

Implicit Representation of Bigranular Rules for Multigranular Data

Stephen J. Hegner
DBMS Research of New Hampshire, USA

M. Andrea Rodríguez
University of Concepción, Chile

DEXA 2018
29th International Conference on
Database and Expert Systems Applications
Regensburg, Germany
04 September 2018

The Idea of Multigranular Attributes

Spatio-temporal
attributes

Thematic
attributes

cmn = comuna/county
prv = provincia/province
rgn = región/region

Place	Time	Births
<i>Puerto_Montt_cmn</i>	<i>Y2017Q1</i>	<i>b₁</i>
<i>Llanquihue_prv</i>	<i>Y2017Q1</i>	<i>b₂</i>
<i>Los_Lagos_rgn</i>	<i>Y2017Q1</i>	<i>b₃</i>
<i>Puerto_Montt_cmn</i>	<i>Y2017</i>	<i>b₄</i>

Granules: The domain values are called *granules*.

Granular order: The granules of spatial and temporal attributes have inherent order structure.

Spatial containment: $Puerto_Montt_cmn \sqsubseteq Llanquihue_prv \sqsubseteq Los_Lagos_rgn$

Temporal interval containment: $Y2017Q1 \sqsubseteq Y2017$

Typical constraints: Functional dependency (FD) $\{Place, Time\} \rightarrow Births$.

- The number of births is monotonic w.r.t. space and time, so

$$b_1 \leq b_2 \leq b_3, b_2 \leq b_4.$$

Lattice-Like Operations on Granules

Place	Time	Births
<i>Osorno_prv</i>	Y2017Q1	b_1
<i>Llanquihue_prv</i>	Y2017Q1	b_2
<i>Chiloé_prv</i>	Y2017Q1	b_3
<i>Palena_prv</i>	Y2017Q1	b_4
<i>Los_Lagos_rgn</i>	Y2017Q1	b_5

Join: The four provinces join to the region.

$$\text{Los_Lagos_rgn} = \sqcup \{ \text{Osorno_prv}, \text{Llanquihue_prv}, \text{Chiloé_prv}, \text{Palena_prv} \}$$

Meet: Distinct provinces are disjoint (six possibilities in all).

$$\sqcap \{ \text{Osorno_prv}, \text{Llanquihue_prv} \} = \perp$$

Disjoint Join: The four provinces join *disjointly* to the region.

$$\text{Los_Lagos_rgn} = \sqcup^{\perp} \{ \text{Osorno_prv}, \text{Llanquihue_prv}, \text{Chiloé_prv}, \text{Palena_prv} \}$$

Consequence: $\sum_{i=1}^4 b_i = b_5$.

Observation: These lattice-like operations are *partial*.

Approaches to Modelling Multigranular Domains

Question: How can multigranular attributes be modelled for effective implementation? **Two possibilities:**

Single-structure model: Widely used in *Geographic Information Systems*.

Fix a domain \mathcal{D} ; semantics of granule g defined by a subset $\text{Sem}(g) \subseteq \mathcal{D}$.

- For a spatial attribute, $\mathcal{D} =$ suitable subset of $\mathbb{R} \times \mathbb{R}$.
- $\text{Sem}(g)$ the geographic region represented by g .

Advantages: Extensive model; well-developed theory and practice.

Disadvantages: Extreme resource demands, both space and time.

Constraint-based model: Work directly with constraints of the form $g_1 \sqsubseteq g_2$, $g \sqsubseteq \bigsqcup\{g_i \mid 1 \leq i \leq k\}$, and $g = \bigsqcup\{g_i \mid 1 \leq i \leq k\}$, among others.

Advantages: Only represent as much information as needed.

Challenges: Constraint inference,

Constraint retrieval (based upon features), Constraint consistency.

A Logic for Multigranular Domains

- A logic for representing knowledge within multigranular attributes has been developed [Hegner & Rodríguez, 2016, 2017].

Granule Expressions (terms): g , $\sqcup\{g \mid g \in S\}$, $\sqcap\{g \mid g \in S\}$, \perp , \top

Granule Rules (sentences): Combine granule expressions using \sqsubseteq (and equality). **Examples:** $g_1 \sqsubseteq g_2$, $\sqcap\{g_1, g_2\} = \perp$.

Semantics/Models: Set semantics. Fix a domain \mathcal{D} .

- A model assigns to each granule a semantics $\text{Sem}(g) \subseteq \mathcal{D}$ in a manner which respects the operations.

Example: $g_1 \sqsubseteq g_2$ iff $\text{Sem}(g_1) \subseteq \text{Sem}(g_2)$.

Satisfiability (of a set of rules): Related to distributivity of order operators.

- Very complex problem in theory.
- Not a problem in practice — axioms are models of “real” things.

Needs for an inference/lookup mechanism: The theory does not provide any such mechanism, beyond raw testing.

Goal of this work: Provide such a mechanism for the *common case*.

Requirements for Join-Rule Lookup

Join rule: One of the following forms:

	$g = \sqcup S$	$g = \sqcup S$
Head	$\textcircled{g} \sqsubseteq \sqcup S$	$g \sqsubseteq \sqcup \textcircled{S}$ Body

Primary Lookup requirements:

Head lookup: Given a granule g' , find all rules with head g' .

Body lookup: Given a set T of granules, find all rules with T contained in the body ($T \subseteq S$).

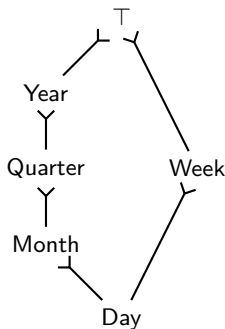
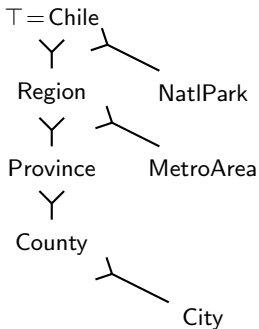
Complication: In the general case, the rules to be found may need to be derived first.

- A *knowledge base* of rules, not just a *database*.

Goal: For rule retrieval, *make the common case fast*.

Question: What is the common case?

Granularities — Organizing Granules



- The granules of each attribute are partitioned into a hierarchy of *granularities*.

Order: $G_1 \leq G_2 \Leftrightarrow ((\forall g_1 \in \text{Granules}\langle G_1 \rangle)(\exists g_2 \in \text{Granules}\langle G_2 \rangle)(g_1 \sqsubseteq g_2))$.

Examples: Every county is contained in a (unique) province.
Every day is contained in a (unique) week.

Disjointness: Distinct granules of the same granularity are disjoint.

Common Cases of Bigranular Rules

A common case: A family of *bigranular* rules between a pair of granules.

Bigranular rule: All granules in the body are of the same granularity.

- The head is necessarily of a different granularity.

Common case 1: Equality-join order property:

$G_1 \trianglelefteq G_2 \stackrel{\text{def}}{=} \text{every granule of } G_2 \text{ is the (disjoint) join of granules of } G_1.$

Example of common case 1: Province \trianglelefteq Region.

$Los_Lagos_rgn = \left[\perp \right] \{ Osorno_prv, Llanquihue_prv, Chiloé_prv, Palena_prv \}.$

$BíoBío_rgn = \left[\perp \right] \{ Arauco_prv, BíoBío_prv, Concepción_prv, ~~Ñuble_prv~~ \}.$

Common case 2: Subsumption-join order property:

$G_1 \otimes G_2 \stackrel{\text{def}}{=} \text{every granule of } G_2 \text{ is contained in the (disjoint) join of granules of } G_1.$

Example of common case 2: MetroArea \otimes Province

$Gran_Puerto_Montt_urb \sqsubseteq \left[\perp \right] \{ Puerto_Montt_cmn, Puerto_Varas_cmn \}$

Resolvability and the Nondisjointness Relation

Context: Let \mathcal{C} denote the set of constraints which hold on the multigranular attribute under consideration.

- Given a rule φ , there are three possible cases.

(a) $\mathcal{C} \models \varphi$

(b) $\mathcal{C} \models \neg\varphi$

(c) Neither of these

Resolvability: The rule φ is *resolvable* from \mathcal{C} if one of (a) or (b) holds.

- Written $\mathcal{C} \stackrel{\pm}{\models} \varphi$.

Context: $\langle G_1, G_2 \rangle$ a pair of granularities.

Full disjointness resolvability: $\langle G_1, G_2 \rangle$ is *fully disjointness resolvable* if $(\forall g_1 \in \text{Granules}\langle G_1 \rangle)(\forall g_2 \in \text{Granules}\langle G_2 \rangle)(\mathcal{C} \stackrel{\pm}{\models} (\bigcap \{g_1, g_2\} = \perp))$.

A compact relational representation under full disjointness resolvability:

$$\text{NRel}_{\langle G_1, G_2 \rangle} = \{ \langle g_1, g_2 \rangle \mid \mathcal{C} \models (\bigcap \{g_1, g_2\} \neq \perp) \}$$

- $\mathcal{C} \models ((\bigcap \{g_1, g_2\} = \perp))$ iff $(g_1, g_2) \notin \text{NRel}_{\langle G_1, G_2 \rangle}$.

Symmetry: Note that $\text{NRel}_{\langle G_1, G_2 \rangle} = \text{NRel}_{\langle G_2, G_1 \rangle}$ always holds.

- Use $\text{NRel}_{\langle G_1, G_2 \rangle}$ rather than $\text{DRel}_{\langle G_1, G_2 \rangle}$ (the corresponding relation for disjointness) because $\text{NRel}_{\langle G_1, G_2 \rangle}$ is usually much smaller.

Main Representation Theorem for \sqsubseteq

Recall: $G_1 \sqsubseteq G_2$ iff

$$(\forall g_2 \in \text{Granules}\langle G_2 \rangle)(\exists S \subseteq_f \text{Granules}\langle G_1 \rangle)(C \models (g_2 = \bigsqcup S)).$$

Main Theorem: If $G_1 \sqsubseteq G_2$ holds, then so do the following.

- (a) $\langle G_1, G_2 \rangle$ is fully disjointness resolvable. In other words, for $G_1 \sqsubseteq G_2$ to hold, there must be complete information about disjointness of the collective granules of G_1 and G_2 .
- (b) In the above “recall” formula, for each $g_2 \in \text{Granules}\langle G_2 \rangle$,
$$S = \{g_1 \in \text{Granules}\langle G_1 \rangle \mid \langle g_1, g_2 \rangle \in \text{NRel}_{\langle G_1, G_2 \rangle}\}.$$
In words, each $g_1 \in \text{Granules}\langle G_2 \rangle$ is the join of those granules in $\text{Granules}\langle G_1 \rangle$ with which it is not disjoint. This is the only possibility.

Fast head-driven lookup: To identify the body of the rule with head

$g_2 \in \text{Granules}\langle G_2 \rangle$, it suffices to find all matches to g_2 in $\text{NRel}_{\langle G_1, G_2 \rangle}$.

Fast body-driven lookup: To identify the head of the rule with

$S' \subseteq \text{Granules}\langle G_1 \rangle$ as part of its body, it suffices to find the unique g_1 which matches every member of S' in $\text{NRel}_{\langle G_1, G_2 \rangle}$.

Main Representation Theorem for \otimes

Recall: $G_1 \otimes G_2$ iff

$$(\forall g_2 \in \text{Granules}\langle G_2 \rangle)(\exists S \subseteq_f \text{Granules}\langle G_1 \rangle)(\mathcal{C} \models (g_2 \sqsubseteq \sqcup S)).$$

- The result is similar to that of \trianglelefteq , with one additional condition necessary.
- With subsumption, a *resolved minimality* condition is necessary.

Motivating example: Consider `MetroArea` \trianglelefteq `Province`.

Trivial “solution”: `Gran_Puerto_Montt_urb` \sqsubseteq `Granules` \langle Province \rangle .

Resolved minimality: In the “recall” formula, for any proper subset $S' \subsetneq S$,
 $(\mathcal{C} \models \neg(g_2 \sqsubseteq \sqcup S'))$.

- In other words, if any element is removed from S , the assertion becomes false (not just fails to be true).
- Under those conditions, a theorem analogous to that for \trianglelefteq holds.

Conclusions and Current Directions

Conclusions:

Representation of common-case join rules:

- Handles both equality- and subsumption-join order.
- Applies in a constraint-based framework with incomplete information.
- Uses only (non)disjointness information about granules.

Current Directions:

Implementation in MGDB: *MGDB* is a PostgreSQL-based multigranular DBMS under development at the University of Concepción.

- Test data of administrative and political subdivisions of Chile provide many instances of equality- and subsumption-join order.
- The ideas of this paper are being applied to the efficient implementation of the associated rules.