

# An Extension of Datalog for Graph Queries<sup>\*</sup>

Mirjana Mazuran<sup>1</sup>, Edoardo Serra<sup>2</sup>, and Carlo Zaniolo<sup>3</sup>

<sup>1</sup> Politecnico di Milano DEI – [mazuran@elet.polimi.it](mailto:mazuran@elet.polimi.it)

<sup>2</sup> University of Calabria DEIS – [eserra@deis.unical.it](mailto:eserra@deis.unical.it)

<sup>3</sup> University of California, Los Angeles – [zaniolo@cs.ucla.edu](mailto:zaniolo@cs.ucla.edu)

**Abstract.** Supporting aggregates in recursive logic rules is a crucial long-standing problem for Datalog. To solve this problem, we propose Datalog<sup>FS</sup> that supports queries and reasoning on the number of distinct occurrences satisfying given goals, or conjunction of goals, in rules. By using a generalized notion of multiplicity called frequency, we show that graph queries can be easily expressed in Datalog<sup>FS</sup>. This simple extension preserves all the desirable semantic and computational properties of logic-based languages, while significantly extending their application range to support efficiently page-rank, and social-network queries.

## 1 Introduction

Due to the emergence of many important application areas we are now experiencing a major resurgence of interest in Datalog for parallel and distributed programming [1] and for expressing and supporting subsets of Description Logic for ontological queries [2]. Other lines of work are exploring execution of recursive queries in the MapReduce framework [3] and in Data Stream Management Systems [4]. The abundance of new applications underscores the need to tackle and solve crucial Datalog problems that remain unsolved and restrict its effectiveness by e.g., disallowing the use of aggregates in recursion. This problem is very challenging since basic aggregates violate the requirement of monotonicity on which the least fixpoint semantics of Datalog is based.

*Related Work* The notion of *stratification* with respect to negation and aggregates is simple for users to master [5, 6]. Unfortunately, stratification (into a finite number of strata) is too restrictive and cannot support the efficient formulation of many graph optimization algorithms, which typically require the use of extrema and counting in recursion [7].

The importance of optimization and graph applications have motivated much research work seeking to solve these problems. These proposals follow three main approaches: *i*) supporting infinite levels of stratifications using Datalog<sub>1S</sub> programs [6]; *ii*) attempting to preserve the fixpoint computation via continuous aggregates and non-deterministic *choice* constructs [8, 9], and *iii*) seeking to achieve monotonicity by using partial orders that are more general than set-containment [10]. These past solutions had limited generality and often required

---

<sup>\*</sup> Extended Abstract

sophisticated users and compilers. We next introduce  $\text{Datalog}^{FS}$  which does not suffer from these problems.

In the next section we introduce  $\text{Datalog}^{FS}$  via simple examples. In Section 3 we introduce constructs that support facts and predicates having multiple occurrences and, in Section 4, we review important graph applications. In Section 5, we show how to implement  $\text{Datalog}^{FS}$  efficiently.

## 2 $\text{Datalog}^{FS}$ by Example

Consider, for instance a database of facts as follows:

`person(adam). person(marc). person(jerry). person(tom).`  
`son(marc, tom). son(marc, jerry). son(tom, eddy). son(tom, adam). son(tom, john).`

The following rule defines fathers with at least two sons:

$$\text{twosons}(X) \leftarrow \text{person}(X), \text{son}(X, Y1), \text{son}(X, Y2), Y2 \neq Y1.$$

$\text{Datalog}^{FS}$  allows the following equivalent expression for our `twosons` rule:

$$\text{twosons}(X) \leftarrow \text{person}(X), 2:[\text{son}(X, Y)].$$

The goal,  $I: [\mathbf{b-expression}]$ , is called a frequency support goal (or *FS-Goal*), and “ $I$ ” is a positive integer, called *Running-FS clause*. The expression in the bracket is called *b-expression*, and can either consists of a single positive predicate or a conjunction of positive predicates [11]. The convenience of FS-goals is clear if we want to find people with many sons:

$$\text{sixsons}(X) \leftarrow \text{person}(X), 6:[\text{son}(X, Y)].$$

will retrieve all persons who have at least six sons. An equivalent rule can be expressed using the  $\neq$  operator. Indeed we can start as follows:

$$\text{sixsons}(X) \leftarrow \text{person}(X), \text{son}(X, Y1), 5:[\text{son}(X, Y2), Y2 \neq Y1].$$

and proceed inductively, and obtain a rule containing six goals `son(X, Yj)`, where  $j = 1, \dots, 6$  and  $6 \times 5$  goals saying, that every  $Y$  be different from every other  $Y$ . If we are interested in links between web pages, which could easily be thousands, it becomes clear that the approach based on  $\neq$  is totally impractical, and without FS-goals we would need a `COUNT` aggregate. Yet, aggregates bring in the curse of non-monotonicity and recursion becomes a problem. At the semantic level, our  $\text{Datalog}^{FS}$  rules can instead be viewed as standard Horn clauses whereby the standard monotonicity-based semantics of negation-free  $\text{Datalog}$  is preserved.

We now clarify the scope of variables in  $\text{Datalog}^{FS}$ . Predicate `friend(X, Y)` denotes that person  $X$  views person  $Y$  as a friend (no assumption of symmetry):

*Example 1.* Pairs of friends (F1, F2) where F1 and F2 have at least 3 friends:

$$\text{popularpair}(X, Y) \leftarrow \text{friend}(X, Y), 3:[\text{friend}(X, V1)], 3:[\text{friend}(Y, V2)].$$

*Example 2.* Pairs of friends (F1, F2) who have at least three friends in common

$$\text{sharethree}(X, Y) \leftarrow \text{friend}(X, Y), 3:[\text{friend}(X, V), \text{friend}(Y, V)].$$

There are two kinds of variables in rules with FS-goals. The first are those, such as  $X$  and  $Y$  in Example 2, that appear in the head of the rule or in some goal outside the b-expression. These will be called *global* variables. They are basically the universally qualified variables of the standard Horn Clauses, and have the whole rule as scope. Variables  $X$  and  $Y$  in Example 1 are global for that rule.

Other variables only appear in b-expressions and their scope is *local* to the b-expression, where they appear (e.g.  $V1$  and  $V2$  in Example 1, and  $V$  in Example 2). Thus,  $K[...]$  can be viewed as an existential declaration of local variables under the following constraint: there exist at least  $K$  assignments of the local variables that satisfy the b-expression. Example 2 states that there exist at least 3 distinct  $V$  occurrences each denoting a person who is a friend to both  $X$  and  $Y$ . The scope of existential variables is local to the b-expression: in Example 1 replacing  $V1$  and  $V2$  with  $V$  would not change the meaning of our rule. Let us now express that an assistant professor to be advanced to associate professor should have an H-index of at least 13:

*Example 3.* Our candidate must have authored at least 13 papers each of which has been referenced at least 13 times. The database table `author(Author, Pno)` lists all papers (co-)authored by a person, while the atom `refer(PnFrom, PnTo)` denotes that paper `PnFrom` contains a reference to paper `PnTo`.

$$\begin{aligned} \text{atleast13(PnTo)} &\leftarrow 13: [\text{refer(PnFrom, PnTo)}]. \\ \text{hindex13(Author)} &\leftarrow 13: [\text{author(Author, Pno), atleast13(Pno)}]. \end{aligned}$$

These simple examples could also be expressed using the count aggregate. Yet, count and other aggregates are non-monotonic with respect to the partial ordering defined by set containment, and cannot be used in recursive rules. Indeed, the meaning and efficient implementation of Datalog programs with recursive rules are based on their least fixpoint semantics<sup>4</sup>, which is only guaranteed to exist when the program rules define monotononic mappings. Now, continuous count is obviously monotonic with respect to set-containment. This is formally proven in [11] by rewriting the running FS-construct with equivalent, although inefficient, Horn clauses (that use list and thus we will avoid in the actual implementation, as shown in Section 5).

**Final-FS construct and Stratification.** The semantics of Datalog<sup>FS</sup> [11] allows to use variables rather than constants in the specification of FS goals. This is useful, for instance, to find the actual number of sons a person has:

*Example 4.* How many sons does a person have?

$$\begin{aligned} \text{csons(PName, N)} &\leftarrow \text{person(PName), N: [son(PName, Sname)], } \neg \text{morethan(PName, N)}. \\ \text{morethan(PName, N)} &\leftarrow N1: [\text{son(PName, _)}], N1 > N. \end{aligned}$$

Thus `csons` must belong to a stratum that is strictly higher than `son`, whereas with respect to `person` it could be in the same stratum or in the one above it. The need to find the maximum value satisfying a running-FS clause is so common that we provide a construct called Final-FS and denoted by the operator `=!`.

<sup>4</sup> Naturally, by “least fixpoint” of a program, we mean “least fixpoint of its immediate consequence operator” [12].

*Example 5.* How many sons does a person have?

$$\text{csons}(\text{Name}, N) \leftarrow \text{person}(\text{Name}), N = ![\text{son}(\text{Name}, \_)].$$

The formal semantics of the *Final-FS* construct is defined as the rewriting of Example 5 into Example 4, which makes use of negation, whereby we will require that our Datalog<sup>FS</sup> programs be stratified w.r.t. Final-FS goals.

**Recursive Datalog<sup>FS</sup>.** Consider the following example:

*Example 6.* Some people will come to the party for sure. Others will also come once they learn that three or more of their friends will come.

$$\begin{aligned} \text{willcome}(X) &\leftarrow \text{sure}(X). \\ \text{willcome}(Y) &\leftarrow 3: [\text{friend}(Y, X), \text{willcome}(X)]. \end{aligned}$$

One person might be more timid than another, and different people could require a different number of friends before they also join the party. Thus, if `requires(Person, PNumber)` denotes the number of friends required by a person, where `PNumber` must be a positive integer (whereas `sure` denotes people who will come even if none of their friends will), we have the following program:

*Example 7.* A person will join the party if a sufficient number of friends join.

$$\begin{aligned} \text{join}(X) &\leftarrow \text{sure}(X). \\ \text{join}(Y) &\leftarrow \text{requires}(Y, K), K: [\text{friend}(Y, X), \text{join}(X)]. \end{aligned}$$

### 3 Multi-Occuring Predicates

As discussed in [10], there are numerous examples where it is desirable that certain predicates are counted as providing a support level greater than one. For instance, we might use the following representation to denote that the paper with DBLP identifier “MousaviZ11” is currently cited in six papers: `ref(“MousaviZ11”):6`. Thus, `Pno = “MousaviZ11”` contributes with count six to the b-expression of the rule:

*Example 8.* Total reference count for an author.

$$\text{tref}(\text{Authr}):N \leftarrow N: [\text{author}(\text{Authr}, \text{Pno}), \text{ref}(\text{Pno})].$$

The clauses “:6” and “:N” used in the above fact and rule head will be called *FS-Assert* clauses. The semantics of programs  $P$  with frequency assert clauses is defined by *expanding* it into its  $\bar{P}$  equivalent, which is obtained as follows: Each rule in  $P$  with head  $q(X_1, \dots, X_n):K \leftarrow \text{Body}$  is replaced by

$$\bar{q}(X_1, \dots, X_n, J) \leftarrow \text{lessthan}(J, K), \text{Body}.$$

where `lessthan(J, K)` is a recursive predicate that generates all positive integers up to  $K$ , included. Thus, we have that, as a result of this expansion, `ref(“MousaviZ11”):6` contributes with six to the reference count of each author of that paper. A property of frequency statements is that, when multiple statements hold for the same fact *only the largest value* is significant, the others are subsumed and can be ignored.

Bill-of-materials (BOM) applications represent a well-known example of the need for recursive queries. Our database might contain records `assbl(Part, Subpart, Qty)` which, for each part number, gives the immediate subparts used in its assembly and their quantity (e.g. a bicycle has 1 frame and 2 wheels as immediate subparts). At the bottom of BOM DAG, we find the basic parts that are purchased from external suppliers and described by `basic(Pno, Days)` denoting the days needed to obtain that basic part. Several interesting BOM applications are naturally expressed by combining aggregates and recursion, as follows:

*Example 9.* How many basic parts does an assembled part contain?

```

cassb(Part, Sub):Qty ← assbl(Part, Sub, Qty).
cbasic(Pno):1 ← basic(Pno, _).
cbasic(Part):K ← K:[cassb(Part, Sub), cbasic(Sub)].
cntbasic(Prt, C) ← C =![cbasic(Prt)].

```

The count of basic parts is not retrieved using the goal `N:[cbasic(frame)]` since this will return all positive integers up to the max value. Goal `N =![cbasic(frame)]` is used instead as this returns the exact count of the basic subparts.

*Example 10.* How many days until delivery?

```

delivery(Pno):Days ← basic(Pno, Days).
delivery(Part):Days ← assb(Part, Sub, _), Days:[delivery(Sub)].
actualDays(Part, CDays) ← CDays =![delivery(Part)].

```

For each assembled part, we find each basic subpart along with the number of days this takes to arrive. Observe that the argument `Pno` is projected out, and only the number of days associated with it is retained, whereby the maximum of the number of days required by any basic part is derived.

## 4 Advanced Graph Applications

We now show some examples that use Datalog<sup>FS</sup> with positive rational numbers. As explained in more details in [11], we can assume that we use rational numbers with the same large denominator, and thus easily derive equivalent Datalog<sup>FS</sup> rules for their numerators.

**Diffusion Models.** The Jackson-Yariv Diffusion Model (JYDM) [13] provides a powerful abstraction on how social structures influence the spread of a behavior and trends in Social Networks. We use the JDYM to understand how a tweet will spread in the Twitter network. Let `followd(X, Y)` indicate that user `X` is followed by user `Y` and `coeff(X, C)` means that `C` is the coefficient of how susceptible to change node `X` is. Predicate `b(X)` denotes, if true, that node `X` will retweet. We assume that an agent, `source(X)`, first posts the tweet and starts its diffusion.

*Example 11.* Modeling a retweet in Datalog<sup>FS</sup>.

```

b(X) ← source(X).
b(X) ← coeff(X, C), K ≥ 1/C, K:[followd(Y, X), b(Y)].

```

The last rule finds each node  $X$  for which the condition  $K \times C \geq 1$  holds by testing the equivalent condition  $K \geq 1/C$ . When this is satisfied  $b(X)$  is set to true. The answer to the query  $?b(Y)$  will thus list all the users that propagated the tweet that originated from the node specified by `source`.

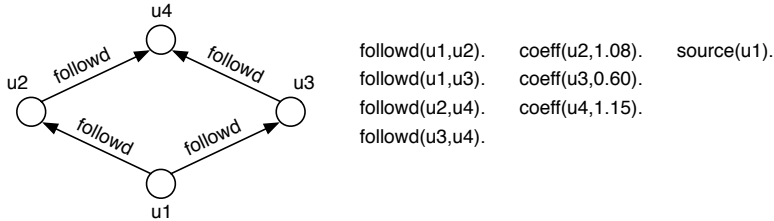


Fig. 1. Retweet modeling in Datalog<sup>FS</sup>.

Then, applying the program in Example 11 to the Twitter network in Figure 1 the following atoms are derived:  $b(u_1)$ ,  $b(u_2)$ ,  $b(u_4)$ .

**Markov Chains and Page Rank.** A Markov chain is represented by the transition matrix  $W$  of  $s \times s$  components where  $w_{ij}$  is the probability to go from state  $i$  to state  $j$  in one step. A Markov chain is *irreducible* if for each pair of states  $i, j$ , the probabilities to go from  $i$  to  $j$  and from  $j$  to  $i$  in one or more steps is greater than zero.

Computing stabilized probabilities of a Markov chain has many real-world applications, such as estimating the distribution of population in a region, and determining the Page Rank of web nodes. Let  $P$  be a vector of stabilized probabilities of cardinality  $s$ , the equilibrium condition in terms of matrices is:  $P = W \cdot P$ .

Computing this fixpoint is far from trivial and irreducible chains can be modeled quite naturally in Datalog<sup>FS</sup>. If  $p\_state(X) : K$  denotes that  $K$  is the probability of node  $X$ ,  $1 \leq X \leq s$ , and  $w\_matrix(Y, X) : W$  denotes that the arc from  $Y$  to  $X$  has weight  $W$ , then we compute the fixpoint as follows:

```
p_state(X):K ← K:[p_state(Y), w_matrix(Y, X)].
w_matrix(1, 1):w11.
w_matrix(1, 2):w12.
⋮
w_matrix(s, s):wss.
```

It is important to notice that each fixpoint of such program is an equilibrium  $P = W \cdot P$  of the Markov Chain represented by matrix  $W$ . In order to find a non trivial fixpoint ( $\neq 0$ ) for program  $P$ , we add baseline facts i.e. a set of facts  $p\_state(1) : 0.1$ .  $p\_state(2) : 0.1$ .  $\dots$   $p\_state(s) : 0.1$ ., that guarantee that the least fixpoint contains facts with predicate  $p\_state$ . Such program is called  $Pbl$  and is a Datalog<sup>FS</sup> program for which we can compute the least fixpoint efficiently. Moreover, every fixpoint of  $Pbl$  is also a fixpoint for  $P$ . Indeed, for any interpretation  $I$  that contains all the baseline facts, the application of either operators produce the same result: i.e.,  $T_P(I) = T_{Pbl}(I)$ . Therefore any fixpoint of  $T_P$  that contains all the baseline facts is also a fixpoint for  $T_{Pbl}$  and vice-versa.

But since, by its very definition, the least model of  $Pbl$  contains all the baseline facts, we have that every fixpoint for  $T_{Pbl}$  is also a fixpoint for  $T_P$ . The opposite of course is not true since the null fixpoint of  $T_P$ , and possibly others, are not fixpoint for  $T_{Pbl}$ . However, if  $T_P$  has a fixpoint that is positive at all nodes, then by multiplying the frequency at all nodes by a large enough finite constant, we obtain a fixpoint for  $T_P$  that contains all the baseline facts of  $T_{Pbl}$ . Since, for each irreducible Markov chain there exists a not trivial fixpoint, also  $T_P$  has one that is not null at every node, then there exists a finite fixpoint for  $T_{Pbl}$ . Therefore, the least fixpoint for  $T_{Pbl}$  is finite. That is:

**Theorem 1.**

- *The least fixpoint of the baseline Datalog<sup>FS</sup> program that models an irreducible Markov chain is finite.*
- *Every non-null solution of an irreducible Markov chain can be obtained by scaling the least fixpoint solution of its baseline Datalog<sup>FS</sup> model.*

In summary, while there has been a significant amount of previous work on Markov chains, the use of Datalog<sup>FS</sup> has provides us with a model and a simple computation algorithm which is valid for all irreducible Markov chains, including periodic ones.

## 5 Efficient Implementation

The greater expressive power of Datalog<sup>FS</sup> combines with its amenability to efficient implementation via the following three optimization steps: (i) differential fixpoint, (ii) Magic Sets, and (iii) Max-optimization. Since (ii) is basically the same as that in Datalog [11], we will discuss here (i) and (iii).

In Datalog<sup>FS</sup> the differential fixpoint step is applied to recursive rules after they are transformed to ensure that every goal in the b-expression also appears outside the bracket. To satisfy this requirement, the second rule in Example 7 is transformed by repeating outside the bracket the two clauses inside the bracket, producing the following rule:

$$\text{join}(Y) \leftarrow \text{requires}(Y, K), \text{friend}(Y, X1), \text{join}(X1), K: [\text{friend}(Y, X), \text{join}(X)].$$

Since we have renamed the local variables ( $X$  for the case at hand), and since  $K \geq 1$ , this transformation does not change the meaning of the rule. However it greatly simplifies its symbolic differentiation since the bracketed expression can now be treated as a constant. Thus the rule becomes linear and its  $\delta$  version is:

$$\delta \text{join}(Y) \leftarrow \text{requires}(Y, K), \text{friend}(Y, X1), \delta \text{join}(X1), K: [\text{friend}(Y, X), \text{join}(X)].$$

The Max-optimization transforms the delta rules so obtained by replacing the running FS-construct with the final FS-construct. For instance, we start by rewriting the delta rule above into the following one that preserves its operational semantics:

$$\begin{aligned} \delta \text{join}(Y) \leftarrow & \text{requires}(Y, K1), \text{friend}(Y, X1), \delta \text{join}(X1), \\ & K: [\text{friend}(Y, X), \text{join}(X)], K \geq K1. \end{aligned}$$

Now, we can replace the running-FS  $K: [\dots]$  by the final-FS  $K =! [\dots]$  and still preserve the operational semantics of the rule, due to the fact that the rule only uses  $K$  in monotonic functions and predicates (e.g., predicates that if they are true for  $K$  they are also true for every value larger than  $K$ ). This optimization would not be possible if the body uses some non-monotonic predicate, e.g., a goal that checks that  $K$  is even. A similar situation occurs in Examples 9, where instead of the running-FS construct in the recursive  $\delta$  rules we can use the final-FS construct. Indeed the values produced by the former satisfy the external goal  $C =! [\text{cbasic}(\text{part})]$  iff the values produced by the latter do. Again this equivalence is due to the monotonicity of the arithmetic and boolean predicates used in the rule. Such monotonicity holds in all examples given in this paper and the many examples of practical interest discussed in [11].

## 6 Conclusion

In this paper, we studied the important problem of allowing aggregates in recursive Datalog rules, and proposed a solution of surprising simplicity for this long-standing challenge. Our Datalog<sup>FS</sup> approach is based on using continuous aggregate-like functions that allow us to query and reason on the frequency with which predicates and conjunctions of predicates occur.

## References

1. J. M. Hellerstein. Datalog redux: experience and conjecture. In *PODS*, pages 1–2, 2010.
2. G. Gottlob, G. Orsi, and A. Pieris. Ontological queries: Rewriting and optimization. In *ICDE*, pages 2–13, 2011.
3. F. N. Afrati, V. R. Borkar, M. J. Carey, N. Polyzotis, and J. D. Ullman. Map-reduce extensions and recursive queries. In *EDBT*, pages 1–8, 2011.
4. C. Zaniolo. The logic of query languages for data streams. In *Logic and Databases 2011. EDBT 2011 Workshops*, pages 1–2, 2011.
5. S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
6. C. Zaniolo, S. Ceri, C. Faloutsos, R. T. Snodgrass, V. S. Subrahmanian, and R. Zicari. *Advanced Database Systems*. Morgan Kaufmann, 1997.
7. I. S. Mumick and O. Shmueli. How expressive is stratified aggregation? *Annals of Mathematics and Artificial Intelligence*, 15:407–435, 1995.
8. S. Greco and C. Zaniolo. Greedy algorithms in datalog. *TPLP*, 1(4):381–407, 2001.
9. F. Arni, K. Ong, S. Tsur, H. Wang, and C. Zaniolo. The deductive database system *ddl++*. *TPLP*, 3(1):61–94, 2003.
10. I. S. Mumick, H. Pirahesh, and R. Ramakrishnan. The magic of duplicates and aggregates. In *VLDB*, pages 264–277, 1990.
11. M. Mazuran, E. Serra, and C. Zaniolo. Graph languages in Datalog<sup>FS</sup>: from abstract semantics to efficient implementation. Technical report, UCLA, 2011.
12. J. W. Lloyd. *Foundations of Logic Programming, 2nd Edition*. Springer, 1987.
13. M. O. Jackson and L. Yariv. Diffusion on social networks. *Economie Publique*, 2005.