

# Relational Similarity-Based Databases II

Vilem Vychodil (Palacky University, Olomouc)

# DAMOL

DATA ANALYSIS AND MODELING LAB

Palacky University, Olomouc, Czech Republic



euROPEAN  
social fund in the  
czech republic



EUROPEAN UNION



MINISTRY OF EDUCATION,  
YOUTH AND SPORTS



OP Education  
for Competitiveness

INVESTMENTS IN EDUCATION DEVELOPMENT

# Preliminaries

**utilized notions** (part I):

- attributes, relation scheme
- tuple, cartesian product
- domain, similarity
- ranked data table (over domains with similarities)

**in case you missed it:** [goto DAMOL webpage](#)

Belohlavek, Vychodil: *Relational Similarity-Based Databases I*

**new notion:** **object constants** = names for individual values that can appear in a database

- $\mathcal{C}$ : denumerable set of all object constants
- object constants denoted by  $c, d, c_1, c_2, \dots$
- note the difference: name (syntactic notion)  $\neq$  value (semantic notion)

# Type Declarations and Types

**motivation:** The need to distinguish between different *types* of attributes (to avoid mixing apples and pears) on the syntactic level.

## Definition (types, type declarations)

A **type declaration** for a set  $Y$  of attributes and a set  $\mathcal{C}$  of object constants is a structure  $\mathbf{\Lambda} = \langle \Lambda, \lambda \rangle$ , where  $\Lambda$  is a nonempty set of elements called **types**, and  $\lambda$  is a map  $\lambda: Y \cup \mathcal{C} \rightarrow \Lambda$  such that for each  $\tau \in \Lambda$  there is an infinite  $Y' \subseteq Y$  such that  $\lambda(Y') = \{\tau\}$ .

### notes:

- $\lambda$  in  $\mathbf{\Lambda}$  assigns to each attribute and object constant  $z$  its type  $\lambda(z)$
- $\lambda(y) = \lambda(z)$  means: *y and z have the same type* (under  $\mathbf{\Lambda} = \langle \Lambda, \lambda \rangle$ )
- for each type there are infinitely many attributes of that type  
( $\Lambda$  is thus at most a denumerable set)

**can be extended:** hierarchy of types (not discussed here)

# Types vs. Domains

**domains** = sets of all possible values of attributes

- $D_y$  is domain of attribute  $y \in Y$ ;
- nonempty (at most) denumerable set of all possible values  $y$

**type**  $\times$  **domain**

- type = syntactic notion (a name),
- domain = semantic notion (a set of permissible values)

**our approach:**

*“Domain is an interpretation of a type.”*

**notes:** different approaches to domains vs. types:

- E. F. Codd: notion of a domain
- C. J. Date: domains = types (type is a synonym for a domain)

# Consistency Between Types and Domains

## motivation:

- The distinction between types and domains has its benefits (various issues can be resolved syntactically based on types and their mutual relationship without considering concrete sets of values)
- there is a need to keep “types and domains in sync”

## Definition (domains with similarities respecting type declaration)

Set  $\{\langle D_y, \approx_y \rangle \mid y \in Y\}$  of domains with similarities **respects** type declaration  $\mathbf{\Lambda} = \langle \Lambda, \lambda \rangle$  if for all  $y_1, y_2 \in Y$ , the following condition is satisfied: If  $\lambda(y_1) = \lambda(y_2)$ , then  $\langle D_{y_1}, \approx_{y_1} \rangle$  coincides with  $\langle D_{y_2}, \approx_{y_2} \rangle$ .

## note:

- for  $\lambda(y_1) \neq \lambda(y_2)$ , we can have  $D_{y_1} \cap D_{y_2} \neq \emptyset$

# Database Scheme

## motivation:

- relational queries involve (sets of) RDTs denoted by names;
- in order to formalize the concept, we introduce database schemes

## Definition (database scheme, relation symbols)

A **database scheme on  $Y$**  is any couple  $\langle \mathbb{R}, \varrho \rangle$  where  $\mathbb{R}$  is a finite nonempty set of **relation symbols** and  $\varrho$  is a map  $\varrho: \mathbb{R} \rightarrow 2^Y$  assigning to each relation symbol  $r \in \mathbb{R}$  a relation scheme  $\varrho(r) \subseteq Y$ .

In terminology of Date, relation symbols correspond to RELVARs (relation variables):

- relation symbol is a *name* that can be *bound* to an RDT
- $\varrho(r)$  prescribes a relation scheme of an RDT that can be bound to  $r$
- RDTs are bound to relation symbols in *database instances*, ...

# Database Instance

## Definition (database instance)

A **database instance**  $\mathcal{D}$  (shortly, an **instance**) of a database scheme  $\langle \mathbb{R}, \varrho \rangle$  and object constants  $\mathfrak{C}$  over domains  $\{\langle D_y, \approx_y \rangle \mid y \in Y\}$  with similarities is a triplet

$\mathcal{D} = \langle U^{\mathcal{D}}, \mathbb{R}^{\mathcal{D}}, \mathfrak{C}^{\mathcal{D}} \rangle$  where

- 1  $U^{\mathcal{D}} = \{\langle D_y, \approx_y \rangle \mid y \in Y\}$  is a set of domains with similarities;
- 2  $\mathbb{R}^{\mathcal{D}}$  is a set of RDTs such that for each relation symbol  $r \in \mathbb{R}$  there is an RDT  $r^{\mathcal{D}} \in \mathbb{R}^{\mathcal{D}}$  on  $\varrho(r)$  over  $U^{\mathcal{D}}$ ;
- 3  $\mathfrak{C}^{\mathcal{D}} = \{\mathfrak{c}^{\mathcal{D}} \in \bigcup_{y \in Y} D_y \mid \mathfrak{c} \in \mathfrak{C}\}$  is a set of values interpreting object constants.

If  $\mathfrak{C}$  is clear from the context, we say that  $\mathcal{D}$  is an instance of  $\langle \mathbb{R}, \varrho \rangle$ . Moreover, we say that  $\mathcal{D} = \langle U^{\mathcal{D}}, \mathbb{R}^{\mathcal{D}}, \mathfrak{C}^{\mathcal{D}} \rangle$  **respects type declaration**  $\Lambda = \langle \Lambda, \lambda \rangle$  if  $U^{\mathcal{D}}$  respects  $\Lambda$  and if for each  $\mathfrak{c} \in \mathfrak{C}$ ,  $\mathfrak{c}^{\mathcal{D}} \in D_y$  whenever  $\lambda(\mathfrak{c}) = \lambda(y)$ .

## notes on database schemes $\times$ database instances:

- $\langle \mathbb{R}, \varrho \rangle$  and  $\mathfrak{C}$  determine a *predicate language* (syntactic notions)
- $\mathcal{D} = \langle U^{\mathcal{D}}, \mathbb{R}^{\mathcal{D}}, \mathfrak{C}^{\mathcal{D}} \rangle$  is a first-order *predicate structure* (semantic notion)

# A General View on Relational Queries

general approach to queries:

## Definition (query)

A *query* on database scheme  $\langle \mathbb{R}, \rho \rangle$  and object constants  $\mathfrak{C}$  is any partial recursive function from the set of all database instances of  $\langle \mathbb{R}, \rho \rangle$  and  $\mathfrak{C}$  to the set of all RDTs.

**typically:**

- queries are induced by evaluating expressions or formulas in database instances

**two (basic) types of query systems in our model:**

- 1 queries are induced by evaluating *relation algebra expressions* ...
- 2 queries are induced by evaluating *formulas of domain relational calculus* ...  
in database instances



# Relation Algebra

## Definition (Relation Algebra Expressions, RA-expressions)

Let  $\langle \mathbb{R}, \varrho \rangle$  be a database scheme on  $Y$  and let  $\Lambda = \langle \Lambda, \lambda \rangle$  be a type declaration for  $Y$  and  $\mathfrak{C}$ . Then, **relational algebra expressions over  $\langle \mathbb{R}, \varrho \rangle$**  (shortly, **RA-expressions**) are defined as follows:

- 1 If  $r \in \mathbb{R}$ , then  $r$  is RA-expression on  $\varrho(r)$ ;
- 2 if  $a \in L$ , then  $\bar{a}_\emptyset$  is RA-expression on  $\emptyset$ ;
- 3 if  $c \in \mathfrak{C}$ ,  $y \in Y$ , and  $\lambda(c) = \lambda(y)$ , then  $[y:c]$  is RA-expression on  $\{y\}$ ;
- 4 if  $E_1$  and  $E_2$  are RA-expressions on  $R$ ,  
then  $(E_1 \cap E_2)$  and  $(E_1 \cup E_2)$  are RA-expressions on  $R$ ,
- 5 if  $E_1$ ,  $E_2$ , and  $E_3$  are RA-expressions on  $R$ , then  $(E_1 \rightarrow^{E_3} E_2)$  is RA-expression on  $R$ ;
- 6 if  $E_1$  is RA-expression on  $R_1$  and  $E_2$  is RA-expression on  $R_2$ ,  
then  $(E_1 \bowtie E_2)$  is RA-expression on  $R_1 \cup R_2$ ;
- ⋮

# Relation Algebra

## Definition (... continuation)

⋮

- 7 if  $E$  is RA-expression on  $T$  and  $R \subseteq T$ , then  $\pi_R(E)$  is RA-expression on  $R$ ;
- 8 if  $E_1$  is RA-expression on  $R$ ,  $E_2$  is RA-expression on  $S \subseteq R$ , and  $E_3$  is RA-expression on  $T = R \setminus S$ , then  $(E_1 \div^{E_3} E_2)$  is RA-expression on  $T$ ;
- 9 if  $E$  is RA-expression on  $R$ ,  $y \in R$ , and  $z \in R \cup \mathfrak{C}$  such that  $\lambda(y) = \lambda(z)$ , then  $\sigma_{y \approx z}(E)$  is RA-expression on  $R$ ;
- 10 if  $E$  is RA-expression on  $R$ , then  $\Delta E$  and  $\nabla E$  are RA-expressions on  $R$ ;
- 11 if  $E$  is RA-expression on  $R$ , and  $h: R \rightarrow Y$  is an injective map such that for each  $y \in R$  we have  $\lambda(h(y)) = \lambda(y)$ , then  $\rho_h(E)$  is RA-expression on  $h(R)$ .

All RA-expressions result by application of 7–11. In addition, if  $E$  is RA-expression on  $R$ , we call  $R$  the **relation scheme of  $E$**  and denote it by  $\text{sch}(E)$ .

## Definition (Results of RA-expressions in Database Instances)

Let  $\mathcal{D} = \langle U^{\mathcal{D}}, \mathbb{R}^{\mathcal{D}}, \mathcal{C}^{\mathcal{D}} \rangle$  be an instance of database scheme  $\langle \mathbb{R}, \varrho \rangle$  which respects a type declaration  $\Lambda = \langle \Lambda, \lambda \rangle$  for  $Y$  and  $\mathcal{C}$ . Then, for each RA-expression  $E$  over  $\langle \mathbb{R}, \varrho \rangle$ , we define an RDT  $E^{\mathcal{D}}$ , called the **result of  $E$  in  $\mathcal{D}$**  as follows:

- 1 If  $E$  is  $r \in \mathbb{R}$ , then  $E^{\mathcal{D}} = r^{\mathcal{D}}$ ;
- 2 if  $E$  is  $\bar{a}_{\emptyset}$ , then  $E^{\mathcal{D}} = a_{\emptyset}$ ;
- 3 if  $E$  is  $[y:c]$ , then  $E^{\mathcal{D}}$  is  $[y:c^{\mathcal{D}}]$ ;
- 4 if  $E$  is  $E_1 \cap E_2$ , then  $E^{\mathcal{D}} = E_1^{\mathcal{D}} \cap E_2^{\mathcal{D}}$ ; if  $E$  is  $E_1 \cup E_2$ , then  $E^{\mathcal{D}} = E_1^{\mathcal{D}} \cup E_2^{\mathcal{D}}$ ;
- 5 if  $E$  is  $E_1 \xrightarrow{E_3} E_2$ , then  $E^{\mathcal{D}} = E_1^{\mathcal{D}} \xrightarrow{E_3^{\mathcal{D}}} E_2^{\mathcal{D}}$ ;
- 6 if  $E$  is  $E_1 \bowtie E_2$ , then  $E^{\mathcal{D}} = E_1^{\mathcal{D}} \bowtie E_2^{\mathcal{D}}$ ;
- 7 if  $E$  is  $\pi_R(F)$ , then  $E^{\mathcal{D}} = \pi_R(F^{\mathcal{D}})$ ;
- 8 if  $E$  is  $E_1 \div^{E_3} E_2$ , then  $E^{\mathcal{D}} = E_1^{\mathcal{D}} \div^{E_3^{\mathcal{D}}} E_2^{\mathcal{D}}$ ;
- 9 if  $E$  is  $\sigma_{y \approx z}(F)$  and  $z \in Y$ , then  $E^{\mathcal{D}} = \sigma_{y \approx z}(F^{\mathcal{D}})$ ; if  $z \in \mathcal{C}$ , then  $E^{\mathcal{D}} = \sigma_{y \approx z^{\mathcal{D}}}(F^{\mathcal{D}})$ ;
- 10 if  $E$  is  $\Delta F$ , then  $E^{\mathcal{D}} = \Delta F^{\mathcal{D}}$ ; if  $E$  is  $\nabla F$ , then  $E^{\mathcal{D}} = \nabla F^{\mathcal{D}}$ ;
- 11 if  $E$  is  $\rho_h(F)$ , then  $E^{\mathcal{D}} = \rho_h(F^{\mathcal{D}})$ .

## Example (Relation Algebra Expressions)

$r$

$\sigma_{y \approx c}(r)$

$0.5\emptyset \bowtie \sigma_{y \approx c}(r)$

$\pi_{\{x,y\}}(0.5\emptyset \bowtie \sigma_{y \approx c}(r))$

$\pi_{\{x,y\}}(0.5\emptyset \bowtie \sigma_{y \approx c}(r)) \bowtie s$

$\pi_{\{x,y\}}(0.5\emptyset \bowtie \sigma_{y \approx c}(r)) \bowtie \Delta s$

$\sigma_{x \approx z}(\pi_{\{x,y\}}(0.5\emptyset \bowtie \sigma_{y \approx c}(r)) \bowtie \Delta s)$

$\vdots$

# Notes on Integrity Constraints

## classic RM:

- $E \subseteq F$  (with  $E, F$  being RA-expressions)
- functional dependencies, keys, referential integrity constraints ...

## similarity-based integrity constraints:

- threshold-based:  $S(E, F) \geq a$ 
  - $S(E, F) \geq a$  is **valid** in  $\mathcal{D}$  if  $S(E^{\mathcal{D}}, F^{\mathcal{D}}) \geq a$
- graded approach  $S(E, F)$ 
  - $S(E^{\mathcal{D}}, F^{\mathcal{D}})$  is a **degree of validity** of  $S(E, F)$  in  $\mathcal{D}$

## Example (similarity-based (soft) constraints)

For  $Q$  (high-quality wines),  $E$  (expensive wines),

$$S(\pi_{\{wine-id\}}(Q), \pi_{\{wine-id\}}(E)) \geq 0.95$$

means “high-quality wines are expensive at least to degree 0.95”.

# Domain Relational Calculus: The Basic Idea

**object variables** = names for (possibly changing) values in database instances

- $X$ : denumerable set of all object variables
- (object) variables from  $X$  are denoted by  $x, x', x_1, x_2, \dots$
- object constraint (has a fixed value in  $\mathcal{D}$ )  $\times$   
object variable (values not determined solely by  $\mathcal{D}$ )

**from formulas to queries:**

- take a suitable first-order formula  $\varphi$ , and evaluate  $\varphi$  in  $\mathcal{D}$
- note: the value of  $\varphi$  in  $\mathcal{D}$  depends on values of *free variables*
- if the free variables are bound to  $y$ -values of a tuple  $t$ , then the value of  $\varphi$  in  $\mathcal{D}$  (under the valuation of free variables) is the *degree to which  $\varphi$  is true in  $\mathcal{D}$*  and can be seen as a rank of  $t$ .
- collect all tuples with nonzero ranks  $\Rightarrow$  result
- issues: domain independence

# Range Component, Range Declaration

**requirement:** finiteness of query results

## Definition (range component, range declaration)

A **range component** (over  $\langle \mathbb{R}, \varrho \rangle$  and  $\mathfrak{C}$ ) is (i) any object constant from  $\mathfrak{C}$  and (ii) any expression  $r(y)$ , where  $r \in \mathbb{R}$ , and  $y \in \varrho(r)$ . A **range declaration** (over  $\langle \mathbb{R}, \varrho \rangle$  and  $\mathfrak{C}$ ) is any finite set  $R$  of range components.

## Definition (type compatibility of range declarations)

A range declaration  $R$  is **compatible** with type declaration  $\mathbf{\Lambda} = \langle \Lambda, \lambda \rangle$  if all attributes and constants appearing in  $R$  have the same type  $\tau \in \Lambda$  which is called the **type of  $R$**  and denoted  $\lambda(R)$ . The set of all range declarations over  $\langle \mathbb{R}, \varrho \rangle$  and  $\mathfrak{C}$  compatible with  $\mathbf{\Lambda}$  is denoted by  $\text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$ . Any map  $rd: X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$  is called a **range declaration for variables from  $X$** .

**note:** value  $rd(x)$  is called a range declaration for  $x$

# Values of Range Declarations in Database Instances

## syntax $\times$ semantics:

- range declaration (syntactic notion),
- value of range declaration (semantic notion) introduced as follows:

### Definition (value of range declaration)

Let  $\mathcal{D} = \langle U^{\mathcal{D}}, \mathbb{R}^{\mathcal{D}}, \mathcal{C}^{\mathcal{D}} \rangle$  be an instance of database scheme  $\langle \mathbb{R}, \varrho \rangle$  and let  $R \in \text{Rd}(\mathbb{R}, \varrho, \mathcal{C}, \Lambda)$  be a range declaration. The following set

$$\mathbb{R}^{\mathcal{D}} = \{c^{\mathcal{D}} \mid c \in R\} \cup \{d \in D_y \mid \text{if } r(y) \in R \text{ and } \pi_{\{y\}}(r^{\mathcal{D}})(\{\langle y, d \rangle\}) > 0\}$$

is called the **value of  $R$  in database instance  $\mathcal{D}$** .

## meaning:

- in  $\mathcal{D}$ , free variable  $x$  can take values only from  $rd(x)^{\mathcal{D}}$
- ranges of bound variables are specified in quantifiers ...



# Formulas

## Definition (DRC-RD formulas)

Let  $\langle \mathbb{R}, \varrho \rangle$  be a database scheme on  $Y$ , let  $\mathbf{\Lambda}$  be a type declaration for  $Y$  and  $\mathfrak{C}$ , and let  $X$  be a set of object variables. Then, **formulas of the domain relational calculus with range declarations** (shortly, **DRC-RD formulas** or **formulas**) over  $\langle \mathbb{R}, \varrho \rangle$ ,  $\mathfrak{C}$ , and  $L$  are defined as follows:

- 1 If  $r \in \mathbb{R}$  is a relation symbol such that  $\varrho(r) = \{y_1, \dots, y_n\}$  and  $x_1, \dots, x_n \in X$  are arbitrary variables, then  $r(y_1:x_1, \dots, y_n:x_n)$  is formula;
- 2 if  $a \in L$ , then  $\bar{a}$  is formula;
- 3 if  $x \in X$  and  $z \in X \cup \mathfrak{C}$ , then  $x \approx z$  is formula;
- 4 if  $\varphi$  and  $\psi$  are formulas then  $(\varphi \otimes \psi)$ ,  $(\varphi \wedge \psi)$ ,  $(\varphi \vee \psi)$ ,  $(\varphi \Rightarrow \psi)$  are formulas;
- 5 if  $\varphi$  is formula, then  $\Delta\varphi$  is formula;
- 6 if  $\varphi$  is formula,  $x \in X$ , and  $R \in \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$ , then  $(\forall x \in R)\varphi$  and  $(\exists x \in R)\varphi$  are formulas and  $\varphi$  is called the scope of the quantifiers  $(\forall x \in R)$  and  $(\exists x \in R)$ , respectively.

All formulas result by application of 1–6.

# Notes on Formulas

## unordered nature of attributes:

- $\mathfrak{r}(y_1:x_1, y_2:x_2)$  is (syntactically) equivalent to  $\mathfrak{r}(y_2:x_2, y_1:x_1)$

**we adopt the usual notions:** *subformula*, *free* and *bound* occurrences of variables

## Example (natural reading of formulas)

Formulas are read the usual way.

For instance,

- $\mathfrak{r}(y_1:x, y_2:z) \otimes \mathfrak{s}(y_1:x, y_2:z)$  reads “the value of  $y_1$  denoted by  $x$  and the value of  $y_2$  denoted by  $z$  are related according to  $\mathfrak{r}$  *and* the value of  $y_1$  denoted by  $x$  and the value of  $y_2$  denoted by  $z$  are related according to  $\mathfrak{s}$ ”.
- $(\exists x \in R) \mathfrak{r}(y_1:x, y_2:z)$  reads “*there is* a value of  $y_1$  denoted by  $x$  (which ranges over  $R$ ) which is according to  $\mathfrak{r}$  related to the value of  $y_2$  denoted by  $z$ ”, ...

# Formulas and Type Safety

## motivation:

- define formulas which together with range declarations for (free) variables “make sense” from the point of view of the type system; (e.g.,  $\mathfrak{r}(y:x)$  makes sense only if the type of  $y$  and the range declaration for  $x$  is the same)

## Definition (type safety)

Formula  $\varphi$  is called **safe with respect to**  $rd: X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$  if the following conditions hold for any variable  $x$  that occurs in  $\varphi$ :

- 1 If the occurrence of  $x$  in a subformula  $\psi$  of  $\varphi$  is free in  $\varphi$  and if
  - $\psi$  is  $\mathfrak{r}(y:x, \dots)$ , then  $\lambda(y) = \lambda(rd(x))$ ; or if
  - $\psi$  is  $x \approx c$ , then  $\lambda(c) = \lambda(rd(x))$ ; or if
  - $\psi$  is  $x \approx y$  and this occurrence of  $y$  is free in  $\varphi$ , then  $\lambda(rd(x)) = \lambda(rd(y))$ .
- 2 For any subformula  $(Q_{x \in \mathbb{R}})\psi$  of  $\varphi$  with  $Q$  being  $\forall$  or  $\exists$ , the formula  $\psi$  is safe with respect to  $rd': X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$ , where  $rd'(y) = rd(y)$  for each  $y \neq x$ , and  $rd'(x) = \mathbb{R}$ .

# Valuation of Object Variables

## values of objects variables $\times$ objects constants:

- *object constants* are assigned values by database instances,
- (*free*) *object variables* are assigned values by valuations (as in predicate logic),
- specific feature: we only allow values from values of range declarations

## Definition (valuation, values of object variables)

Let  $\mathcal{D} = \langle U^{\mathcal{D}}, \mathbb{R}^{\mathcal{D}}, \mathcal{C}^{\mathcal{D}} \rangle$  be a database instance of  $\langle \mathbb{R}, \varrho \rangle$  where  $U^{\mathcal{D}} = \{ \langle D_y, \approx_y \rangle \mid y \in Y \}$  and let  $rd: X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathcal{C}, \mathbf{\Lambda})$  be a range declaration for variables. A map  $v: X \rightarrow \bigcup_{y \in Y} D_y$  such that  $v(x) \in rd(x)^{\mathcal{D}}$  is called a  **$\mathcal{D}$ -valuation** (of variables from  $X$  with respect to  $rd$ ) and  $v(x)$  is called the **value of  $x$  under  $v$** .

## notation:

- for two  $\mathcal{D}$ -valuations  $w$  and  $v$  we write  $w =_x v$  whenever  $w(x') = v(x')$  for each  $x' \in X$  different from  $x$

# Interpretation of Formulas

## Definition (interpretation of formulas)

Let  $\mathcal{D} = \langle U^{\mathcal{D}}, \mathbb{R}^{\mathcal{D}}, \mathfrak{C}^{\mathcal{D}} \rangle$  be an instance of database scheme  $\langle \mathbb{R}, \varrho \rangle$  which respects a type declaration  $\Lambda = \langle \Lambda, \lambda \rangle$  for  $Y$  and  $\mathfrak{C}$ , let  $\mathbf{v}$  be a  $\mathcal{D}$ -valuation, and let  $\varphi$  be a formula which is safe with respect to range declaration  $rd: X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \Lambda)$  for variables from  $X$ . Then, we define a degree  $\|\varphi\|_{\mathcal{D}, \mathbf{v}}^{rd} \in L$ , called the **degree to which  $\varphi$  is true in  $\mathcal{D}$  under  $\mathbf{v}$  and  $rd$**  as follows:

- 1 If  $\varphi$  is  $\mathfrak{r}(y_1:\mathfrak{x}_1, \dots, y_n:\mathfrak{x}_n)$ , then  $\|\varphi\|_{\mathcal{D}, \mathbf{v}}^{rd} = \mathfrak{r}^{\mathcal{D}}(t)$ ,  
where  $t$  is a tuple on  $\varrho(\mathfrak{r})$  such that  $t(y_i) = \mathbf{v}(\mathfrak{x}_i)$  for all  $i = 1, \dots, n$ ;
- 2 if  $\varphi$  is  $\bar{a}$ , then  $\|\varphi\|_{\mathcal{D}, \mathbf{v}}^{rd} = a$ ;
- 3 if  $\varphi$  is  $\mathfrak{x} \approx \mathfrak{z}$ , then  $\|\varphi\|_{\mathcal{D}, \mathbf{v}}^{rd} = \mathbf{v}(\mathfrak{x}) \approx_y \mathbf{v}(\mathfrak{z})$  where  $\lambda(y) = \lambda(rd(\mathfrak{x}))$ ;  
if  $\varphi$  is  $\mathfrak{x} \approx \mathfrak{c}$ , then  $\|\varphi\|_{\mathcal{D}, \mathbf{v}}^{rd} = \mathbf{v}(\mathfrak{x}) \approx_y \mathfrak{c}^{\mathcal{D}}$  where  $\lambda(y) = \lambda(rd(\mathfrak{x}))$ ;
- ⋮

# Interpretation of Formulas

## Definition (... continuation)

⋮

4 if  $\varphi$  is  $\psi \otimes \chi$ , then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \|\psi\|_{\mathcal{D},v}^{rd} \otimes \|\chi\|_{\mathcal{D},v}^{rd}$ ;

if  $\varphi$  is  $\psi \wedge \chi$ , then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \|\psi\|_{\mathcal{D},v}^{rd} \wedge \|\chi\|_{\mathcal{D},v}^{rd}$ ;

if  $\varphi$  is  $\psi \vee \chi$ , then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \|\psi\|_{\mathcal{D},v}^{rd} \vee \|\chi\|_{\mathcal{D},v}^{rd}$ ;

if  $\varphi$  is  $\psi \Rightarrow \chi$ , then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \|\psi\|_{\mathcal{D},v}^{rd} \rightarrow \|\chi\|_{\mathcal{D},v}^{rd}$ ;

5 if  $\varphi$  is  $\Delta\psi$ , then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \begin{cases} 1, & \|\psi\|_{\mathcal{D},v}^{rd} = 1, \\ 0, & \text{otherwise;} \end{cases}$

6 if  $\varphi$  is  $(\forall x \in R)\psi$  then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \bigwedge \{ \|\psi\|_{\mathcal{D},w}^{rd'} \mid w =_x v \text{ and } w(x) \in R^{\mathcal{D}} \}$ ;

if  $\varphi$  is  $(\exists x \in R)\psi$  then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \bigvee \{ \|\psi\|_{\mathcal{D},w}^{rd'} \mid w =_x v \text{ and } w(x) \in R^{\mathcal{D}} \}$ ;

where in both cases  $rd' : X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$  is a range declaration for variables such that  $rd'(y) = rd(y)$  for each  $y \neq x$ , and  $rd'(x) = R$ .

# Results of Queries in DRC-RD

## Definition (target component, target set, query result)

Let  $\mathcal{D}$  be an instance of  $\langle \mathbb{R}, \rho \rangle$  which respects type declaration  $\Lambda = \langle \Lambda, \lambda \rangle$  for  $Y$  and  $\mathcal{C}$  and let  $\varphi$  be a formula which is safe with respect to range declaration  $rd: X \rightarrow \text{Rd}(\mathbb{R}, \rho, \mathcal{C}, \Lambda)$  for variables from  $X$ .

A **target component** (based on  $rd$ ) is any expression  $y:z$  such that  $y \in Y$ ,  $z \in X$ , and  $\lambda(y) = \lambda(rd(z))$ . A **target set** (based on  $rd$ )

$$\mathcal{T} = \{y_1:z_1, \dots, y_n:z_n\}$$

is any finite set of target components (all  $y_1, \dots, y_n$  and  $z_1, \dots, z_n$  pairwise distinct).

For  $\mathcal{T}$  and  $\varphi$ , we define an RDT  $\mathcal{T}\varphi^{rd, \mathcal{D}}$  on the relation scheme  $R = \{y_1, \dots, y_n\}$ , called the **result of query  $\mathcal{T}\varphi$  in  $\mathcal{D}$  under  $rd$**  as follows:

$$(\mathcal{T}\varphi^{rd, \mathcal{D}})(r) = \begin{cases} \bigvee \{ \|\varphi\|_{\mathcal{D}, v}^{rd} \mid \text{for each } i: v(z_i) = r(y_i) \}, & \text{if } r \in \prod_{y_i \in R} rd(z_i)^{\mathcal{D}}, \\ 0, & \text{otherwise.} \end{cases}$$

Moreover, if  $v(z_i) = r(y_i)$  for all  $i = 1, \dots, n$ , we say that  $v$  is *induced by  $r$  and  $\mathcal{T}$* .

## Theorem

The following are true for each safe formula  $\varphi$ :

- 1 If for each  $x$  which is free in  $\varphi$ , we have  $v(x) = w(x)$ , then  $\|\varphi\|_{\mathcal{D},v}^{rd} = \|\varphi\|_{\mathcal{D},w}^{rd}$ .
- 2 If for each  $x$  which is free in  $\varphi$ , we have  $rd(x) = rd'(x)$ , then for any  $v$  which is a  $\mathcal{D}$ -valuation with respect to both  $rd$  and  $rd'$ , we have  $\|\varphi\|_{\mathcal{D},v}^{rd} = \|\varphi\|_{\mathcal{D},v}^{rd'}$ .
- 3 For each target set  $\mathcal{T}$ , there is safe formula  $\psi$  so that the variables in  $\mathcal{T}$  are exactly the free variables in  $\psi$  and for each  $\mathcal{D}$ , we have  $\mathcal{T}\varphi^{rd,\mathcal{D}} = \mathcal{T}\psi^{rd,\mathcal{D}}$  and

$$(\mathcal{T}\psi^{rd,\mathcal{D}})(r) = \|\varphi\|_{\mathcal{D},v}^{rd}$$

for all  $\mathcal{D}$ -valuations  $v$  induced by  $r$  and  $\mathcal{T}$ .

**consequence:** simplified notation of queries:

$$\{y_1:x_1 \in R_1, \dots, y_n:x_n \in R_n\}\varphi,$$

meaning we consider  $\{y_1:x_1, \dots, y_n:x_n\}\varphi$  under a range declaration of variables  $rd: X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$  such that  $rd(x_i) = R_i$  for all  $i = 1, \dots, n$  and all free variables in  $\varphi$  are exactly  $x_1, \dots, x_n$ .



# Generality of Similarity-Based Databases

## Theorem (on generalization of the ordinary DRC-RD)

*If  $\mathbf{L}$  is a two-element Boolean algebra then each  $\mathcal{T}_\varphi^{\mathcal{D}}$  can be obtained as a result of a query in the ordinary DRC-RD calculus provided that all similarity relations in  $\mathcal{D}$  are identities.* □

## Theorem (crisp case of DRC-RD)

*If  $\mathbf{L}$  is arbitrary complete residuated lattice, all  $r^{\mathcal{D}}$  in  $\mathcal{D}$  are nonranked, all similarities in  $\mathcal{D}$  are identities, and the only constants of truth degrees appearing in  $\varphi$  are  $\bar{0}$  or  $\bar{1}$ , then  $\mathcal{T}_\varphi^{\mathcal{D}}$  is nonranked and a corresponding ordinary relation can be obtained as a result of a query in the ordinary DRC-RD calculus.* □

**meaning:** the similarity-based model (and its DRC-RD) is a (nontrivial) generalization of the RM (and DRC-RD)

# Relational Completeness

## Theorem (RA is more general than DRC-RD)

Let  $\varphi$  be a safe formula with respect to  $rd: X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$  and let  $\mathcal{T}$  be a target set based on  $rd$ . Then, there is an RA-expression  $E_\varphi$  over  $\langle \mathbb{R}, \varrho \rangle$  such that  $\mathcal{T}_\varphi^{rd, \mathcal{D}} = E_\varphi^{\mathcal{D}}$  for any instance  $\mathcal{D}$  of  $\langle \mathbb{R}, \varrho \rangle$  which respects  $\mathbf{\Lambda}$ .  $\square$

## Theorem (DRC-RD is more general than RA)

Let  $E$  be an RA-expression over  $\langle \mathbb{R}, \varrho \rangle$ . Then, there are a range declaration  $rd_E: X \rightarrow \text{Rd}(\mathbb{R}, \varrho, \mathfrak{C}, \mathbf{\Lambda})$ , a formula  $\varphi_E$  which is safe with respect to  $rd_E$ , and a target set  $\mathcal{T}_E$  based on  $rd_E$  such that  $E^{\mathcal{D}} = \mathcal{T}_E^{\varphi_E^{rd_E, \mathcal{D}}}$  holds true for any instance  $\mathcal{D}$  of  $\langle \mathbb{R}, \varrho \rangle$  which respects  $\mathbf{\Lambda}$ .  $\square$

### details in:



BELOHLAVEK, R. AND VYCHODIL, V.

Relational Similarity-Based Databases, Part I: Query Systems (in preparation)

## Derived Operations: Similarity-Based $\theta$ -restrictions

**idea:** for any RA-expression  $E$  and safe formula  $\theta$  whose free variables are in a correspondence with the attributes from  $\text{sch}(E)$ , we can introduce a (**similarity-based**)  $\theta$ -restriction  $\sigma_\theta(E^{\mathcal{D}})$  of RDT  $E^{\mathcal{D}}$  by  $\theta$  as follows:

$$\sigma_\theta(E^{\mathcal{D}}) = \mathcal{T}_E(\varphi_E \otimes \theta)^{rd_E, \mathcal{D}},$$

where  $\varphi_E$ ,  $rd_E$ , and  $\mathcal{T}_E$  correspond to  $E$

**generalization of similarity-based restrictions:** take  $y \approx c$  for  $\theta$

Example (threshold similarity-based restriction)

Example of a derived restriction: take  $\bar{a} \Rightarrow y \approx c$  for  $\theta$ , i.e.

$$(\sigma_{\bar{a} \Rightarrow y \approx c}(\mathcal{D}))(r) = \mathcal{D}(r) \otimes (a \rightarrow r(y) \approx_y c^{\mathcal{D}})$$

for all  $r \in \text{Tupl}(R)$ .

# Derived Operations: Similarity-Based Semijoins and Closures

**tuple similarity:** for any tuples  $r_1, r_2 \in \text{Tupl}(R)$ , we put

$$r_1 \approx_R r_2 = \bigwedge_{y \in R} r_1(y) \approx_y r_2(y)$$

and call  $r_1 \approx_R r_2$  the **degree to which  $r_1$  and  $r_2$  are similar**

**similarity-based semijoins:** for any RDTs  $\mathcal{D}_1$  and  $\mathcal{D}_2$  on  $R \cup S$  and  $S \cup T$  such that  $R \cap S = R \cap T = S \cap T = \emptyset$ , define a (**natural**) **similarity-based semijoin**  $\mathcal{D}_1 \ltimes^{\approx} \mathcal{D}_2$  of  $\mathcal{D}_1$  and  $\mathcal{D}_2$  (in this order) by putting

$$(\mathcal{D}_1 \ltimes^{\approx} \mathcal{D}_2)(rs) = \mathcal{D}_1(rs) \otimes \bigvee_{t \in \text{Tupl}(T)} \bigvee_{s' \in \text{Tupl}(S)} (\mathcal{D}_2(s't) \otimes s' \approx_S s)$$

for all  $r \in \text{Tupl}(R)$  and  $s \in \text{Tupl}(S)$

**similarity-based closures:** particular case for  $R = T = \emptyset$ , non-ranked  $\mathcal{D}_1$ , and  $\mathcal{D}_2 \subseteq \mathcal{D}_1$

$$(\mathcal{D}_1 \ltimes^{\approx} \mathcal{D}_2)(s) = \mathcal{D}_1(s) \otimes \bigvee_{s' \in \text{Tupl}(S)} (\mathcal{D}_2(s') \otimes s' \approx_S s)$$

## Derived Operations: Further Join-like Operations

**similarity-based  $\theta$ -join:** using cross joins and  $\theta$ -restrictions, for  $\mathcal{D}_1$  and  $\mathcal{D}_2$  on disjoint relation schemes, define  **$\theta$ -join** by

$$\mathcal{D}_1 \bowtie_{\theta} \mathcal{D}_2 = \sigma_{\theta}(\mathcal{D}_1 \bowtie \mathcal{D}_2).$$

**similarity-based join:** for  $\mathcal{D}_1$  and  $\mathcal{D}_2$  on  $R \cup S$  and  $S \cup T$  such that  $R \cap S = R \cap T = S \cap T = \emptyset$  and  $\mathcal{D}_3$  on  $S$ , define a (**natural**) **similarity-based join**  $\mathcal{D}_1 \bowtie_{\tilde{\mathcal{D}}_3} \mathcal{D}_2$  of  $\mathcal{D}_1$  and  $\mathcal{D}_2$  by putting

$$(\mathcal{D}_1 \bowtie_{\tilde{\mathcal{D}}_3} \mathcal{D}_2)(rst) = \mathcal{D}_3(s) \otimes \bigvee_{s', s''} (\mathcal{D}_1(rs') \otimes s' \approx_S s \otimes s \approx_S s'' \otimes \mathcal{D}_2(s''t)),$$

where  $s', s''$  range over  $\text{Tupl}(S)$ .

# Similarity-Based Functional Dependencies

## Definition (tuple similarity)

For an RDT  $\mathcal{D}$  on  $R$  over  $\{\langle D_y, \approx_y \rangle \mid y \in R\}$ , tuples  $r_1, r_2 \in \text{Tupl}(R)$  and  $C \in L^R$ , we introduce a **degree**  $r_1(C) \approx_{\mathcal{D}} r_2(C)$  to which  $r_1$  and  $r_2$  from  $\mathcal{D}$  have similar values on attributes from  $C$  by

$$r_1(C) \approx_{\mathcal{D}} r_2(C) = (\mathcal{D}(r_1) \otimes \mathcal{D}(r_2)) \rightarrow \bigwedge_{y \in R} (C(y) \rightarrow r_1(y) \approx_y r_2(y)).$$

**similarity-based functional dependencies:**  $A \Rightarrow B$  ( $A, B \in L^R$ )

## Definition (semantics of SBFDs)

Let  $\mathcal{D}$  be an RDT over  $\{\langle D_y, \approx_y \rangle \mid y \in R\}$  and let  $A, B \in L^R$ . A **degree**  $\|A \Rightarrow B\|_{\mathcal{D}}$  to which  $A \Rightarrow B$  is true in  $\mathcal{D}$  is defined by

$$\|A \Rightarrow B\|_{\mathcal{D}} = \bigwedge_{r_1, r_2 \in \text{Tupl}(R)} ((r_1(A) \approx_{\mathcal{D}} r_2(A))^* \rightarrow r_1(B) \approx_{\mathcal{D}} r_2(B)).$$

For  $A \Rightarrow B$ , there is RA-expression  $E$  such that  $\|A \Rightarrow B\|_{r, \mathcal{D}} = E^{\mathcal{D}}(\emptyset)$  for any  $\mathcal{D}$ .

## Further Extensions

### general comparators:

- other comparators than similarity
- *domains with general comparators*  $\langle D_y, \approx_y, \lesssim_y, \dots \rangle$

### additional connectives:

- linguistic hedges (nontrivial for general  $\mathbf{L}$ , important for dependencies)
- difference-like operations (and adjoint) non-idempotent disjunctions



$$a \ominus 0 = a, (a \ominus b) \ominus c = (a \ominus c) \ominus b, \bigvee_{i \in I} a_i \ominus b = \bigvee_{i \in I} (a_i \ominus b)$$

$$a \oplus b = \bigvee \{c \in L \mid c \ominus a \leq b\}$$

### additional quantifiers:

- can be done following Hájek's approach (with some adjustments)
- graded quantifiers like "many" (cf. approach by Badia)

# References

-  BADIA, A. 2009.  
*Quantifiers in Action*. Advances in Database Systems, vol. 37.  
Springer.
-  BELOHLAVEK, R. 2002.  
*Fuzzy Relational Systems: Foundations and Principles*.  
Kluwer Academic Publishers, Norwell, MA, USA.
-  HÁJEK, P. 1998.  
*Metamathematics of Fuzzy Logic*.  
Kluwer Academic Publishers, Dordrecht, The Netherlands.
-  LACROIX, M. AND PIROTTE, A. 1977.  
Domain-oriented relational languages.  
In *VLDB*. IEEE Computer Society, 370–378.
-  MAIER, D. 1983.  
*Theory of Relational Databases*.  
Computer Science Pr, Rockville, MD, USA.



# Concluding Remarks

## consequences of RA = DRC-RD:

- formal background for:
  - new derived (relational) operations
  - integrity constraints (similarity-based functional dependencies and more)
  - various types of data decomposition

## further results/directions:

- dependency theory
- efficiency issues
  - algorithms for query execution
  - algorithms for reasoning with dependencies
- design of (implementable) query languages

## details in:



BELOHLAVEK, R. AND VYCHODIL, V.

Relational Similarity-Based Databases, Part I: Query Systems

Relational Similarity-Based Databases, Part II: Dependencies in Data