



RÈGLES EXISTENTIELLES: CHASE

THÉORIE DES BASES DE CONNAISSANCES
HMIN312M
COURS DE ML MUGNIER

EXISTENTIAL RULES

$$\forall x \forall y (\text{Body } [x,y] \rightarrow \exists z \text{ Head } [x,z])$$

$x, y, z :$
sets of variables

any **positive conjunction** (without functional symbols)

$$\forall x (\text{actor}(x) \rightarrow \exists z \text{ play}(x,z))$$

$$\forall x \forall y (\text{siblingOf}(x,y) \rightarrow \exists z (\text{parentOf}(z,x) \wedge \text{parentOf}(z,y)))$$

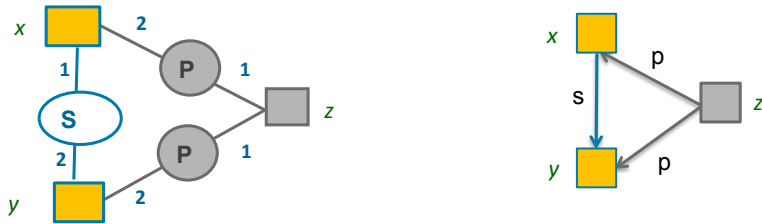
Key point: ability to assert **the existence of unknown entities**

Crucial for representing ontological knowledge in **open domains**

GRAPH VIEW OF (EXISTENTIAL) RULES

$$\forall x \forall y (\underbrace{\text{Body}[x,y]}_{\text{graph}} \rightarrow \exists z \underbrace{\text{Head}[x,z]}_{\text{graph}})$$

$$\forall x \forall y (\text{siblingOf}(x,y) \rightarrow \exists z (\text{parentOf}(z,x) \wedge \text{parentOf}(z,y)))$$



The rule head has 2 kinds of variables (or unlabelled term nodes):

- **frontier**: shared with the body
- **existential**

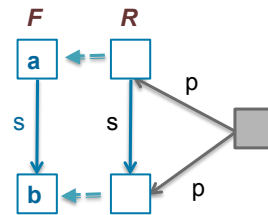
GENERATION OF FRESH (UNKNOWN) INDIVIDUALS

$$R = \forall x \forall y (\text{siblingOf}(x,y) \rightarrow \exists z (\text{parentOf}(z,x) \wedge \text{parentOf}(z,y)))$$

$$F = \text{siblingOf}(a,b)$$

R is **applicable** to F if there is a **homomorphism** h from $\text{body}(R)$ to F

$$\begin{matrix} x \rightarrow a \\ y \rightarrow b \end{matrix}$$



Applying R to F w.r.t. h produces $F \cup h(\text{head}(R))$

where a **fresh variable** ($a \ll \text{null}$) is created for each existential variable in R

$$F' = \exists z_0 (\text{siblingOf}(a,b) \wedge \text{parentOf}(z_0,a) \wedge \text{parentOf}(z_0,b))$$

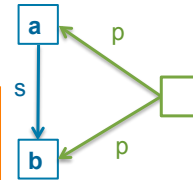
$$\text{Notation (when needed)} : \alpha(F,R,h) = F \cup h^{\text{safe}}(\text{head}(R))$$

where h^{safe} is a substitution of $\text{var}(\text{head}(R))$

such that:

$$h^{\text{safe}}(x) = h(x) \text{ if } x \text{ in } \text{frontier}(R)$$

otherwise $h^{\text{safe}}(x)$ is a fresh variable (a null)



KNOWLEDGE BASES WITH EXISTENTIAL RULES

$\mathcal{K} = (F, \mathcal{R})$ where

\mathcal{R} is a set of **existential rules**

F is a set of **facts** (rules with an empty body): existential conjunctions of atoms

Forward chaining called « **chase** » (let us still denote by F^* the result of the chase)

Main change with respect to Datalog rules: F^* can be **infinite**

$R = \text{person}(x) \rightarrow \exists y \text{ hasParent}(x,y) \wedge \text{person}(y)$

$F = \text{person}(a)$

$\wedge \text{person}(y_0) \wedge \text{hasParent}(a, y_0)$

$\wedge \text{person}(y_1) \wedge \text{hasParent}(y_0, y_1)$

Etc.

but it remains a **universal model**

Hence, for BCQs: $\mathcal{K} \models q$ iff $q \geq F^*$

Other changes: F^* is **not unique** (but all F^* we will see are logically equivalent)

DERIVATION

◦ **Trigger** for a set of facts F : $(R, h) \mid h$ homomorphism from $\text{body}(R)$ to F

◦ **Derivation**: $(F_0 = F) (R_1, h_1) F_1 (R_2, h_2) F_2, \dots$
 where for all i , (h_i, R_i) trigger for F_{i-1}
 and $F_i = \alpha(F_{i-1}, R_i, h_i)$

When the triggers are not needed, we note $(F_0 = F), F_1, F_2, \dots$

◦ The notion of derivation can be generalized by
 taking F_i **equivalent** to $\alpha(F_{i-1}, R_i, h_i)$ instead of **equal**

◦ **Different chase variants** with their own rule application criteria

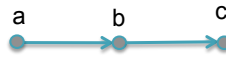
→ notion of *legal* rule application (or: *active* trigger (R_i, h_i)):

A chase variant considers only derivations with **active** triggers

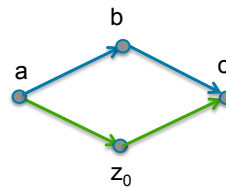
DIFFERENT VARIANTS OF THE CHASE

All variants we will see compute **universal models** of the KB
but they differ on how they handle **redundancies** possibly caused by nulls

$p(a,b), p(b,c)$



$p(a,b), p(b,c),$
 $\exists z_0 p(a,z_0) p(z_0,c)$



$z_0 \mapsto b$

Core: set of atoms without homomorphism to one of its strict subsets

14

OBLIVIOUS CHASE

Oblivious (or naive): « performs all rule applications according to all new triggers »

A trigger (R,h) to F_i is *active on F_i* iff this trigger has *not* already been used
in the derivation from F_0 to F_{i-1}

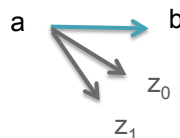
$R = p(x,y) \rightarrow \exists z p(x,z)$

$F = p(a,b)$

$p(a,z_0)$

$p(a,z_1)$

...



*stupid rules to keep
examples simple!*

infinite
derivation

RuleML+RR

15

OBLIVIOUS / SEMI-OBLIVIOUS = SKOLEM CHASE

Semi-oblivious: consider only homomorphisms that differ on the **rule frontier** (x)

A trigger (R,h) to F_i is *active on F_i* iff there is **no** trigger (R,h') such that $h'(x) = h(x)$ for **all** x in frontier(R) in the derivation from F_0 to F_{i-1}

$$F = p(a,b) \quad R = p(x,y) \rightarrow \exists z p(x,z)$$



Skolem chase: similar behavior

- (1) skolemize rules: in R, replace each existential variable z by a function $f_R^z(\text{frontier}(R))$
- (2) perform the oblivious chase on skolemized rules

$$R = p(x,y) \rightarrow p(x,f(x))$$

$$p(a,b)$$

$$p(a,f(a))$$

Skolemization can be seen as a way of naming existential variables and « tracking » the nulls created during the semi-oblivious chase

RuleML+RR

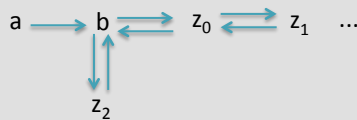
16

RESTRICTED (AKA STANDARD) CHASE

Restricted: do not perform a rule application that brings *only* redundant information

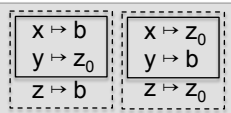
A trigger (R,h) to F_i is *active on F_i* iff h **cannot** be extended to homomorphism h' : body \cup head $\rightarrow F_i$

$$F : p(a,b) \quad R : p(x,y) \rightarrow \exists z p(y,z), p(z,y)$$



(semi-) oblivious chase:
infinite

$$a \rightarrow b \rightarrow z_0$$



restricted chase:
halts after one rule application

RuleML+RR

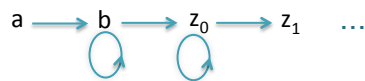
17

RESTRICTED CHASE: NATURAL BUT TRICKY

- For the same KB, some derivations may halt while others may not

$$F : p(a,b) \quad R_1: p(x,y) \rightarrow \exists z p(y,z) \\ R_2: p(x,y) \rightarrow p(y,y)$$

If R_1 is always applied before R_2 for a given homomorphism of $p(x,y)$:



If R_2 is applied first:



RuleML+RR

18

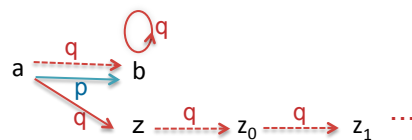
CORE CHASE

Iterate:

- perform a finite number of rule applications as in the restricted chase
- compute the **core** of the result

$$F : p(a,b), q(b,b), q(a,z) \quad \text{where } z \text{ is a variable}$$

$$R_1: p(x,y) \rightarrow q(x,y) \\ R_2: q(x,y) \rightarrow \exists z q(y,z)$$



The **restricted** chase only checks redundancy of **newly** added atoms
 \Rightarrow infinite here

The **core** chase outputs $\{ p(a,b), q(b,b), q(a,b) \}$

The **core** chase allows to detect **global** redundancies

19

WHEN DOES A CHASE HALT?

- o **Terminating** derivation:
(1) finite and (2) there is no active trigger on the last factbase
- o A chase derivation has to be **fair**: no active trigger is indefinitely delayed
Formally: if (R, h) is an active trigger on F_i
then there is F_j with $j > i$ such that F_j is obtained by applying (R, h)
or (R, h) is not active anymore on F_j

Terminating = finite and fair

$R_1: p(x, y) \rightarrow \exists z p(y, z)$

$R_2: p(x, y) \rightarrow p(y, y)$

$F = p(a, b)$

unfair infinite derivation: apply only $R_1 \dots$

(semi-) oblivious: all fair derivations are infinite

restricted: some terminating derivations,
some infinite fair derivations

core: all fair derivations are terminating

For a chase variant C , C **halts** on a KB K if all **fair** derivations on K are finite

IN SHORT

All previous chase variants compute **universal** models of a KB

They can be **strictly ordered wrt termination**:

oblivious < semi-oblivious = skolem < restricted < core

$[X < Y]$ means that: for any KB K , if X -chase halts on K then Y -chase halts on K
and there is a KB on which Y -chase halts but not X -chase]

Only the **core** chase halts if and only if the KB admits a **finite universal model**
but it is **costly** (involves homomorphisms from the whole factbase)

The **O**, **S-O** and **core** chases yield a **unique** result (up to the name of nulls):
all fair derivations for a given chase variant yield the same result on a given KB
but not the **R chase**: we can even have finite and infinite fair derivations

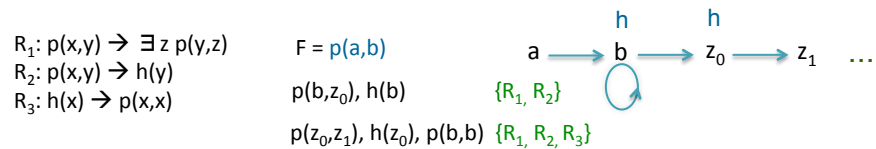
The **R chase** seems to achieve a **good tradeoff**
redundancy elimination / efficiency of computation (when it stops)
but its behavior is **difficult to control**

TRICKY RESTRICTED CHASE

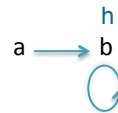
Open question:

is there an ordering strategy that terminates more often than the others?

- **Breadth-first ordering** is a natural candidate (iterate:
 - (1) compute all rule body homomorphisms to the current factbase,
 - (2) apply all active triggers according to these homomorphisms)
- however, it is **not optimal for restricted chase** termination



Optimal order: apply R_2 then R_3 (ie delay application of R_1)



- Usual heuristic: at each step, first saturate with all **datalog rules**, then apply an existential rule
- would be optimal on this example, is it always the case?