

Rekonstrukce databázového modelu na základě nepřesných dat

Martin Řimnáč

ITAT 2004



Ústav informatiky
Akademie věd ČR



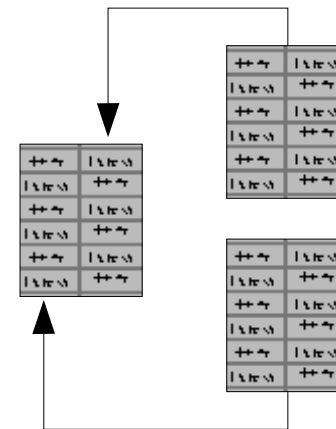
Katedra kybernetiky
České vysoké učení technické

Dekompozice modelu

Nalezení takového databázového modelu, který odpovídá

- Požadované normální formě
- Vstupním datům (složená relace)

Ařes	++ *	řřes	++ *
++ *	řřes	++ *	řřes
Ařes	++ *	řřes	++ *
++ *	řřes	++ *	řřes
Ařes	++ *	řřes	++ *
++ *	řřes	++ *	řřes
Ařes	++ *	řřes	++ *
++ *	řřes	++ *	řřes
Ařes	++ *	řřes	++ *
++ *	řřes	++ *	řřes
Ařes	++ *	řřes	++ *
++ *	řřes	++ *	řřes
Ařes	++ *	řřes	++ *
++ *	řřes	++ *	řřes



Metoda zahrnuje:

- Detekci funkčních závislostí
- Odstranění redundantních závislostí
- Generování databázového modelu

Funkční závislosti

Atribut Y je funkčně závislý na atributu X ($X \rightarrow Y$), pokud vzoru (X) odpovídá právě jeden obraz (Y).

X	Y
+	+
+	+
+	+
+	+
+	+
+	+

Vlastnosti:

- Transitivita
- Hierarchie

$X \rightarrow Y$ a $Y \rightarrow Z$, pak $X \rightarrow Z$
 $X \rightarrow Y$, pak i $\{X, Z\} \rightarrow Y$

Značení:

- Není funkčně závislý
- Vzájemně závislý

$X \not\rightarrow Y$

$X \leftrightarrow Y$

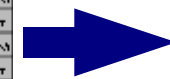
$X \rightarrow Y$ a $Y \rightarrow X$

Detekce funkční závislosti

Detekce funkční závislosti $X \rightarrow Y$:

- Seřazení záznamů podle X, sekundárně Y
- Testování všech záznamů podle definice

X	Y	Z	W
1	1	1	1
2	1	1	1
2	1	1	1
3	2	1	1
4	2	1	1
1	1	1	1
2	1	1	1
2	1	1	1
3	2	1	1
4	2	1	1
1	1	1	1
2	1	1	1
2	1	1	1
3	2	1	1
4	2	1	1



- $X \rightarrow Y$
- $X \rightarrow Z$
- $X \rightarrow W$
- $W \leftrightarrow Z$

Složitost:

- Rekonstrukce tabulky $o(1)$ až $o(NC)$
- Seřazení záznamů $o(N \log(N))$
- Testování rovnosti $o(N)$

X	Y
1	1
2	1
2	1
3	2
4	2

$X \rightarrow Y$

X	Y
1	1
2	1
2	1
2	2
3	2

$X \not\rightarrow Y$

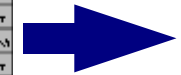
Jednoatributové funkční závislosti: $C(C-1)$ testů.

Detekce funkční závislosti

Detekce funkční závislosti $X \rightarrow Y$:

- Seřazení záznamů podle X , sekundárně Y
- Testování všech záznamů podle definice

X	Y	Z	W
1	1	1	1
1	1	2	1
1	1	3	1
1	1	4	1
1	2	1	1
1	2	2	1
1	2	3	1
1	2	4	1
2	1	1	1
2	1	2	1
2	1	3	1
2	1	4	1
2	2	1	1
2	2	2	1
2	2	3	1
2	2	4	1
2	3	1	1
2	3	2	1
2	3	3	1
2	3	4	1
2	4	1	1
2	4	2	1
2	4	3	1
2	4	4	1



- $X \rightarrow Y$
- $X \rightarrow Z$
- $X \rightarrow W$
- $W \leftrightarrow Z$

Složitost:

- Rekonstrukce tabulky $o(1)$ až $o(NC)$
- Seřazení záznamů $o(N \log(N))$
- Testování rovnosti $o(N)$

X	Y
1	1
2	1
2	1
3	2
4	2

$X \rightarrow Y$

X	Y
1	1
2	1
2	1
2	2
3	2

$X \not\rightarrow Y$

Jediný chybný záznam rozhodne o neplatnosti funkční závislosti

Jednoatributové funkční závislosti: $C(C-1)$ testů.

Matice závislostí

Popis funkčních závislostí

- Pomocí matice závislostí M
- Omezení pouze na jednoatributové závislosti

Matice závislostí

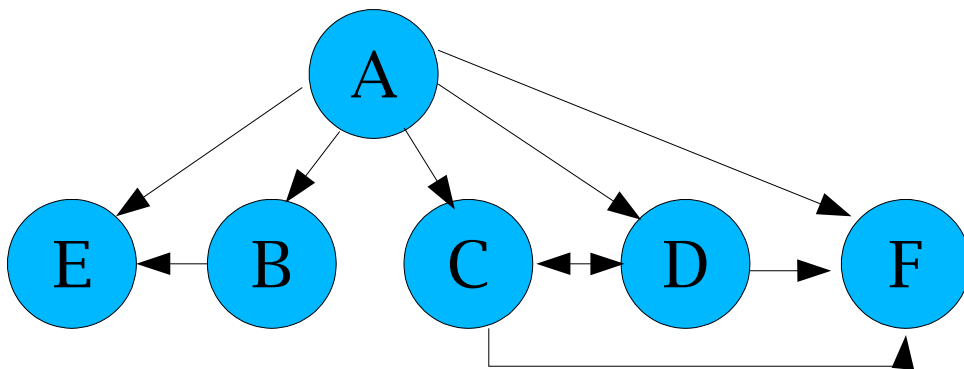
- $M_{ij} = +1$ pokud $A_i \rightarrow A_j$ a $A_j \not\rightarrow A_i$
- $M_{ij} = -1$ pokud $A_i \not\rightarrow A_j$ a $A_j \rightarrow A_i$
- $M_{ij} = 0$ jinak

Ilustrační příklad

A	B	E	C	D	F
Výrobek	Podkategorie	Kategorie	Cena	Cena+DPH	Cena-Akce
...
...

Detekované funkční závislosti:

- $A \rightarrow B$, $A \rightarrow E$, $A \rightarrow C$, $A \rightarrow D$, $A \rightarrow F$, $B \rightarrow E$, $C \leftrightarrow D$, $C \rightarrow F$



	A	B	E	C	D	F
A	0	1	1	1	1	1
B	-1	0	1	0	0	0
E	-1	-1	0	0	0	0
C	-1	0	0	0	0	1
D	-1	0	0	0	0	1
F	-1	0	0	-1	-1	0

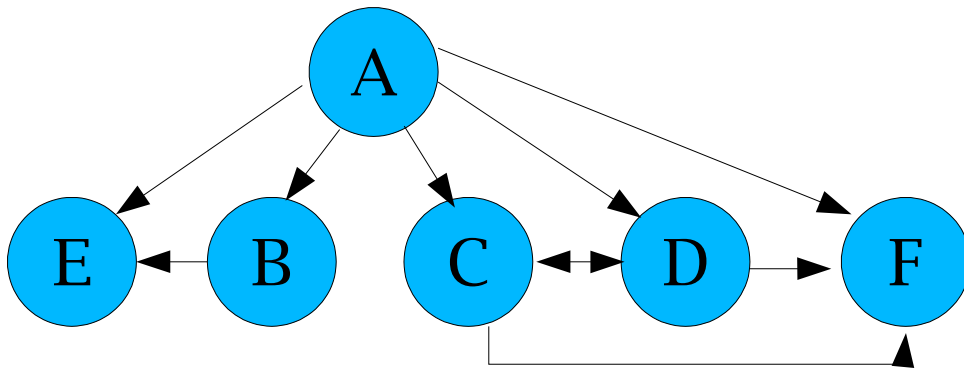
Uspořádání na matici závislostí

Definujme uspořádání atributů na základě tak, že

$$\sum_k m_{ik} > \sum_k m_{jk}, \text{ pak } i < j$$

Vlastnost (díky tranzitivitě):

Pokud $A_i \rightarrow A_j$, pak $i < j$.



	A	B	E	C	D	F
A	0	1	1	1	1	1
B	-1	0	1	0	0	0
E	-1	-1	0	0	0	0
C	-1	0	0	0	0	1
D	-1	0	0	0	0	1
F	-1	0	0	-1	-1	0
Σ	-5	0	2	1	1	3

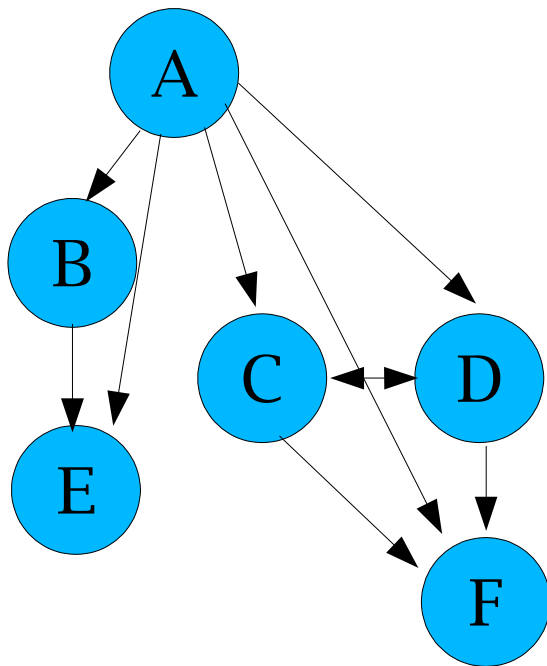
Uspořádání na matici závislostí

Definujme uspořádání atributů na základě tak, že

$$\sum_k m_{ik} > \sum_k m_{jk}, \text{ pak } i < j$$

Vlastnost (díky tranzitivitě):

Pokud $A_i \rightarrow A_j$, pak $i < j$.



	A	B	C	D	E	F
A	0	1	1	1	1	1
B	-1	0	0	0	1	0
C	-1	0	0	0	0	1
D	-1	0	0	0	0	1
E	-1	-1	0	0	0	0
F	-1	0	-1	-1	0	0
Σ	-5	0	1	1	2	3

Vzájemné závislosti

Vzájemná závislost atributů $A_i \leftrightarrow A_j$

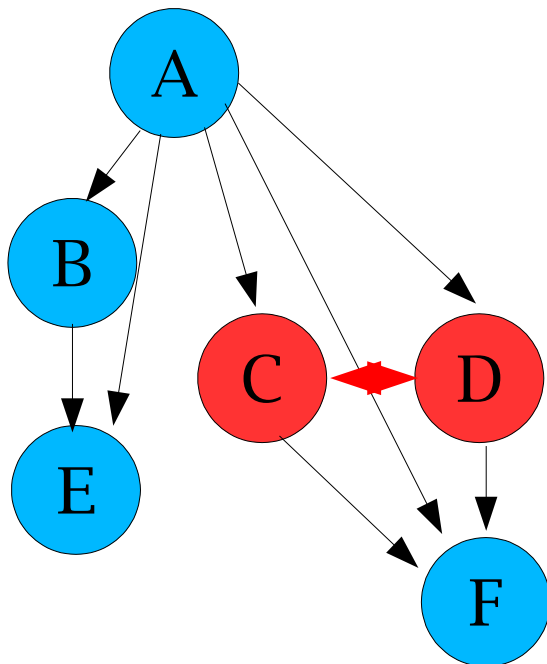
$A_i \rightarrow A_j$ a zároveň $A_j \rightarrow A_i$

Platí:

$$\sum_k m_{ik} = \sum_k m_{jk}$$

$$m_{ik} = m_{jk}$$

$$m_{ki} = m_{kj}$$



	A	B	C	D	E	F
A	0	1	1	1	1	1
B	-1	0	0	0	1	0
C	-1	0	0	0	0	1
D	-1	0	0	0	0	1
E	-1	-1	0	0	0	0
F	-1	0	-1	-1	0	0

Vzájemné závislosti

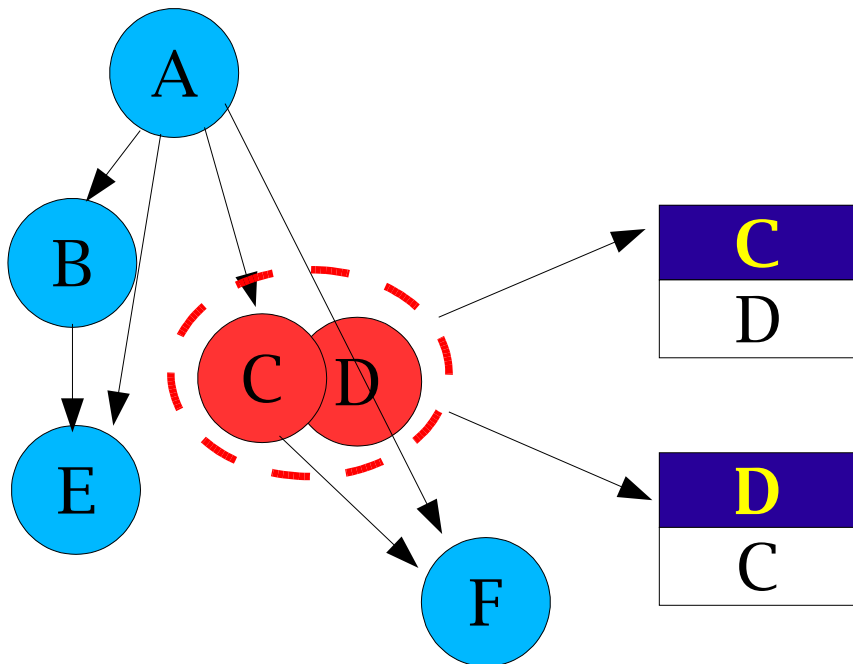
Vzájemná závislost atributů $A_i \leftrightarrow A_j$

Vytvoření vzájemné relace

Neurčitost primárního klíče – výběr libovolného atributu

Redukce matice (omezení stavového prostoru)

Ve schématu zůstane pouze atribut představující primární klíč



	A	B	CD	E	F
A	0	1	1	1	1
B	-1	0	0	1	0
CD	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Vzájemné závislosti

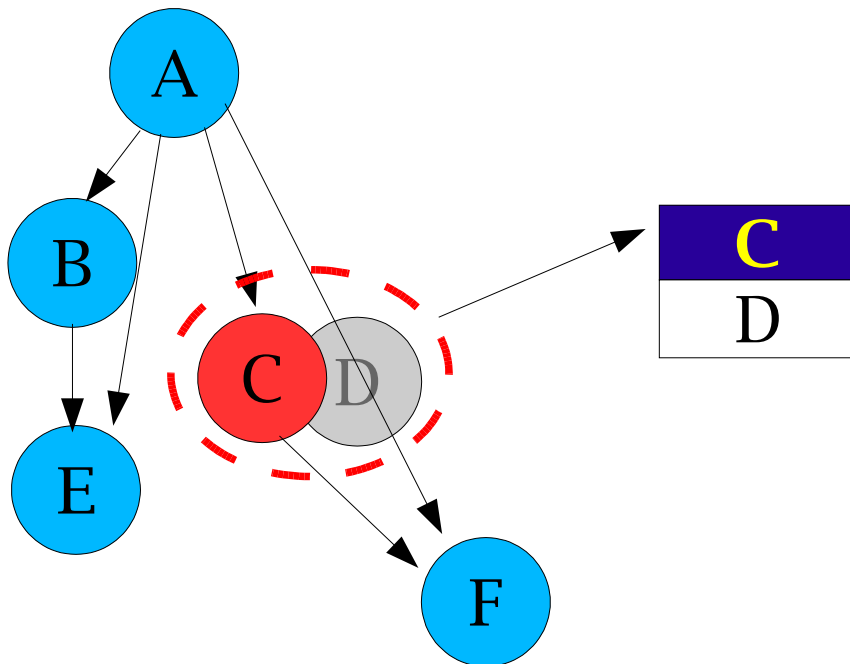
Vzájemná závislost atributů $A_i \leftrightarrow A_j$

Vytvoření vzájemné relace

Neurčitost primárního klíče – výběr libovolného atributu

Redukce matice (omezení stavového prostoru)

Ve schématu zůstane pouze atribut představující primární klíč



	A	B	C	E	F
A	0	1	1	1	1
B	-1	0	0	1	0
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Antisymetrie matice

Matice obsahuje pouze jednosměrné funkční závislosti

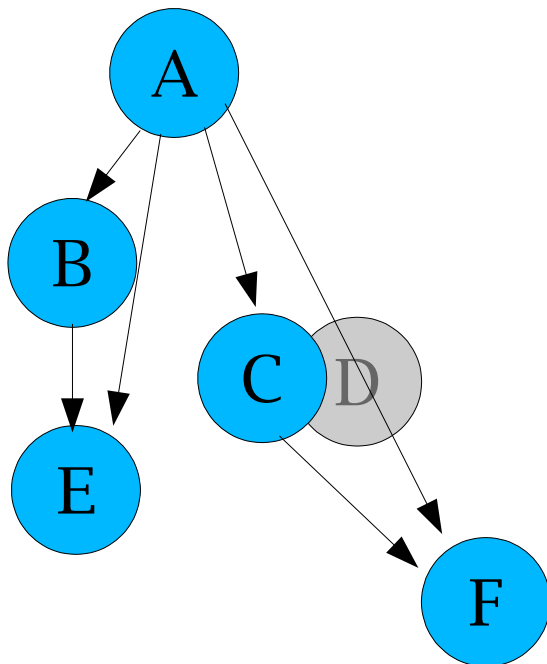
$$A_i \rightarrow A_j \text{ a } A_j \not\rightarrow A_i$$

Antisymetrická matice

$$m_{ij} = -m_{ji}$$

Kladné prvky nad diagonálou, nulová diagonála

Redukce

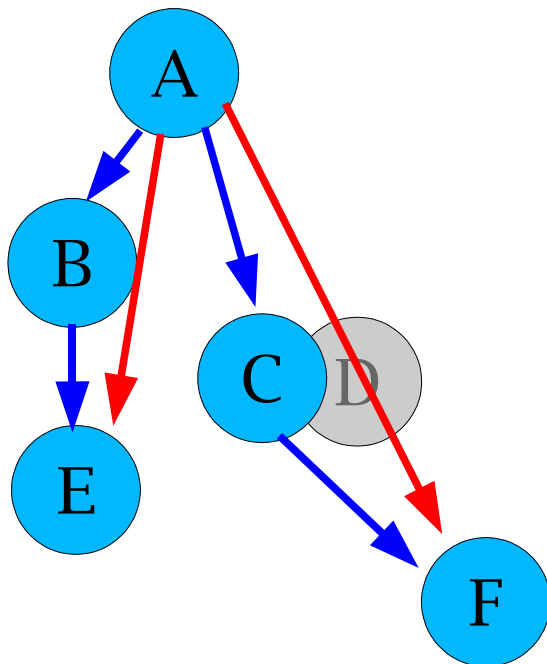


C
D

	A	B	C	E	F
A	0	1	1	1	1
B	-1	0	0	1	0
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Tranzitivita

Redundance díky tranzitivitě
 $X \rightarrow Y$ a $Y \rightarrow Z$, pak $X \rightarrow Z$



C
D

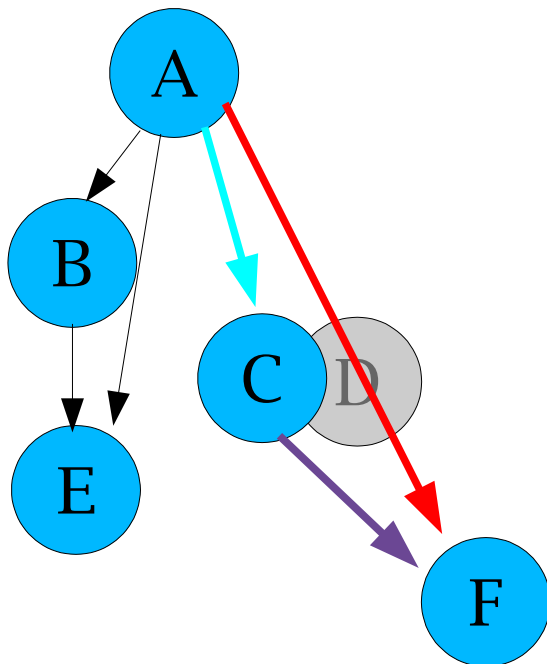
	A	B	C	E	F
A	0	1	1	1	1
B	-1	0	0	1	0
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Odstranění redundance

Pokud $m_{ij} = +1$,

pak v i -tém řádku nuluj všechny $+1$,

pokud ve stejném sloupci je v j -tém řádku $+1$.



C
D

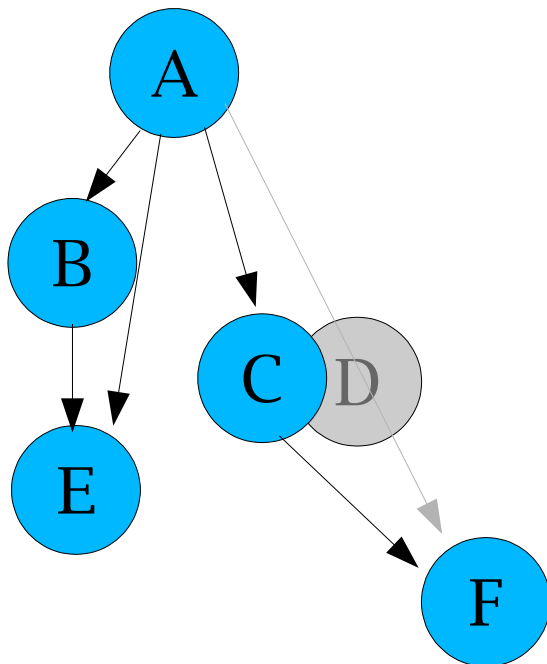
	A	B	C	E	F
A	0	1	1	1	1
B	-1	0	0	1	1
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Odstranění redundance

Pokud $m_{ij} = +1$,

pak v i -tém řádku nuluj všechny $+1$,

pokud ve stejném sloupci je v j -tém řádku $+1$.



C
D

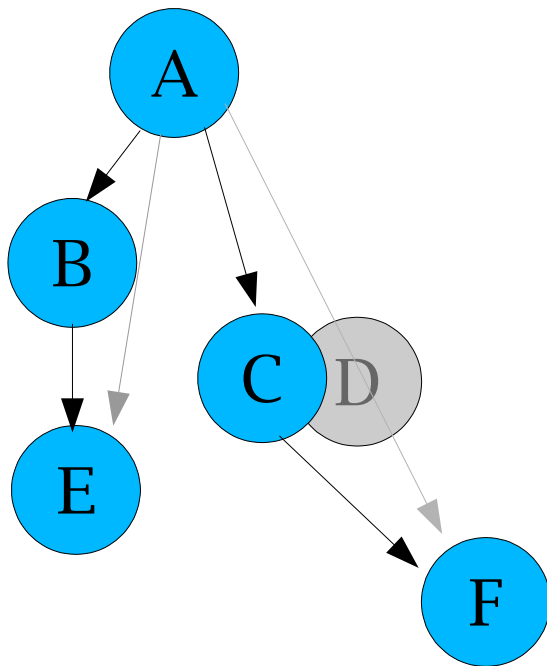
	A	B	C	E	F
A	0	1	1	1	0
B	-1	0	0	1	0
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Odstranění redundance

Pokud $m_{ij} = +1$,

pak v i -tém řádku nuluj všechny $+1$,

pokud ve stejném sloupci je v j -tém řádku $+1$.



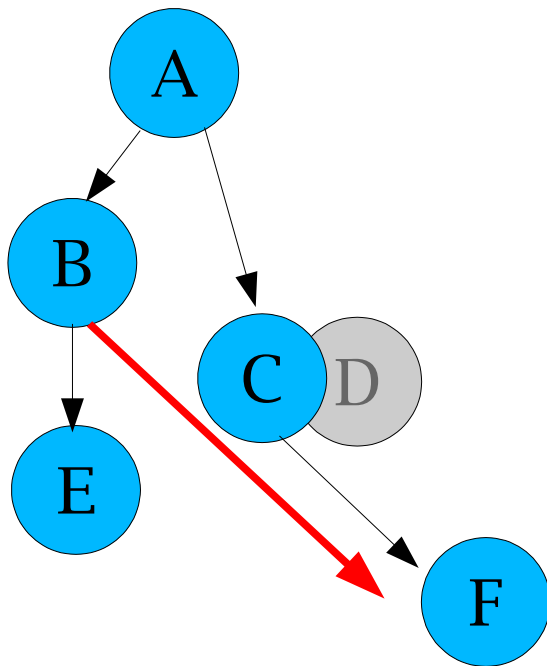
C
D

	A	B	C	E	F
A	0	1	1	0	0
B	-1	0	0	1	0
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Odstranění redundance

Můžeme považovat za redundantní všechny závislosti, mezi nimiž a diagonálou je alespoň jedna +1.

Pokud do uzlu vstupuje více hran, může se jednat o speciální případ vícehodnotové závislosti.



C
D

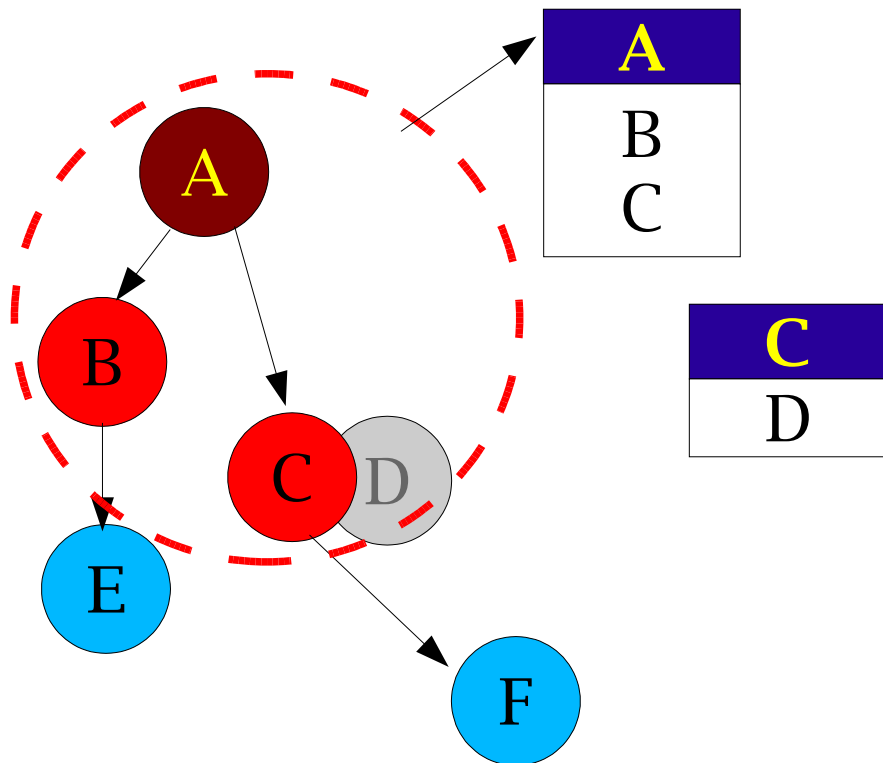
	A	B	C	E	F
A	0	1	1	0	0
B	-1	0	0	1	1
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Rekonstrukce relací

Každý řádek představuje relaci.

Atribut řádku je primárním klíčem relace

Relaci dále tvoří všechny atributy $s + 1$.



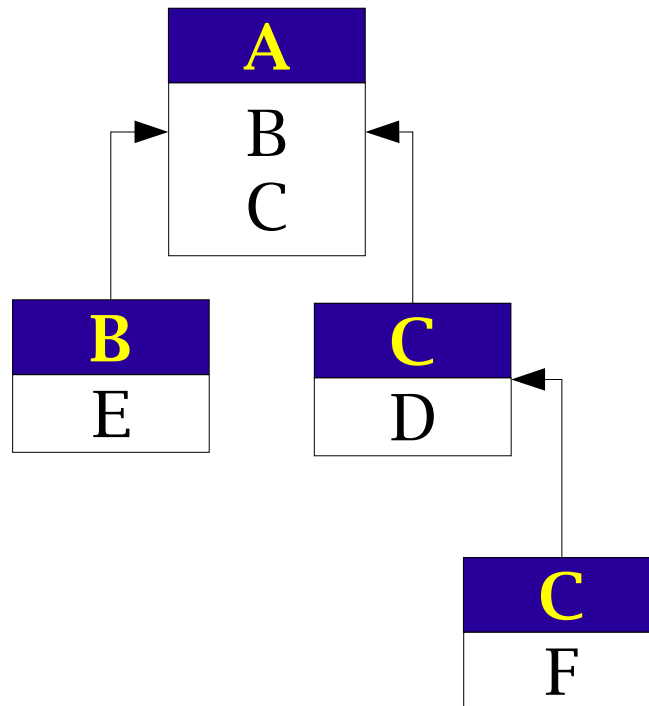
	A	B	C	E	F
A	0	1	1	0	0
B	-1	0	0	1	0
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Rekonstrukce relací

Každý řádek představuje relaci.

Atribut řádku je primárním klíčem relace

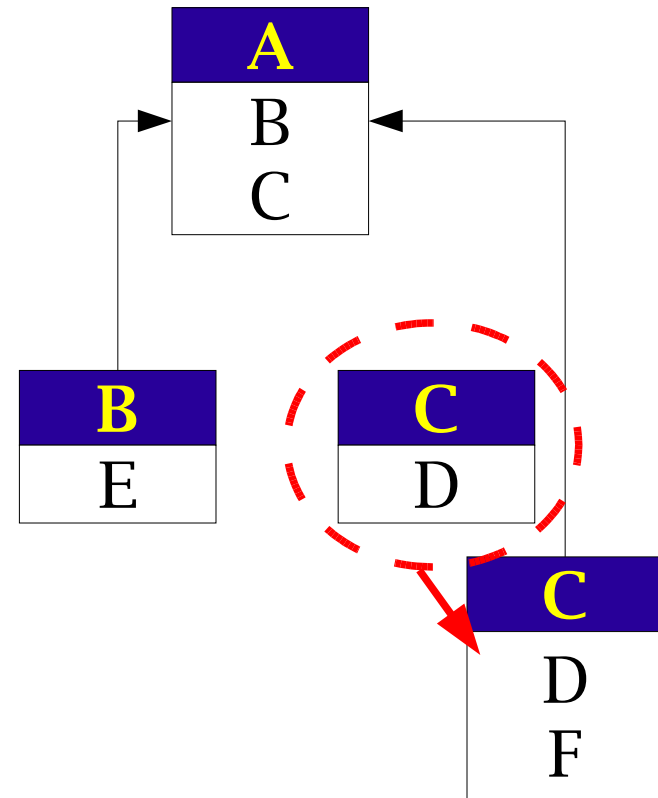
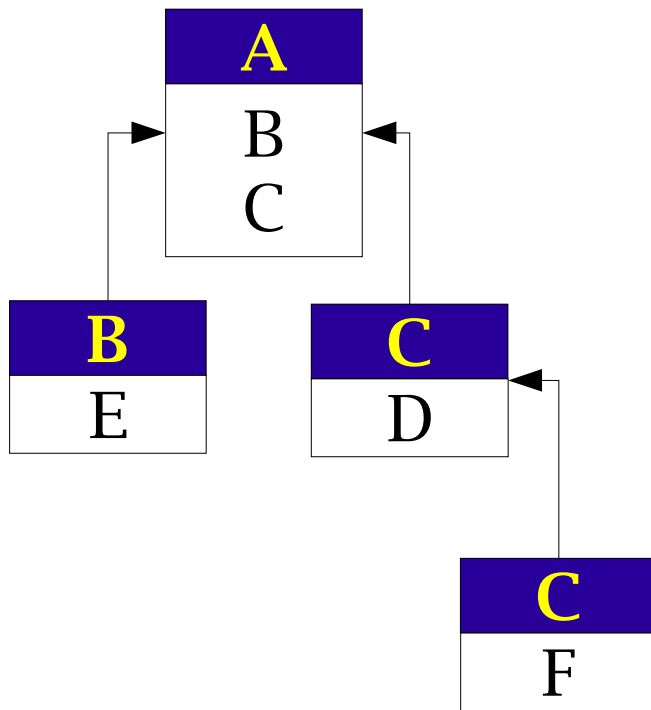
Relaci dále tvoří všechny atributy s +1.



	A	B	C	E	F
A	0	1	1	0	0
B	-1	0	0	1	0
C	-1	0	0	0	1
E	-1	-1	0	0	0
F	-1	0	-1	0	0

Rekonstrukce relací

Relaci vzniklou ze vzájemných relací je možné zakomponovat do relace o úroveň výše nebo níže.



Oba modely splňují kritéria 3. normální formy.

Motivace pro fuzzy

Pokud rekonstruujeme nekompletní data, klasická definice funkční závislosti je příliš restriktivní.

X	Y
1	1
2	1
2	1
2	2
3	2

Jediný chybný záznam rozhodne o neplatnosti funkční závislosti

Klasický přístup
 $X \not\rightarrow Y$

4/5 záznamů funkční závislosti vyhovují

Fuzzy přístup
 $X \rightarrow Y$ (@80%)

Matice fuzzy závislosti

$$M_{ij} = c_{ij} / C - c_{ji} / C$$

c_{ij} počet záznamů splňujících $A_i \rightarrow A_j$

c_{ji} počet záznamů splňujících $A_j \rightarrow A_i$

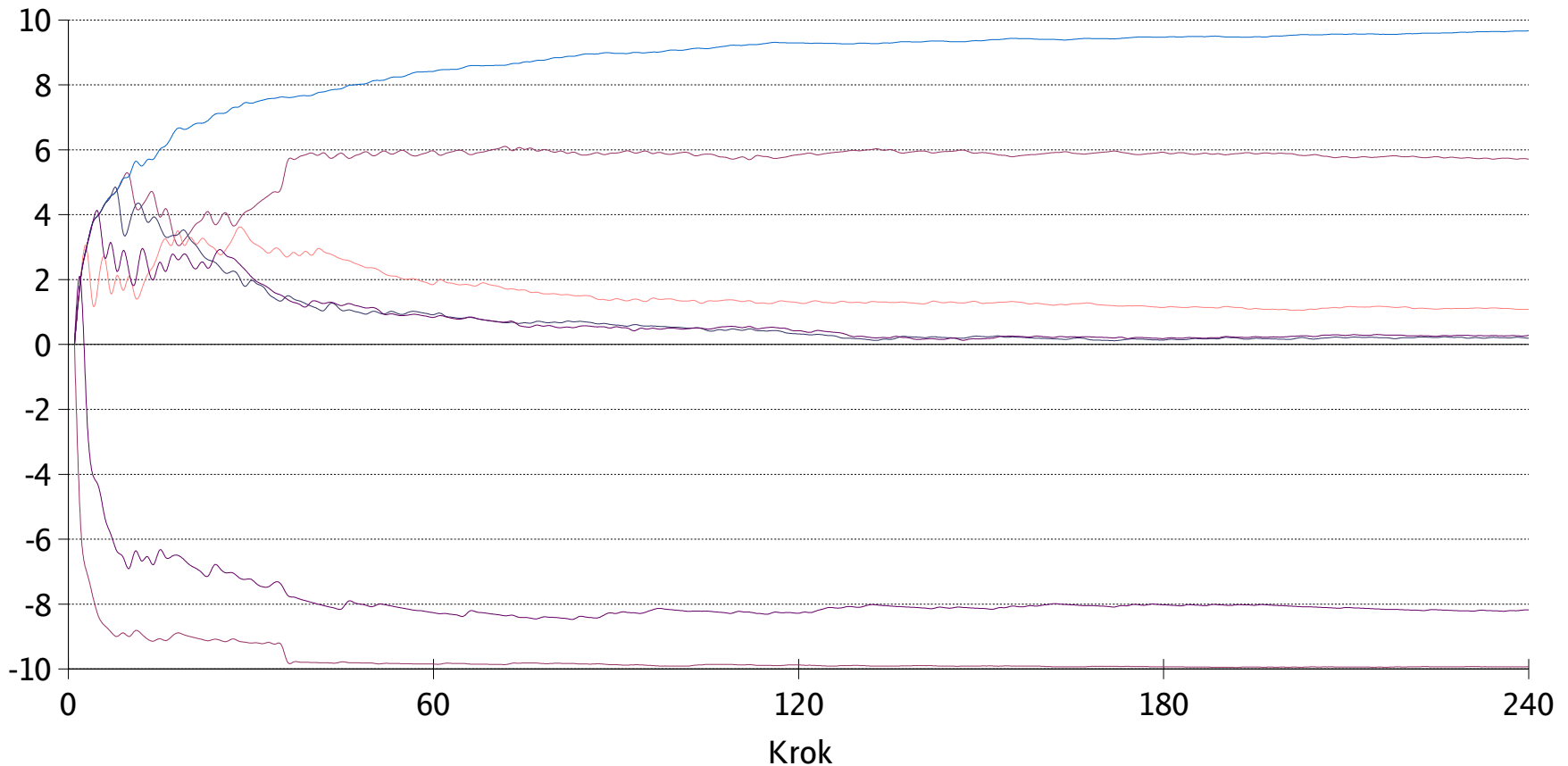
C celkový počet záznamů

Uspořádání atributů pro fuzzy

Kritérium pro uspořádání shodné s nefuzzy verzí:

$$\sum_j m_{ij}$$

Stabilita součtu fuzzy závislostí



Nefuzzy závislosti

- + Implementovaná metoda s ohledem na možnost fuzzyfikace
- + Extenzivní přístup k datům
- + Ohled na výpočetní složitost, redukce stavového prostoru
- Nevhodné pro nekompletní data
- Výsledek: 1 model (náhodně vybraný z více)

Fuzzy varianta

- + Intenzivní přístup k datům
- + Vychází z nefuzzy metody
- + Možnost odhadu extenzivního řešení - stabilita
- + Výsledek: Ohodnocená množina modelů
(vhodné kvůli nejednoznačnosti)