

Data Structure Estimation for RDF Oriented Repository Building (Extended Abstract)¹⁾

Martin Řimnác

Institute of Computer Science, Czech Republic, Europe.

Email: `rimnacm@cs.cas.cz`

Abstract:

Mechanisms for accessing and training the data repository using a binary matrix formalism are presented. The repository is designed for a data storage through corresponding instances of simple attribute functional dependencies, which can be seen as similar to the binary predicate formalism being used by the RDF semantic web format.

Two mechanisms for querying a repository, the generalisation and the specialisation, are given. Furthermore, the incremental repository training mechanism with no extra requirements on the input data form is described: The extensional functional dependency system is used as a generalised view on the stored data; the algorithm is inspired by the functional dependency discovery approach.

1 Introduction

The aim of the paper is to propose a system, which enables unknown structured data to be stored and consequently to be accessed through search queries; these queries are based on the specialisation (which objects satisfy a condition given by a query) and the generalisation (which fact are valid for all objects satisfying a query). The generalised relationships generated by the features of instances in the repository will be found very important for the internal organisation; a generalised description by a functional dependency system for a data structure is chosen.

The proposed system uses the binary matrices for the internal representation and is restricted to relationships between single attributes (similar to the RDF semantic web format [1, 2]). The proposal is also inspired by results of the functional dependency discovery approach [3, 4, 5, 6, 7], which estimates a data structure by an extensional functional dependency system.

¹⁾The work was supported by the project 1ET100300419 of the Program Information Society (of the Thematic Program II of the National Research Program of the Czech Republic) "Intelligent Models, Algorithms, Methods and Tools for the Semantic Web Realization" and partly by the Institutional Research Plan AV0Z10300504 "Computer Science for the Information Society: Models, Algorithms, Applications".

Note, the returned model $M(R')$ always principally depends on its input relation R' and can be seen only as an estimation and may be changed upon the relation R' will be extended by new tuples. The extensional model can be generated using the extensional dependency definition

$$(A_i \rightarrow_{R'} A_j) \in M_{R'} \Leftrightarrow \exists \mathcal{P}' : \mathcal{D}_\alpha^{R'}(A_i) \rightarrow \mathcal{D}_\alpha^{R'}(A_j) \quad (1)$$

2 Repository Definition and Retrieval

The repository is a system enabling data storage and retrieval. Let a binary matrix notion be used for the internal representation: Denote by \mathcal{T} a set of tuples representing a part of universum. Further, let \mathcal{A} be the set of attributes corresponding to the set of tuples \mathcal{T} , $t_i(A)$ a value of the attribute $A \in \mathcal{A}$ in the tuple $t_i \in \mathcal{T}$ and \mathcal{D} a set of all values. The set \mathcal{A} , respectively \mathcal{D} is also accessible through the index $I_{\mathcal{A}}$, respectively $I_{\mathcal{D}}$ given by the mapping $\mathcal{A} \rightarrow N$, respectively $\mathcal{D} \rightarrow N$. The basic element in the repository will be then an attribute-value pair, which will be indexed by the corresponding projection $I_{\mathcal{E}} : \mathcal{A} \times \mathcal{D} \rightarrow N$.

The **active domain matrix** is a binary matrix $\Delta_{\mathcal{T}} = [\delta_{ij}]$ indicating for each element (A, v) with the index $i = I_{\mathcal{E}}(A, v)$ and each attribute $A \in \mathcal{A}$ with the index $j = I_{\mathcal{A}}(A)$ the fact $v \in \mathcal{D}_\alpha(A)$:

$$\Delta_{\mathcal{T}} = \{\delta_{ij}\} : \delta_{ij} = \begin{cases} 1 & \text{if } \exists t \in \mathcal{T} : i = I_{\mathcal{E}}(A, t(A)), j = I_{\mathcal{A}}(A) \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

The **extensional functional dependency matrix** $\Omega_{\mathcal{T}} = \{\omega_{ij}\}$ is a binary square matrix representing a projection $I_{\mathcal{A}} \times I_{\mathcal{A}} \rightarrow \{0, 1\}$:

$$\Omega_{\mathcal{T}} = \{\omega_{ij}\} : \omega_{ij} = \begin{cases} 1 & \text{if } (A_j \rightarrow A_i) \in M_{\mathcal{T}} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

Analogically, the **instances** $\Phi_{\mathcal{T}} = \{\phi_{ij}\}$ is organised as a binary square matrix, where

$$\Phi_{\mathcal{T}} = \{\phi_{ij}\} : \phi_{ij} = \begin{cases} 1 & \text{if } j = I_{\mathcal{E}}(A_1, v_1) \text{ implies } i = I_{\mathcal{E}}(A_2, v_2), \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

The instance matrix $\Phi_{\mathcal{T}}$ generates the general description of the relationships – the extensional functional dependency matrix $\Omega_{\mathcal{T}}$ – using

$$\Omega_{\mathcal{T}} = \Delta_{\mathcal{T}}^T \cdot \Phi_{\mathcal{T}} \cdot \Delta_{\mathcal{T}} \quad (5)$$

Further in the inverse way, the functional dependency matrix $\Omega_{\mathcal{T}}$ implies the mask of positions in the instance matrix $\Phi_{\mathcal{T}}^M$, which can be set; the rest of positions can not be activated due to the fact that there is no corresponding functional dependency in the $\Omega_{\mathcal{T}}$ matrix.

$$\Phi_{\mathcal{T}}^M = \Delta_{\mathcal{T}} \cdot \Omega_{\mathcal{T}} \cdot \Delta_{\mathcal{T}}^T \quad (6)$$

2.1 Repository Searching

Let \vec{x}_0 be a vector representing elements to be searched for. There are two querying mechanisms describing one step:

- The **generalisation** – the vector \vec{x}_k is extended by element(s), which are implied by any element in the vector.

$$\vec{x}_{k+1}^G = \Phi_{\mathcal{T}} \cdot \vec{x}_k \quad (7)$$

- The **specialisation** – the result gives elements, which can imply any element(s) in the vector. Note, the result may correspond to several different objects (and if required, these objects have to be separated after each step is done).

$$\vec{x}_{k+1}^S = \Phi_{\mathcal{T}}^T \cdot \vec{x}_k \quad (8)$$

On the condition each tuple $t_i \in \mathcal{T}$ contains the key attribute A_k , only instances of functional dependencies between single attributes can be considered. In such a case, complex attribute relationships can be expressed as extra-added meta-information, the generalisation is extended by (γ_L, γ_R are the left/right sides of the functional dependency)¹⁾:

$$\vec{x}'_{k+1} = \Phi \cdot \vec{x}_k + \Phi((\Phi^T \cdot ((\Delta^T \cdot \gamma_L) \odot \vec{x}_k)) == 1) \odot (\Delta^T \cdot \gamma_R) \quad (9)$$

2.2 Repository Matrix Forms

The matrix $\Phi_{\mathcal{T}}$, respective $\Omega_{\mathcal{T}}$ can be expressed in one of these forms:

- **Minimal:** This form does not contain any redundant items (optimal for an effective repository storage), but the result can be reached generally in several steps (κ^G/κ^S for the $\Phi_{\mathcal{T}}$ generalisation/specialisation, when Φ reasoning is monotonic):

$$1 \leq \kappa^G \leq |\mathcal{A}| \quad (10)$$

$$1 \leq \kappa^S \leq |I_{\mathcal{E}}| \quad (11)$$

- **Full:** The matrix may contain a lot of redundant items, but the result is reachable in one step (optimal for the reasoning).

The minimal form of the instance matrix $\Phi'_{\mathcal{T}}$ for an effective repository storage can be handled by the incremental model skeleton $\Omega'_{\mathcal{T}}$ finding algorithm [8] followed by the corresponding position mask (6) restriction.

¹⁾The operator \odot definition: $C = A \odot B$ when $\forall i, j : c_{ij} = a_{ij}b_{ij}$

3 Repository Training from Input Data Set

This section describes a way the repository including the instance matrix $\Phi_{\mathcal{T}}$ and the functional dependency matrix $\Omega_{\mathcal{T}}$ can be created from the set of tuples \mathcal{T} . The algorithm stores (in k -th step) the instance matrix Φ_k and the corrupted functional dependency matrix \mathcal{U}_k ; the functional dependency matrix Ω_k can be derived by (5) and the active domain matrix Δ is known²).

The incremental algorithm is based on the fact that the functional dependencies and their instances can be easily generated for a model covering only one tuple (in such a case, each attribute value can be derived for any other attribute value) - the tuple corresponds to the square block of 1.

The algorithm input is a set of matrices $\{\Phi_i^\Delta = \vec{x}_i^T \cdot \vec{x}_i\}$, $\vec{x}_i \in \mathcal{T}$, the output is represented by the corrupted functional dependency matrix $\mathcal{U}_{\mathcal{T}}$ and the instance matrix $\Phi_{\mathcal{T}}$. The algorithm without determining the complex attribute functional dependencies has the following structure:

- Merging the current state instance matrix Φ_k and the k -th input tuple matrix Φ_k^Δ :

$$\Phi'_{k+1} = \Phi_k + (\Delta^T \cdot (1 - \mathcal{U}_k) \cdot \Delta) \odot \Phi_k^\Delta \quad (12)$$

- Detecting corrupted functional dependencies (by inconsistent items implying more than one element per attribute from \vec{x}_k):

$$\mathcal{U}_{k+1}^\Delta = \Delta((\Phi_{k+1}^T \cdot \Delta^T) \succ 1) \quad (13)$$

- Removing the inconsistent items by (6):

$$\Phi_{k+1} = \Phi'_{k+1} \odot (\Delta^T \cdot (1 - \mathcal{U}_{k+1}^\Delta) \cdot \Delta) \quad (14)$$

- Actualising the corrupted functional dependency matrix:

$$\mathcal{U}_{k+1} = \mathcal{U}_k + \mathcal{U}_k^\Delta \quad (15)$$

- Continue with next k while all input matrices are not gathered. Finally:

$$\mathcal{U}_{\mathcal{T}} = \mathcal{U}_k, \quad \Phi_{\mathcal{T}} = \Phi_k \quad (16)$$

²In this proposal, Δ is designed as static, in the real algorithm, the matrix has to be updated before new tuple gathering.

4 Conclusion

The paper dealt with a binary matrix formalism usage for handling the repository described by relationship meta-description and instances expressing the binary predicates, i.e. to the relationships between single attributes. The main paper idea is to make the introduction at first to the basic mechanisms to operate in the repository (specialisation and generalisation mechanisms) and secondary to describe a way the input data \mathcal{I} can be stored in the repository.

The consideration of a structure for a global relationship description by the extensional functional dependency system was found useful; The relationships are driven by input data and reversibly affect the repository during next tuple insertion.

Finally, the repository can be seen as triples - a RDF semantic web document, because the same expression formalism - binary predicates - is used. In this situation, the data structure defines/estimates the semantics in an explicit way.

Generally, the proposed mechanisms can be used for storing and querying data from the repository. This incremental proposal being useful for handling data with an a-priori unknown structure is especially oriented towards to gathering data from web page content extractors.

The future work aim is to precise the complexity issue by consideration of sparse matrices and to extend and to test this algorithm for tree formatted documents (XHTML).

References

- [1] G. Antoniou, F. v Harmelen. "A Semantic Web Primer". MIT Press, 2004. ISBN: 0-262-01210-3.
- [2] E. Miller, R. Swick, D. Brickley. "Resource Description Framework". <http://www.w3.org/RDF/> [on-line]. 2004.
- [3] H. Mannila, K.J. Räihä. "Design by Example: An Applications of Armstrong Relations." *Journal of computer and system sciences* 33, pp. 129-141. Academic Press. 1986.
- [4] H. Mannila, K.J. Räihä "Dependency Inference". In *Proc. of VLDB*. pp. 155-158. ISBN: 0-934613-46-X. 1987.
- [5] H. Mannila, K.J. Räihä. "Algorithms for Inferring Functional Dependencies from Relations." In *Data & Knowledge Engineering* 12, pp 83-99. Elsevier. 1994.
- [6] J. Kivinen, H. Mannila "Approximate Inference of Functional Dependencies from Relations". In *Proc. of 4. int. conf. on Database Theory*, Berlin, Germany. pp. 129-149. ISSN: 0304-3975. 1995.
- [7] P.A.Flach, I.Savnik. "Database Dependency Discovery: A Machine Learning Approach". In *AI Communications*, Volume 12/3, pp. 139-160. 1999.
- [8] M. Rimnac "Transforming Current Web Sources for Semantic Web Usage" In *Proc. of SOFSEM 2006. Volume 2. pp 155-165 (ISBN: 80-903298-4-5). 2006.*