

Extracting Fuzzy Summaries from NoSQL Graph Databases

Arnaud Castelltort and Anne Laurent

Abstract Linguistic summaries have been studied for many years and allow to sum up large volumes of data in a very intuitive manner. They have been studied over several types of data. However, few works have been led on graph databases. Graph databases are becoming popular tools and have recently gained significant recognition with the emergence of the so-called NoSQL graph databases. These databases allow users to handle huge volumes of data (e.g., scientific data, social networks). There are several ways to consider graph summaries. In this paper, we detail the specificities of NoSQL graph databases and we discuss how to summarize them by introducing several types of linguistic summaries, namely structure summaries, data structure summaries and fuzzy summaries. We present extraction methods that have been tested over synthetic and real database experimentations.

Keywords Linguistic summaries · Graph databases · NoSQL · Fuzzy graph mining

1 Introduction

Representing data with graphs is now a proven practice. Graphs are used in many applications ranging from linguistics to chemistry and social networks. For instance, graphs allow one to intuitively illustrate the relationships between people, as well as those between people and the organizations they belong to.

Graphs are recognized to play an important role within the pattern recognition field [8]; indeed, they are a key technology for retrieving relevant information, such as in fraud detection [18] or social/biological interactions. Relevant information can be retrieved either via data mining methods or predefined queries [1].

A. Castelltort · A. Laurent(✉)
LIRMM, University of Montpellier, Montpellier, France
e-mail: {arnaud.castelltort,anne.laurent}@lirmm.fr
<http://www.lirmm.fr>

Even if graphs are ubiquitous, their representation and use have taken many forms in the literature and in real applications [3, 17]. Theoretical works have attempted to formalize the representations and treatments of graphs. More recently, they have been used in the semantic web framework, especially with ontologies. However, with the emergence of big data, the increasing volume and complexity of data and treatments, researchers have realized that robust database management systems are required. Additionally, classical relational database engines are not the solution, as shown in the performance comparisons done on this topic [7]. NoSQL is reputed for its suitability to handle big data [12]. NoSQL graph databases have thus been proposed as the best alternative for managing huge volumes of graph data and complex queries. There are several NoSQL engines; among them, Neo4j is one of the most popular.

In the NoSQL graph database model, the objects considered are nodes and relationships. Complex information on both nodes and relationships are managed as properties with $(key, value)$ pairs. In addition to properties, nodes and relationships may be labeled with types. For instance, Fig. 1 shows the relationships between people and the places they live in. The relationship types depict whether people *own* or *rent* their housing. The type of housing can be an apartment or a house: we define this value as the node type. Node information (e.g., age of people) and relationships (e.g., monthly rental fees) can be provided as node and relationship properties (e.g., $key = age$).

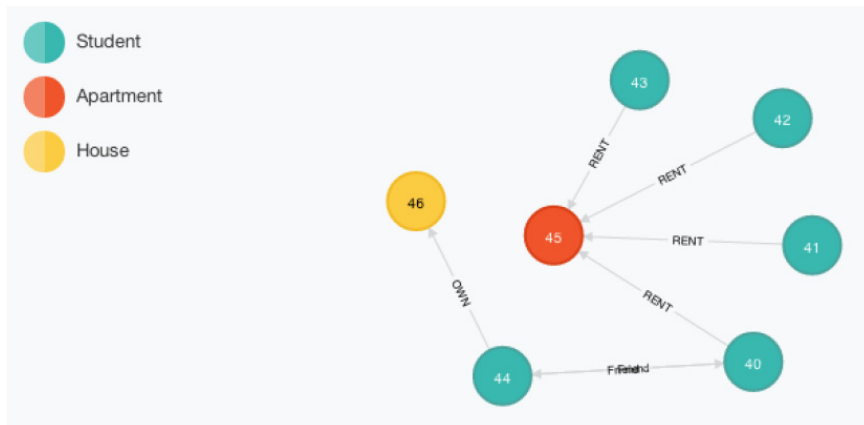


Fig. 1 Example of a Graph Database

As the amount and the complexity of information increase, the need for summaries increases as well. The literature is rich with propositions on linguistic summaries of relational databases [5, 10, 14]. Linguistic summaries are based on protoforms, the first one being $Qy \text{ are } P$ where Q stands for a fuzzy quantifier, y are the objects to be summarized and P is a possible value, such as in *Most students are young*. Linguistic summaries have been extended in many works, for instance to handle time series [2, 13]. However, these works cannot be easily applied to NoSQL

graph databases. The summaries we aim to discover must indeed be transposable to linguistic summaries. Moreover such databases combine several criteria that have never been considered altogether: node and relationship types, complex information contained in the *(key, value)* properties.

In this paper, we thus propose several types of linguistic summaries of NoSQL graph databases. The paper is organized as follows. Section 2 reports existing work on NoSQL graph databases and linguistic summarization. Section 3 introduces our definitions for linguistic summaries of NoSQL graph databases. In Section 4, we introduce the queries used to extract the summaries, by highlighting the power of NoSQL graph databases using the Cypher language. Section 5 concludes the paper and discusses perspectives for future work.

2 Background

Our work is strongly related to NoSQL graph databases and data summarization. We therefore review the basics of these two topics below.

2.1 Graphs

General Concepts. Graphs have been studied for a long time by mathematicians and computer scientists. A graph can be directed or not, labeled or not. It is defined as follows.

Def. 1 (Graph). *A graph G is given by a pair (V, E) where V stands for a set of vertices and E stands for a set of edges with $E \subseteq \{V \times V\}$.*

Def. 2 (Directed Graph). *A directed graph G is given by a pair (V, E) where V stands for a set of vertices and E stands for a set of edges with $E \subseteq (V \times V)$. That is E is a subset of all ordered pairs of V .*

When used in real world applications, graphs need to be provided with the capacity to label nodes and relations, thus leading to the so-called labeled graphs, or property graphs.

Def. 3 (Labeled Directed Graph). *A labeled oriented graph G , also known as oriented property graph, is given by a quadruplet (V, E, α, β) where V stands for a set of vertices and E stands for a set of edges with $E \subseteq (V \times V)$, α stands for the set of attributes defined over the nodes, and β the set of attributes defined over the relations.*

NoSQL graph databases [17] are based on these concepts, attributes and values over the attributes being stored thanks to the *(key, value)* paradigm which is very common in NoSQL databases. Fig. 3 shows a graph and its structure in *(key, value)* pairs.

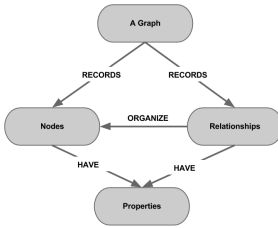


Fig. 2 Labeled Graph

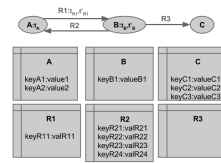


Fig. 3 Node and Relation Properties

Studies have shown that these technologies present good performances, much better than classical relational databases for representing and querying such large graph databases.

Graph Summarization. Data Summarization has been extensively studied in the last decades to produce linguistic sentences, such as *Most of the students are young* [20]. These approaches are based on the so-called *protoforms* (e.g., *Qy are P*) where *Q* is a fuzzy quantifier, *y* are the objects to summarize and *P* is a (fuzzy) predicate. They focus on relational data where source data are represented in the form of tuples defined over a schema. For instance, the tuples (John, 23, 45000), (Mary, 32, 60000) and (Bill, 38, 55000) are three tuples defined on the schema (Name, Age, Salary). The fuzzy quantifiers are defined over the [0, 1] universe of proportions. We may for instance consider two quantifiers *Few* and *Most* which membership functions are displayed by Fig. 4.

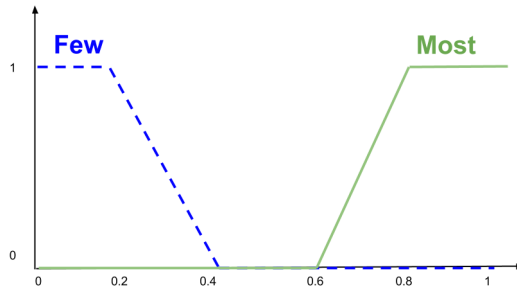


Fig. 4 Example of Fuzzy Quantifier Membership Functions

The quality of linguistic summaries can be assessed by many measures, the seminal one being *T*, the *degree of truth* that can be simply computed with a σ -count:

$$T(Qy's are P) = \mu_Q \left(\frac{1}{n} \sum_{i=1}^n \mu_P(y_i) \right)$$

where n is the number of objects (y_i) that are summarized, and μ_P , and μ_Q are the membership functions of the summarizer and quantifier, respectively.

There are various ways to examine summaries: Researchers have focused on the design of protoforms, quality measures, efficient algorithms, etc [5].

We will not recall here all the literature on fuzzy linguistic summaries which has been amply presented in previous works. We will focus instead on the subject of graph data. In this framework, two main characteristics have to be highlighted. First, graph databases are not provided with a strict and given schema such as relational data. In fact, they are close to semi-structured data. Second, graph databases focus on relationships.

Summarizing graph data has been considered for many years, aiming for instance at compressing such data with the use of supernodes as shown by Fig. 5 from [16].

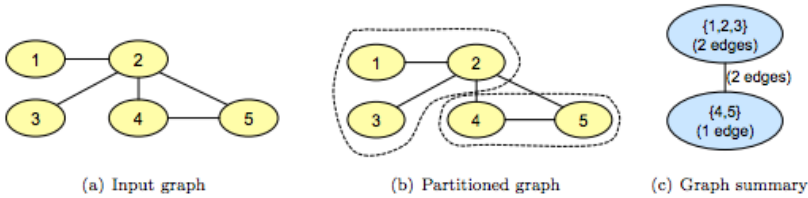


Fig. 5 Graph Summarization using Supernodes [16]

To some extent, graph summarization can be assimilated to graph mining. Graph (and tree) mining is seen as the problem of extracting frequent patterns (subgraphs/subtrees) from a large graph. It is often presented as an extension of the so-called itemset mining methods. Such methods have been successfully applied to large graphs by considering efficient approaches [9, 15, 21].

Several works in the literature have focused on schema extraction in the context of semi-structured graph data, i.e., XML data. The schema extraction problem consists in identifying a schema S from a given set of XML data documents D , such that S captures the structural information of the documents in D in the most minimal way. The schema extraction process is also referred to as schema inference [4]. The underlying structure of a given collection of XML documents can be described using Document Type Definitions (DTD), XML Schema, or via a more general representation such as tree or a graph. The structure extraction techniques in the literature aim to infer three kinds of representations: tree or graph summaries, DTD or XML Schema.

2.2 NoSQL Graph Databases

NoSQL graph databases [17] are based on graph concepts with the following additional points:

- Nodes and relationships are labeled with types;
- Properties are defined over the nodes and relationships stored according to the (key, value) paradigm, which is very common in NoSQL databases.

It should be noted that types are distinguished from properties, as in NoSQL engines such as Neo4j. These types appear in Fig. 1 as colors for nodes (e.g., Student, House) and as labels for relationships (e.g., Owns).

Generalizing the definition of labeled directed graphs, we propose a formal definition of NoSQL graph databases hereafter:

Def. 4 (NoSQL Graph Database). *A NoSQL graph database G is given by a tuple $(V, E, \theta, \tau, \alpha, \beta)$ where*

- α stands for the set of node properties defined by (key:value) pairs;
- β stands for the set of edge properties defined by (key:value) pairs;
- θ stands for the set of node types;
- τ stands for the set of edge types;
- V stands for a set of vertices with $\forall v \in V, v = (id_v, t_v, \kappa_v)$ s.t. $t_v \subseteq \theta$ stands for the types of v , $\kappa_v \subseteq \alpha$ stands for the properties of v and id_v is the vertice identifier;
- E stands for a set of edges with $\forall e \in E, e = (id_e, (v_e^1, v_e^2), t_e, \lambda_e)$ s.t. $(v_e^1, v_e^2) \in \{V \times V\}$, $t_e \subseteq \tau$ stands for the types of e , $\lambda_e \subseteq \beta$ stands for the properties of e and id_e is the edge identifier.

The set of properties of a node v is denoted by α_v , the set of types is denoted by θ_v . The set of properties of a relation e is denoted by β_e , the set of types is denoted by τ_e .

Fig. 3 shows a graph and its structure in (key, value) pairs. In this figure, we have:

- $\alpha = \{(keyA1, valueA1), \dots\}$
- $\beta = \{(keyR11, valueR11), \dots\}$
- $V = \{(A, \{t_A\}, \{(keyA1, valueA1), (keyA2, valueA2)\}), (B, \{t_B, t'_B\}, \{(keyB1, valueB1)\}), (C, \emptyset, \{(keyC1, valueC1), (keyC2, valueC2), (keyC3, valueC3)\})\}$
- $E = \{(R1, (A, B), \{t_{R1}, t'_{R1}\}, \{(keyR11, valueR11)\}), \dots\}$

There exist several NoSQL graph database engines (OrientDB, Neo4J, Hyper-GraphDB, etc.) [3]. Neo4J is considered the best performer [19]. All NoSQL graph databases require developers and users to use graph concepts to query data. Queries are called traversals, referring to the action of visiting elements, i.e. nodes and relationships. There are three main ways to traverse a graph:

- Programmatically: using an API;
- By functional traversal: using a traversal based on a sequence of functions applied to a graph;
- By declarative traversal: explicitly expressing what we want to do and not how we want to do it. The database engine then defines the best way to achieve this goal.

In this paper, we focus on declarative queries over a NoSQL graph database. The Neo4j language is called Cypher.

For instance, in Fig. 6, a query to return the customers who have visited the “Ritz” hotel is displayed. Those customers are both displayed in the list and circled in red in the graph.



Fig. 6 Displaying the Result of a Cypher Query

Cypher clauses are similar to SQL ones. It is based on a “ASCII art” way of writing graph elements. For example, directed relations are written using the `-[]->` symbol. Types and labels are written after a semi-column (:).

More specifically, queries in Cypher have the following syntax¹:

Listing 1.1 Query example on a Graph

```

1 [START]
2 [MATCH]
3 [OPTIONAL MATCH WHERE]
4 [WITH [ORDER BY] [SKIP] [LIMIT]]
5 RETURN [ORDER BY] [SKIP] [LIMIT]
```

As shown above, Cypher is comprised of several distinct clauses which are listed below in Listing 1.2.

Listing 1.2 Cypher Clauses

- 1 **START**: Starting points in the graph, obtained via index lookups or by element IDs.
- 2 **MATCH**: The graph pattern to match, bound to the starting points in **START**.
- 3 **WHERE**: Filtering criteria.
- 4 **RETURN**: What to return.
- 5 **CREATE**: Creates nodes and relationships.
- 6 **DELETE**: Removes nodes, relationships and properties.
- 7 **SET**: Set values to properties.
- 8 **FOREACH**: Performs updating actions once per element in a list.
- 9 **WITH**: Divides a query into multiple, distinct parts.

These operations can even be extended to fuzzy queries [6]. In this work, they are used for computing the linguistic summaries introduced below.

¹ <http://docs.neo4j.org/refcard/2.0/>
<http://docs.neo4j.org/chunked/milestone/cypher-query-lang.html>

3 Building Fuzzy NoSQL Graph Summaries

Below, we define the protoforms for summing up NoSQL graph databases.

3.1 Structure Summaries

Structure summaries are meant to retrieve the structure of the graph embedded in element types, which could somehow be associated with relational database schema. Such summaries could thus be associated with with schema mining in the literature.

Def. 5 (Structure Summary). Let $G = (V, E, \theta, \tau, \alpha, \beta)$ be a NoSQL graph database. A structure summary S of G is defined as $S = (a \text{ } \text{--}[r]\text{ } \text{-->} b, Q)$ where $a, b \in \theta$ (node types), $r \in \tau$ (relation type) and Q is a fuzzy quantifier.

The structure summary can be expressed in a linguistic form as follows: In G , Q of the a r b

Example 1. In the toy example, $(\text{Student } \text{--}[rent]\text{ } \text{-->} \text{apartment}, \text{Most})$, expressed as “Most of the students rent an apartment”, is a structure summary.

3.2 DataStructure Summaries

Data structure summaries are meant to refine structure summaries. They allow one to differentiate cases when schema depend on the value of properties. For instance, depending on their salary, employees may rent apartments instead of houses.

Def. 6 (Data Structure Summary). Let $G = (V, E, \theta, \tau, \alpha, \beta)$ be a NoSQL graph database. A Data Structure Summary S is defined as $S = (a.X \text{ } \text{--}[r.Z]\text{ } \text{-->} b.Y, Q)$ with $a, b \in \theta$ (node types), $r \in \tau$ (relation type), $X, Y \subseteq \alpha$ (node properties), $Z \subseteq \beta$ (relation properties) and Q a fuzzy quantifier.

Example 2. In the toy example, $(\text{Student}(\text{Age} : 28) \text{ } \text{--}[rent(\text{fees} : 1200)]\text{ } \text{-->} \text{apartment}, \text{Few})$ is a data structure summary.

Such summaries are extended in order to allow fuzzy linguistic labels in the refinement. Indeed, it would be both difficult and useless to define summaries on single values such as “the age is 28, as in fuzzy data mining for fuzzy association rule mining. Using fuzzy linguistic labels makes it possible to retrieve fuzzy linguistic summaries where young students and low rental fees are considered.

Def. 7 (Fuzzy Data Structure Summary). Let $G = (V, E, \theta, \tau, \alpha, \beta)$ be a NoSQL graph database. Let F_α and F_β be sets of fuzzy properties. A Fuzzy Data Structure Summary S is defined as $S = (a.X \text{ } \text{--}[r.Z]\text{ } \text{-->} b.Y, Q)$ with $a, b \in \theta$ (node types), $r \in \tau$ (relation type), $X, Y \subseteq \alpha \cup F_\alpha$ (node properties and node fuzzy properties), $Z \subseteq \beta \cup F_\beta$ (relation properties and relation fuzzy properties) and Q a fuzzy quantifier.

Example 3. In the toy example, $(\text{Student}(\text{Age} : \text{young}) \text{ } \text{--}[rent(\text{fees} : \text{low})]\text{ } \text{-->} \text{apartment}, \text{Most})$ is a fuzzy structure summary.

For all these summaries, it is important to discuss the way to assess them. For this purpose, we propose to rely on the extension of the degree of truth.

3.3 Degree of Truth

In our framework, the degree of truth determines the extent to which the relationship appearing in the summary is truthful regarding the fuzzy quantifier. For instance, if the summary mentions that *most of the students rent an apartment*, then the degree of truth describes to which extent a high proportion of students rent an apartment.

There are two ways of calculating this degree. Indeed, it may consist in computing the proportion of students who rent an apartment over the whole student population, or the proportion of students who rent an apartment over the number of students who rent their housing. These two definitions are provided below.

Def. 8 (Degree of Truth of NoSQL Graph Summaries). *Given a graph database G and a summary $S = a -[r]->b$, Q , the degrees of Truth of S in G are defined as:*

$$Truth_1(S) = \mu_Q \left(\frac{\text{count}(\text{distinct}(S))}{\text{count}(\text{distinct}(a))} \right)$$

$$Truth_2(S) = \mu_Q \left(\frac{\text{count}(\text{distinct}(S))}{\text{count}(\text{distinct}(a -[r]-> (?)))} \right)$$

The second type of degree of truth is also called the *diversity of target source* denoted by D_T later on in this paper.

Example 4. *In the toy example from Fig. 1, the degree of truth of the summary " $S = \text{Student} -[\text{rent}]->\text{apartment}$, $Most$ " is given by the membership degree to the fuzzy quantifier² of the ratio between the number of times a relationship appears between a Student and an Apartment (s)he rents over the number of relations of type Rents starting from a Student node. In this example, we thus have $Truth_2(S) = \mu_{Most} \left(\frac{4}{5} \right) = \mu_{Most}(0.8)$*

The ratio appearing in the definition of the degree of truth can be compared to the *confidence* in association rule mining which acts as a conditional probability.

4 Extracting the Summaries

Below we detail how to extract the above-defined summaries from NoSQL graph databases and some first results.

² We do not mention here the detailed membership function of the *Most* quantifier which can be defined in a very classical manner as done in the literature of fuzzy quantifiers and fuzzy summaries.

4.1 Queries

Most of the treatments we propose can be performed effectively by defining queries over the NoSQL graph database, so as to obtain a more declarative than procedural way to extract summaries. This property is based on the fact that NoSQL graph databases provide powerful pattern-matching features, as intended in inductive relational databases [11].

The structure summaries can be retrieved by considering Cypher queries such as:

```
1 MATCH (a)-[r]->(m)
2 RETURN DISTINCT labels(a), type(r), labels(m), count(r)
```

In the above query, all the structures are retrieved together, along with the number of times they appear.

The query from Listing 1.3 extracts the summaries from a graph and calculates the degree of truth for every summary.

Listing 1.3 Retrieving Structure Summaries

```
1 MATCH (a)-[r]->(b)
2 WITH DISTINCT labels(a) AS labelsA, type(r) AS typeR, labels(b) AS labelsC, toFloat(count(*) ←
  ) AS countS
3 MATCH (a1)-[r2]->(m)
4 WHERE labels(a1)= labelsA AND type(r2)= typeR
5 WITH DISTINCT labelsA, typeR, labelsC, countS, labels(a1) AS labelsA1, type(r2) AS typeR2, ←
  count(*) AS count2
6 RETURN labelsA, typeR, labelsC,
7 toFloat(countS)/ count2 AS Truth,
8 MuMost(countS/count2) as TruthMost,
9 MuFew(countS/count2) as TruthFew,
10 MuVeryFew(countS/count2) as TruthVeryFew
```

Where μ_{Most} , μ_{Few} and $\mu_{VeryFew}$ respectively stand for the membership function of the fuzzy subsets *Most*, *Few* and *VeryFew* defined on the $[0, 1]$ universe which are implemented as the *MuMost*, *MuFew* and *MuVeryFew* functions in Cypher.

4.2 Experimental Results

Experiments have been run on synthetic and real databases with a Java implementation run on an Intel Core i5 2.4 Ghz. The Movie real database deals with Directors, Movies and Actors³. It contains about 12,000 movies and 50,000 actors. Below is an example of the results from the synthetic database:

Summary	$TruthMost$	$TruthFew$	$TruthVeryFew$
Student-Rents-Apartment	$\mu_{Most}(4/5)$	$\mu_{Few}(4/5)$	$\mu_{VeryFew}(4/5)$
Student-Friend-Student	$\mu_{Most}(2/5)$	$\mu_{Few}(2/5)$	$\mu_{VeryFew}(2/5)$
Student-Owens-House	$\mu_{Most}(1/5)$	$\mu_{Few}(1/5)$	$\mu_{VeryFew}(1/5)$

³ <http://neo4j.com/developer/example-data>

The Listing 1.4 displays some examples of fuzzy data structure summaries.

Listing 1.4 Fuzzy Data Structure Summaries Extracted from the Movie Database

```

1 DS12: (Person(olD) -[DIRECTED]->Movie.Most); Truth(DS12)=1.0
2 DS16: (Person(young) -[DIRECTED]->Movie.VeryFew); Truth(DS16)=0.21
3 DS17: (Person(middleAge) -[DIRECTED]->Movie.VeryFew); Truth(DS17)=0.47
4 DS18: (Person(olD) -[DIRECTED]->Movie.VeryFew); Truth(DS18)=0.0
    
```

5 Conclusion and Perspectives

In this paper, we have presented an approach to summarize NoSQL graph databases in the form of linguistic summaries. These databases are quite specific with respect to existing work, as they combine several difficulties: a focus on relationships, node and relationship types management, management of complex (key:value) pair information management, etc.

Several types of linguistic summaries are proposed. They can be easily extracted using existing declarative query languages, making it possible to deploy them on every commercial tool. Experiments have been made on Neo4j using the Cypher query language, and have demonstrated the interest of our proposal.

Our work opens several perspectives. First, we plan to run more experiments to test scalability. The propositions could also be extended to several graphs. Moreover, we might try and integrate other types of summaries for managing time and space information. Another perspective is to consider extended linguistic summaries containing several relationships.

References

1. Aggarwal, C.C., Wang, H. (eds.): *Managing and Mining Graph Data*. Advances in Database Systems, vol. 40. Springer (2010)
2. Almeida, R.J., Lesot, M., Bouchon-Meunier, B., Kaymak, U., Moysse, G.: Linguistic summaries of categorical series for septic shock patient data. In: *Proceedings of the IEEE International Conference on Fuzzy Systems, FUZZ-IEEE 2013, Hyderabad, India, July 7–10, pp. 1–8*. IEEE (2013). <http://dx.doi.org/10.1109/FUZZ-IEEE.2013.6622581>
3. Angles, R., Gutiérrez, C.: Survey of graph database models. *ACM Comput. Surv.* 40(1) (2008)
4. Bex, G.J., Neven, F., Vansummeren, S.: Inferring XML schema definitions from XML data. In: Koch, C., Gehrke, J., Garofalakis, M.N., Srivastava, D., Aberer, K., Deshpande, A., Florescu, D., Chan, C.Y., Ganti, V., Kanne, C., Klas, W., Neuhold, E.J. (eds.) *Proceedings of the 33rd International Conference on Very Large Data Bases*, University of Vienna, Austria, September 23–27, pp. 998–1009. ACM (2007)
5. Bouchon-Meunier, B., Moysse, G.: Fuzzy linguistic summaries: where are we, where can we go ? In: *CIFer 2012 IEEE Conf. on Computational Intelligence for Financial Engineering & Economics (CIFer)*, pp. 317–324. IEEE (2012)

6. Castellort, A., Laurent, A.: Fuzzy queries over NoSQL graph databases: perspectives for extending the cypher language. In: Laurent, A., Strauss, O., Bouchon-Meunier, B., Yager, R.R. (eds.) IPMU 2014, Part III. CCIS, vol. 444, pp. 384–395. Springer, Heidelberg (2014)
7. Cattell, R.: Scalable SQL and NoSQL data stores. *SIGMOD Record* **39**(4), 12–27 (2010)
8. Conte, D., Foggia, P., Sansone, C., Vento, M.: Thirty Years Of Graph Matching In Pattern Recognition. *International Journal of Pattern Recognition and Artificial Intelligence* (2004)
9. Cook, D.J., Holder, L.B.: *Mining Graph Data*. John Wiley & Sons (2006)
10. Rasmussen, D., Yager, R.R.: Finding fuzzy and gradual functional dependencies with summary SQL. *Fuzzy Sets and Systems* **106**, 131–142 (1999)
11. De Raedt, L.: A perspective on inductive databases. *SIGKDD Explor. Newsl.* **4**(2), 69–77 (2002)
12. Han, J., Haihong, E., Le, G., Du, J.: Survey on noSQL database. In: Proc. of the 6th International Conference on Pervasive Computing and Applications (ICPCA), pp. 363–366 (2011)
13. Kacprzyk, J., Wilbik, A., Zadrozny, S.: An approach to the linguistic summarization of time series using a fuzzy quantifier driven aggregation. *Int. J. Intell. Syst.* **25**(5), 411–439 (2010). doi:[10.1002/int.20405](https://doi.org/10.1002/int.20405)
14. Kacprzyk, J., Zadrozny, S.: Linguistic database summaries and their protoforms: towards natural language based knowledge discovery tools. *Inf. Sci.* **173**(4), 281–304 (2005). doi:[10.1016/j.ins.2005.03.002](https://doi.org/10.1016/j.ins.2005.03.002)
15. Kuramochi, M., Karypis, G.: Frequent subgraph discovery. In: Proceedings of the 2001 IEEE International Conference on Data Mining, ICDM 2001, pp. 313–320. IEEE Computer Society, Washington, DC (2001). <http://dl.acm.org/citation.cfm?id=645496.658027>
16. LeFevre, K., Terzi, E.: Grass: Graph structure summarization. In: Proceedings of the SIAM International Conference on Data Mining, SDM 2010, April 29 - May 1, Columbus, Ohio, USA, pp. 454–465. SIAM (2010)
17. Robinson, I., Webber, J., Eifrem, E.: *Graph Databases*. O’Reilly (2013)
18. Sadowski, G., Rathle, P.: Fraud detection: Discovering connections with graph databases. In: White Paper - Neo Technology - Graphs are Everywhere (2014)
19. ThoughtWorks: Technology advisory board, May 2013. <http://thoughtworks.fileburst.com/assets/technology-radar-may-2013.pdf>
20. Yager, R.R.: A new approach to the summarization of data. *Information Sciences* **28**(1), 69–86 (1982)
21. Yan, X., Han, J.: gSpan: Graph-based substructure pattern mining. In: Proceedings of the 2002 IEEE International Conference on Data Mining, ICDM 2002, pp. 721–724. IEEE Computer Society, Washington, DC (2002). <http://dl.acm.org/citation.cfm?id=844380.844811>