

Complete Mining of Frequent Patterns from Graphs: Mining Graph Data



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Seminar aus Maschinellem Lernen
Hongtao Yan

Übersicht



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Motivation
- Wiederholung Warenkorbanalyse (Apriori Algorithmus)
- Darstellung der graphstrukturierten Daten
 - Adjacency matrix
 - Induced Subgraph
 - Vorverarbeitung von labeled Links und self-Looped Knoten im Graph
 - Matrixcode
- Extraktion von häufigen Graphs
 - Join Operation
 - Kanonische Form
- Leistungsfähigkeit beurteilen
- Anwendung
 - Web Browsing Pattern Analyse
 - Chemische Karzinogenese Analyse

- Erweiterung des Apriori-Algorithmus zum Itemset Mining auf Graphen
- Darstellung der Graphen als "Matrix"
- Mining häufige "induced" Graphs

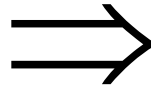
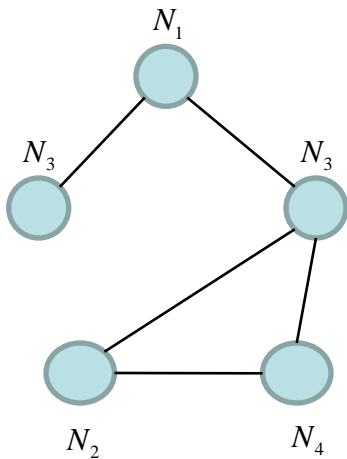
Anwendung

- Web browsing analyse
- Chemische Analyse

Warenkorbanalyse (Apriori Algorithmus)

- Es sei $I = \{i_1, i_2, \dots, i_n\}$ eine Menge von Objekten (Items)
- D sei eine Menge von Transaktionen, wobei jede Transaktion eine Menge von Objekten ist.
- Eine Assoziationsregel ist eine Implikation der Form $X \Rightarrow Y$, wobei X und Y Untermengen von I sind. (Und X und Y keine gemeinsamen Elemente haben)
- Eine Regel $X \Rightarrow Y$ hat den Konfidenzwert c , falls $c\%$ der Transaktionen aus D , die X enthalten auch Y enthalten.
- Eine Regel $X \Rightarrow Y$ hat den Support s , wenn $s\%$ der Transaktionen aus D , X vereinigt Y enthalten.

Adjacency matrix



$$\begin{matrix} & N_1 & N_2 & N_3 & N_3 & N_4 \\ \begin{matrix} N_1 \\ N_2 \\ N_3 \\ N_3 \\ N_4 \end{matrix} & \begin{pmatrix} 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{pmatrix} \end{matrix}$$

$$\text{sup}(G) = \frac{\text{number of transaction including an induced subgraph } G}{\text{total number of transactions}}$$

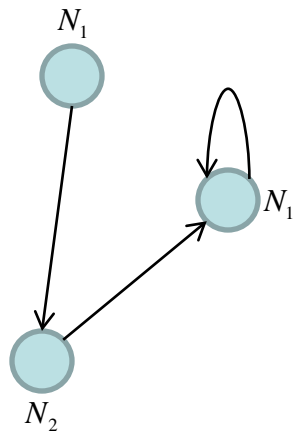
$$\text{sup}(B \Rightarrow H) = \text{sup}(B \cup H) \quad \text{conf}(B \Rightarrow H) = \frac{\text{sup}(B \cup H)}{\text{sup}(B)}$$

Induced Subgraph

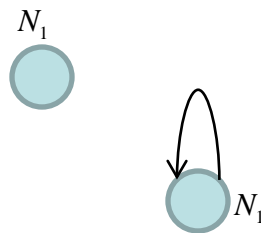
- Formale Definition von „induced Graph“

$$V(G') \subseteq V(G), \quad E(G') \subseteq E(G)$$

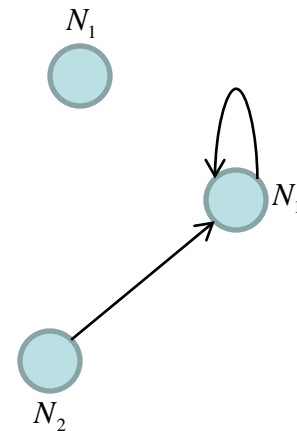
$$\forall u, v \in V(G') \{u, v\} \in E(G) \Leftrightarrow \{u, v\} \in E(G')$$



(a)



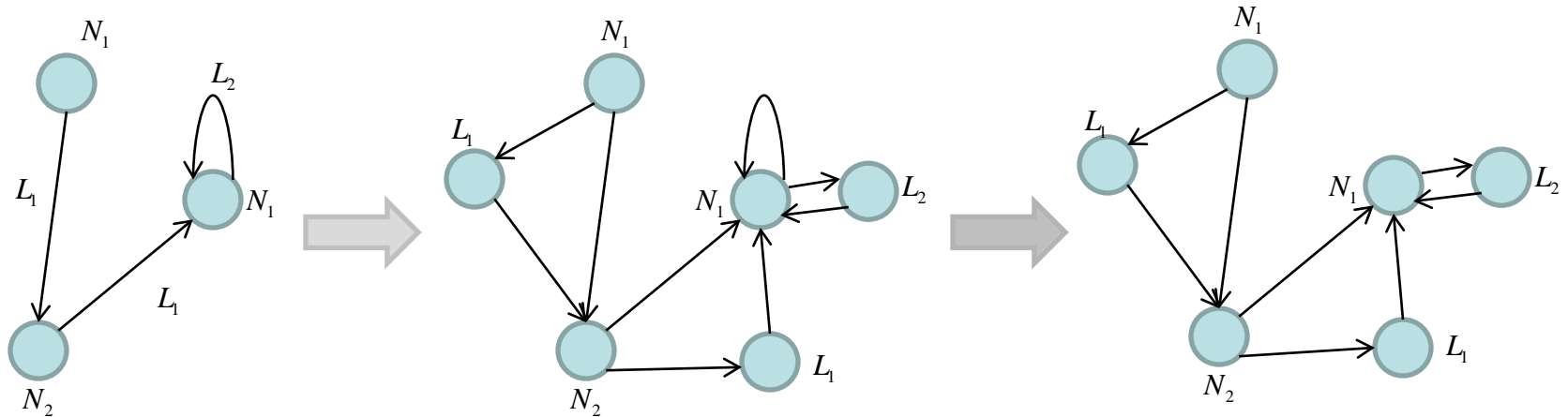
(b)



(c)

Labeled Links und self-Loop

Umwandlung der labeled Links und self-looped Knoten



$$\begin{matrix} N_1 & N_1 & N_2 \\ N_1 \begin{pmatrix} L_2 & 0 & 0 \\ 0 & 0 & L_1 \\ L_1 & 0 & 0 \end{pmatrix} \\ N_2 \end{matrix} \Rightarrow \begin{matrix} N_1 & N_1 & N_2 & L_1 & L_1 & L_2 \\ N_1 \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \\ N_2 \\ L_1 \\ L_1 \\ L_2 \end{matrix} \Rightarrow \begin{matrix} N_1' & N_1 & N_2 & L_1 & L_1 & L_2 \\ N_1' \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \\ N_1 \\ N_2 \\ L_1 \\ L_1 \\ L_2 \end{matrix}$$

Matrixcode

$$X_k = \begin{pmatrix} 0 & x_{1,2} \downarrow & x_{1,3} \downarrow & \cdots & x_{1,k} \downarrow \\ x_{2,1} & 0 & x_{2,3} \downarrow & \cdots & x_{2,k} \\ x_{3,1} & x_{3,2} & 0 & \cdots & x_{3,k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{k,1} & x_{k,2} & x_{k,3} & \cdots & 0 \end{pmatrix}$$

Ungerichtet Graph:

$$code(X_k) = x_{1,2}x_{1,3}x_{2,3}x_{1,4} \cdots x_{k-2,k}x_{k-1,k}$$

Bsp.:

$$\begin{matrix} & N_1 & N_2 & N_3 & N_3 & N_4 \\ \begin{matrix} N_1 \\ N_2 \\ N_3 \\ N_3 \\ N_4 \end{matrix} & \begin{pmatrix} 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{pmatrix} \end{matrix}$$

$$code(X_k) = 0101100101$$

Gerichtet Graph:

$$code(X_k) = c_{1,2}c_{1,3}c_{2,3}c_{1,4} \cdots c_{k-2,k}c_{k-1,k}$$

wobei

$$c_{i,j} = 2x_{j,i} + x_{i,j}$$

Extraktion häufigen induced Subgraph

- Häufig Subgraph mining ist ähnlich zum Itemset Mining
Subgraph G mit k Knoten kann nur frequent sein, wenn alle Subgraphen von G mit $k-1$ Knoten auch frequent sind
- Zwei Subgraph können nur kombinieren, wenn die folgend drei Bedingung erfüllt sind
- Constraint 1:

$$X_K = \begin{pmatrix} X_{K-1} & x_1 \\ x_2^T & 0 \end{pmatrix} \quad Y_K = \begin{pmatrix} X_{K-1} & y_1 \\ y_2^T & 0 \end{pmatrix}$$

$$Z_{K+1} = \begin{pmatrix} X_{K-1} & x_1 & y_1 \\ x_2^T & 0 & z_{k,k+1} \\ y_2^T & z_{K+1,k} & 0 \end{pmatrix} = \left(\begin{array}{c|c} X_k & y_1 \\ \hline y_2^T & z_{k,k+1} \\ \hline z_{K+1,k} & 0 \end{array} \right)$$

Constraint 2:

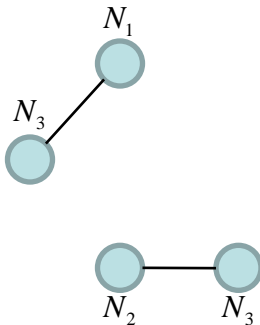
$N(X_k, i)$: Label von i -ten Knoten in Matrix X_k .

$$N(X_k, i) = N(Y_k, i) = N(Z_{k+1}, i) \text{ and } N(X_k, i) \leq N(X_k, i+1) \text{ for } i = 1, \dots, k-1$$

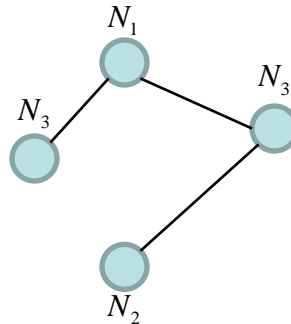
$$N(X_k, k) = N(Z_{k+1}, k), \quad N(Y_k, k) = N(Z_{k+1}, k+1), \text{ and } N(X_k, k) \leq N(Y_k, k)$$

Constraint

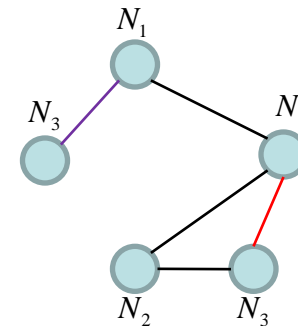
- Constraint 3: code(the first matrix) \leq code(the second matrix)
- Die drei Constraint heißt „join operation“
- Matrix, die durch „join operation“ erzeugt ist, „normal Form“
z.B.:



$$X_4 = \begin{matrix} & N_1 & N_2 & N_3 & N_3 \\ N_1 & \boxed{0} & \boxed{0} & \boxed{1} & \boxed{0} \\ N_2 & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{1} \\ N_3 & \boxed{1} & \boxed{0} & \boxed{0} & \boxed{0} \\ N_3 & \boxed{0} & \boxed{1} & \boxed{0} & \boxed{0} \end{matrix}$$



$$Y_4 = \begin{matrix} & N_1 & N_2 & N_3 & N_3 \\ N_1 & \boxed{0} & \boxed{0} & \boxed{1} & \boxed{1} \\ N_2 & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{1} \\ N_3 & \boxed{1} & \boxed{0} & \boxed{0} & \boxed{0} \\ N_3 & \boxed{1} & \boxed{1} & \boxed{0} & \boxed{0} \end{matrix}$$



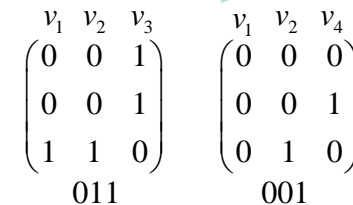
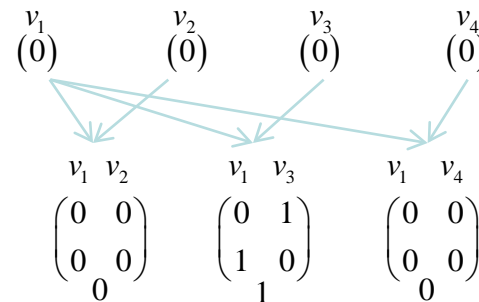
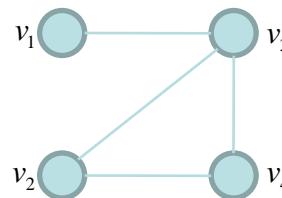
$$Z_5 = \begin{matrix} & N_1 & N_2 & N_3 & N_3 & N_3 \\ N_1 & \boxed{0} & \boxed{0} & \boxed{1} & \boxed{0} & \boxed{1} \\ N_2 & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{1} & \boxed{1} \\ N_3 & \boxed{1} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} \\ N_3 & \boxed{0} & \boxed{1} & \boxed{0} & \boxed{0} & \boxed{?} \\ N_3 & \boxed{1} & \boxed{1} & \boxed{0} & \boxed{?} & \boxed{0} \end{matrix}$$

Normalform

- Direktes Vergrößern der Graphen um jeweils
 - Finde alle Graphen mit k Knoten und erweitere diese um einen weiteren Knoten \Rightarrow Kandidaten für $k+1$ elementare Graphen

$$\begin{matrix}
 v_1 & v_2 & v_3 & v_4 \\
 v_1 & \begin{pmatrix} 0 & 0 & 1 & 0 \end{pmatrix} \\
 v_2 & \begin{pmatrix} 0 & 0 & 1 & 1 \end{pmatrix} \\
 v_3 & \begin{pmatrix} 1 & 1 & 0 & 1 \end{pmatrix} \\
 v_4 & \begin{pmatrix} 0 & 1 & 1 & 0 \end{pmatrix}
 \end{matrix}$$

011011
 X_4



$$\begin{matrix}
 v_1 & v_2 & v_4 & v_3 \\
 \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \\
 001111
 \end{matrix}$$

- Generierung durch entfernen dem i -th Knoten \Rightarrow non-normal form Matrix
- Ein Methode wandelt non-normal form zu normal Form um

Kanonisch Form

- Nach alle Kandidaten Graph erzeugt werden, überprüfen ihre Supportwert
- Problem: mehre Normalform bilden identisch Graph ab (z.B.)

$$X_5 = \begin{matrix} & N_1 & N_1 & N_1 & N_1 & N_1 \\ N_1 & \begin{pmatrix} 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 \end{pmatrix} \end{matrix} \quad Y_5 = \begin{matrix} & N_1 & N_1 & N_1 & N_1 & N_1 \\ N_1 & \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 \end{pmatrix} \end{matrix}$$

- Falsche Ergebnis von Supportwert
- Lösung: Kanonisch Form

$$\Gamma(G) = \{X_k \mid X_k \text{ is an normal form matrix and } G \equiv G(X_k)\}$$

$$C_k \text{ w.r.t } code(C_k) = \min_{X_k \in \Gamma(G)} code(X_k)$$

- Erzeugung durch Permutation der Zeile und Spalten der normalen Form

$$C_k = W_k^T X_k W_k$$

- Annahme: Permutationsmatrix für jeden häufigen induced $k-1$ Supgraph sind bekannt

- Theorem 1: Die erste Matrix einem Kanonischen Form Matrix ist auch ein Kanonischen Form Matrix

$$C_k = \begin{pmatrix} C_{k-1} & C_1 \\ C_2^T & 0 \end{pmatrix}$$

- Theorem 2:

$$\text{code}(C_{k-1}) \leq \text{code}(\text{can}(X_{k-1}^m)) \quad (1 \leq m \leq k, N(X_k, m) = N(X_k, k))$$

- Definition: „pseudo- canonical form“

$$C_k^p = \begin{pmatrix} C_{k-1} & C_1^p \\ C_2^{pT} & 0 \end{pmatrix}$$

$$G(C_k^p) \equiv G(C_k) \equiv G(X_k) \quad \text{code}(C_k) \leq \text{code}(C_k^p)$$

- X_{k-1}^m ($1 \leq m \leq k$) wird zum normale Form umgewandelt.

Permutationsmatrix: T_{k-1}^m .

- Zum Kanonische Form umgewandelt. Permutationsmatrix: S_{k-1}^m

- T_k^m und S_k^m wird wie folgende Gleichung definiert.

$$s_{i,j} = \begin{cases} s_{i,j}^m & 0 \leq i \leq k-1 \text{ and } 0 \leq j \leq k-1 \\ 1 & i = k \text{ and } j = k \\ 0 & \text{otherwise} \end{cases} \quad t_{i,j} = \begin{cases} t_{i,j}^m & i < m \text{ and } j \neq k \\ t_{i-1,j}^m & i > m \text{ and } j \neq k \\ 1 & i = m \text{ and } j \neq k \\ 0 & \text{otherwise} \end{cases}$$

- nach Theorem 3: ein „pseudo-canonical form“

$$C_k^p \text{ w.r.t } code(C_k^p) = \min_m code((T_k^m S_k^m)^T X_k (T_k^m S_k^m))$$

- Und ein Kanonische Form wie folgend

$$C_k \text{ w.r.t } code(C_k) = \min_{U_k \in \Lambda(C_k^p)} code((U_k^T C_k^p U_k))$$

Beispiel



$$X_5 = N_1 \begin{pmatrix} N_1 & N_1 & N_1 & N_1 & N_1 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 \end{pmatrix}$$

$$Y_5 = N_1 \begin{pmatrix} N_1 & N_1 & N_1 & N_1 & N_1 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 \end{pmatrix}$$

$$X_4^m \text{ for } m=1,2,3,4,5 \quad X_4^1 = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} X_4^2 = \begin{pmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} X_4^3 = \begin{pmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} X_4^4 = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{pmatrix} X_4^5 = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix}$$

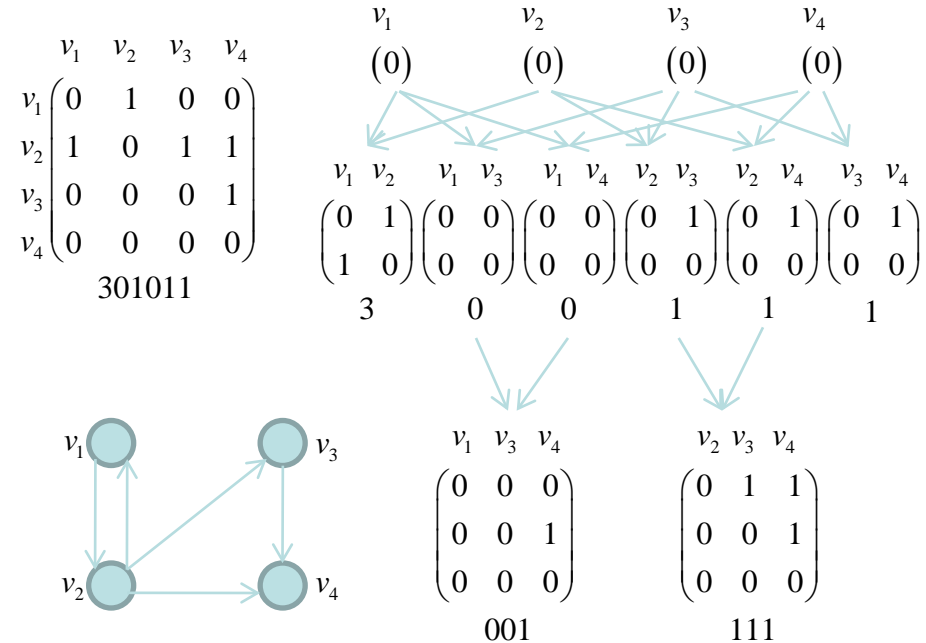
$$T_4^1 = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{pmatrix} \text{ and } S_4^1 = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{pmatrix} \Rightarrow T_5^1 = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}, S_5^1 = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

- Für $m=1,2,3,4$ und 5 $code((T_5^m S_5^m)^T X_5 (T_5^m S_5^m)) = 0011010011, 0011110010, 0111110010, 0000111011, 0000111101$
- $m=4$ ist Minimum C_5^p ist gefunden
- Y_5 ist kanonische Form

- 1) **forall** X_k in a set of candidate frequent induced subgraphs
- 2) $X'_k = X_{k^m}$ $m \leq 1$
- 3) **while** $m \leq k$ **do begin**
- 4) **if** $(N(X_{k^m}, m) = N(X_{k^m}, k))$ **then do begin**
- 5) **if** $(code(X_k) > code((T_k^m S_k^m)^T X_k (T_k^m S_k^m))$ **then do begin**
- 6) $X'_k = (T_k^m S_k^m)^T X_k (T_k^m S_k^m)$;
- 7) **if** (the transformation matrix of X'_k to the canonicalform is known) **then do begin**
- 8) $X'_K = S_K^T X'_k S_K$;
 //where S_k is the matrix to transform X'_k in r.h.s. to its canonicalform./
- 9) **break**;
- 10) **end**
- 11) **end**
- 12) **end**
- 13) $m=m+1$
- 14) **end**
- 15) **if** (the canonicalform of X_K has not been derived in the step 8)
- 16) $X'_k = \text{permutation}(X'_k)$;
- 17) **end**
- 18) Canonicalform of X_k is X'_k ;
- 19) **end**

Generierung häufigen Subgraphs

- Annahme: 2 sitz Subgraph, die code 3 ist, ist nicht häufig Subgraph
- Kombination wenn first Matrix gleich ist
- Kombination wenn Knotenlabel in erst Matrix gleich sind.

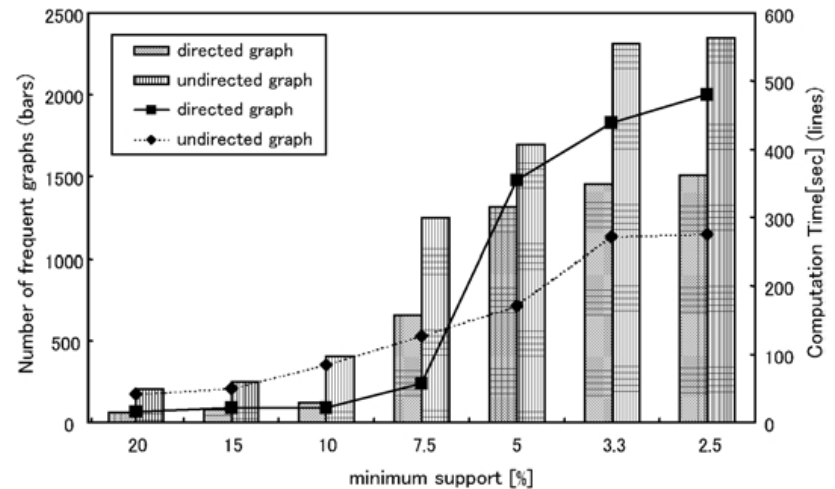
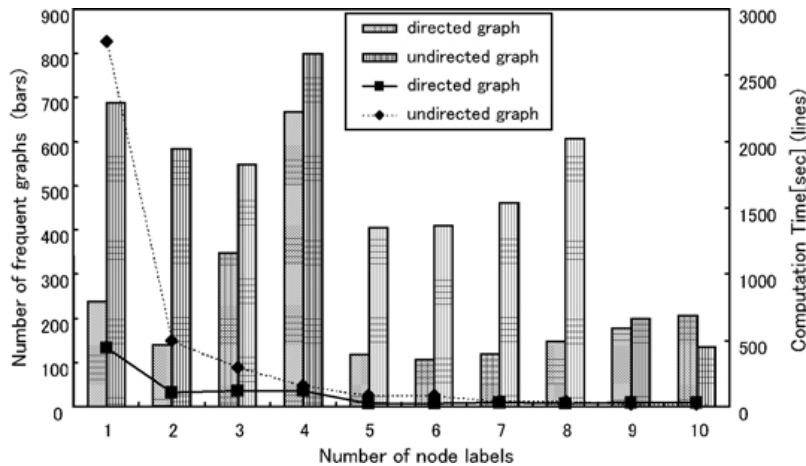
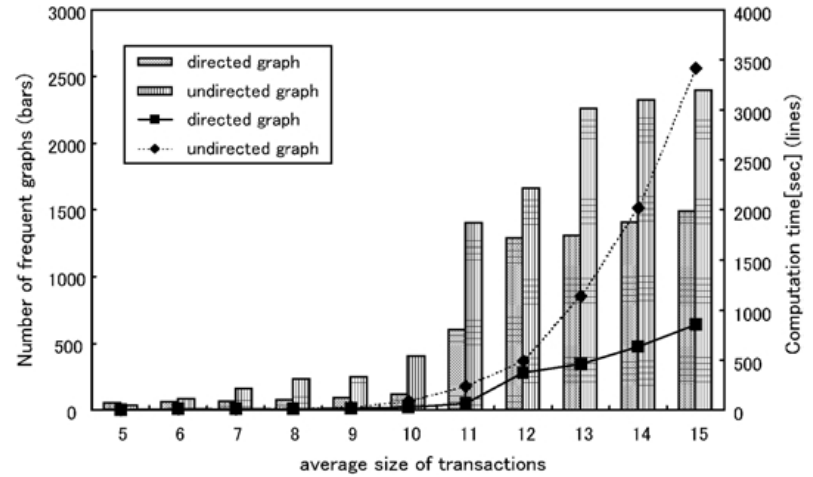
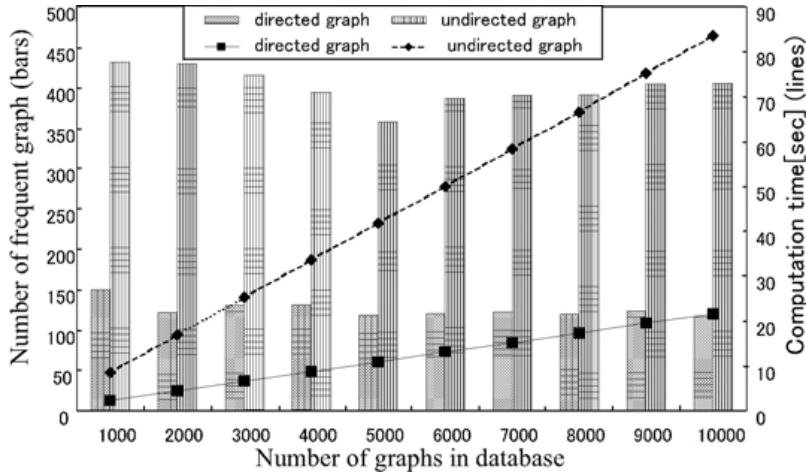


Leistungsbewertung

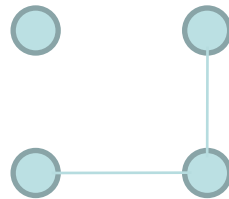
Pentium CPU 400MHz, Ram 128MB

Parameter	Definition	Default value
D	Number of transaction	10,000
T	Average transaction size	10
L	Number of basic patterns	10
I	Average basic patterns size	4
N	Number of node labels	5
p	Link existence probability	50%
minsup	Minimum support	10%

Leistungsbewertung

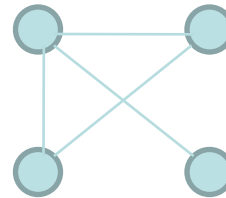


- Ein komplett Graph ist ein Graph, in dem ein Link für jeder Knoten existiert.



000011
000101
000110
010100
100010

(a)



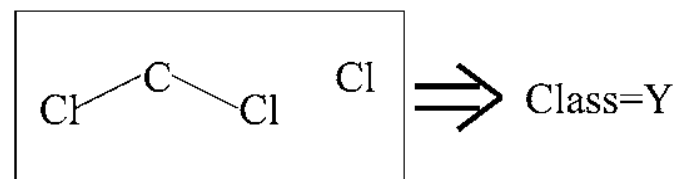
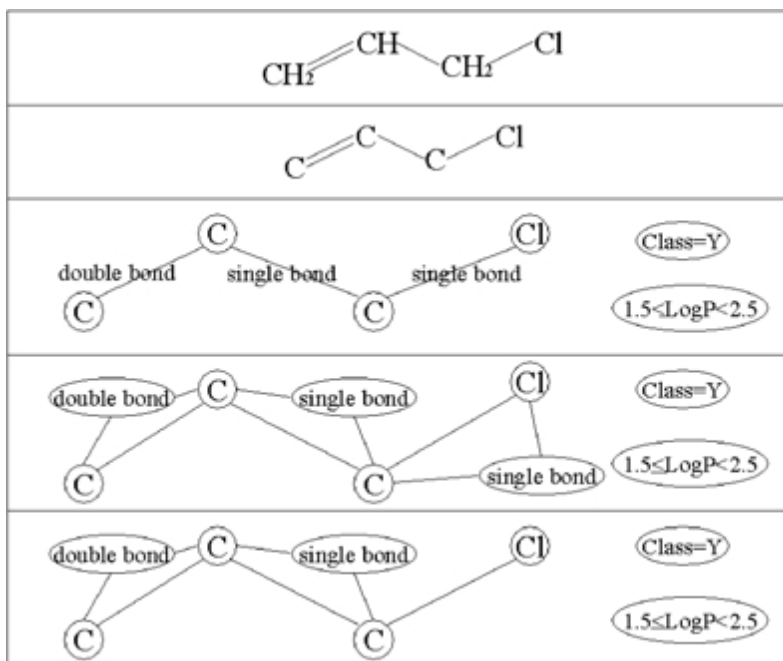
111001
110110
010111
001111

(b)

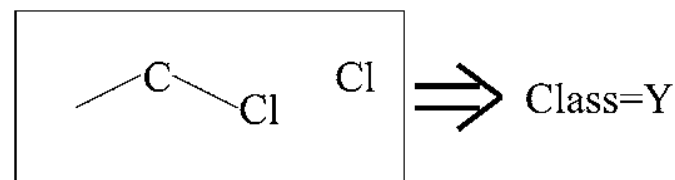
Komplementär

Anwendung

- Web Browsing Analyse
- Chemische Karzinogenes Analyse



Ex1. support=31.7%, confidence=86.7%



Ex2. support=36.6%, confidence=83.3%



Danke für Ihre Aufmerksamkeit