

BIZONYTALAN INFORMÁCIÓK KEZELÉSE LOGIKAI ADATMODELLEKBEN

Achs Ágnes, achs@pmmf.hu

Janus Pannonius Tudományegyetem Pollack Mihály Főiskolai Kara, Pécs

Abstract

In this lecture we deal with the concept of logical data models and give an extension to handling uncertain information. We define the fuzzy Datalog programs as a set of Horn formulae with degrees and give their meaning by defining the deterministic and nondeterministic semantics.

1. Alapfogalmak

Term-nek nevezünk egy változót, konstansot vagy egy $f(t_1, \dots, t_n)$ alakú kifejezést, ahol f függvénszimbólum és t_1, \dots, t_n termek. *Atom* egy $p(t)$ alakú formula, ahol p egy n változós predikátumszimbólum és t egy n elemű termsorozat. *Literál*-nak nevezünk egy atomot (pozitív literál) vagy a negáltját (negatív literál). Literálok diszjunkciója a *klóz*. A változó nélküli termet (atomot, literált, klózt) *alaptermnek* (alapatom, alapliterál, alapklóz) nevezzük.

2. Logikai adatmodell

Egy logikai adatmodell - vagy más szóval tudásbázis - tényekből és következtetési szabályokból épül fel. A tények bizonyos ismerete ket reprezentálnak, melyekből a szabályok segítségével következtethetünk újabb ismeretekre. A deduktív adatbázisok elméletében közismertek a Datalog-szerű nyelvek.

A Datalog Horn-klózek, azaz

$$A \leftarrow B_1, \dots, B_n$$

alakú szabályok halmaza, ahol A, B_i ($i = 1, \dots, n$) pozitív literálok.

Azok a predikátumok, melyekhez tartozó relációkat az adatbázisban tároljuk, az EDB predikátumok, a megfelelő relációk az EDB relációk. Ezeket szokták tényeknek is nevezni. (Alakja: $A \leftarrow$)

Azokat a predikátumokat (relációkat), melyeket logikai szabályok definiálnak, IDB predikátumoknak (relációknak) nevezzük.

Minden Datalog programhoz készíthető egy úgynevezett megelőzési gráf. A gráf csúcsai a programban szereplő közönséges predikátumok. A p predikátumból a q predikátumba akkor vezet él, ha van olyan szabály, melynek feje q és törzsében szerepel p . Ha a gráf kört tartalmaz, akkor a program rekurzív.

A Datalog szabályok kiértékelésekor két utat követhetünk. Az egyik módszer szerint minden IDB predikátumhoz meghatározunk egy relációt az EDB predikátumokhoz tartozó relációk segítségével. Rekurzív programok esetén egy Datalog egyenletrendszerhez jutunk, melyet iteratív módon oldhatunk meg. Ezt a módszert preferálja [1]. [2] a másik utat részesíti előnyben, amely szerint egy rákövetkezési operátort definiálhatunk, s ez alapján juthatunk új tényekhez. Mindkét megoldási algoritmus egy fixpontkeresés, és a

legkisebb fixpontot eredményezi. Ez a fixpont egyúttal a Datalog program egyértelmű minimális modellje. (Vagyis azon tények minimális halmaza, melyek kielégítik a program összes szabályát.)

Nem ilyen egyszerű a helyzet, ha negációt is megengedünk, hiszen a negatív információk kezelése nehéz feladat.

A negatív adatok tömege jóval nagyobb, mint a pozitívaké, ezért csak a pozitívakat tároljuk, s a negatívakra következtetünk. Többféle következtetési mód is ismert, de mindegyik kapcsolódik valamilyen módon a zártvilág-feltételezéshez (CWA). A CWA lényege, hogy ha egy tény logikailag nem következik egy klózalmazból, akkor arra következtetünk, hogy a tény negáltja igaz.

Ha a CWA elvet a közönséges Datalogra alkalmazzuk, akkor is tudunk negatív tényekre következtetni, de ezekből a negatív tényekből már nem határozhatunk meg újabbakat. Hogy tovább tudjunk belőlük következtetni, meg kell engednünk negatív literálok használatát a szabálytörzsben. Így jutunk a Datalog nyelvhez.

A negáció miatt előfordulhat, hogy nem kapunk legkisebb fixpontot, hanem több minimálisat. (Ez hasonlóan több minimális modellt is jelent.) Kérdés, melyiket válasszuk közülük, azaz milyen szemantikát rendelünk a Datalog nyelvhez.

Kétféle szemantikát szoktak értelmezni, az egyik a rétegzés, a másik az inflációs szemantika.

Rétegzéskor egyfajta rendezést vezetünk be a predikátumok között, s a programokat ebben a sorrendben értékeljük ki, vagyis a deklaratív szabályokhoz procedurális kiértékelés társul. Ha ebben a sorrendben végezzük el a kiértékelést, akkor egy negált predikátumra már csak akkor kerül sor, ha előbb minden pozitív előfordulási helyén kiértékeljük.

Az inflációs szemantika esetén nem kívánjuk meg a modell minimalitását, megelégszünk egy inflációs operátor (amely segítségével tényekből újabb tényekre következtetünk) legkisebb fixpontjával, s ezt a fixpontot definiáljuk szemantikaként. Ez a szemantika "hatásosabb", mint a rétegzés. Ez azt jelenti, hogy található olyan lekérdezés, amely ki fejezhető az inflációs Datalog-ban, de nem fejezhető ki a rétegzettben. ([2])

A bizonytalan információk kezelésére azonban nem alkalmas a klasszikus Datalog. A további fejezetekben a bizonytalanság kezelésének egy lehetséges módját adjuk meg. A közönséges Datalog szabályokat kiegészítjük egy bizonytalansági szintet jelző értékkel és egy implikációs operátorral, amely segítségével következtetni tudunk a szabályfej bizonytalansági szintjére. A szabályokhoz és ily módon a predikátumokhoz rendelt bizonytalansági szint azt mutatja meg, hogy az illető predikátum legalább milyen szinten teljesül. Az így kapott program nyelvet fuzzy Datalognak, vagy röviden fDATALOG-nak nevezzük.

2. A fuzzy elmélet alapfogalmai

A nyelv definiálása előtt ismerkedjünk meg a fuzzy halmazok elméletének néhány alapvető fogalmával!

Legyen D egy halmaz! A D fölötti F fuzzy halmaz egy $F : D \rightarrow [0,1]$ függvény. Jelölje $\mathcal{F}(D)$ az összes D fölötti fuzzy halmazt. Ekkor $F \in \mathcal{F}(D)$. A fuzzy halmazokon a metszet és unió szokásos értelmezése a következő:

$$F \cup G (d) \stackrel{\text{def}}{=} \max(F(d), G(d))$$

$$F \cap G (d) \stackrel{\text{def}}{=} \min(F(d), G(d))$$

Fuzzy halmazokra vonatkozó rendezési relációt is értelmezhetünk: $F \leq G$ akkor és csak akkor, ha $F(d) \leq G(d) \forall d \in D$. Mivel $\mathcal{F}(D)$ minden részhalmazának létezik legkisebb felső, illetve legnagyobb alsó korlátja, ezért $(\mathcal{F}(D), \leq)$ teljes háló.

A fuzzy halmazokat gyakran $F = \bigcup_{d \in D} (d, \alpha_d)$ módon jelölik, ahol $(d, \alpha_d) \in D \times [0, 1]$.

A következtetések elvégzéséhez szükségünk van az *implikációs operátor* fogalmára. Az implikációs operátorok megválasztásával és tulajdonságaival kapcsolatban sok vizsgálat folyik. Ezekről a vizsgálatokról ad összefoglaló képet [3]. A külön böző implikációs operátorokat a normák és conormák segítségével értelmezik. Ezek a metszet és az unió műveletének a fuzzy halmazokra való kiterjesztése. A következő táblázatban részletezés nélkül összefoglaljuk a leggyakoribb implikációs operátorokat:

jelölés	név	formula
$I_1(x, y)$	Gödel	1 ha $x \leq y$ y egyébként
$I_2(x, y)$	Lukasiewicz	1 ha $x \leq y$ $1 - x + y$ egyébként
$I_3(x, y)$	Goguen	1 ha $x \leq y$ y/x egyébként
$I_4(x, y)$	Kleene-Dienes	$\max(1-x, y)$
$I_5(x, y)$	Reichenbach	$1 - x + xy$
$I_6(x, y)$	Zadeh	$\max(1-x, \min(x,y))$
$I_7(x, y)$	Gaines-Rescher	1 ha $x \leq y$ 0 egyébként

3. A fuzzy Datalog szintaktikája

1. Definíció:

fDATALOG szabály egy $(r ; I ; \beta)$ hármas, ahol r egy

$$Q \leftarrow Q_1, \dots, Q_n \quad (n \geq 0)$$

alakú formula, Q atom (a szabály feje), Q_1, \dots, Q_n literálok (a szabály törzse), I egy implikációs operátor és $\beta \in (0, 1]$ (a szabály bizonytalansági foka vagy szintje).

Az fDATALOG szabály biztonságos, ha

- a fejben előforduló összes változó szerepel a törzsben is;
- az összes olyan változó, amely negatív literálban szerepel, előfordul pozitív literálban is.

Egy fDATALOG program biztonságos fDATALOG szabályok véges halmaza.

Az $(A \leftarrow ; I ; \beta)$ alakú szabályokat, ahol A alapatom, tényeknek nevezzük.

Egy P program *Herbrand univerzuma* (H_P) a P-ben előforduló konstansokból és függvény-szimbólumokból képzett összes lehetséges alapatom halmaza. A P *Herbrand bázisa* (B_P) a P-ben előforduló predikátumszimbólumokból képzett összes lehetséges olyan alapatom halmaza, melynek argumentumai H_P elemei. Egy $(r; I; \beta)$ szabály P-beli *alapelőfordulása* egy olyan szabály, amelyet úgy kapunk r-ből, hogy az összes r-beli X változót $\Phi(X)$ -szel helyettesítjük, ahol Φ az r-ben előforduló változók H_P -be való leképezése.

Az $(r; I; \beta)$ szabály összes alapelőfordulásának halmazát $(\text{ground}(r); I; \beta)$ -val jelöljük. A P program alapelőfordulása:

$$\text{ground}(P) = \cup_{(r; I; \beta) \in P} (\text{ground}(r); I; \beta)$$

2. Definíció:

Egy P program interpretációja B_p egy fuzzy részhalmaza:

$$N_p \in \mathcal{F}(B_p),$$

azaz $N_p = \bigcup_{A \in B_p} (A, \alpha_A)$

A konjunkció és negáció szintjét a szokásos módon értelmezzük, azaz tetszőleges A_1, \dots, A_n alapatomok esetén :

$$\alpha_{A_1 \wedge \dots \wedge A_n} \stackrel{\text{def}}{=} \min(\alpha_{A_1}, \dots, \alpha_{A_n})$$

$$\alpha_{\neg A} \stackrel{\text{def}}{=} 1 - \alpha_A$$

3. Definíció:

Egy interpretáció a P program modellje, ha minden

$$(A \leftarrow A_1, \dots, A_n; I; \beta) \in \text{ground}(P)$$

esetén

$$I(\alpha_{A_1 \wedge \dots \wedge A_n}, \alpha_A) \geq \beta.$$

Az M legkisebb modell, ha bármely N modell esetén $M \leq N$. M minimális modell, ha nincs olyan $N \neq M$ modell, hogy $N \leq M$ teljesüljön.

Az egyszerűség kedvéért $\alpha_{A_1 \wedge \dots \wedge A_n}$ -t $\alpha_{\text{törzs}}$ -szel és α_A -t α_{fej} -jel jelöljük.

Megjegyzés:

Az fDATALOG a közönséges Datalog általánosításának tekinthető, hiszen ha minden szabályhoz az I_{γ} -tel jelölt, ún. Gaines-Rescher implikációs operátort és $\beta = 1$ bizonytalansági szintet rendelünk, akkor Datalog szabályok halmazát kapjuk.

4. A fuzzy DATALOG szemantikája

Két rákövetkezési transzformációt értelmezzünk, egy determinisztikus és egy nem determinisztikus transzformációt. A determinisztikus transzformáció segítségével párhuzamosan értékeljük ki a program szabályait, míg nem determinisztikus esetben a szabályokat tetszőleges sorrendben, egyenként egymás után véve végezzük a következtetést. A két transzformációnak megfelelően kétféle szemantikát rendelhetünk az fDATALOG programokhoz. Ezek a szemantikák pozitív programok esetén megegyeznek, de negációt tartalmazó programoknál különböző eredményt adhatnak.

4. Definíció

A $DT_p : \mathcal{F}(B_p) \rightarrow \mathcal{F}(B_p)$ és $NT_p : \mathcal{F}(B_p) \rightarrow \mathcal{F}(B_p)$ rákövetkezési transzformációkat a következő módon értelmezzük:

$$DT_P(X) = \{ \cup \{ (A, \alpha_A) \} \mid (A \leftarrow A_1, \dots, A_n; I; \beta) \in \text{ground}(P), (|A_i|, \alpha_{A_i}) \in X, 1 \leq i \leq n, \alpha_A = \max(0, \min\{ \gamma \mid I(\alpha_{\text{törzs}}, \gamma) \geq \beta \}) \} \cup X$$

és

$$NT_P(X) = \{ (A, \alpha_A) \} \cup X$$

ahol $(A \leftarrow A_1, \dots, A_n; I; \beta) \in \text{ground}(P), (|A_i|, \alpha_{A_i}) \in X, 1 \leq i \leq n,$
 $\alpha_A = \max(0, \min\{ \gamma \mid I(\alpha_{\text{törzs}}, \gamma) \geq \beta \})$

Itt $|A_i|$ - val azt a $p(\underline{c})$ atomot jelöljük, melyre $A_i = p(\underline{c})$ vagy $A_i = \neg p(\underline{c})$, p egy k argumentumú predikátumszimbólum és \underline{c} k alaptermből álló lista.

A vizsgált implikációs operátorok körét le kell szűkítenünk, ugyanis nem minden I esetén létezik a bizonytalansági-szint függvény (bizonytalansági-szint függvénynek nevezzük az $f(I, \alpha, \beta) = \min(\{ \gamma \mid I(\alpha, \gamma) \geq \beta \})$ függvényt), azaz nem határozható meg tetszőleges $\alpha_{\text{törzs}}, \beta$ - hoz a kívánt α_A érték. Könnyen belátható, hogy bármely $I \in \{ I_1, I_2, I_3, I_4, I_5, I_7 \}$ esetén tetszőleges $\alpha_{\text{törzs}}, \beta$ értékhez létezik az $f(I_i, \alpha_{\text{törzs}}, \beta)$ ($i = 1, 2, 3, 4, 5, 7$) érték, míg $I = I_6$ esetén nem értelmezett minden α, β esetén az $f(I_6, \alpha, \beta)$ függvény.

A tényhalmazból kiindulva a transzformációk egymás utáni alkalmazásával (vagyis a transzformáció hatványainak meghatározásával) értelmezünk egy halmaz sorozatot. Ezen sorozatok fixpontjaként értelmezzük majd az fDATALOG szemantikáját.

Belátható a következő két állítás:

1. Állítás:

A DT_P transzformációnak van fixpontja, vagyis létezik olyan $X \in \mathcal{F}(B_P)$, melyre $DT_P(X) = X$. Az NT_P transzformációnak van fixpontja, vagyis létezik olyan $X \in \mathcal{F}(B_P)$, melyre az X halmaz bármely szabályát is választva $NT_P(X) = X$. Ha P pozitív, akkor a két fixpont megegyezik, és ez a közös érték a legkisebb fixpont.

A fixpontokat $\text{lfp}(DT_P)$, illetve $\text{lfp}(NT_P)$ -vel jelöljük.

2. Állítás:

$\text{lfp}(DT_P)$ és $\text{lfp}(NT_P)$ a P modelljei.

Az előző két állítás igazsága teszi lehetővé, hogy az fDATALOG szemantikáját a következő módon definiáljuk:

5. Definíció:

$\text{lfp}(DT_P)$ az fDATALOG P program determinisztikus szemantikája.

$\text{lfp}(NT_P)$ az fDATALOG P program nondeterminisztikus szemantikája.

Pozitív programok esetén a két szemantika megegyezik, s belátható, hogy ez a fixpont legkisebb modell is.

Belátható az is, hogy az I_1, I_2, I_3, I_4, I_7 operátorok, esetén létezik algoritmus a legkisebb modell meghatározására, azaz a fixpontszámítási eljárás véges sok lépésben befejeződik. Az I_5 operátor alkalmazásakor előfordulhat, hogy a program nem terminál véges lépésszámon belül.

Példa:

1. $r(a) \leftarrow ; I_1; 0.8$

2. $p(x) \leftarrow r(x), \neg q(x); I_1; 0.6$
3. $q(x) \leftarrow r(x); I_1; 0.5$
4. $p(x) \leftarrow q(x); I_1; 0.8$

Ekkor

$$\text{lfp}(\text{DT}_p) = \{(r(a), 0.8); (p(a), 0.6); (q(a), 0.5)\}$$

Nemdeterminisztikus kiértékeléskor különböző megoldásokat kaphatunk, hiszen az függ a kiértékelendő szabályok sorrendjétől.

Ha a kiértékelendő szabályok sorrendje 1.,2.,3.,4. akkor

$$\text{lfp}(\text{NT}_p) = \{(r(a), 0.8); (p(a), 0.6); (q(a), 0.5)\},$$

míg 1.,3.,2.,4. -es sorrendben

$$\text{lfp}(\text{NT}_p) = \{(r(a), 0.8); (p(a), 0.5); (q(a), 0.5)\}.$$

Mint a példából is láthatjuk, $\text{lfp}(\text{DT}_p)$ nem mindig minimális modell, azaz a determinisztikus szemantikát alkalmazva a kapott fixpont nem minden esetben adja a kívánt minimális modellt. Nem determinisztikus esetben azonban bizonyos feltételek mellett biztosítható a minimalitás. Ez a feltétel a rétegzés. A rétegzés meghatároz egy kiértékelési sorrendet, amelyben a negatív literálokat értékeljük ki először, s így módon minimális modellt kapunk.

A rétegzéshez előbb a függőségi gráf fogalmát kell megadnunk. Ez egy olyan irányított gráf, melynek csúcsai a P predikátumai. Egy p predikátumból akkor vezet él egy q predikátumba, ha P -nek van olyan szabálya, melynek törzsében p vagy $\neg p$ szerepel, és amelynek fej-predikátuma q .

Egy program rétegzett, ha az összes olyan esetben, amikor egy szabály fej-predikátuma p , és a törzs negált literálja $\neg q$, a függőségi gráfban nincs út p -ből q -ba.

Egy P program rétegzése a P predikátumszimbólumainak olyan P_1, \dots, P_n részhalmazokra való partíciója, melyre teljesülnek a következő feltételek:

a/ ha $p \in P_i, q \in P_j$ és nincs él q -ból p -be, akkor $i \geq j$.

b/ ha $p \in P_i, q \in P_j$ és van olyan p fejpredikátumú szabály, melynek törzse $\neg q$, akkor $i > j$.

A P_1, \dots, P_n halmazokat rétegeknek nevezzük.

Legyen P egy rétegzett FDATALOG program P_1, \dots, P_n rétegzéssel. Jelölje P_i^* a P_i réteghez tartozó összes P -beli szabály halmazát, azaz az összes olyan szabály halmazát, melynek fej-predikátuma P_i -ben van.

Legyen

$$L_1 = \text{lfp}(\text{NT}_{P_1^*}),$$

ahol a fixpontoszámolási eljárás kiindulópontja a program tényeinek halmaza.

Legyen

$$L_2 = \text{lfp}(\text{NT}_{P_2^*}),$$

ahol a számolás kiindulópontja az előbb kapott L_1 ,

$$L_n = \text{lfp} (NT_{P_n^*}),$$

ahol a kiindulópont L_{n-1} .

Más szóval, először kiszámítjuk a P első rétegéhez tartozó L_1 fixpontot. Ezen fixpont meg határozása után léphetünk a következő rétegekre.

Megjegyzés:

$$\text{lfp} (NT_{P_i^*}) = \text{lfp} (DT_{P_i^*})$$

Belátható, hogy L_n a P minimális fixpontja és egyúttal minimális modellje.

Ha P rétegzett fDATALOG program, akkor L_n a P minimális fixpontja.

4. Konklúzió

Az előadásban a Datalog-szerű nyelvek egy lehetséges fuzzy kiterjesztését adtuk meg úgy, hogy a Datalog szabályokat kiegészítettük egy bizonytalansági szintet jelző értékkel és egy implikációs operátorral. Bizonyos rákövetkezési transzformációk fixpontjaként definiáltuk az fDATALOG determinisztikus és nem-determinisztikus szemantikáját és beláttuk, hogy ez a fixpont pozitív programok esetén legkisebb, negácót is tartalmazó programok esetén pedig bizonyos feltételek mellett minimális modell.

Az fDATALOG fuzzy adatokra is kiterjeszhető. A hasonlósági relációk felhasználásával ilyen irányban indul el [4].

Irodalomjegyzék

- [1] J.D.Ullman : Principles of database and knowledge-base systems
Computer Science Press, Rockville, 1988.
- [2] S.Ceri - G.Gottlob - L.Tanca : Logic Programming and Databases
Springer-Verlag Berlin, 1990
- [3] Didier Dubois - Henri Prade: Fuzzy sets in approximate reasoning, Part 1:
Inference with possibility distributions
Fuzzy Sets and Systems 40 (1991) (143-202)
- [4] Ágnes Achs - Attila Kiss : Fuzzy extension of Datalog
Acta Cybernetica 12 (1995) (153-166), Szeged
- [5] Ágnes Achs - Attila Kiss : Fixpoint query in fuzzy Datalog
megjelenőben : Annales, Bp
- [6] Attila Kiss : On the least models of fuzzy Datalog programs
International Conference on Information Processing and Management of Uncertainty in
Knowledge-based system, Mallorca, 1993. (465-471)
- [7] Vilém Novák : Fuzzy sets and their applications
Adam Hilger Bristol and Philadelphia, 1989.