# Discovering Informative Connection Subgraphs in Multi-relational Graphs

Cartic Ramakrishnan[*]
cartic@cs.uga.edu

William H. Milnor
milnor@cs.uga.edu

Matthew Perry[*]
mperry@cs.uga.edu

Amit P. Sheth[*]
amit@cs.uga.edu

[*]LSDIS Lab, University of Georgia, Athens GA

## ABSTRACT

Discovering patterns in graphs has long been an area of interest. In most approaches to such pattern discovery either quantitative anomalies, frequency of substructure or maximum flow is used to measure the interestingness of a pattern. In this paper we introduce heuristics that guide a subgraph discovery algorithm away from banal paths towards more "informative" ones. Given an RDF graph a user might pose a question of the form: *"What are the most relevant ways in which entity X is related to entity Y?"* the response to which is a subgraph connecting *X* to *Y*. We use our heuristics to discover informative subgraphs within RDF graphs. Our heuristics are based on weighting mechanisms derived from edge semantics suggested by the RDF schema. We present an analysis of the quality of the subgraphs generated with respect to path ranking metrics. We then conclude presenting intuitions about which of our weighting schemes and heuristics produce higher quality subgraphs.

## Keywords

Subgraph Discovery, Multi-Relational graphs, Semantic Pattern discovery in RDF Graphs

## 1. INTRODUCTION

*"I keep six honest serving-men (They taught me all I knew); Their names are What and Why and When And How and Where and Who."*— (Rudyard Kipling, from "The Elephant's Child" in Just So Stories 1902). The six questions in this quote by Rudyard Kipling are often tools we as humans use in an attempt to gain knowledge. *How* two entities are related is arguably the most crucial question among these. Discovering relevant sequences of relationships between two entities answers this question. We envision a system, which supports its users in discovering ways in which a pair of entities are related. It is very likely that semantic search engines [1] of the future will need to support such a discovery process. Applications for such a search paradigm can be found in areas such as conflict of interest detection [2] and financial risk analysis [3]. To this end, we investigate techniques that provide users with a chain of relationships between entities in response to queries of the following kind: *"What are the most relevant ways in which entity X is related to entity Y?"* The notion of relevance is critical to the definition of such a query. This becomes clear when one considers the small-world phenomenon [4, 5]. Given a knowledgebase and any two entities X and Y there could be a myriad of relatively short chains (i.e. six degrees) of relationships linking the two. Hence the need for some way of semantically constraining the discovery of possible ways in which X and Y could be related. Faloutsos et.al. [6] address this issue by developing an algorithm to extract relatively small connection subgraphs. They define the *Connection Subgraph Problem* as follows:

**Given:** an edge-weighted undirected graph *G*, vertices *s* and *t* from *G* and an integer budget *b*

**Find:** a connected subgraph *H* containing *s* and *t* and at most *b* other vertices that maximizes a "*goodness*" function *g(H)*.

Faloutsos et.al. [6] applied their techniques to a graph where nodes represented famous people and the edges between these nodes represented strength of acquaintance between them. These connection strengths were derived from name co-occurrences in Web pages. All edges in their dataset therefore have exactly the same interpretation.

Clearly this weighting scheme will not work for finding relevant subgraphs in RDF [7] graphs. Also, naively using a uniform weight on each edge is insufficient, as the semantics of each property type (edge) in RDF is different. Therefore a systematic way of weighting edges based on the semantics conveyed by the ontology represented using RDF schema [8] is needed.

To adapt the approach in [6] to the more general case of an RDF graph:

- We propose heuristics for edge weighting that depend indirectly on the semantics of entity and property types in the ontology and on characteristics of instance data. More specifically, we define *Class* and *Property Specificity, Instance Participation Selectivity* and a *Span Heuristic*.

- We evaluate the generated subgraphs using path ranking schemes suggested in [9-11].

- We present empirical evidence that our weighting schemes do indeed help identify "informative" patterns in the output subgraphs.

- We present results that support the electricity based [6] model for RDF graph relevance.

Section 2 presents related work. In sections 3 and 4 we discuss our algorithms and heuristics respectively. This is followed by a discussion of the dataset for our experiments in Section 5. Section 6 presents our results and evaluations thereof. We conclude in Section 7 with a look at future research directions.

## 2. Related Work

Reasoning and knowledge discovery over graph data models has been studied in the Graph Mining community and more recently in the context of the Semantic Web. The remainder of this section highlights work which is most relevant to ours.

The work most directly related to graph-based knowledge discovery and reasoning for the Semantic Web is that of Semantic Associations which were first introduced in [12]. Semantic Associations (termed ρ-operators) represent meaningful directed paths in an RDF meta-base. To the best of our knowledge this is the only existing work of this type. Anyanwu and Sheth define the *ρ-path* operator among others. Two entities *X* and *Y* are said to be

1

*ρ-path* associated if there exists a sequence of properties (relationships) starting at *X* connecting intermediate entities and ending at *Y*. The nature of web data [4] often leads to an overwhelming number of associations between two entities. To combat this problem, [9, 10] propose to rank Semantic Associations. As an alternative approach, the method in [13] filters the search space before computing associations. They adapt the HITS algorithm [14] to compute importance of Semantic Web resources and then only consider nodes with importance greater than some threshold when computing Semantic Associations. Their preprocessing step based on importance thresholds is likely to discount those paths that contain even a single unimportant node. Our approach to this information overload problem is fundamentally different from these two. We try to find the 'best' set of associations which contain a visually comprehendible number of resources.

There has been a considerable amount of work done in the field of Graph Mining to detect patterns in graphs. Patterns discovered are characterized either by their anomalous nature or frequent occurrence, among other things. Efficient algorithms have been developed for many variations of the frequent subgraph discovery problem [15-17]. Community and group detection is another well-studied graph mining problem which attempts to discover communities and groups based on link analysis. The problem has been studied on both the web graph [18, 19] and other data sets [20]. These graph mining problems however focus on graphs with single node types and single edge types. For the Semantic Web and Link Mining we need algorithms which take into account the semantics of different node and edge types. Community detection and mining in multi-relational networks has recently received a lot of attention [21]. *Novel Link Discovery* was introduced in [11] and involves finding *novel* paths between entities, *novel* loops and significantly connected nodes. The methodology used in this work considers different node and edge types but differs from ours in that importance is determined purely from rarity. Also the paths examined are considerably shorter than the ones we examine.

## 3. ALGORITHMS

Our method for finding a connection subgraph between two RDF resources is based on the algorithms from [6]. The authors present an algorithm for extracting a so-called *candidate graph* from an input graph. They also propose an algorithm based on electrical circuits to extract a display graph from the candidate for a given budget *b*. For our purposes we refer to these as *Candidate ρ-graph* and *Display ρ-graph*. We assume that the properties (edges) in the RDF graph are bidirectional (*i.e.* every relationship has a corresponding inverse relationship). This assumption is necessary because two resources may not be connected by a directed path but by a path which contains inverse relations. Ignoring this path could exclude vital information about the connections between the entities.

### Candidate ρ-graph Generation Algorithm

The *Candidate ρ-graph* generation algorithm is used to prune the search space in very large graphs. It is based on a notion of distance between two nodes. The algorithm grows a set *S* around the source node *s* and a set *T* around the sink node *t* (*s* and *t* are referred to as the roots of their respective sets) until a certain threshold is met: a maximum number of expanded nodes or maximum number of cut edges between *S* and *T*. At each iteration, a pending list is maintained for each of these sets which

consists of those nodes $n \notin S$ and $n \notin T$ and adjacent to some node $k \in S$ or $k' \in T$. The sets *S* and *T* are expanded by choosing from the pending list the node with shortest distance to either *s* or *t*. Let *u'* be the predecessor of *u* (the node adjacent to *u* on the shortest path to its root). For an edge *(u, v)* the distance between *u* and *v* is given by:

$$ distance \quad (u,v) = log \left( \frac{(degree \ (u) + degree \ (v))^2}{w(u,v) * \beta_{u' \to u \to v}} \right) $$

$w(u,v)$ and $\beta_{u' \to u \to v}$ are a function of the three heuristics explained in Section 4 *viz.* Class and Property Specificity, Instance Participation Selectivity and Span. The length of a path is the sum of the length of its edges. The aim of our initial experiments is to determine the quality of the *Candidate ρ-graph* in terms of its ability to capture the best paths between the query endpoints.

### Display ρ-graph Generation Algorithm

The *Display ρ-graph* generation algorithm extracts a small connection subgraph from the input graph. In [6] the authors present a rather elegant solution to this by modeling the graph as an electrical circuit where the edge weights represent the conductance values in the circuit. They use the fact that current flows from high voltage to low voltage to impose direction on an otherwise undirected graph. Using Ohm's law and Kirchoff's law, a system of linear equations is created with voltages at each node as a variable in these equations. Solving this system of equations gives voltages at each node. This step takes $\Theta(n^3)$ time, which motivates the need for the *Candidate ρ-graph* generation. The greedy *Display ρ-graph* generation algorithm attempts to find a graph of at most *b* nodes (set to a maximum of 100 in our experiments) which maximizes the amount of total current delivered from the start node to the end node. Starting with an empty subgraph, this algorithm iteratively adds paths until meeting the budget *b*. At each of the iterations, a dynamic programming algorithm is used to find the path which has the maximum ratio of delivered current to number of new nodes added to the subgraph. This choice may not be globally optimal, hence the greedy nature of the algorithm. In our experiments we test this model based on current flow used to compute these display ρ-graphs.

## 4. HEURISTICS

RDFS vocabulary allows users to represent classes and relationships (properties) connecting them thereby indirectly imposing meaning on resources that are instances of these classes. We define three quantities (Class and Property Specificity, Instance Participation Selectivity, and Span) indirectly based on semantics and RDF statement types and frequencies. Our aim in doing this is to use semantics suggested by the schema to systematically convert an arbitrary un-weighted RDF graph into an edge-weighted graph appropriate as input to the algorithms described previously.

We define a schema *S* as the union of the set of classes (*C*) and property types (*P*). Further, we define an RDF data store $R = \langle \Pi, I \rangle$ where $\Pi = \bigcup S$ and *I* is the set of class and property instances corresponding to the schemas. A single entity could be an instance of multiple classes belonging to different schemas in $\Pi$. We assume that such an entity instance is uniquely

2

identified by one URI. In other words, no data integration operation is required.

## Class and Property Specificity (CS and PS)

Intuitively more specific resources (entity instances and properties) participating in a path, convey more information than general ones. For instance, it is more informative if one knows that Michael Jordan was a *basketball player* as opposed to knowing that he is a *person*. Similarly, knowing that Rudy Giuliani was an *employee of* New York City is less informative than the fact that he was *mayor of* New York City.

As a result of the rdfs:subClassOf and rdfs:subPropertyOf properties provided by RDF schema it is possible to impose a partial ordering of properties and classes in the schema resulting in a wellformed hierarchy of classes and properties. For a given property p, let $H(p)$ be the length of the longest path in the hierarchy tree that contains p, and for a given class c, let $H'(c)$ be the length of the longest path in the hierarchy tree from the root to c. Properties and classes at the root of their respective hierarchy trees in the schema are considered most general while those at the leaves of these trees are considered most specific. Therefore a measure of specificity can be associated with each class or property commensurate with its position in its hierarchy. Let the depth of an arbitrary property in its property hierarchy be $d(p_i)$ and the depth of an arbitrary class in its class hierarchy be $d(c_j)$. Therefore, the specificity of property $p_i$ and class $c_j$ are given by $\mu(p_i) = \frac{d(p_i)}{H(p)}$  $\mu(c_j) = \frac{d(c_j)}{H'(c)}$. Every resource that is an instance of the class $c_j$ is assigned the weight $\mu(c_j)$. If a resource $r$ is an instance of $k$ distinct classes it is assigned the value $\mu(r) = \max_{1 \leq m \leq k} \mu(c_m)$, since we want the most specific nodes and properties to be in the output subgraph. To convert this node weight into an edge weight, the value of each resource weight is equally distributed among all edges incident on the resource *r*. This weighting scheme favors nodes with lower degree since the node specificity is divided equally among its incident edges, therefore edges incident on nodes with high degree will get a lower weight.

## Instance Participation Selectivity (IPS)

Another guideline we use is that rarer facts are typically more informative that frequently occurring ones [22]. Consider the example shown in Fig.1. The example shows two relationships *lives_in* and *council_member_of* defined on the classes *Person* and *City*. The instances $p_1, p_{2...} p_m$ of the class *Person* are members of the council of *City* $c_1$, hence the relationship *council_member_of* between each $p_1, p_{2...} p_m$ to $c_1$. Instances of class *Person* $p_{m+1}, p_{m+2}, \ldots p_{k-2}, p_{k-1}, p_k$ represent people who live in *City* $c_1$ and therefore are related to $c_1$ by the relationship *lives_in*. From the perspective of the node $c_1$, following an edge labeled *lives_in* will lead to one node among *k-m* possible nodes. In contrast, following an edge labeled *council_member_of* will lead to one node among *m* nodes. Given that rarer paths are considered more informative, the amount of information gained by choosing to traverse the *council_member_of* relationship to a node in the set $\{p_1, p_{2...} p_m\}$ is more than the gain achieved by choosing to traverse the *lives_in* relationship to a node in the set $\{p_{m+1}, p_{m+2}, \ldots p_{k-2}, p_{k-1}, p_k\}$.
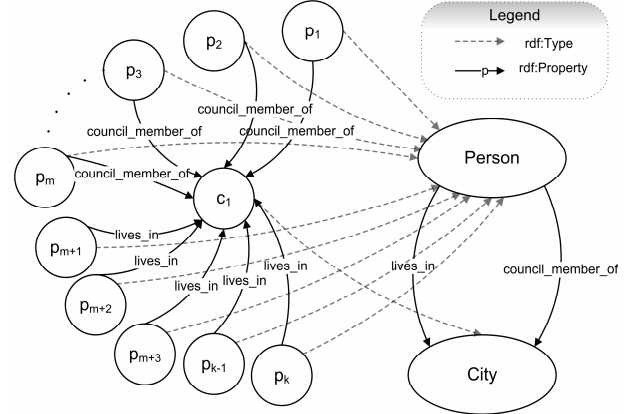


**Figure 1 Illustrative example for Instance Participation Selectivity**

This is akin to choosing the hop with maximum information gain. To define this heuristic formally, we introduce the notion of the type of an RDF statement. The type of an RDF statement $\langle s,p,o \rangle$ is defined as the triple $\pi = \langle C_i, p, C_j \rangle$ where $typeOf(s) = C_i$ and $typeOf(o) = C_j$. Further, $| \pi |$ is thus the number of statements of type $\pi$ in a given RDF instance base. We therefore define *Instance Participation Selectivity* for each RDF statement as $\sigma_\pi = 1/| \pi |$. Going back to Figure 1, let $\pi = \langle$ *Person, lives_in, City* $\rangle$ and $\pi' = \langle$ *Person, council_member_of, City* $\rangle$. According to this example, $\sigma_\pi = 1/(k-m)$ and $\sigma_{\pi'} = 1/m$ and if $k > m$ then $\sigma_{\pi'} > \sigma_\pi$.

## The Span Heuristic (SPAN)

In [9] the authors define a ranking metric known as *Refraction*. Given a path of the form $v_1, e_1, v_2, e_2 \ldots e_{n-2}, v_{n-1}, e_{n-1}, v_n$ from $v_1$ to $v_n$, where $v_i \in$ *Resources* and $e_i \in$ *Properties* $\forall i$ $1 \leq i \leq n$, this path is said to *refract* if there exists at least a pair of statements $\langle v_i, e_i, v_{i+1} \rangle$, $\langle v_{i+1}, e_{i+1}, v_{i+2} \rangle$ such that $\neg \exists (S|e_i, e_{i+1} \in S)$. In other words, this path passes through an instance of classes that belong to more than one schema. The number of such occurrences measures the extent to which a given path conforms to a schema. We consider resources that are instances of classes belonging to different schemas as being indicative of informative paths between the given entities, since they tie different domains together. What makes such paths interesting is the fact that these paths represent a deviation from the expected paths suggested by the schemas. For example, in our scenario in Figure 4 an instance of the class *Person* may be classified as both an instance of class *Actor* in the *Entertainment* domain and an instance of class *SpokesPerson* in the *Business* domain. Such an instance serves to link different schemas.

As a heuristic that favors the inclusion of such *refracting* paths in output subgraphs, we define a Span weight for an edge based on the class types of its two endpoints. Let us consider the example in Figure 2. For every node *v* in a given RDF graph we can define a set called *SchemaCover,* which is the set of schemas to which the classes (types) of *v* belong. Formally,

$$SchemaCover = \{S | \exists C \in S \land typeOf(v) = C\}$$

The *SchemaCover* for each of the nodes in the set $\{u', u, v_1, v_2, v_3, v_4, v_5\}$ is shown adjacent to the respective node in Figure 2. To

favor paths that span as many schemas as possible the search algorithm favors nodes that are classified under as many "new" schemas as possible at each step. By "new" we mean schemas that have been least recently encountered along a particular path. Let *SDiff(u,u')* represent the number of new schemas seen as a result of traversing the edge *(u,u')*, where the value of *SDiff(u,u')* = |*SchemaCover(u')-SchemaCover(u)*|. The idea behind *SDiff* is to ensure that the discovery algorithm chooses the edge with maximum schema variety to traverse next. Using the absolute *SchemaCover* difference for one hop only, could lead to the following situation: The hop from *u* to *u'* (based on high *SDiff(u, u')*) is chosen. Subsequently, $v_i$ is chosen such that it maximizes *SDiff(u',v_i)*. But *u'* and $v_i$ are covered by exactly the same schemas (as in node *u'* and $v_5$ in Figure 2). Therefore *SDiff* alone does not ensure that search will continue along paths that have nodes classified under more diverse schemas. To reduce the chance of this problem we define the *Cumulative Schema Difference* $CSDiff(u,u',v_i)=1+SDiff(u,v_i)+SDiff(u',v_i)$ for $v_i \in adj[u]-\{u'\}$. We normalize this *Cumulative Schema Difference* (*CSDiff*) measure to compute a factor $\beta_{u'\to u\to v_i}$ ;

$$\beta_{u'\to u\to v_i} = \frac{CSDiff}{1+2(m-1)}, \text{ where } m \text{ is the number of schemas}$$

We then obtain the adjusted weight given by;

$$w'(u,v_i)=\beta_{u'\to u\to v_i}*w(u,v_i)$$



SDiff (u,v₁)  = | {A, B, C, D} - {A, B} | = 2
SDiff (u,v₂)  = | {B} - {A, B} | = 0
SDiff (u,v₃)  = | {A, B, C} - {A, B} | = 1
SDiff (u,v₄)  = | {A, C, D} - {A, B} | = 2
SDiff (u,v₅)  = | {A} - {A, B} | = 0
SDiff (u',v₁) = | {A, B, C, D} - {A} | = 3
SDiff (u',v₂) = | {B} - {A} | = 1
SDiff (u',v₃) = | {A, B, C} - {A} | = 2
SDiff (u,v₄)  = | {A, C, D} - {A} | = 2
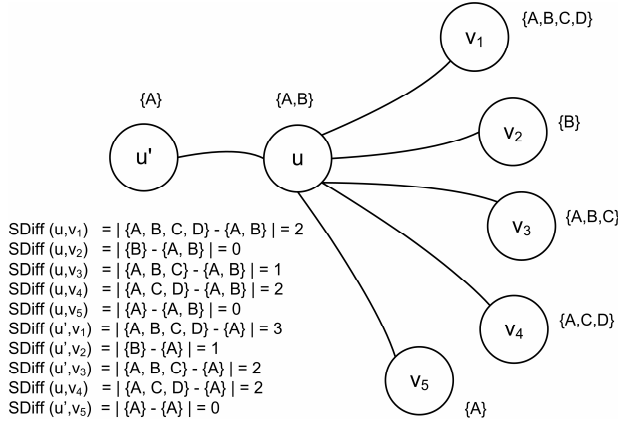SDiff (u',v₅) = | {A} - {A} | = 0

**Figure 2 Example of Span metric computation**

The effect of the factor *β* is to bias edge weights in the following way. Successor nodes that are instances of classes belonging to schemas other than those of the current and previous node are more likely to be visited, quantified by the two *SDiff* terms of *CSDiff*. More specifically, in the case of the example in Figure 2, a partial ordering is induced by the adjusted weights *w'(u, v_i)*, on the nodes as follows $v_1 \succ v_4 \succ v_3 \succ v_2 \succ v_5$. The node $v_1$ is therefore visited next. However, the measure *β* is not sufficient to distinguish between nodes in all cases. Consider the example in Figure 3. Nodes $v_1$ and $v_2$ have the same value of $\beta_{u'\to u\to v_i}$; but $v_1$ should be more desirable than $v_2$ because it has a larger *SchemaCover* value. For such cases, we define a factor called *SchemaCoverFactor* $\alpha(u, v)$:

$$\alpha(u,v)=\frac{1}{2}\left(\frac{|SchemaCover(u)|+|SchemaCover(v)|}{m}\right),$$
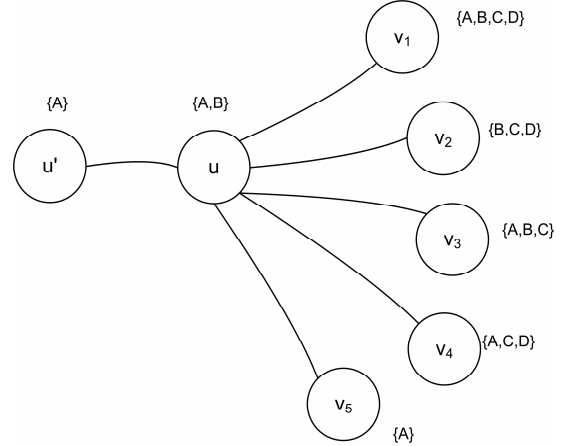
where *m* is the number of schemas.



**Figure 3 Influence of the Schema Cover Factor α**

SDiff (u,v₁)  = | {A, B, C, D} - {A, B} | = 2
SDiff (u',v₁) = | {A, B, C, D} - {A} | = 3
CSDiff        = 1+ SDiff (u,v₁)+SDiff (u',v₁)

$$\beta_{u'\to u\to v_1} = \frac{CSDiff}{1+2(m-1)}$$

$$\beta_{u'\to u\to v_1} = \frac{1+3+2}{1+2(4-1)} = 0.8571$$

$$\beta_{u'\to u\to v_2} = \frac{CSDiff}{1+2(m-1)}$$

$$\beta_{u'\to u\to v_1} = \frac{1+3+2}{1+2(4-1)} = 0.8571$$

SDiff (u,v₂)  = | {B,C,D} - {A, B} | = 2
SDiff (u',v₂) = | {B,C,D} - {A} | = 3
CSDiff        = 1+ SDiff (u,v₂)+SDiff (u',v₂)

$$\alpha(u,v_1)=\frac{1}{2}\left(\frac{|SchemaCover(u)|+|SchemaCover(v_1)|}{m}\right)$$

$$\alpha(u,v_2)=\frac{1}{2}\left(\frac{|SchemaCover(u)|+|SchemaCover(v_2)|}{m}\right)$$

$$\alpha(u,v_1)=0.75 \text{ and } \alpha(u,v_2)=0.625$$

Note how the values of α for the two pairs of nodes that are in consideration are different even though the β values are the same. We therefore use α as a tie breaker in such cases. As per the calculations shown in Figure 3 this factor treats the node $v_1$ preferentially over node $v_2$ i.e. $v_1 \succ v_2$ thus resolving the ambiguity. Note that since the factor β is a property on a 3 node sequence it is query dependent and therefore cannot be pre-computed. The value of $\beta_{u'\to u\to v_i}$ is therefore computed on-the-fly during the candidate generation process.

## Combining three heuristics

The values of all the produces using all the heuristics discussed above lie in the interval (0,1]. Although different weighting could be assigned to each heuristic we give them all equal weights in our experiments. The initial weight of an edge is given by the following formula:

$$w(u,v)=\frac{\mu(p_{u\to v})+\frac{1}{2}\left(\frac{\mu(u)}{degree(u)}+\frac{\mu(v)}{degree(v)}\right)+\sigma_\pi+\alpha(u,v)}{4}$$

where $p_{u\to v}$ is the property connecting the resource node *u* and *v*, and π is the type of the statement $\langle u, p_{u\to v}, v\rangle$. $\beta_{u'\to u\to v_i}$ is then used to adjust the weights $w(u, v_i) \ \forall i\, 1\le i\le n$ as follows: $\beta_{u'\to u\to v_i}* w(u, v_i)$. This is done during the candidate generation phase when the path leading to a given edge is known.

## 5. DATASET AND SCENARIO

We used a synthetic dataset for our experiments since we needed control over characteristics of the data. This helps us ensure that our results are not unduly affected by unknown aspects i.e. connectivity, relative instance distribution etc. of the dataset.

4

Collection of real world data follows an almost opportunistic approach since availability often dictates design. As a result there is room for skew in instance data population. This skew may not always reflect real-world distributions, as was observed in our experience with SWETO [23]. To circumvent this we built a utility [24] that takes as input a set of schemas and a properties file specifying relative distributions of instances of classes and properties that would be expected in the real world. For example, consider two classes in the *Business* ontology (Appendix Fig. A2.): *Trustee* and *Employee.* It would be reasonable to assume that if there are 5,000 instances of the class *Employee* then there are **unlikely** to be 1,000 instances of the class *Trustee.* Instances of the class *Trustee* are more likely to be approximately 10. These numbers are domain specific. Our method for assigning values to these relative distributions is empirical and a discussion of this issue is beyond the scope of this paper. The result of running this utility is an RDF graph that contains nodes and edges that are instances of classes and property types belonging to any or all of the classes in the given schemas. The graph for our experiments contains 30,000 nodes and 45,000 edges.
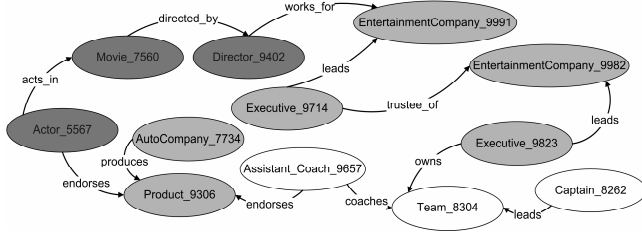


**Figure 4 Example snippet of a subgraph returned for the query ρ(Actor_5567, Captain_8262) on our synthetic dataset– (Nodes in the above graph are color-coded according to the schemas their class belongs to. White nodes for Sports schema classes, Light Grey for Business and Dark Grey for Entertainment)**

As a motivation for the domains used in our dataset consider the following example. A fraud investigator with the Securities and Exchange Commission (SEC) receives the following piece of information about a week after the stock prices for *EntertainmentCompany_9982* plummet. *Actor_5567* sold 70% of his shares of *EntertainmentCompany_9982* one week after *Capt_8262* sold all of his shares in the same company. Both transactions took place two weeks before the prices plummeted. The example subgraph shown in Fig.4, might help an investigator visualize the connections between the resources: *Actor_5567* and *Captain_8262*.

## 6. RESULTS AND EVALUATION

We recognize the fact that the notion "best" subgraph is very subjective and dependent on the user's perspective. It is however desirable to have an objective measure that could be used to quantify the quality of a generated subgraph. The issue of judging relevance of paths i.e. path ranking has been addressed in [9] and [10]. In [11] the authors use *rarity* of the path as a measure of its interestingness. To the best of our knowledge these are the only three efforts that measure path relevance. We therefore use these path ranking mechanisms to evaluate the quality of both the *Candidate ρ-graph* and the *Display ρ-graph.* In our experiments

the *Candidate ρ-graphs* generated contained 3000[1] nodes and the *Display ρ-graphs* were restricted to a maximum of 100 nodes making them easy to visualize.

## Evaluation using Path Ranks

In our data set there are over 60 million paths of length 13 between the two endpoints used in Fig. 4. Paths of this length are unlikely to be of much interest to the user. To evaluate our subgraphs, we run an exhaustive *k-hop* limited Depth-First Search (DFS) on the input graph between the two entities. We use a depth limit of 9 hops for our experiments for feasibility of path enumeration for ranking. Note that both the *Candidate ρ-graph* and *Display ρ-graph* generated do contain arbitrary length paths, but we only consider paths of length at most 9 for fairness of comparison. We represent the paths returned by the *k-hop* DFS as the set *FGPaths₉ (paths of up to 9 hops in the full graph).* There are therefore 30 distinct *FGPaths₉* sets, one for each query in our experiments. We rank the paths in each of the *FGPaths₉* sets using the ranking mechanisms proposed in [9] and [10] in addition to what we call *Rarity Rank* based on the method suggested in [11]. While *Rarity Rank* ranks paths based only on the rarity of the edges that constitute the path, SemRank [9] uses a combination of Rarity, Specificity and Refraction (measure of conformity to schema) to rank paths. The mechanism proposed in [10] is the most flexible since it allows one to incorporate context, trust in the statements that constitute the path, path length, specificity and rarity.

The rank of a path *p* based on the *Rarity Rank* scheme is given by the inverse of the number of paths that share the same type as path *p*. Each of the ranking mechanisms applied to the set *FGPaths₉* results in a list of ranked paths. Let us assume that this leads to ranking from 1→M where M is the rank of the least relevant path. Let this set of ranked paths be represented as *FGRankedPaths₉.* We therefore have three distinct scales (*FGRankedPaths₉* sets) against which the quality of both *Candidate ρ-graph* and the *Display ρ-graph* can be measured. In all of the graphs, a table shown below the x-axis represents the 16 possible combinations of the four heuristics we use *viz. Class* and *Property Specificity (CS and PS), Instance Participation Selectivity (IPS)* and *The Span Heuristic (SPAN).*

### *Measuring Candidate ρ-graph quality*

To measure *Candidate ρ-graph* we compare the best paths in the entire graph to those in the *Candidate ρ-graph.* Let *CGPaths₉* represent the set of paths in the *Candidate ρ-graph* with maximum length 9. For each path $p_{candidate} \in CGPaths_9$ we count the number of paths $p \in FGRankedPaths_9$ such that *rank(p) > rank(p_{candidate}).* This gives us the rank of each path in the *Candidate ρ-graph* with respect to all paths in the set *FGRankedPaths₉.* The score of a path is given by:

$$score(p_{candidate}) = \left| FGRankedPaths_9 \right| - rank(p_{candidate})$$

The quality of the *Candidate ρ-graph* is therefore given by:

---

[1] This was the observed number of nodes in the *Candidate ρ-graph* for all the 30 queries used in our experiments. Further investigation revealed that this was an artifact of the connectivity of our dataset.

$$Q(CGPaths_9) = \frac{\sum\limits_{p_{candidate}\in CGPaths_9}\left(score(p_{candidate})\right)}{\frac{|CGPaths_9|}{\sum\limits_{r=1}\left(|FGRankedPaths_9|-r\right)}}$$
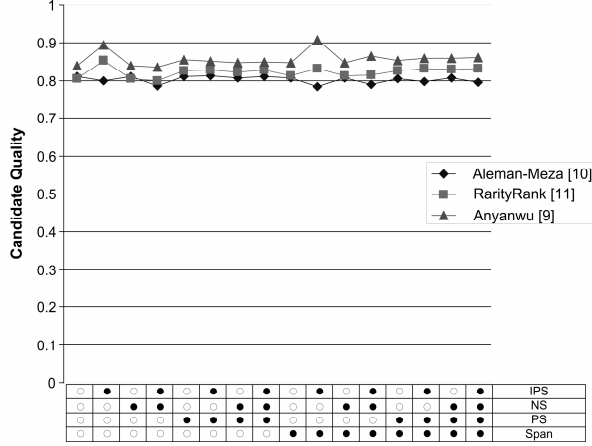


**Figure 5 Quality of the Candidate ρ-graph as percentage of maximum score for averaged over 30 queries**

Figure 5 shows that the Candidate ρ-graph containing *k* paths obtained using our edge weighting schemes achieves between 80—90% of the score that can be achieved by choosing the *top-k* ranked paths from the full graph (entire dataset of 30,000 nodes and 45,000 edges). The Candidate ρ-graphs in our results typically contain 30-40% of the paths in the entire graph between the endpoints yet are 80-90% as *"good"* as the top paths in the entire graph between the two endpoints.

*Measuring Display ρ-graph quality*

Similar to *Candidate ρ-graph* quality, we compare the paths in the *Display ρ-graph* to the best paths in the entire graph. Let the set *DGPaths* represent the paths in the *Display ρ-graph*. The rank of a path in the *Display ρ-graph* is computed exactly the same way the rank of a path in the *Candidate ρ-graph* is computed, as is the score.

$$score(p_{display}) = |FGRankedPaths_9| - rank(p_{display})$$

The quality of a display ρ-graph is computed by comparing its cumulative score to the best possible display that could be obtained from the ranked set of paths in the full graph. We refer to this best possible display as *Pseudo-Display*. In our experiments we use a *budget* of 100 nodes for our *Display ρ-graphs*. Starting with an empty *Pseudo-Display* graph and the path with rank 1 in the set *FGRankedPaths_9* we add paths to the *Pseudo-Display* until 100 nodes have been added. The cumulative score of the *Pseudo-Display* is then computed as the sum of the scores of the paths. The quality of a *Display ρ-graph* is therefore given by:

$$Q(DGPaths) = \frac{\sum\limits_{p_{display}\in DGPaths}score(p_{display})}{\sum\limits_{p_{pseudo}\in Pseudo-Display}score(p_{pseudo})}$$

Figure 6 shows that starting with the Candidate ρ-graphs with 80—90% quality the Display ρ-graphs computed capture a maximum of 84% of the score that can be obtained by taking the best paths in the full graph. Our results show the quality of Display ρ-graphs with respect to SemRank [9] to be surprisingly low – 43%. Further investigation of the methods used revealed that the difference between the ranking scheme in [10] and that in [9] is that in the former instance node degrees affect the rank of a path (nodes of lower degree being favored) whereas in the latter rank of path is determined purely by properties in the path. Our heuristics favor lower degree nodes and hence the observed trend. A personal communication with the authors of [9] revealed that extending SemRank to include the effect of nodes is an intended follow up to this work.
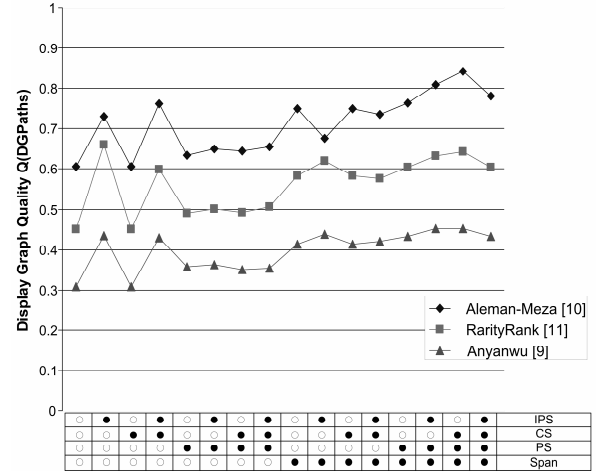


**Figure 6 Quality of the Display ρ-graph – Note that all weighting heuristics turned off results in poor graph quality in contrast with all heuristics turned on**

*Successive Display ρ-graph quality*

With the intention of validating the current flow model for subgraph relevance [6] we conducted the following experiment. We computed what we term as *Successive Display ρ-graphs*. To construct these displays we successively run the *Display ρ-graph* generation algorithm on the candidate ρ-graph. At each successive run we discount the paths used in previous displays. This results in the next best Display ρ-graph at every successive run. This process is repeated five times in our experiments to obtain five Display ρ-graphs. The current flow in each of these Display ρ-graphs is plotted relative to the current flow in the first Display ρ-graphs on a *log* scale in Figure 7. The quality of these Display ρ-graphs is plotted relative to the quality of the first Display ρ-graph in Figure 8. There is a large difference both in the current flow and the display quality between the first display and the next display. This confirms that there is a correspondence between current flow in the Display ρ-graphs and their quality. This in turn supports the electricity based model for RDF graph relevance. Note that the plots below are averages of the relative differences of successive displays over all ranking schemes.
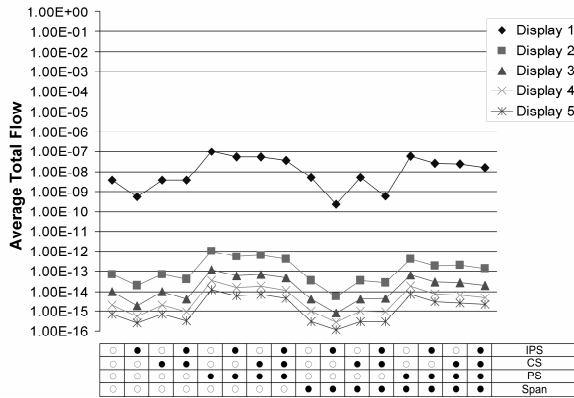
6

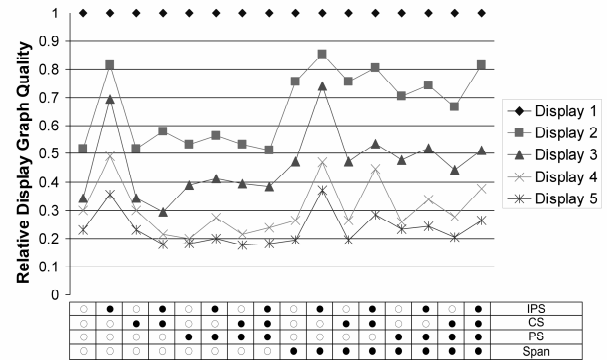**Figure 7 Current Flow in 5 Successive Display ρ-graphs**



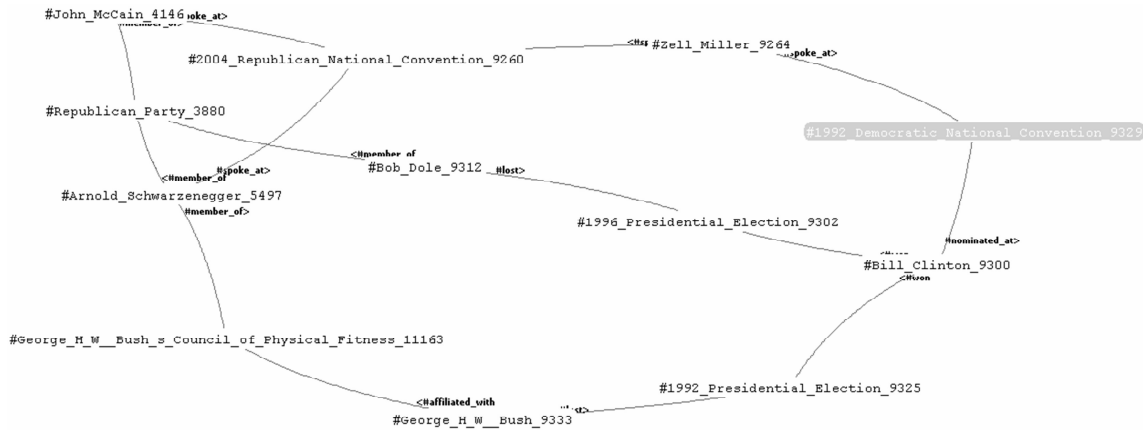**Figure 8 Quality of 5 Successive Displays relative to the best**



**Figure 9 Display graph with budget =12 WITHOUT using edge weighting scheme**
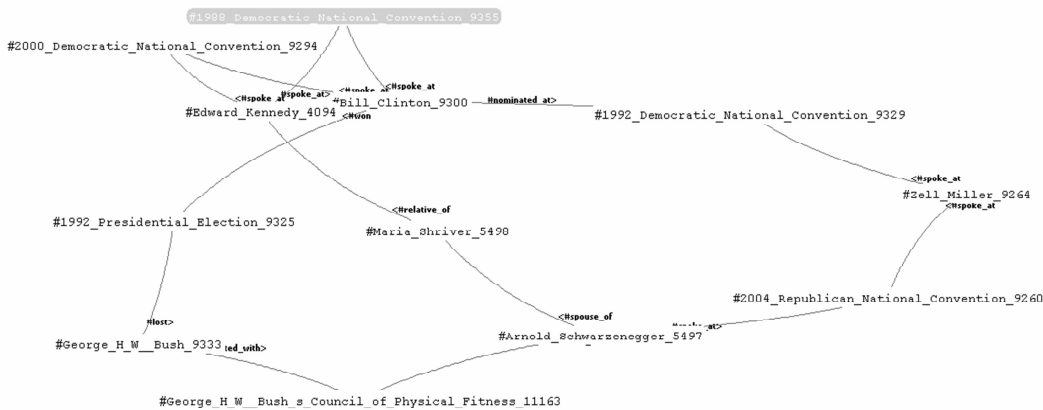


**Figure 10 Display graph with budget =12 using edge weighting scheme**

We conducted another experiment on factual data extracted from public access government websites. For this particular experiment we picked two persons of interest at random. Figures 9 and 10 compare Display ρ-graphs generated for end points: *Edward Kennedy* and *Bill Clinton* both spoke at consecutive Democratic National Conventions. This indicated a strong similarity in their political ideologies. A strong familial relationship between *Edward Kennedy* and *Arnold Schwarzenegger* via *Maria Shriver* is also seen in the output in Figure 10. In comparison to these paths the paths in Figure 9 are rather uninteresting. This provides empirical evidence that our

*Arnold Schwarzenegger* and *Bill Clinton*. In both cases the budget is set to 12 nodes. Figure 10 contains a very interesting path that connects *Arnold Schwarzenegger* to *Bill Clinton* via *Maria Shriver* followed by *Edward Kennedy*. weighting schemes are useful. We believe that further investigation of semantic methods for subgraph discovery is warranted.

## 7. CONCLUSIONS AND FUTURE WORK

Our results suggest that using edge weights generated by our weighting scheme results in highly relevant *Candidate ρ-graphs*,

7

where relevance is judged using established path ranking metrics. Further evidence supporting this claim can be seen from quality of the *Display ρ-graphs*. The ranking metrics proposed by Aleman-Meza et.al. [10] in our experiments show that the quality of the *Display ρ-graphs* are best when using Class Specificity (CS), Instance Participation Selectivity (IPS) and Span together. Results for the Successive Displays serve to support the electricity flow based model for RDF subgraph relevance, besides validating our edge weighting schemes. Results presented in this paper seem very promising for application domains like Ontology based Scientific Discovery where the ability to visualize relevant relationships between metadata entities is crucial. As a follow up to this work we plan to apply our techniques to develop tools for finding correlations between Glycosylation patterns and patterns of gene expression within a cell line in the Glycomics [25] domain. We further propose to develop algorithms to support queries involving $k$ endpoints for RDF graphs. Another interesting direction involves formalizing the notion of *Context* and investigating *Context-Aware Subgraph Discovery* algorithms.

# 8. ACKNOWLEDGEMENTS

# 9. REFERENCES

[1] Guha, R.V. and E.M. Rob McCool. *Semantic search*. in *The Twelfth International World Wide Web Conference*. 2003. Budapest, Hungary.

[2] http://lsdis.cs.uga.edu/projects/semdis/coi.

[3] http://lsdis.cs.uga.edu/projects/semdis/HS-brief.pdf.

[4] Albert, R. and A.-L. Barabasi, *Statistical mechanics of complex networks.* Reviews of Modern Physics, 2002. **74**(1): p. 47 - 97.

[5] Milgram, S., *The Small World Problem.* Psychology Today, 1967. May 1967: p. 60-67.

[6] Faloutsos, C., K.S. McCurley, and A. Tomkins. *Fast discovery of connection subgraphs*. In *Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2004. Seattle, Washington.

[7] Lassila, O. and R.R. Swick. *Resource Description Framework (RDF) Model and Syntax Specification, W3C Recommendation*. 1999

[8] RDFS, *http://www.w3.org/TR/rdf-schema/*.

[9] Anyanwu, K., A. Maduko, and A. Sheth. *SemRank: Ranking Complex Relationship Search Results on the Semantic Web*. In *The Fourteenth International World Wide Web Conference*. 2005. Chiba, Japan.

[10] Aleman-Meza, B., et al., *Ranking Complex Relationships on the Semantic Web.* IEEE Internet Computing, Special Issue: Information Discovery: Needles and Haystacks, 2005. **9**(3): p. 37 - 44.

[11] Lin, S.-d. and H. Chalupsky. *Unsupervised Link Discovery in Multi-relational Data via Rarity Analysis*. In *Third IEEE International Conference on Data Mining (ICDM 2003)*. 2003. Melbourne, Florida.

[12] Anyanwu, K. and A. Sheth. ρ-*Queries: Enabling Querying for Semantic Associations on the Semantic Web*. In *The Twelfth International World Wide Web Conference*. 2003. Budapest, Hungary.

[13] Mukherjea, S. and B. Bamba. *BioPatentMiner: An Information Retrieval System for BioMedical Patents*. In *Thirtieth International Conference on Very Large Data Bases*. 2004. Toronto, Canada.

[14] Kleinberg, J.M., *Authoritative Sources in a Hyperlinked Environment.* Journal of the ACM, 1999. **46**(5): p. 604 - 632.

[15] Yan, X. and J. Han. *CloseGraph: Mining Closed Frequent Graph Patterns*. in *Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2003. Washington, DC, USA.

[16] Huan, L., W. Wang, and J. Prins. *SPIN: Mining Maximal Frequent Subgraphs from Graph Databases*. In *Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2004. Seattle, WA.

[17] Kuramochi, M. and G. Karypis, *GREW: A Scalable Frequent Subgraph Discovery Algorithm*, in *Fourth IEEE International Conference on Data Mining (ICDM 2004)*. 2004: Brighton, UK. p. 439 - 442.

[18] Flake, G.W., et al., *Self-organization of the web and identification of communities.* IEEE Computer, 2002. **35**(3): p. 66-71.

[19] Gibson, D., J. Kleinberg, and P. Raghavan. *Inferring Web Communities from Link Topology*. in *9th ACM Conference on Hypertext and Hypermedia (HyperText 98)*. 1998. Pittsburgh, PA.

[20] Adibi, J., et al. *The KOJAK Group Finder: Connecting the Dots via Integrated Knowledge-Based and Statistical Reasoning*. in *AAAI*. 2004. San Jose, California, USA.

[21] Cai, D., et al. *Mining Hidden Community in Heterogeneous Social Networks*. In *ACM-SIGKDD Workshop on Link Discovery: Issues, Approaches and Applications (LinkKDD'05)*. 2005. Chicago, IL.

[22] Shannon, C.E., *A Mathematical Theory of Communication.* Bell System Technical Journal, July and October, 1948. **27**: p. 379-423; 623-656.

[23] Aleman-Meza, B., et al. *SWETO: Large-Scale Semantic Web Test-bed*. In *16th International Conference on Software Engineering & Knowledge Engineering (SEKE2004): Workshop on Ontology in Action*. 2004. Banff, Canada.

[24] Perry, M. (2005) *TOntoGen: A Synthetic Data Set Generator for Semantic Web Applications*. AIS SIGSEMIS Bulletin **Volume**, 46-48

[25] Sheth, A., et al. *Semantic Web technology in support of Bioinformatics for Glycan Expression*. In *W3C Workshop on Semantic Web for Life Sciences*. October 27-28, 2004. Cambridge, Massachusetts.

8