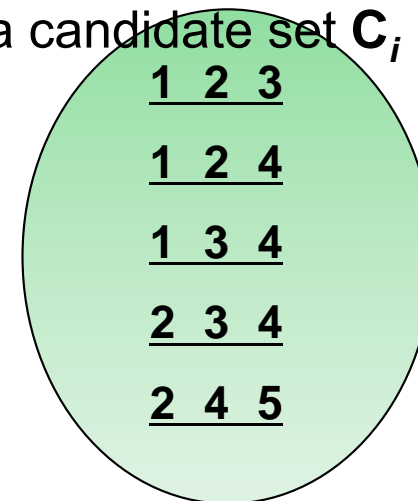


	1	2	3	4	5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

Horizontal layout

- Scanning the database tuple by tuple
- For each tuple t , **increment** the count of all candidate itemset that are subsets of t .

0) Given a candidate set C_i



1) Scan the database sequentially tuple by tuple

	items				
t1	1	3			
t2	2	4			
t3	1	2	4	5	..
t4	2	3			
t5	1	2	4		



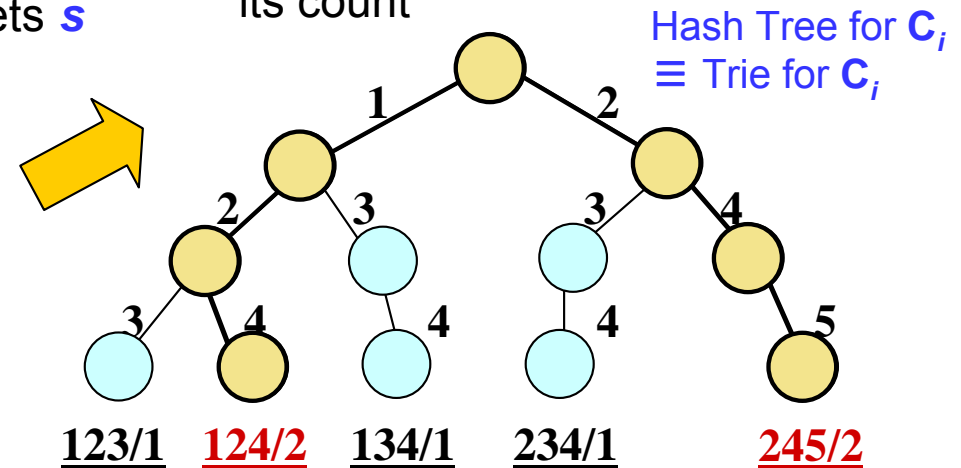
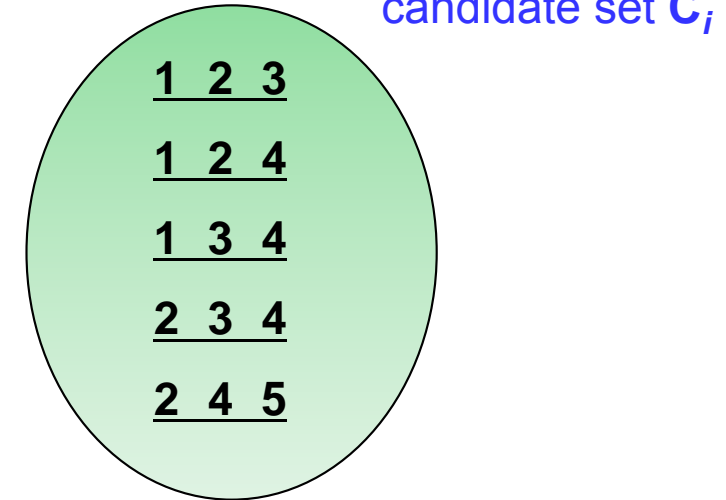
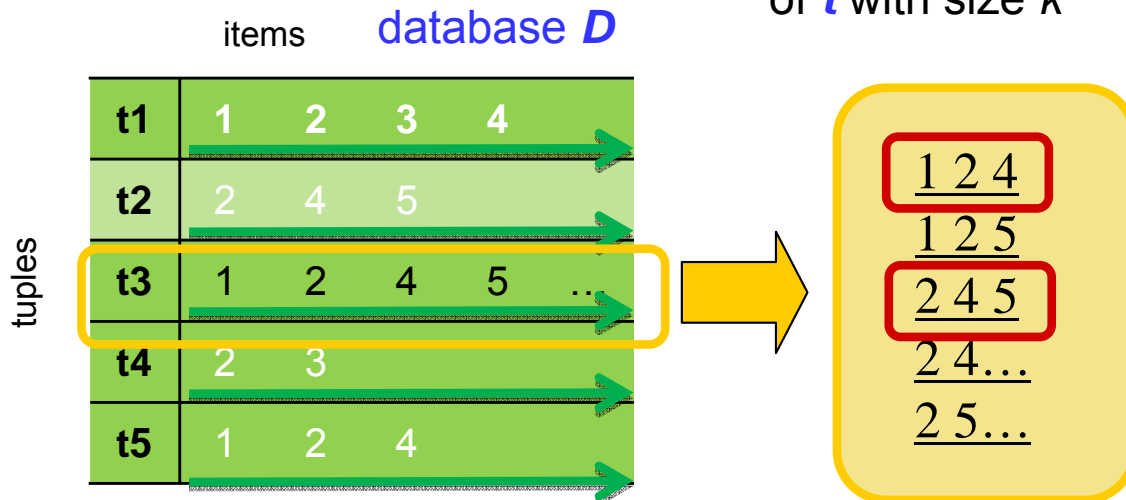
Horizontal layout

- Compute the frequencies of all candidate sets in C_i by scanning the database tuple by tuple
- For each tuple t , **increment** the count of all candidate itemset that are subsets of t .

1) Scan the database sequentially tuple by tuple

2) For each tuple t , enumerate all subsets s of t with size k

3) Lookup the subset s in the hash tree and then increment its count



● **アルゴリズム Apriori**(T : データベース, $0 \leq \sigma \leq |T|$: minsup):

出力: 全頻出集合の集合 F .

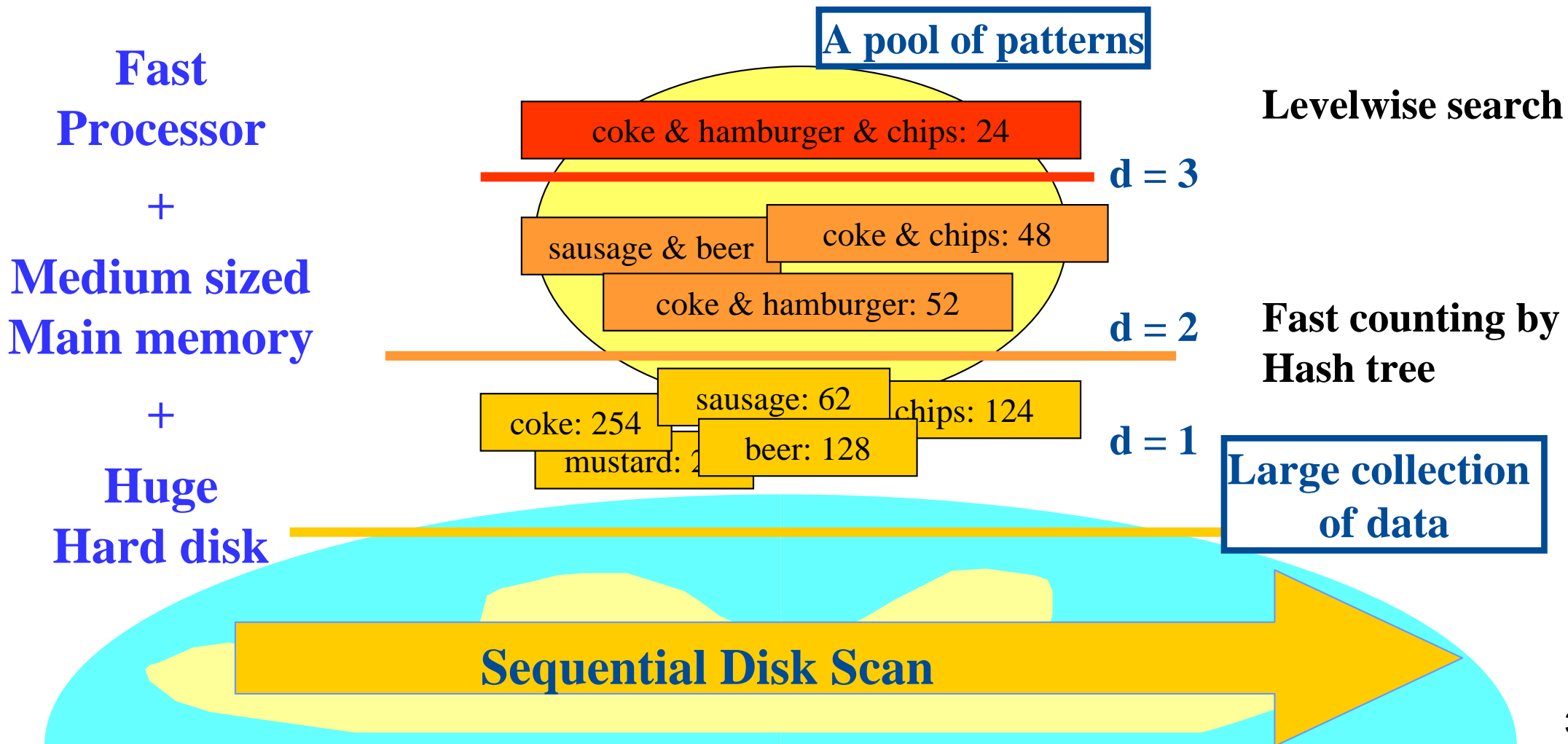
● サイズ1の頻出集合の全体 F_1 を計算する. $i = 2$.

● **while** ($F_{i-1} \neq \emptyset$) **do** // **ステージ i** :

- STEP1: F_{i-1} に含まれるサイズ $(i-1)$ の頻出集合同士を組み合わせて, サイズ i の候補集合の集まり C_i を計算する.
- STEP2: データベースを一度だけ走査し, 各候補集合 $X \in C_i$ の頻度 $\text{Freq}(X)$ を一度に計算する. $F_i = \emptyset$.
- STEP3: C_i から頻度 σ 以上のすべての候補集合を取り出し, F_i に加える.
- STEP4: $i = i + 1$.

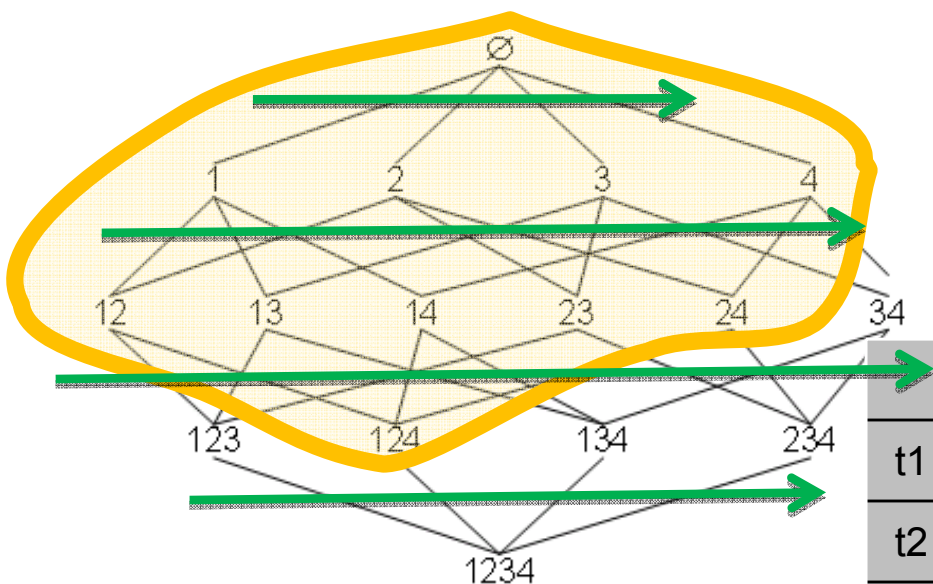
● **return** $F = F_1 \cup \dots \cup F_i$ を出力する.

The state-of-the-art algorithm for association rules
Clever use of the present computer technology



Apriori algorithm [1994]

- Breadth-first search (BFS)
- Horizontal data layout



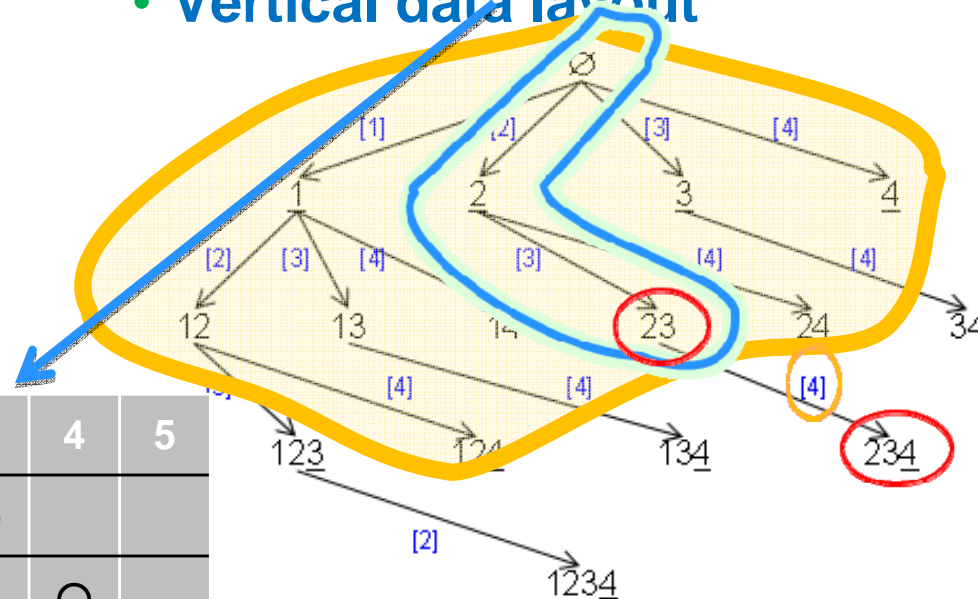
- 1st generation
- External memory algorithm

	1	2	3	4	5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

database

Backtrack algorithm [1997-1998]

- Depth-first search (DFS)
- Vertical data layout



- 2nd generation
- In-core algorithm
- Space efficient



Backtracking Algorithm

- The 2nd generation frequent itemset miners
- Eclat [Zaki et al., KDD'97], [Morishita DS'98], [Bayardo SIGMOD'98]

DFS (Depth-first search)

- **The set enumeration tree $T = (F, E_{pa}, \phi)$** for the class of frequent sets
 - Starting from ϕ , search **F** from smaller to larger
- **Pattern growth approach:**
 - Grow each itemset by attaching a new item, *the tail element*.
- **Implementation:** The search structure is implicitly implemented by a recursive procedure.

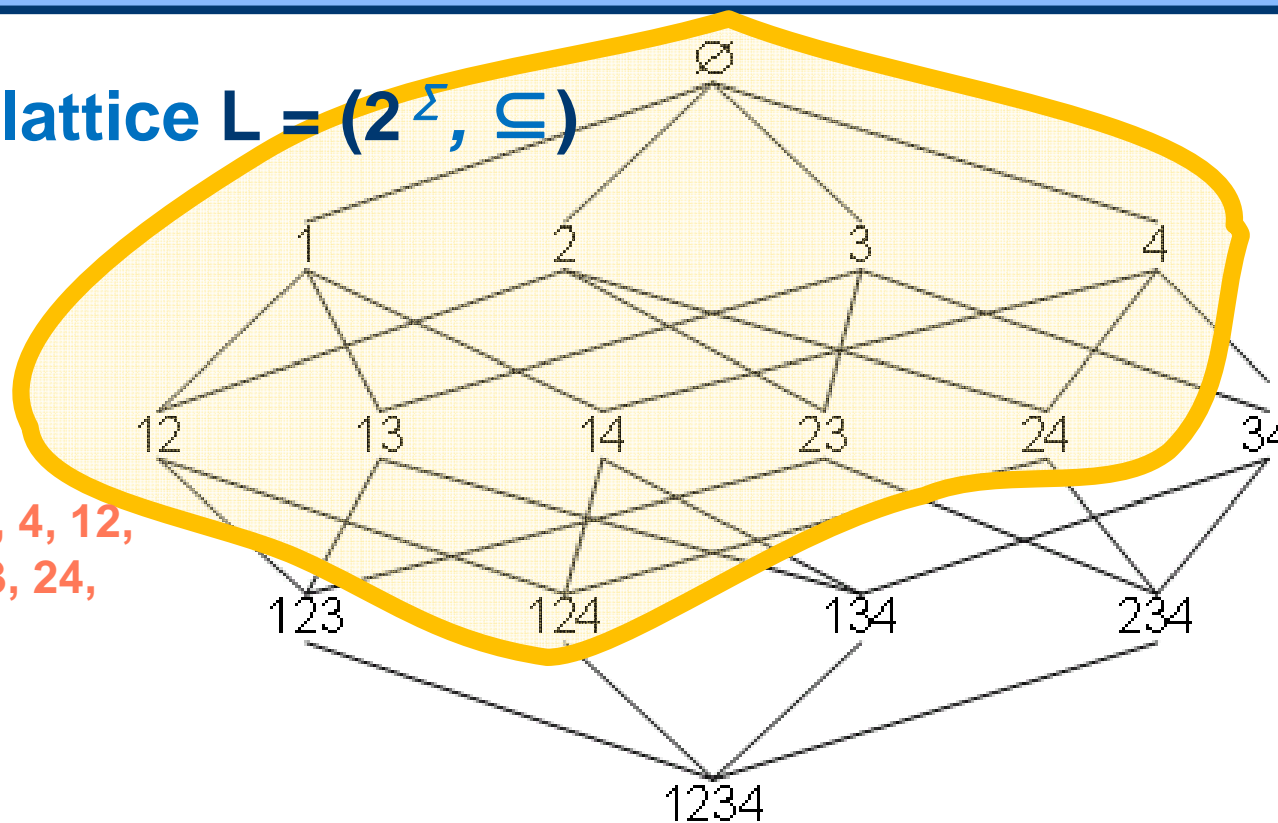
Vertical Data Layout

- For each enumerated itemset X , maintain its occurrence list $\text{Occ}(X)$ to compute the frequencies.
- Incremental update of $\text{Occ}(X)$ is possible: “*Downward closure*”

DFS algorithm

- How to generate all subsets $X \subseteq \Sigma$ in depth-first search?
- Use enumeration of subsets! (岡本先生の講義の「部分集合列挙」)

The set lattice $L = (2^\Sigma, \subseteq)$



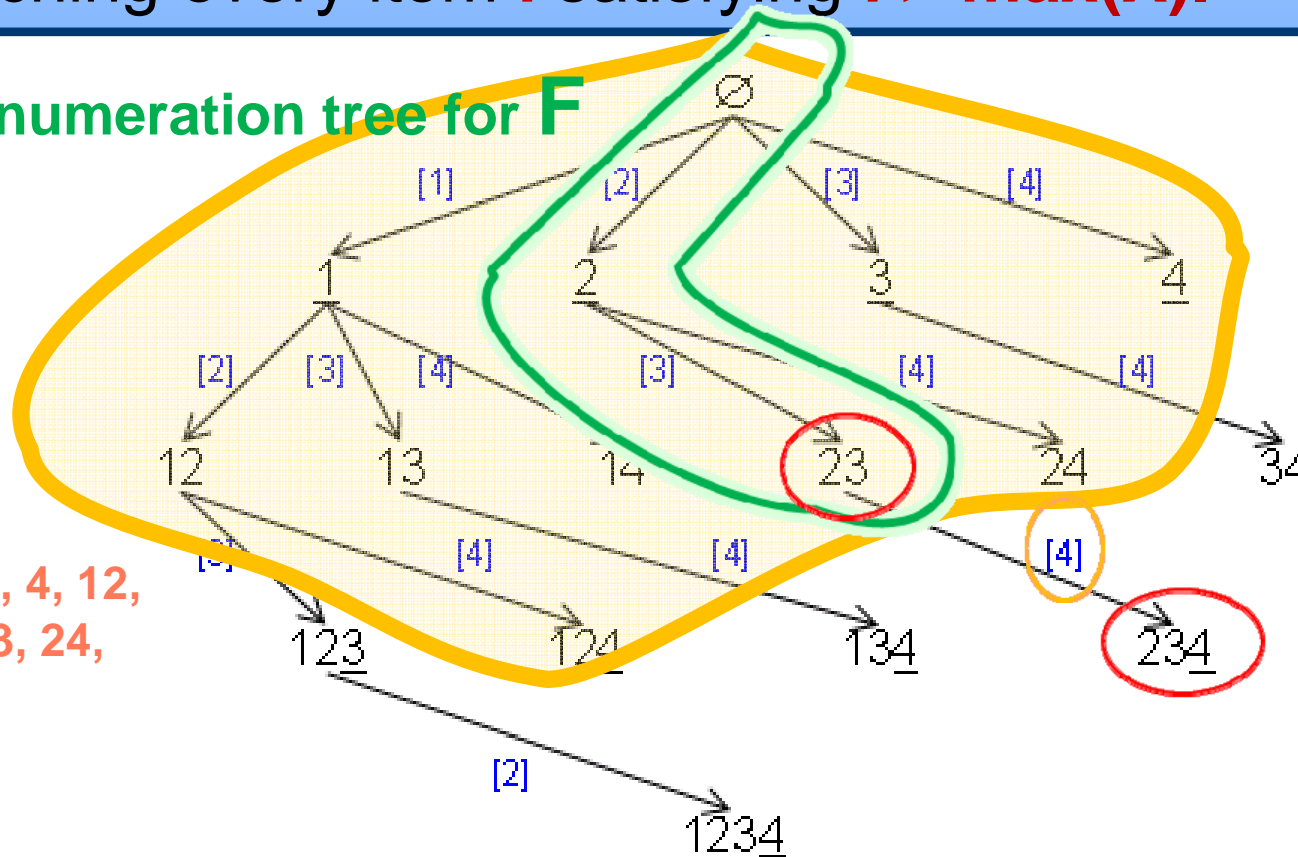
The set lattice for
 $\Sigma = \{1,2,3,4\}$

$$\sigma = 2$$

$F = \{\emptyset, 1, 2, 3, 4, 12, 13, 14, 23, 24, 124\}$

- A **spanning tree** for all frequent itemsets
- Assign the **unique parent** $\text{Pa}(Y) = Y - \{ \max(Y) \} =: X$ to each non-empty frequent set Y .
- Given a parent X , we can generate **its children** $Y = X \cup \{ i \}$ by attaching every item i satisfying $i > \max(X)$.

The set enumeration tree for F



The set lattice for $\Sigma = \{1,2,3,4\}$

$\sigma = 2$
 $F = \{ \emptyset, 1, 2, 3, 4, 12, 13, 14, 23, 24, 124 \}$

Algorithm Backtrack

- BacktrackExpand(\emptyset , Occ(\emptyset), 0, n).

procedure BacktrackExpand(X , Occ(X), k , n)

- Input: X : itemset, Occ(X), k : tail, n : maximum item;
- if $\text{Freq}(X) = |\text{Occ}(X)| < \sigma$ then,
 - return. //Backtrack
- Output a frequent set X .
- for $i = k+1, \dots, n$ do:
 - 出現リスト $\text{Occ}(X \cup \{i\})$ を計算する.
 - BacktrackExpand($X \cup \{i\}$, Occ($X \cup \{i\}$), i , n).

Properties:

- For any $X1, X2, Occ(X \cup Y) = Occ(X) \cap Occ(Y)$. (basic)
- For any set X , item a , $Occ(X \cup \{a\}) = \{ t \in Occ(X) : a \in t \}$

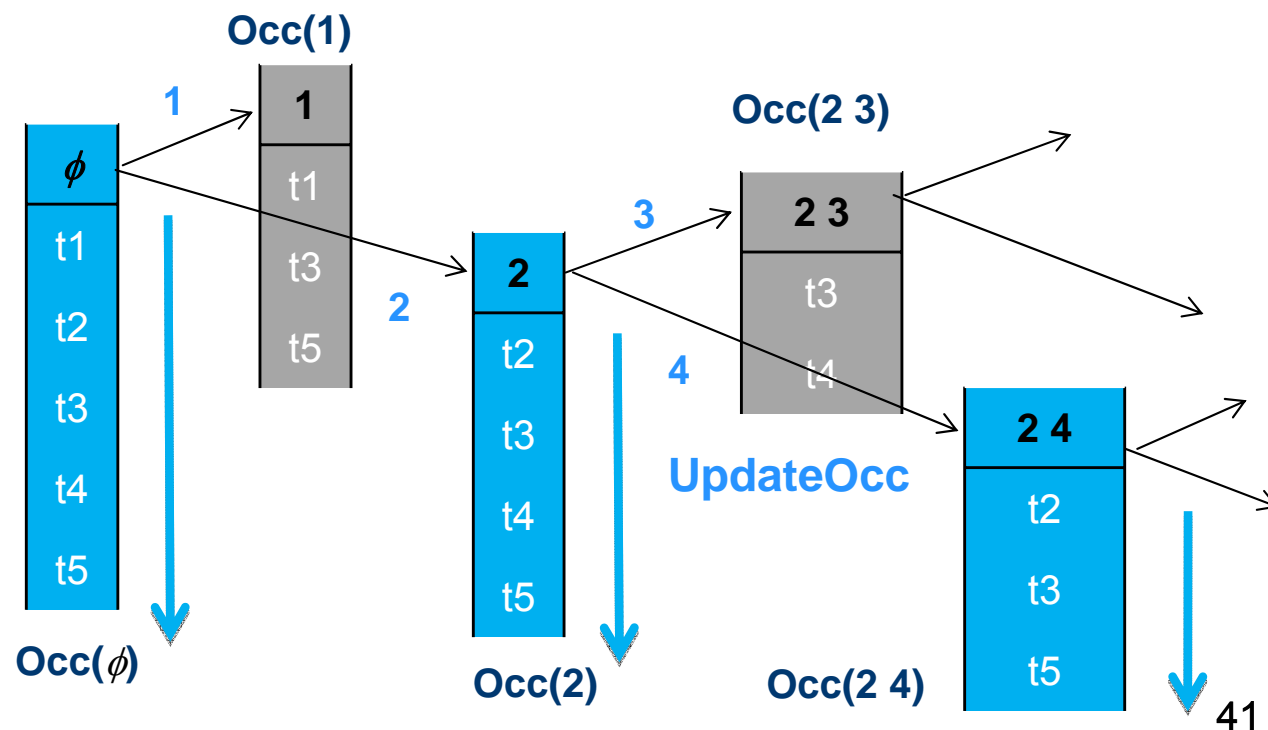
	items				
t1	1	2	3	4	
t2	2	4	5		
t3	1	2	4	5	...
t4	2	3			
t5	1	2	4		

Vertical layout:
The occurrence list representation

	1	2	3	4	5
t1	t1	t2	t1	t2	t4
t3	t3	t3	t3	t3	
t5	t5	t4	t4	t5	
		t5			

database D

Update of occurrence lists:
Downward closure [Zaki et al '97]



DFS over the set enumeration tree

Properties:

- For any $X1, X2$, $\text{Occ}(X \cup Y) = \text{Occ}(X) \cap \text{Occ}(Y)$. (basic)
- For any set X , item a , $\text{Occ}(X \cup \{a\}) = \{ t \in \text{Occ}(X) : a \in t \}$

手続き $\text{UpdateOcc}(X, a, O(X))$

Input: $X \subseteq \Sigma$ と, $a \in \Sigma$, $O(X)$.

Output: $O(X \cup \{a\})$.

- $Y = X \cup \{a\}$; $\text{Occ}(Y) = \emptyset$;
- **foreach** $t \in O(X)$ **do**
 - **if** $a \in t$ **then** $\text{Occ}(Y) = \text{Occ}(Y) \cup \{t\}$;
- **return** $\text{Occ}(Y)$;

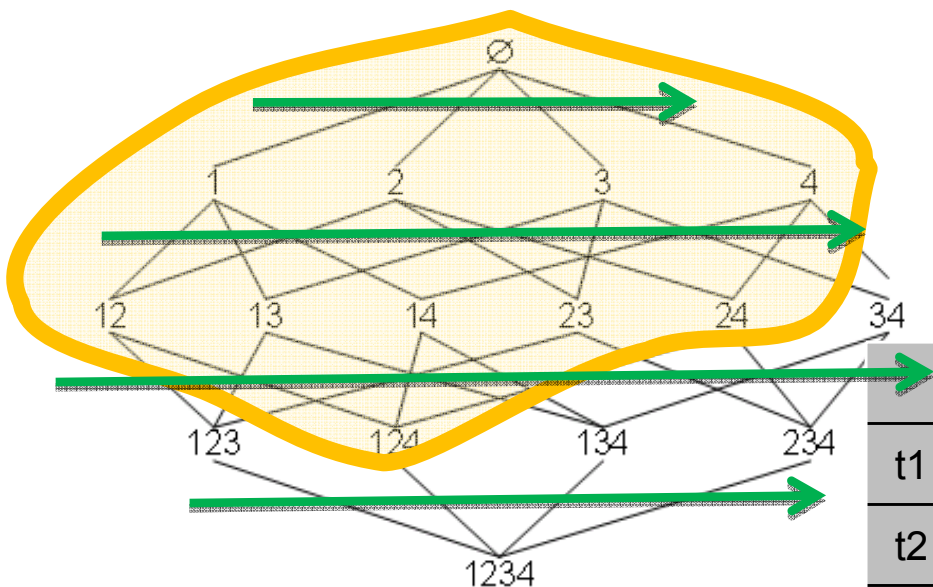
Downward closure [Zaki et al '97]

Theorem

- The algorithm **Backtrack** can be implemented to enumerates **all σ -frequent itemsets X** in a given transaction database T
 - in $O(l \cdot |\text{Occ}(X)|)$ time per frequent itemset
 - using $O(l \cdot |\text{Occ}(X)|)$ space
 - where l is the maximum length of tuples in T and $|\text{Occ}(X)|$ is the number of occurrences of X .
 - Space and time efficient (**poly-delay & poly-space**)
-
- On the other hand, the algorithm **Apriori** requires $O(l \cdot |\text{Occ}(X)|)$ time per frequent itemset and $O(\max_i \|F_i\|)$ space.

Apriori algorithm [1994]

- Breadth-first search (BFS)
- Horizontal data layout



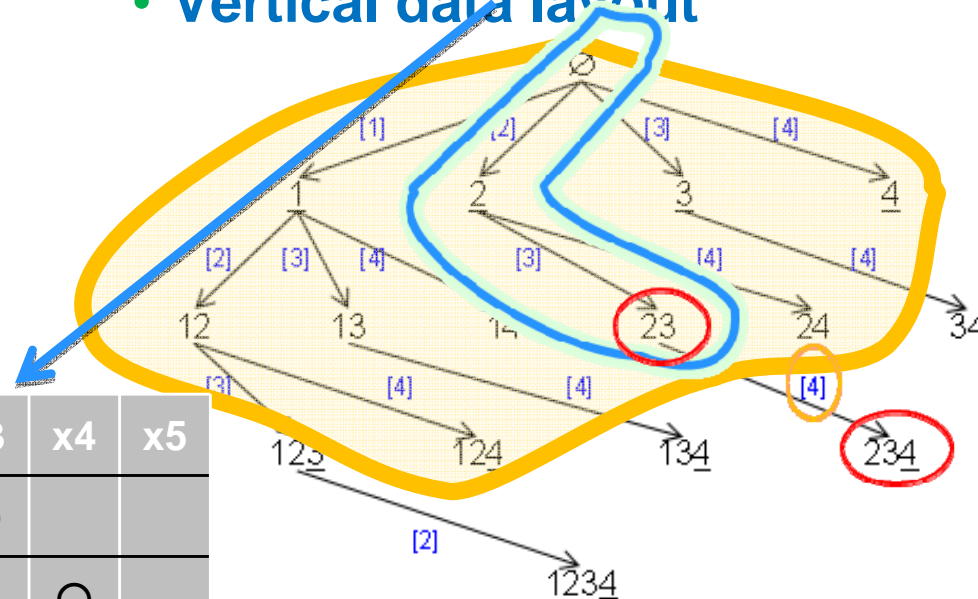
- 1st generation
- External memory algorithm

	x1	x2	x3	x4	x5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

database

Backtrack algorithm [1997-1998]

- Depth-first search (DFS)
- Vertical data layout



- 2nd generation
- Space efficient
- In-core algorithm



Summary: Horizontal & Vertical Layout

2008年2月29日
列挙学校

Horizontal layout (Apriori)

t1	x1	x3		
t2	x2	x4		
t3	x1	x2	x3	x4
t4	x2	x3		
t5	x1	x2	x4	

tuples

items

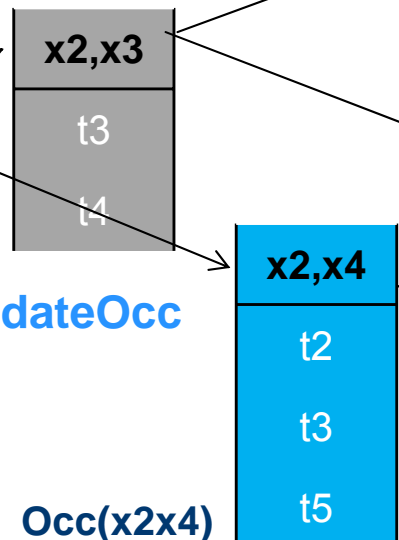
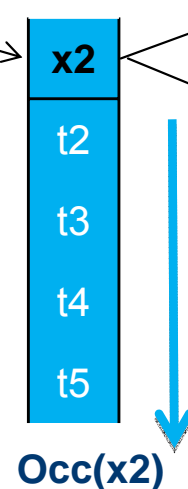
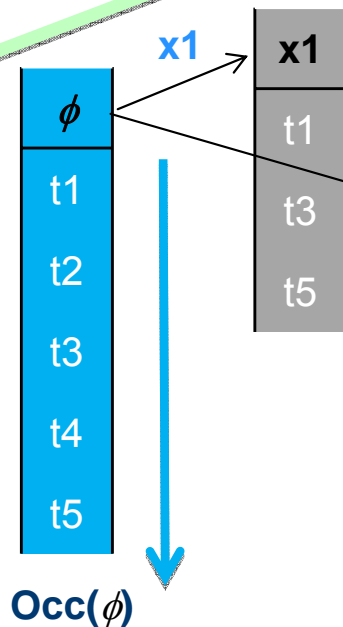
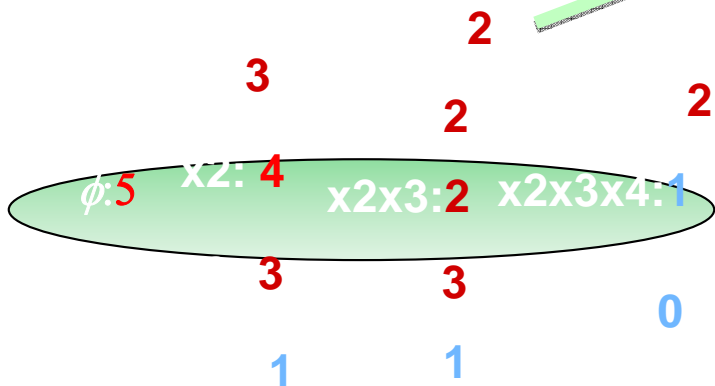
	x1	x2	x3	x4	x5
t1	○		○		
t2		○		○	
t3	○	○	○	○	
t4		○	○		○
t5	○	○		○	

Vertical layout (Backtrack)

x1	x2	x3	x4	x5
t1	t2	t1	t2	t4
t3	t3	t3	t3	
t5	t4	t4	t5	
	t5			

Scan & Count

Hash tree for (itemset, frequency)



[Han, Pei, Yin, SIGMOD'00]

- パターンとデータを**トライ**(Trie/Prefix tree/FP-Tree)に圧縮して格納
- DFS+出現リスト方式の変種
- 実際のデータに対して高速といわれている

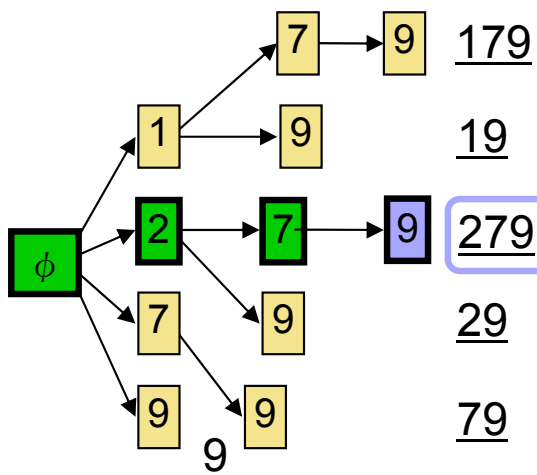
Example) Expand pattern **27** into **279** and then update its occurrence list

Pattern **279** appears 4 times in the database!

頻出集合

Φ
 {1} {2} {7} {9}
 {1,7} {1,9} {2,7}
 {2,9} {7,9}
 {1,7,9} {2,7,9}

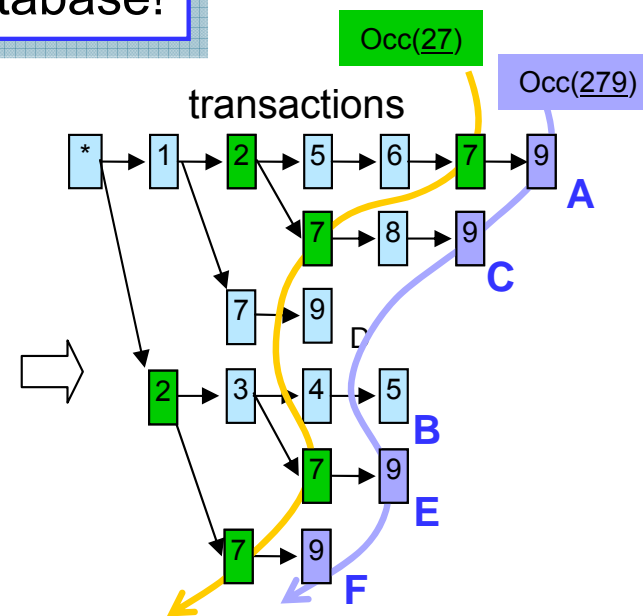
patterns



データベース

A: 1,2,5,6,7,9
 B: 2,3,4,5
 C: 1,2,7,8,9
 D: 1,7,9
 E: 2,3,7,9
 F: 2,7,9

transactions



頻出パターン木 (FP-Tree)

= トライを用いて圧縮されたパターン集合

トライを用いて圧縮されたデータベース

■ LCM (Linear-time Closed Itemset Miner)

- 超高速代表元集合マイニングアルゴリズム
- 理論的性能保証: 出力線形時間アルゴリズム
- 公開 & 応用

くわしくは宇野毅明の
講義で...

高速な代表元パターンマイニング

(宇野・有村, DS'04, FIMI'04)

ニュース速報 (Nov. 1, 2004, Brighton, UK)

第2回 FIMI Workshop で宇野・清見・
有村組(LCM ver.2) が優勝！
FIMI'04 Best Implementation Award

- FIMI Workshop:
「頻出アイテム集合マイニング実装(FIMI)」に関する
データマイニング分野のプログラミングコンテスト
今年は, 8+5 実装が最終エントリー (問合せ80件超)

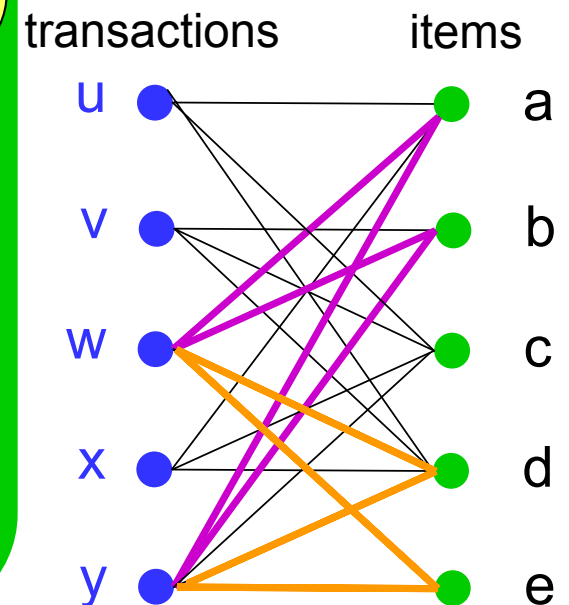
- LCMアルゴリズム: PPC拡張による出力多項式時間の 宇野毅明先生
閉アイテム集合列挙アルゴリズム [宇野・有村 DS'04] (NII, A01班)



今年は
勝ちま
した

極大2部クリーク C

$(\{w, y\}, \{a, b, d\})$



- LCM, 宇野先生HP, program codes: <http://research.nii.ac.jp/~uno/codes.htm>
- FIMI Frequent Itemset Mining Implementations Repository: <http://fimi.cs.helsinki.fi/> 47

.Frequent Itemset Mining

- Finding all "frequent" sets of elements appearing in a given transaction data base.
- Introduced by Agrawal and Srikant [VLDB'94]
- One of the most basic data mining problem

.Algorithms

- BFS algorithm — Apriori [Agrawal et al, '94]
- DFS algorithm — Backtrack [Zaki et al. '97; Morishita'98]

.Applications

- Feature extraction for SVM & Boosting
- Closed and Maximal Frequent Itemset Mining
- Sequences and Graphs

パート1: データマイニングと頻出集合発見

- [AIS 93] R. Agrawal, T. Imielinski, A. N. Swami: Mining Association Rules between Sets of Items in Large Databases, Proc. SIGMOD Conference 1993, pp. 207-216, 1993.
- [AMST 96] R. Agrawal, H. Mannila, R. Srikant, H. Toivonen, A. I. Verkamo: Fast Discovery of Association Rules, Advances in Knowledge Discovery and Data Mining, pp. 307-328, 1996.
- [Agrawal & Srikant, VLDB'94] R. Agrawal, R. Srikant: Fast Algorithms for Mining Association Rules in Large Databases. Proc. VLDB 1994, pp. 487-499, 1994. (Apriori)
- [Bayardo, SIGMOD'98] R. J. Bayardo Jr.: Efficiently Mining Long Patterns from Databases, Proc. SIGMOD Conference 1998: pp. 85-93, 1998.
- [Han et al. SIGMOD'00] J. Han, J. Pei, Y. Yin, Mining Frequent Patterns without Candidate Generation, Proc. SIGMOD Conference 2000, pp. 1-12, 2000. (FP-growth)
- [Morishita & Sese PODS'00] S. Morishita, J. Sese: Traversing Itemset Lattice with Statistical Metric Pruning, Proc. PODS 2000, pp. 226-236, 2000.
- [Zaki et al., KDD'97] M. J. Zaki, S. Parthasarathy, M. Ogihara, W. Li, New algorithms for fast discovery of association rules, Proc. KDD 1997, pp. 283-286, 1997. (Eclat)
- [宇野, 有村] 宇野毅明, 有村博紀, 「データインテンシブコンピューティング その2 - 頻出アイテム集合発見アルゴリズム -」, 人工知能学会誌, レクチャーシリーズ「知能コンピューティングとその周辺」, (編)西田豊明, 第2回, Vol.22, No.3, 2007年5月. (PDF: <http://www-ikn.ist.hokudai.ac.jp/~arim/jtalks.html>)

列挙学校：第3コマ

パート2: 系列とグラフのマイニング

有村 博紀

北海道大学大学院 情報科学研究科 コンピュータサイエンス専攻

email: {arim,kida}@ist.hokudai.ac.jp

<http://www-ikn.ist.hokudai.ac.jp/~arim>

パート1: データマイニングと頻出集合発見

- データマイニング
- 頻出集合マイニング

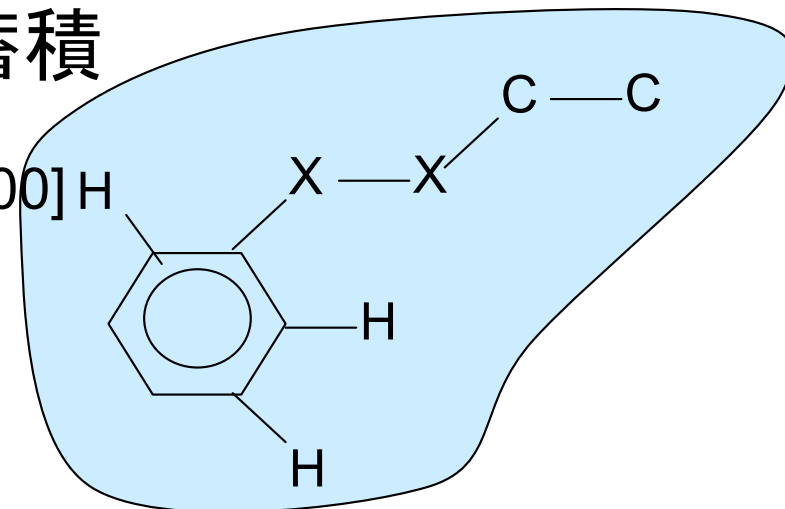
パート2: 系列とグラフのマイニング

- 系列とグラフのマイニング概観
- 列挙によるパターンのマイニング
- 列挙とマイニングに関連した話題

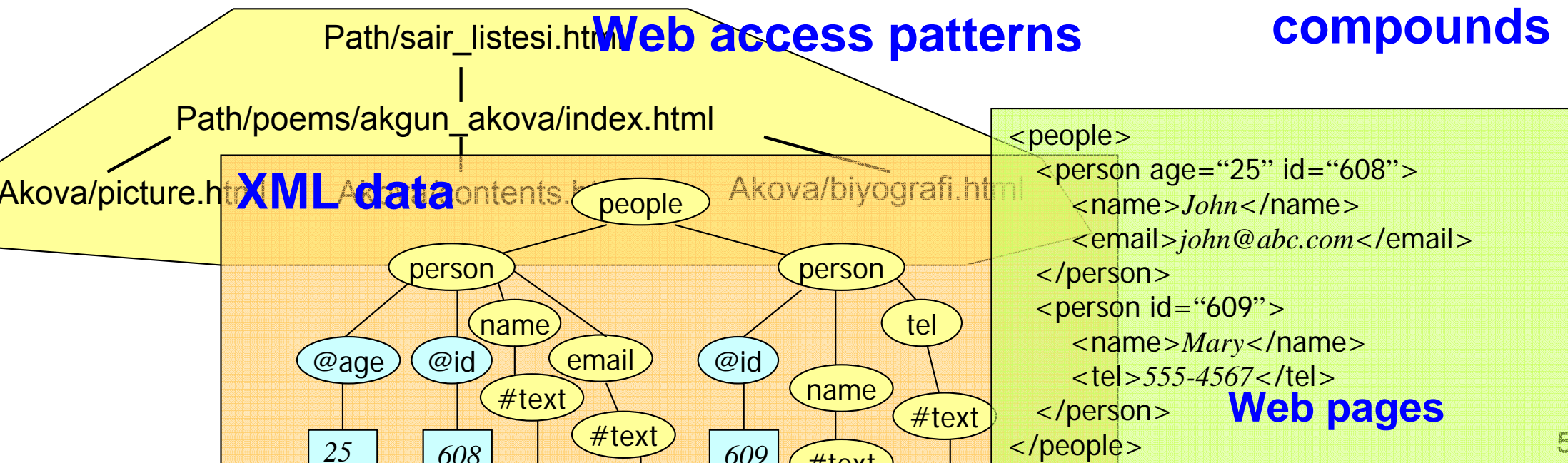
Mining Semi-Structured Data

- Sequences
- Trees
- Graphs

- 大量のウェブページやXML文書が蓄積
- **半構造データ** [Abiteboul, Buneman, Suciu 00]
- ウェブページとXML [2nd ed, W3C 00],
- ゲノムデータベース [Entrez, NCBI]
- 化合物データや, CADデータ
- 自然言語データも



Chemical compounds





History of Sequence Mining

~ 1995		HMM (Hidden Markov Model]*1) [Haussler, Krogh et al., 26 th HICSS 1993]
1996		[Durvin, Eddy, Krogh, Mitchison, Cambridge, 1998]
1997	Mining frequent episodes in an event sequence (BFS) [Mannila, Toivonen, Verkamo, DMKD, 1(3), 1997]	MOTIF program [Smith, Annau <i>et al.</i> , PNAS, 87 , 1990]
1998		PRATT program [Jonassen et al., Protein Sci. 4 , 1995]
1999		TEIRESIAS program [Rigoutsos, Floratos, Bioinformatics 14 , 1998]
2000	Mining frequent itemset sequence (DFS) Spade [Zaki, 1998; Mach. Learn. 2000]	eMOTIF program [Nevill-Manning et al., PNAS 95 , 1998]
2001	Mining frequent sequences (DFS + Projected DB) PrefixSpan [Pei, Han, et al., ICDE 2001]	
2002	Mining frequent closed frequent sequences (DFS) CloSpan [Yan, Han, et al., SDM 2003]	PROJECTION (Random projection) [Buhler, Tompa, J. Comp. Bio. 9(2) , 2002]
2003	Mining frequent closed sequences (DFS) BIDE [Wang & Han, ICDE 2004]	
2004	Mining frequent closed sequences with wildcards (DFS) AU [Arimura, Uno, LNCS 3827, ISAAC2005]	a PTAS algorithm for concensus motif problem in Hamming distance [Li, Ma, Wang, STOC 1999]
2005	Mining frequent closed sequences with gaps (DFS) AU [Arimura, Uno, LNAI 4914, LLLL2007]	
2006		
2007		
2008		



History of Tree & Graph Mining

~1995

Algorithm for finding subgraphs by MDL principle

1996

Subdue [Holder *et al.* (KDD'94)]

1997

Finding frequent paths [Wang and Liu (KDD'97)]

1998

Finding Semi-structured Schema
[Nestrov, Abiteboul *et al.* (SIGMOD'98)]

1999

Finding frequent subgraphs

2000

AGM [Inokuchi, Wahio, Motoda (PKDD'00, MLJ. 2003)]

2001

FSG [Kuramochi *et al.* (ICDM'01)]

2002

Finding frequent / optimal ordered trees

FREQT [Asai *et al.* (SDM'02)], **Treeminer** [Zaki (KDD'02)]

2003

Finding frequent subgraphs

gSpan [Yan and Han (ICDM'02)], **VG** [Venetik, *et al.* (ICDM'02)]

Finding frequent unordered trees

UNOT [Asai, Uno, Nakano, Arimura (SDM'03)], **NK** [Nijssen, Kok (MGTS'03)]

Frequent Maximal/Closed Trees & Graphs mining

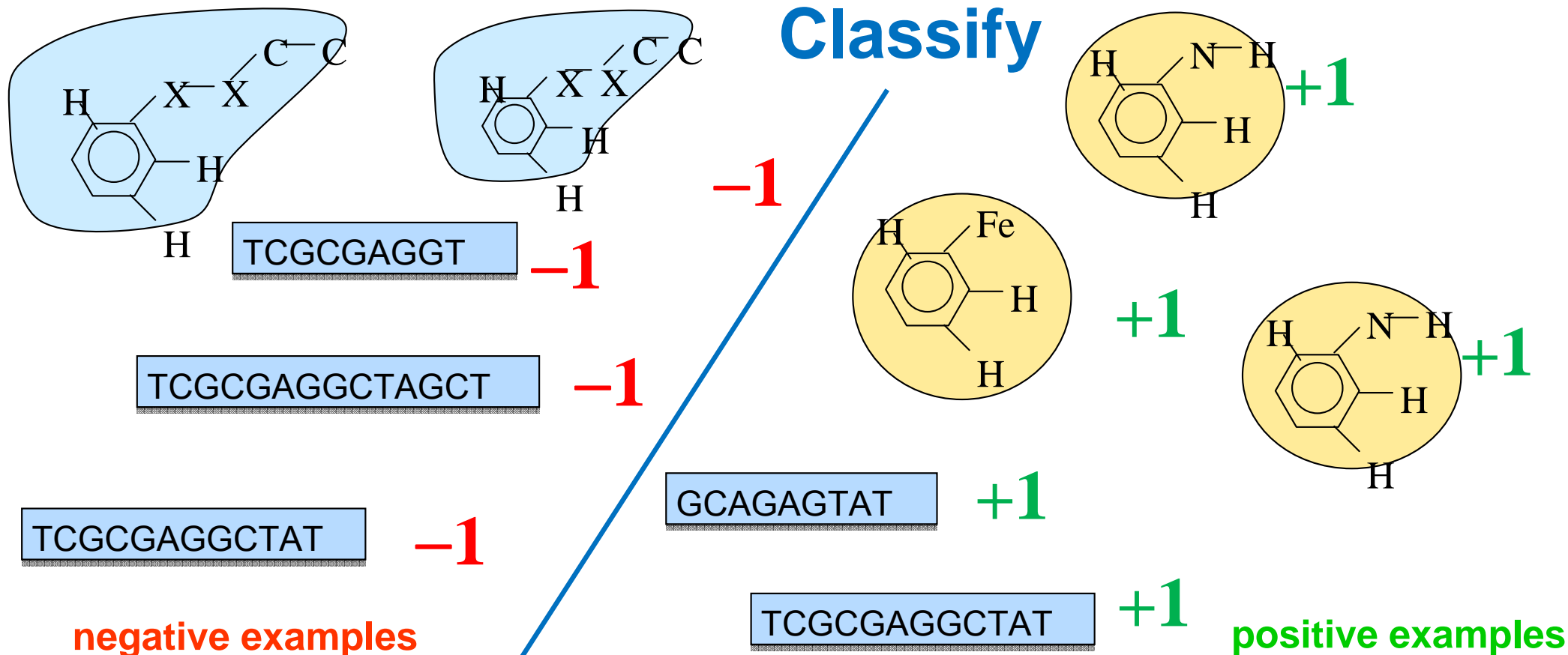
[Yan&Han '03; Termier *et al.*'04] and many algorithms in 2000s

プラン

- ここでは、時間が許す限り、系列、木、グラフのパターンに関する具体的なデータマイニングアルゴリズムと最近の話題を紹介する予定
- 機械学習で用いられるカーネル法と列挙に関する話題にも触れたい

Learning Graphs from examples

- unknown function $f: \text{Graphs} \rightarrow \{+1, -1\}$



Kernel methods and Enumeration

- Kernel methods for learning structured data
- #P-hardness and enumeration
- Hardness of learning and prediction

パート1: データマイニングと頻出集合発見

- データマイニング
- 頻出集合マイニング

パート2: 系列とグラフのマイニング

- 系列とグラフのマイニング概観
- 列挙によるパターンのマイニング
- 列挙とマイニングに関連した話題

最新版PPTは次からダウンロードしてください: <http://www-ikn.ist.hokudai.ac.jp/~arim/jtalks.html>

Hiroki Arimura, Hokkaido University, 2008