

# 機械学習 第12回 パターンマイニング

立命館大学 情報理工学部

福森 隆寛

Beyond Borders

# 講義スケジュール

(第1～4回、第14回) (第5～13回、第15回)

□ 担当教員：村上 陽平先生・福森 隆寛

1	機械学習とは、機械学習の分類
2	機械学習の基本的な手順
3	識別（１）
4	識別（２）
5	識別（３）
6	回帰
7	サポートベクトルマシン
8	ニューラルネットワーク

9	深層学習
10	アンサンブル学習
11	モデル推定
12	パターンマイニング
13	系列データの識別
14	強化学習
15	半教師あり学習

□ 担当教員：叶 昕辰先生（第16回の講義を担当）

# 今回の講義内容

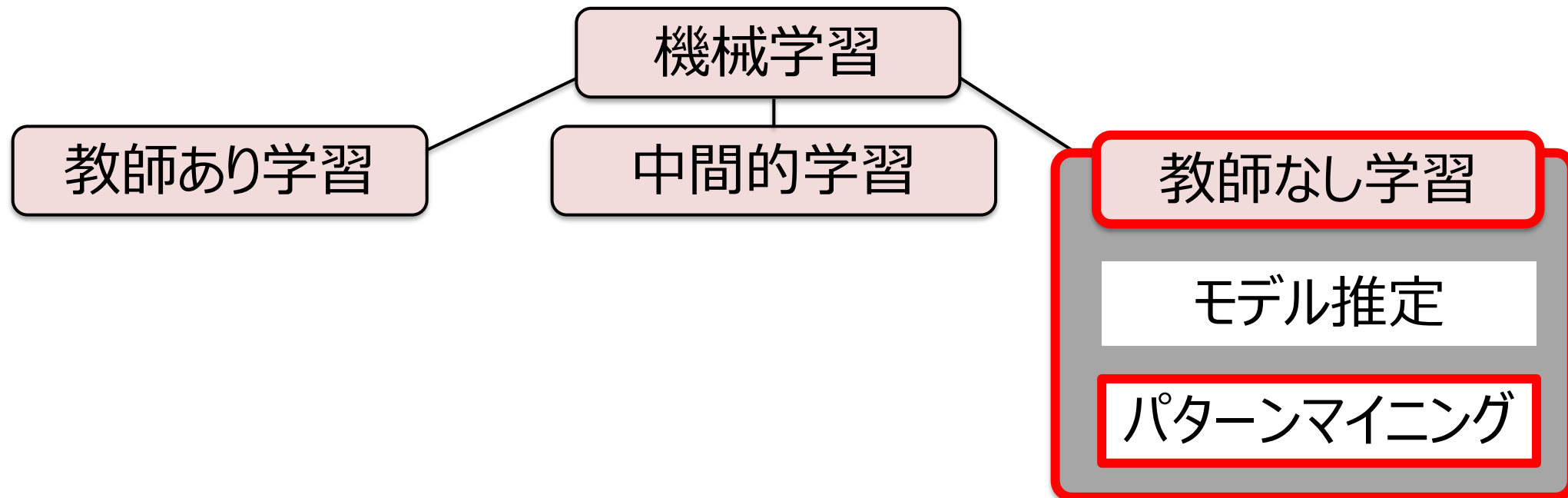
- 取り扱う問題の定義
- パターンマイニング
- Aprioriアルゴリズム
  - 頻出項目抽出<sup>ひんしゅつ</sup>
  - 連想規則抽出<sup>れんそう</sup>
- FP-Growthアルゴリズム
- 推薦システムにおける学習
  - 協調フィルタリング<sup>きょうちよう</sup>
  - Matrix Factorization
- 演習問題

# 取り扱う問題の定義：教師なし・パターンマイニング

□ **カテゴリ形式**の特徴ベクトルを入力して、  
そのデータに潜んでいる有用なパターンを見つける

※ 教師なし学習の問題での学習データは、以下で構成される

入力データの特徴ベクトル  $\leftarrow \{x_i\}, i = 1, 2, \dots, N \rightarrow$  学習データの総数  
(カテゴリ形式) ※ 正解情報は与えられていない



# パターンマイニング

## □ パターンマイニング

- データ中に何度も出現するパターンを抽出したり（**頻出項目抽出**）、そのパターンに基づいた規則を発見する（**連想規則抽出**）手法

## □ パターンマイニングの応用例

- ネットショッピングサイトなどでのお勧め商品の提示<sup>ていじ</sup>
  - 商品Aと商品Bを購入している客が多くいるとき、商品Aのみを購入した客に対して商品Bを勧める
- データからの連想規則の抽出による新たな知見の獲得
  - 血液の生化学検査項目<sup>けつえき せいかがくけんさ</sup>の値から腫瘍マーカー<sup>しゅよう</sup>の値の<sup>こうてい</sup>高低が推定できれば、高価な腫瘍マーカー検査の回数を減らせる

# 頻出項目抽出

## □ パターンマイニングで扱うデータの単位<sup>た ん い</sup>

### ■ トランザクション：1個のデータ

- 例：スーパーマーケットの売り上げ記録の場合
  - ある人の1回分の<sup>かい</sup> <sup>もの</sup>買い物で同時に買われた物の集合がトランザクション

## □ 支持度<sup>し じ ど</sup> (support)

### ■ 全データに対して、ある項目集合が出現する割合

$$\text{support}(\text{items}) = \frac{T_{\text{items}}}{T}$$

- $T$ ：全トランザクション<sup>けんすう</sup>件数
- $T_{\text{items}}$ ：項目集合itemsが出現するトランザクション件数

# 演習問題12-1（10分間）

□ 以下のスーパーマーケットの売り上げ記録がある

■ 商品点数：{ミルク、パン、バター、雑誌<sup>ざっし</sup>}の4点

■ トランザクション：6件（6回分の購入記録）

1. このデータは「疎らなデータ」、「密なデータ」のどちらであるか考えよ

2. 項目集合{ミルク、パン}の支持度を示す  $\text{support}(\{\text{ミルク、パン}\})$  を求めなさい

3. 商品点数が1000種類であったとき、可能な項目集合の数はいくつか？  
ヒント：商品は「買った」「買っていない」の2種類に分類できると考える

No.	ミルク	パン	バター	雑誌
1	t	t		
2		t		
3				t
4		t	t	
5	t	t	t	
6	t	t		

スーパーマーケットの売り上げ記録  
(tはその商品が買われたことを意味する)

# 頻出項目抽出：Aprioriアルゴリズム

## □ Aprioriアルゴリズム

- 「a prioriな原理」の対偶<sup>たいぐう</sup>を用いて、  
小さな項目集合から支持度の計算をはじめ、  
項目集合を大きくする際に、頻出でない項目集合を  
それ以上拡張<sup>かくちょう</sup>させず、調べる項目集合を減らす方法<sup>しら</sup>

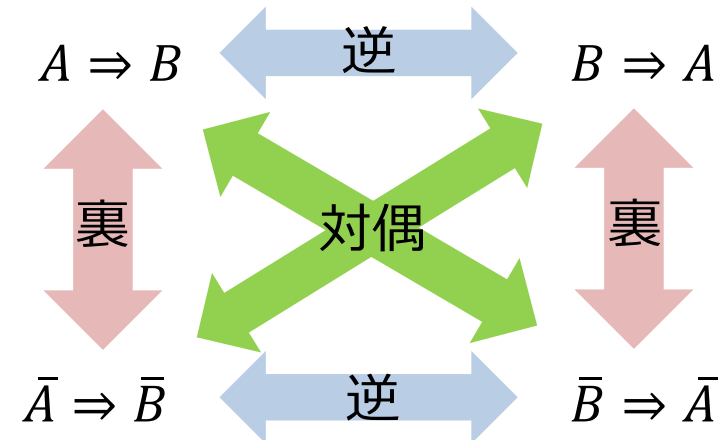
### ■ a prioriな原理

- “a priori” は、ラテン語で  
「経験<sup>けいけん</sup>的認識<sup>しんしき</sup>に先立<sup>さきだ</sup>つ先天<sup>せんてんてき</sup>的、自明<sup>じめいてき</sup>的な認識や概念」
- 今回の場合、「学習データ（経験<sup>けいけん</sup>的認識<sup>しんしき</sup>）に関係なく、  
当たり前<sup>あたりまえ</sup>に成り立つこと」という意味



# 頻出項目抽出：Aprioriアルゴリズム

- 【補足】<sup>めいだい</sup>命題「AならばB」について
- 逆：BならばA
  - <sup>うら</sup>裏：AでないならばBでない
  - **対偶**（逆の裏）：BでないならばAでない
  - 「AならばB」が成立するとき
    - 逆や裏：必ずしも成り立つわけではない
    - **対偶**：必ず成り立つ



# 頻出項目抽出：Aprioriアルゴリズム

## □ 今回の「a prioriな原理」

### 命題

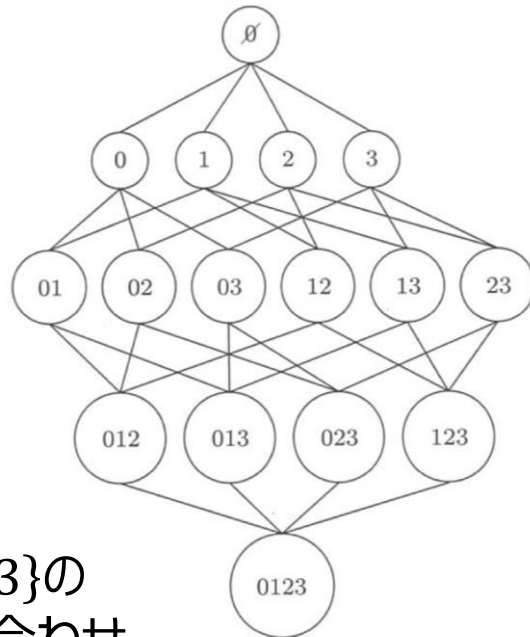
ある項目集合が頻出ならば  
その部分集合も頻出である

### 対偶

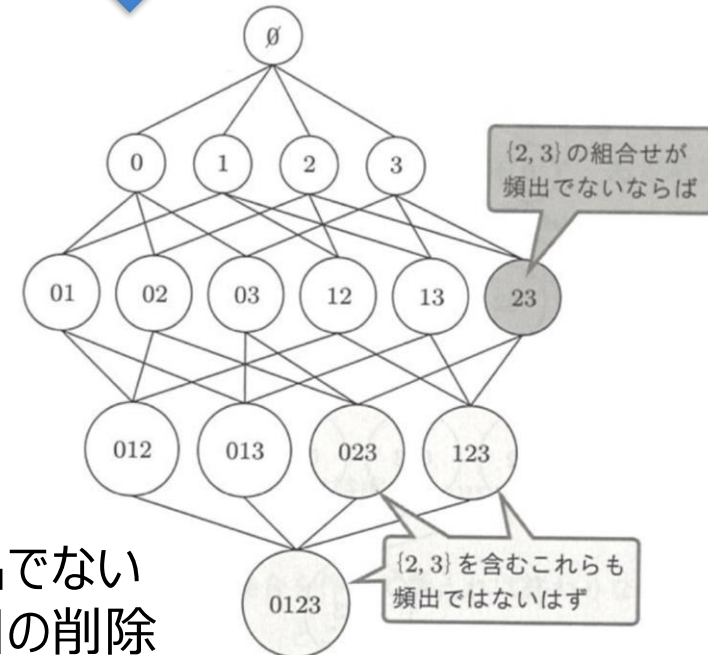
ある項目集合が頻出でないならば  
その項目集合を含む集合も頻出ではない



これを図で表現すると



項目{0,1,2,3}の  
全ての組み合わせ



頻出でない  
項目の削除

# 頻出項目抽出

## □ Aprioriアルゴリズム（頻出項目抽出）

```
入力：正解なしデータ  $D$   
出力：頻出項目集合  
 $F_1 \leftarrow$  要素数1の頻出項目集合  
 $k = 2$   
while  $F_1 \neq \emptyset$  do  
     $C_k \leftarrow F_{k-1}$  の各要素の組み合わせ  
    for all  $x \in D$  do  
        for all  $c \in C_k$  do  
            if  $c \subset x$  then  
                 $c.count \leftarrow c.count + 1$   
            end if  
        end for  
         $F_k \leftarrow \{c \in C_k | c.count > \text{閾値}\}$   
    end for  
     $k \leftarrow k + 1$   
end while  
return  $\bigcup_k F_k$ 
```

$F_k$  : 要素数 $k$ の頻出項目集合

$C_k$  :  $F_{k-1}$ の要素を組み合わせ  
作られる頻出項目集合の候補

※ 予め抽出する項目集合の支持度の  
閾値を決めておき、それを超えるものを  
頻出項目集合とする

# 連想規則抽出：Aprioriアルゴリズム

- 規則の学習では、どの特徴（またはその組み合わせ）が結論部になるかわからない
  - 「商品Aを購入したならば、商品Bを購入することが多い」、「商品Cを購入したならば、商品DとEを購入することが多い」のように得られた規則のそれぞれが、異なった条件部・結論部を持つことが多くなる
- 連想規則の作成手順
  1. Aprioriアルゴリズムで頻出項目を抽出
  2. 頻出項目の要素を、条件部と結論部に分けて可能な規則集合を生成
  3. **規則の有用性**を評価し、役立ちそうなものを絞り込む

# 連想規則抽出：Aprioriアルゴリズム

## □ 連想規則の作成手順（つづき）

### ■ 規則の有用性を評価する基準

#### • 確信度（confidence）

- 規則の条件部が起こったときに結論部が起こる割合
- 確信度が高いほど、この規則に当てはまる事例が多い

$$\text{confidence}(A \Rightarrow B) = \frac{\text{support}(A \cup B)}{\text{support}(A)}$$

#### • リフト値（lift）

- 規則の結論部だけが<sup>たんだく</sup>単独で起こる割合と条件部が起こったときに結論部で起こる割合との比
- リフト値が高いほど、得られる情報の多い規則である

$$\text{lift}(A \Rightarrow B) = \frac{\text{confidence}(A \Rightarrow B)}{\text{support}(B)}$$

# 演習問題12-2（5分間）

- 「ハム→<sup>たまご</sup>卵」という規則について、以下の情報が得られたとする
  - 項目集合 {ハム、卵} の支持度が0.1
  - 「ハム→卵」という規則の確信度が0.7、リフト値：5
  
- 上記の条件において、以下の項目を求めなさい
  1. ハムと卵を同時に購入している客の割合
  2. ハムを購入した客の中で卵も購入する客の割合
  3. 「ハムを既<sup>すで</sup>に買った客が卵を買う確率」は「任意の客が卵を買う確率」の何倍か？

# 連想規則抽出：Aprioriアルゴリズム

□ 「a priori原理」から、評価値の高い規則を絞り込む

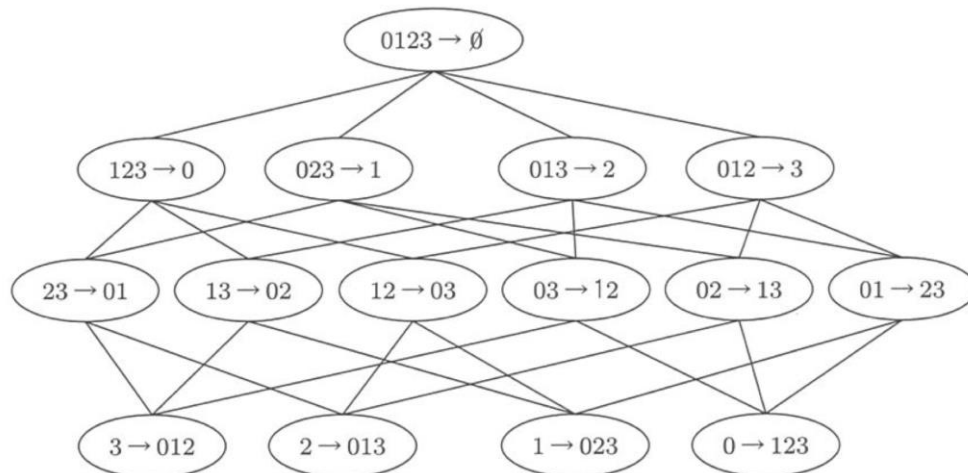
## 命題

ある項目集合を結論とする規則の確信度が高ければ、その部分集合を結論とする規則の確信度も高い

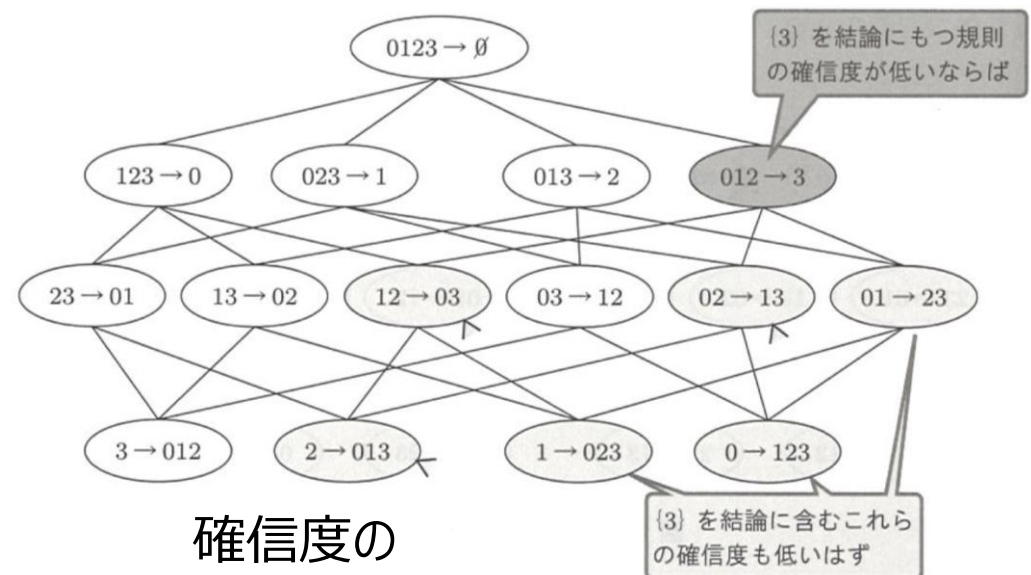
## 対偶

ある項目集合を結論とする規則の確信度が低ければ、その項目集合を含む項目集合を結論とする規則の確信度も低い

↑ 「商品Aを購入したならば、商品B・Cを購入する」という規則の確信度が高ければ  
「商品Aを購入したならば、商品Bを購入する」という規則の確信度も高い...ということ



可能な連想規則集合  
(頻出項目集合が {0,1,2,3} の場合)



確信度の  
低い規則の削除

# 連想規則抽出

## □ Aprioriアルゴリズム

$F_k$  : 要素数 $k$ の頻出項目集合  
 $H_m$  : 要素数 $m$ の結論部の集合

入力 : 頻出項目集合  $F_k$  ( $k \geq 2$ )

出力 : 連想規則集合

**for all**  $f_k \in F_k$  **do**

$H_1 \leftarrow \emptyset$

$A \leftarrow \{a_{k-1} \mid a_{k-1} \subset f_k\}$

**if** confidence( $a_{k-1} \Rightarrow f_k - a_{k-1}$ ) > 閾値 **then**

規則  $a_{k-1} \Rightarrow f_k - a_{k-1}$  を出力

$H_1$  に  $\{f_k - a_{k-1}\}$  を加える

**end if**

ap-genrules( $f_k, H_1$ ) /\* 次スライドのアルゴリズムを参照<sup>さんしょう</sup> \*/

**end for**

※ 予め抽出する連想規則の確信度の閾値を決めておく



# 連想規則抽出

## □ Aprioriアルゴリズム (つづき)

### ■ ap-genrules( $f_k, H_m$ )

**if**  $k > m + 1$  **then**

$H_m$ の要素を組み合わせて $H_{m+1}$ を作成

**for all**  $h_{m+1} \in H_{m+1}$  **do**

**if** confidence( $f_k - h_{m+1} \Rightarrow h_{m+1}$ ) > 閾値 **then**

規則 $f_k - h_{m+1} \Rightarrow h_{m+1}$ を出力

**else**

$H_{m+1}$ から  $h_{m+1}$  を削除

**end if**

**end for**

ap-genrules( $f_k, H_{m+1}$ )

**end if**

# FP-Growthアルゴリズム

## □ FP-Growthアルゴリズム

- Aprioriアルゴリズムを高速化する手法
- トランザクションデータをコンパクトな情報に変換して、そのコンパクトな情報に対してパターンマイニングする
- トランザクションデータをコンパクトにする手順
  1. 特徴を頻出頻度順に並べ替える
  2. 頻度の高い順からその情報をまとめると  
商品Aの購入が100件で、その内、商品Bの同時購入が40件、  
商品Cの同時購入が30件...  
のように多数のトランザクションの情報を<sup>てみじか</sup>手短かに表現
  3. 2. の情報を木構造で表現

# FP-Growthアルゴリズム

## □ FP-Growthアルゴリズムの例（※ 頻出項目抽出の場合）

### ■ 以下のトランザクション集合を学習データとして考える

#### 1. 学習データに対する特徴の出現頻度の計算とフィルタリング

```
1 {r,z,h,j,p}
2 {z,y,x,w,v,u,t,s}
3 {z}
4 {r,x,n,o,s}
5 {y,r,x,z,q,t,p}
6 {y,z,x,e,q,s,t,m}
```

#### 学習データ

出現する文字は特徴名を表す  
(出現したときは、その特徴の値がtとなる)

- ・特徴を出現頻度順にソート
- ・出現頻度が低い特徴を  
フィルタにかけて消去  
(今回は2回以下しか  
現れないものを消去)

```
1 {z,r}
2 {z,x,y,s,t}
3 {z}
4 {x,s,r}
5 {z,x,y,r,t}
6 {z,x,y,s,t}
```

#### ソーティングと フィルタリング後の 学習データ

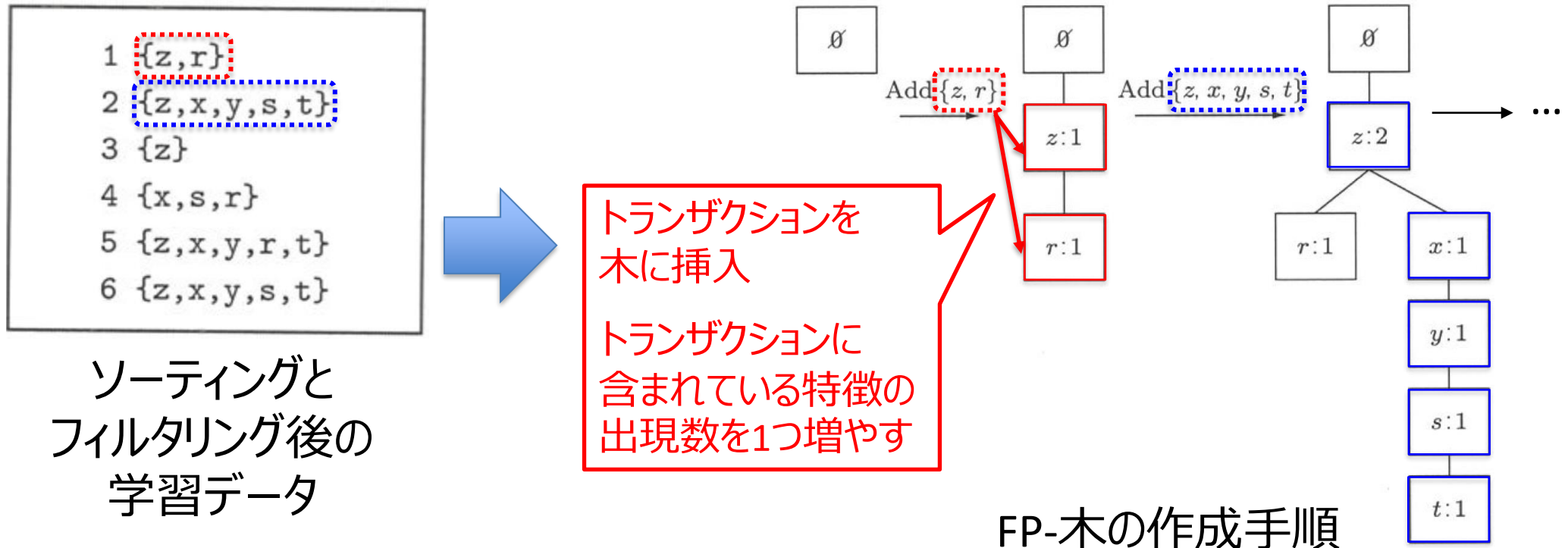
# FP-Growthアルゴリズム

## □ FP-Growthアルゴリズムの例（つづき）

### 2. FP-木（Frequent Pattern Tree）の作成

- 最初にnullというカテゴリを付けた根ノード（ルート）を用意し、トランザクションを順にその木に挿入する

そうにゅう



# FP-Growthアルゴリズム

## □ FP-木挿入アルゴリズム

■  $T$  : トランザクション、 $FP$  : FP-木

トランザクション $T$ の先頭要素 $t$ を取り出す

せんとう

**if**  $t \in FP$  **then**

$t$ に対応するノード $N$ のカウントを1増やす

**else**

ノード $N$ を作成し、カウントを1として $FP$ につなぐ

**end if**

**if**  $T$  に残りの要素がある **then**

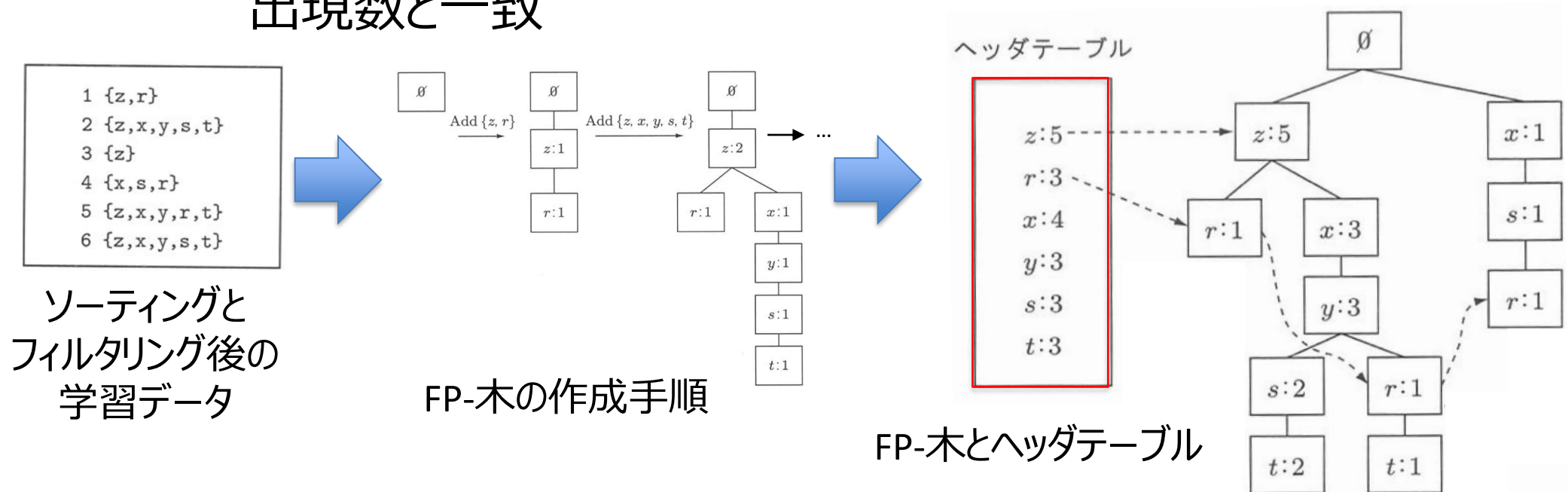
FP-木挿入( $T$ の残りの要素,  $N$ をルートとするFP-木)

**end if**

# FP-Growthアルゴリズム

## □ FP-Growthアルゴリズムの例（つづき）

3. <sup>かんせい</sup>完成したFP-木に対して特徴を見出しとする <sup>みだ</sup>ヘッダテーブルを作成し、その頻度を記録する
4. FP-木に出現する同じ要素をリンクで結ぶ
  - リンクを辿って集めた出現数は、全体のトランザクション集合での出現数と一致



# FP-Growthアルゴリズム

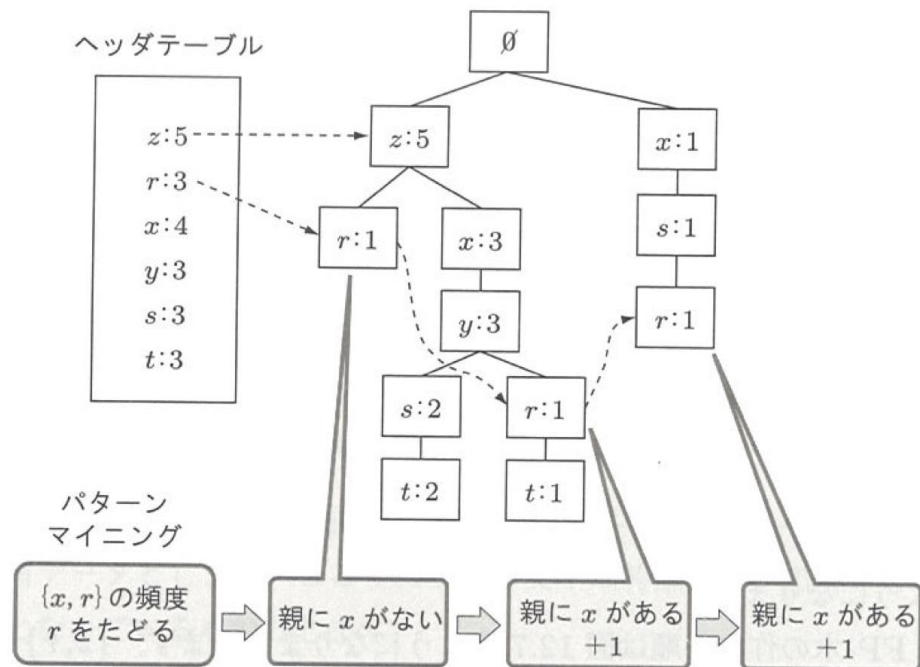
## □ FP-Growthアルゴリズムの例（つづき）

### 5. FP-木に対して、パターンマイニングを行う

例えば

**頻出項目抽出で{x,r}の頻度を求める**

1. ヘッダテーブルから頻度の少ない方を選ぶ
  - rが3回、xが4回なので、rのリンクを辿りながら頻度を計算
2. 最初のrから親を辿る
  - xが無いので、{x,r}が共起していないと判断
3. 次のrのリンクを辿り、その親を辿る
  - xがあったので、rの頻度をカウントに加える
4. 親にxが出現するパスになるrの頻度を足し続けて、最終的な{x,r}の頻度を得る



FP-木のマイニング

# 推薦システムにおける学習

## □ 個人こじんに対して推薦すいせんを行うシステムの構築

- トランザクションデータ（1個のデータが、1件の売りに相当）を個人に対応付けてまとめると、どの個人がどの商品を購入しているのかがわかる

## □ 協調フィルタリング

- 推薦システムにおける学習手法の1つ
  - 新規しんきユーザがある商品を購入した時、購入パターンが似ているユーザを探し、「そのユーザが購入していて、かつ新規ユーザが購入していない商品」を推薦するというのが基本的な考え方
- 購入データを低次元の行列に分解し、ユーザ・商品の特徴を低次元のベクトルで抽出する



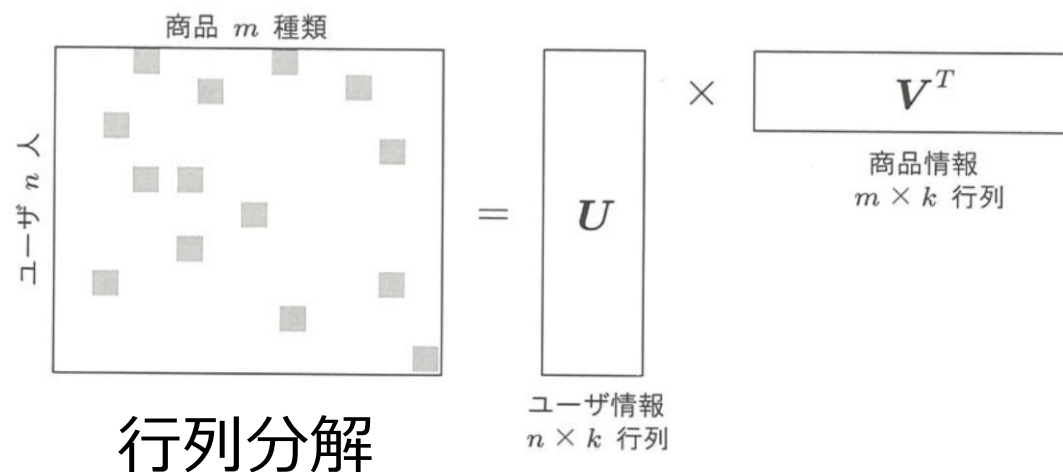
# 推薦システムにおける学習

## □ 協調フィルタリング（つづき）

- 購入データを低次元の行列に分解し、ユーザ・商品の特徴を低次元のベクトルで抽出する

## ■ Matrix Factorization

- 疎らなデータを低次元行列の積に分解
  - Alternating Least Squares手法：値のある要素のみを使って行列分解
  - Non-negative Matrix Factorization：分解した行列の要素が全て非負 ひふ



# 演習問題12-3（10分間）

---

- 実際に協調フィルタリングが用いられている事例を調べなさい