

Logic and Relational Databases

bioactivity

mol	ACE	SSRI
m1	0	0
m2	1	0
m3	0	1
⋮	⋮	⋮

atoms

mol	name	type	X	Y	Z
m1	a1	C	2.1	-1.3	3.4
m1	a2	O	3.0	-1.0	3.0
⋮	⋮	⋮	⋮	⋮	⋮

bonds

mol	atom1	atom2	type
m1	a1	a2	2
m1	a1	a3	1
⋮	⋮	⋮	⋮

bioactivity (m1, 0, 0).
bioactivity (m2, 1, 0).
bioactivity (m3, 0, 1).
⋮

atom (m1, a1, c, 2.1, -1.3, 3.4).
atom (m1, a2, o, 3.0, -1.0, 3.0).
⋮

bond (m1, a1, a2, 2).
bond (m1, a1, a3, 1).
⋮

"Extensional Database"

Datalog

Definite program containing no non-nullary function symbols.

$\text{grandparent}(X, Z) \leftarrow \text{parent}(X, Y), \text{parent}(Y, Z).$
 $\text{parent}(X, Y) \leftarrow \text{father}(X, Y).$
 $\text{parent}(X, Y) \leftarrow \text{mother}(X, Y).$

$\text{father}(\text{adam}, \text{bill}).$ $\text{mother}(\text{anne}, \text{bill}).$
 $\text{father}(\text{adam}, \text{beth}).$ $\text{mother}(\text{anne}, \text{beth}).$
 $\text{father}(\text{bill}, \text{cathy}).$ $\text{mother}(\text{cathy}, \text{donald}).$
 $\text{father}(\text{donald}, \text{eric}).$ $\text{mother}(\text{diana}, \text{eric}).$

mother and father are defined extensionally.
grandparent and parent are defined intensionally.

Some Interesting Notes

- With function symbols, even the "simplest" logic program is undecidable — more specifically there exist choices of literals $L_1, L_2,$ & L_3 such that it is undecidable

if $L_1 \models L_4$
 $L_2 \leftarrow L_3.$

for some literals L_4 ,

- Moreover, there exist 2-literal clauses

- Moreover, there exist 2-literal clauses such that the other 2-literal clauses entailing/entailed by them are undecidable.
- But Datalog is decidable. (Not using SLD-resolution but using other approaches such as Magic Sets – Ramakrishnan, Ch.15.)

What are we still missing?

- Field names.
Syntactic sugar – will live without.
- Types for fields.
Just unary predicates:
male(tom). person(x) ← male(x).
- Query language.
Coming next... will show that Datalog with negation as finite failure has at least the computational power of relational algebra.

Primitive Operations of Relational Algebra

- Union
- Set Difference
- Cartesian (cross-) Product
- Projection
- Selection

Union

$$R_1 \cup R_2 = \{ \langle x_1, \dots, x_n \rangle \mid \langle x_1, \dots, x_n \rangle \in R_1 \vee \langle x_1, \dots, x_n \rangle \in R_2 \}$$

In Datalog, use predicate symbols r_1/n & r_2/n for R_1 & R_2 , respectively. Use r/n for their

union :

$$r(x_1, \dots, x_n) \leftarrow r_1(x_1, \dots, x_n).$$

$$r(x_1, \dots, x_n) \leftarrow r_2(x_1, \dots, x_n).$$

Example: $\text{parent}(X, Y) \leftarrow \text{father}(X, Y).$
 $\text{parent}(X, Y) \leftarrow \text{mother}(X, Y).$

Difference

$$R_1 \setminus R_2 = \{ \langle x_1, \dots, x_n \rangle \mid \langle x_1, \dots, x_n \rangle \in R_1 \wedge \langle x_1, \dots, x_n \rangle \notin R_2 \}$$

In Datalog:

$$r(x_1, \dots, x_n) \leftarrow r_1(x_1, \dots, x_n), \text{ not } r_2(x_1, \dots, x_n).$$

Example:

$$\text{father}(X, Y) \leftarrow \text{parent}(X, Y), \text{ not } \text{mother}(X, Y).$$

Cartesian Product

$$R_1 \times R_2 = \{ \langle x_1, \dots, x_m, y_1, \dots, y_n \rangle \mid \langle x_1, \dots, x_m \rangle \in R_1 \wedge \langle y_1, \dots, y_n \rangle \in R_2 \}$$

In Datalog:

$$r(x_1, \dots, x_m, y_1, \dots, y_n) \leftarrow r_1(x_1, \dots, x_m), r_2(y_1, \dots, y_n).$$

Example:

$$\text{mfcouple}(X, Y) \leftarrow \text{male}(X), \text{female}(Y).$$

Projection

$$\pi_{i_1, \dots, i_m} R = \{ \langle x_{i_1}, \dots, x_{i_m} \rangle \mid \langle x_1, \dots, x_n \rangle \in R, \\ 1 \leq i_1, \dots, i_m \leq n \}$$

(Projection deletes and/or rearranges one or more columns from R .)

In Datalog:

$$pr(x_{i_1}, \dots, x_{i_m}) \leftarrow r(x_1, \dots, x_n).$$

Example: $father(X) \leftarrow father(X, Y)$.

Selection

$$\delta_F(R) = \{ \langle x_1, \dots, x_n \rangle \mid \langle x_1, \dots, x_n \rangle \in R \text{ and formula } \\ F \text{ is true of } \langle x_1, \dots, x_n \rangle \}$$

F is allowed to contain comparisons ($=, \geq$, etc.) among variables and constants, joined by \wedge, \vee, \neg .
field names values

Datalog translation depends on F . Where

$D(F)$ is the Datalog representation of F :

$\delta_F R$ is represented as

$$rs(x_1, \dots, x_n) \leftarrow r(x_1, \dots, x_n), D(F).$$

Example: $\delta_{Y \geq 1,000,000} INCOME(X, Y)$ is

$millionaire(X, Y) \leftarrow income(X, Y), Y \geq 1,000,000.$

Notes on Selection

- Replace field names in F by corresponding variable names during translation.
- Using $;$ for \vee and not for \wedge , translate F directly.

$$R$$

f_1	f_2	f_3	f_4	f_5
\vdots	\vdots	\vdots	\vdots	\vdots

$\delta_{(f_2=f_3) \vee (f_4 > f_5 \wedge f_5 \neq 10)}$ R becomes

$sr(X_1, \dots, X_5) \leftarrow r(X_1, \dots, X_5),$
 $((X_2 = X_3) ;$
 $(X_4 > X_5, X_5 \neq 10)).$

Can build other relational operators from those relational algebra primitives.

Example: Natural Join

$$R \bowtie S = \pi_A \delta_{R.T_1=S.T_1 \wedge \dots \wedge R.T_k=S.T_k} (R \times S)$$

where T_1, \dots, T_k are the attributes (field names) appearing in both R and S , and A is the list of all attributes except $S.T_1, \dots, S.T_k$.

$father(X, Y) \bowtie parent(Y, Z):$

$cross(X, Y_1, Y_2, Z) \leftarrow father(X, Y_1), parent(Y_2, Z).$

$sel(X, Y_1, Y_2, Z) \leftarrow cross(X, Y_1, Y_2, Z), Y_1 = Y_2.$

$proj(X, Y, Z) \leftarrow sel(X, Y_1, Z).$

Could do in one step:

$gf(X, Y, Z) \leftarrow father(X, Y), parent(Y, Z).$

Missing Values

f1	f2	f3
a	1	?
b	?	g
c	3	f

Like existential variables: we know some value must go there, but we don't know what.

Datalog: no existential quantifiers. Use skolem constants. Every missing value is represented by a distinct skolem constant.

$p(a, 1, sk-1).$
 $p(b, sk-2, g).$
 $p(c, 3, f).$