

Relational Similarity Machines

Ryan A. Rossi
Palo Alto Research Center
rrossi@parc.com

Rong Zhou
Palo Alto Research Center
rzhou@parc.com

Nesreen K. Ahmed
Intel Labs
nesreen.k.ahmed@intel.com

ABSTRACT

This paper proposes Relational Similarity Machines (RSM): a fast, accurate, and flexible relational learning framework for supervised and semi-supervised learning tasks. Despite the importance of relational learning, most existing methods are unable to handle large noisy attributed networks with low or even modest levels of relational autocorrelation. Furthermore, they also have issues with efficiency, accuracy, scalability, and flexibility. For instance, many existing methods perform poorly for multi-class classification problems, graphs that are sparsely labeled, or network data with low relational autocorrelation. In contrast, the proposed relational learning framework is designed to be (i) fast for learning and inference at real-time interactive rates, and (ii) flexible for a variety of learning settings, constraints, and application domains. The experiments demonstrate the effectiveness of RSM for a variety of tasks and data.

Keywords

Statistical relational learning; collective classification; semi-supervised learning (SSL); multi-class; node classification; interactive machine learning

1. INTRODUCTION

Networks (relational data, graphs) encode dependencies between entities (people, computers, proteins) and allow us to study phenomena across social [2], technological [26], and biological domains [9]. Recently, relational machine learning (RML) methods were developed to leverage relational dependencies [17, 13, 40] between nodes to improve predictive performance [24, 16, 31, 29, 13, 36].

Relational classifiers can sometimes outperform traditional *iid* ML techniques by exploiting dependencies between class labels (attributes) of related nodes. However, the performance of RML methods can degrade when there are few labeled instances (majority of neighboring instances are also unlabeled). Collective Classification (CC) aims to solve this problem by iteratively predicting labels and propagating them to related instances [45]. Unfortunately, the performance of CC methods may also degrade when there are very few labels available

(e.g., label density < 0.01) [31]. In both situations, if not careful, the performance of RML methods may degrade to a point where *iid* techniques perform better.

Despite the fundamental importance of these techniques, the vast majority of RML methods rely on a significant amount of relational autocorrelation (homophily [33]) existing in the data. It has been noted that RML techniques may perform worse than *iid* methods when there is low or even modest relational autocorrelation. In addition, existing methods also have difficulty learning with graph data that is large, noisy, probabilistic, sparsely labeled, attributed, and are sensitive to many other issues and data characteristics that often arise in practice. Moreover, obtaining labeled data is also expensive, and thus RML methods should be robust to learning with few labeled instances. Furthermore, relational representation and transformations of the nodes, links, and/or features can dramatically affect the capabilities and results of such algorithms [40]. In general, RML methods are sensitive to many other important and fundamental issues that arise in practice (and thus unique to interconnected and dependent relational data).

To solve these problems, this paper introduces relational similarity machines (RSM) — an efficient, flexible, and highly accurate relational learning framework. RSM is based on the fundamental notion of similarity and is well-suited for classification and regression tasks in arbitrary networked data. It also generalizes to both graph-based supervised and semi-supervised learning (SSL) [48] and gives rise to a variety of methods for both settings. In particular, RSM is extremely fast, accurate, and flexible with many interchangeable components. RSM learning and inference is fast, space-efficient, and highly scalable for *real-time interactive learning*. In addition, RSM is designed to be intuitive and easy to adapt and encode application constraints. Moreover, RSM has many other attractive properties including its robustness to noisy links, nodes, and attributes, as well as graph data with only a few labeled instances and attributes.

Existing methods perform poorly for data that exhibits low relational autocorrelation, which occurs frequently in practice. In contrast, RSM naturally handles relational data with varying levels of autocorrelation by adjusting a simple hyperparameter. In particular, the key contributions of this work as well as the advantages of RSM over existing methods are as follows:

- A general principled relational learning framework is proposed for large-scale supervised and semi-supervised learning (SSL) in multi-dimensional networks that are noisy, sparsely labeled, and contain only modest or even low levels of

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

MLG KDD 2016 San Francisco, CA USA

© 2016 Copyright held by the owner/author(s).

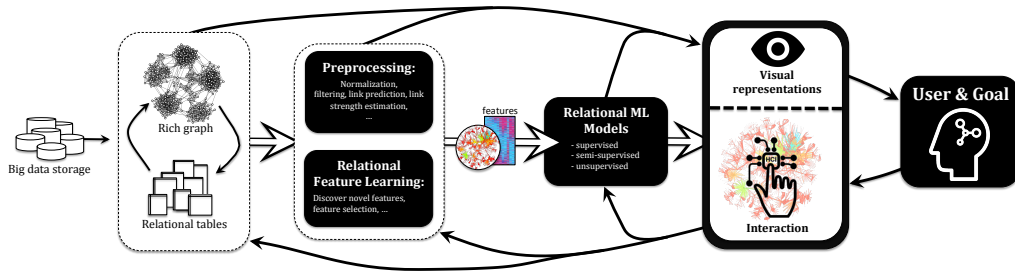


Figure 1: Interactive Relational Machine Learning (iRML) Paradigm. Visual encoding of the results from the various components are indicated via the top arrows, whereas user interactions are represented by the arrows at the bottom.

relational autocorrelation.

- Naturally generalizes for large multi-class problems
- Flexible with many interchangeable components (*e.g.*, similarity function, collective inference)
- Fast learning and inference methods (for applications requiring real-time performance).
- Easily parallelizable for shared- and distributed-memory architectures with speedups that are nearly linear.
- The contribution of the relational information can be learned from the data directly (or quickly tuned by the user).
- Effective for sparsely labeled graph data as well as situations where the labels of the training data are skewed towards one or more class (*e.g.*, represented disproportionately).
- RSM generalizes across the space of potential classifiers including simple approaches that leverage non-relational *iid* data, to those that leverage the graph topology (and possibly a set of attributes), as well as classifiers that leverage both. In fact, these are all special cases that arise under certain conditions (hyper-parameters combinations). Further, a collective inference approach is derived for the space of classifiers that leverage graph topology in some fashion (where G may be given as input or perhaps created, *e.g.* using SSL).
- RSM is robust to networks with modest or even low levels of relational autocorrelation
- Excellent predictive accuracy and performs significantly better than existing approaches with comparable time and space constraints.
- Amenable for streaming and online settings with efficient updates, learning, and inference.

Note that RSM is also useful for heterogeneous networks and leverages both link and node attributes.

2. INTERACTIVE RML

Relational Machine Learning (RML) [17] methods exploit the relational dependencies between nodes to improve predictive performance [24]. However, these approaches often fail in practice due to low relational autocorrelation, noisy links, sparsely labeled graphs, and data representation [40].

2.1 iRML Paradigm

To overcome these problems, the Interactive Relational Machine Learning (*iRML*) paradigm was envisioned [42] in which

users interactively specify relational models and data representation (via transformation techniques for the graph structure and features), as well as perform evaluation, analyze errors, and make adjustments and refinements in a closed-loop [42]. In this work, humans interact with relational learning algorithms by providing input (in the form of labels, similarity/kernel function, hyper-parameters, priors, confidence/uncertainty about particular instances, learning rate, corrections, rankings, probabilities, evaluation) while observing the output (in the form of predictions, uncertainty, feedback, and any useful visual representation of the data).

Desiderata for the *iRML* paradigm are as follows:

- **Immediate visual feedback.** *iRML* methods should be optimized for the way humans learn [1, 46, 5]. Thus, they should provide *immediate* and *continuous visual feedback* upon every interaction (*e.g.*, change of a slider for filtering uncertain or misclassified nodes, selection of a subgraph for modeling, or correcting the class label of a node). Further, interactive queries need to be rapid, incremental and reversible with immediate visual feedback.
- **Flexibility and generality.** Methods should be useful for a wide variety of data, constraints, and learning scenarios. They need to also be robust for learning sparsely labeled graphs, noisy relational data, low and varying levels of relational autocorrelation, and other problems that frequently arise in practice.
- **Effectiveness.** Models must have good predictive quality with low error, variance, and bias.
- **Scalable methods.** Fast time- and space-efficient learning methods capable of interactive rates is an important and key requirement. The requirement of rapid and interactive model updates often dictates trading off accuracy with speed. Thus, network sampling methods may be used to balance speed and accuracy.
- **Accessibility and simplicity.** To be accessible to domain but non-ML experts, *iRML* methods must be carefully designed to be simple, intuitive, and easy-to-use. Whenever possible, assistance and guidance from the system is desired.
- **Principled models.** Another challenge is the design of intuitive learning and inference methods to facilitate interactive reasoning, understanding, and derivation of theoretical behavior and guarantees. This enables quick understanding and refinement by the user, while also providing a means to backtrack if warranted to understand a specific outcome or anomaly.
- **Unified & expressive models.** A unifying family of

relational learning methods that express a large and multifaceted space of relational models. These models must perform well across a variety of different data, characteristics, and assumptions. They must also generalize to a variety of learning settings (e.g., relational active learning, online/incremental learning).

2.2 Visual Representation and Interaction

Visual analytic techniques are also developed that combine a wide variety of visual representations and interactive techniques to exploit human capabilities for seeing, exploring, and understanding large amounts of information. A few of the features/advantages are given below:

- After each graph manipulation or user interaction (e.g., via visual interactions, change of a slider), the *i*RML method is updated immediately with visual feedback after each change. For instance, suppose a user finds that a node is labeled incorrectly via the visual analytic techniques, and simply corrects the error immediately by simply clicking the node and updating its field. Immediately after the correction is made, all features, models, and visual representations that depend on that value are immediately updated on the fly. Hence, after each graph manipulation (or user interaction) such as inserting, filtering, or permanently deleting nodes and links.
- Select or hover over nodes and edges to analyze their class label, uncertainty, estimated class label probabilities, as well as other properties, attributes, relational features, as well as topology features such as graphlet statistics [4, 3].
- Real-time interactive visual graph querying and filtering capabilities, e.g., filter all uncertain nodes above a user-specified threshold, or select all misclassified nodes for further analysis.
- Color and size of nodes and edges can be visually encoded as a function of measures derived from learning and inference such as node (edge) labels, uncertainty, or a learned feature or meta-feature, among other possibilities. Non-relational attributes and network properties derived from the topology may also be used.
- Multiple potentially disconnected subgraphs may be selected for learning, e.g., by brushing over interesting regions of the network visually or by sliders that control filtering/selection, etc.
- Nodes (or edges) can be labeled easily, e.g., double-clicking a node of interest (or set of nodes after selection).

Nearly all visualizations are interactive and support brushing, linking, zooming, panning, tooltips, etc. Multiple visual representations of the graph data are supported, including the multi-level graph properties (e.g., interactive scatter plot matrix, and other statistical plots). Dynamic network analysis and visualization tools to understand the temporal dependencies to better leverage them in learning. Network may also be searched via textual query (e.g., node name). There are many other features including full customization of the visualization (color, size, opacity, background, fonts, etc), text annotation, graph layouts, collision detection, fish eye, and many others. See [42] for more details. Note RSM was implemented in the interactive graph mining and visualization platform proposed in [5].

3. RELATIONAL SIMILARITY MACHINES

Given an attributed graph $G = (V, E, \mathbf{X}, Y, C)$ where V is a set

of $n = |V|$ nodes (e.g., IoT devices in a sensor network), E is a set of $m = |E|$ edges (e.g., communications between devices). We define $\mathbf{X} \in \mathbb{R}^{n \times d}$ as a matrix where rows represent nodes and columns are features. Further, let \mathbf{x}_i be the feature vector for $v_i \in V$. For clarity, we also use $\mathbf{X}_{i:}$ for the i^{th} row of \mathbf{X} and $\mathbf{X}_{:,k}$ to denote the k^{th} column of \mathbf{X} .

DEFINITION 3.1. *Given a graph G , a known set of node labels $Y^\ell = \{y_i | v_i \in V^\ell\}$ for nodes $V^\ell \subset V$, the within-network classification task is to infer Y^u — the set $V^u = V \setminus V^\ell$ of remaining vertices with unknown labels.*

The set of class labels is denoted as C and $k = |C|$ is the number of unique labels. Let $\mathbf{A} = [a_{ij}]$ be the adjacency matrix of G where $A_{ij} = 1$ if there exists $(v_i, v_j) \in E$ and $A_{ij} = 0$ otherwise. Furthermore, \mathbf{A} may encode an edge attribute, that is, $A_{ij} = \Phi_{ij}$.

A general computational framework for RSM is given in Alg. 1. The RSM framework gives rise to a large space of potential relational learning methods due to RSM’s powerful representation and flexibility, e.g., as many of the learning components in Alg. 1 are naturally interchangeable. It is straightforward to learn these from data directly, similar to how hyperparameters of a particular instantiation of RSM are learned via k-fold cross-validation (CV).

3.1 Similarity Functions

We define a few parameterized similarity functions. Note that one may also interpret $\mathbf{X}\mathbf{Z}^\top$ as a similarity matrix with the dot product as similarity measure. Thus, RSM classifies a test point according to a weighted sum of its similarities with the training points.

Radial Basis Function (RBF): Given vectors \mathbf{x}_i and \mathbf{z}_j of length d , the RBF similarity function is:

$$\mathbf{S}(\mathbf{x}_i, \mathbf{z}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{z}_j\|_2^2}{2\sigma^2}\right)$$

where the radius of the RBF function is controlled by choice of σ (i.e., tightness of the similarity measure).

Polynomial functions: Polynomial functions offer another representative similarity measure for vectors of uniform length:

$$\mathbf{S}(\mathbf{X}, \mathbf{Z}) = \|\mathbf{X}, \mathbf{Z}\|^q$$

3.2 Class Prior Estimation

Let $\mathbf{P} = [P_{ik}] = [\mathbf{p}_1^\top \cdots \mathbf{p}_i^\top \cdots]$ where each \mathbf{p}_i is a k -dimensional row vector. Given n training nodes V^ℓ with known labels, informally, we learn a model which is used as an initial estimate for the nodes with unknown class labels, that is, \mathbf{p}_j , $\forall j = 1, \dots, m$ and $v_j \in V^u$. For simplicity, the experiments in Section 5 use the following approach:

- S1** Estimate the class-prior probability $\hat{\mathbb{P}}(y)$ from the training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ as $\hat{\mathbb{P}}(y) = n_y/n$ and set $\mathbf{p}_j^{(0)} = \hat{\mathbb{P}}(y) : \forall j = 1, \dots, m$ where n_y is the number of training nodes in V^ℓ in class y .
- S2** Estimate \mathbf{p}_j , $\forall v_j \in V^u$ from training data $\{(\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in \mathbb{R}^d, y_i \in \{1, 2, \dots, k\}\}_{i=1}^n$ using RSM-iid for simplicity, however, any efficient and accurate approach will suffice.
- S3** Iteratively compute new estimates \mathbf{p}_j , $\forall v_j \in V^u$ using a fast linear-time approach based on belief propagation.

Algorithm 1 Relational Similarity Machines (RSM). A general and flexible relational learning framework based on the notion of maximum similarity.

```

1: Normalize all data (if needed)
2: Generate an ordering of  $V^u$ 
3: Estimate the class priors  $\mathbf{P} = [\dots \mathbf{p}_i \dots]^\top$  where  $\mathbf{p}_i \in \mathbb{R}^k$ 
   is the estimated prior for  $v_i$  (see Section 3.2).
4: Compute graph topology features using  $G = (V, E)$  (e.g.,
   based on k-graphlets where  $k = \{3, 4, \dots\}$  [3], among other
   important and efficient graph properties) and append
   these to  $\mathbf{X}$  and  $\mathbf{Z}$ 
5: Set  $\tau \leftarrow 1$ 
6: repeat                                 $\triangleright$  outer iterations  $\tau = 1, 2, \dots, \tau_{\max}$ 
7:   Compute relational features based on neighbor classes
8:   Compute relational features based on neighbor attributes
9:   Append the relational features from Line 7-8 to the cur-
   rent set of features and renormalize if needed.
10:  for each  $v_i \in V^u$  in order           $\triangleright$  next test instance
11:    Set  $\mathbf{w}_i^R, \mathbf{w}_i^I$  to  $\mathbf{0} = [0 \dots 0] \in \mathbb{R}^k$ 
12:    parallel for each  $v_j \in V^\ell$            $\triangleright$  Supervised
13:      Set  $s_{ij}$  to be  $\mathbf{S}(\mathbf{z}_i, \mathbf{x}_j)$ 
14:      Let  $k$  be the index of the class label  $y_j$  for  $v_j$ 
15:      if  $v_j \in \Gamma_h(v_i)$  then
16:        Update  $w_{ik}^R \leftarrow w_{ik}^R + p_{ik} \cdot s_{ij}$ 
17:      else Update  $w_{ik}^I \leftarrow w_{ik}^I + p_{ik} \cdot s_{ij}$ 
18:    end parallel
19:    parallel for each  $v_j \in V^u$            $\triangleright$  Semi-supervised
20:      Set  $s_{ij}$  to be  $\mathbf{S}(\mathbf{z}_i, \mathbf{z}_j)$ 
21:      for each class  $k \in \mathcal{C}$  do
22:        if  $v_j \in \Gamma_h(v_i)$  then
23:           $w_{ik}^R \leftarrow w_{ik}^R + p_{ik} p_{jk} \cdot s_{ij}$ 
24:        else  $w_{ik}^I \leftarrow w_{ik}^I + p_{ik} p_{jk} \cdot s_{ij}$ 
25:      end parallel
26:      Normalize  $\mathbf{w}_i^R$  and  $\mathbf{w}_i^I$  s.t.  $\sum w_{ik}^R = \sum w_{ik}^I = 1$ 
27:    end for
28:    Update  $\mathbf{p}_i$  via Eq. (1),  $\forall v_i \in V^u$ 
29:    Compute confidence  $\mathbf{c}$  and assign predictions for top- $k$ 
    most confident
30:    Include  $\mathbf{P}, \mathbf{W}^R, \mathbf{W}^I$ , and/or  $\mathbf{c}$  as features (if warranted)
31:    Set  $\tau \leftarrow \tau + 1$  and renormalize data if needed
32:  until stopping criterion is reached or  $\tau > \tau_{\max}$ 
33:  parallel for each  $v_i \in V^u$  do
34:     $y_i \leftarrow \arg \max_{k \in \mathcal{C}} P_{ik}$            $\triangleright$   $k$  is the class label

```

This corresponds to a simple collective learning approach that essentially meshes the current probability distribution vector \mathbf{p}_j for $v_j \in V^u$ with the estimated probability vectors from the neighbors of v_j denoted $\Gamma(v_j)$ s.t. probability vectors for the neighbors $\bar{\mathbf{p}}_i = \mathbf{e}^\top \mathbf{P}_s / r$ where $\mathbf{P}_s = [\dots \mathbf{p}_i^\top \dots] \in \mathbb{R}^{r \times k}$, \mathbf{e}^\top is a vector of all ones, and $r = |\Gamma_h(v_j)|$. Thus, $\bar{\mathbf{p}}_i = \mathbf{e}^\top \mathbf{P}_s / r$ is the “centroid” of estimates from v_j ’s neighborhood, and is used to compute $\mathbf{p}_j^{(\tau+1)}$ along with the previous local estimate of v_j . This is repeated for all vertices across a number of iterations until the estimates converge.

This approach differs from existing work that uses a “local model” to obtain an estimate for each node [30, 45, 25].

However, a key advantage of RSM is its flexibility, and as such, other semi-supervised estimation methods are also applicable and may lead to further improvements [14].

3.3 Update Equations

A key contribution of this work is the proposed relational learning framework that learns from iid and *relational data* that is both *labeled* and *unlabeled*. To the best of our knowledge, this work is the first to leverage and investigate using both *labeled* and *unlabeled* relational and *iid* data in learning. This leads to four different learning scenarios: (a) labeled neighbors, (b) unlabeled neighbors, and (c) labeled and (d) unlabeled nodes that are not neighbors. Thus, RSM exploits the estimates from labeled and unlabeled *iid* and relational data simultaneously (and in a collective fashion) to improve predictive performance. The update equation used in this paper is simply:

$$\mathbf{p}_i^{(\tau+1)} = \underbrace{\alpha \mathbf{w}_i^R}_{\text{Neighbor}} + \underbrace{(1 - \alpha) \mathbf{w}_i^I}_{\text{Non-neighbor}} + \underbrace{\omega \mathbf{p}_i^{(\tau)}}_{\text{previous}} \quad (1)$$

where α is a hyperparameter which satisfies $0 \leq \alpha \leq 1$, and $\mathbf{w}_i \in \mathbb{R}^k$ is the non-negative vector with $W_{ik} \geq 0$. Further, ω determines the influence given to the previous estimate $\mathbf{p}_i^{(\tau)}$ and τ is the iteration.

3.4 Semi-Supervised Learning

Semi-supervised learning (SSL) lies between unsupervised and supervised learning as it uses both unlabeled and labeled data for deciding the class of a node. Partially labeled data is often found in practice, due to it being costly and time-consuming to obtain labels. Note RSM is also flexible for the case where there are very few labeled instances and supports a number of intuitive and principled semi-supervised learning strategies. For instance, at each iteration (global update), we estimate uncertainty (for each unlabeled instance), and predict the labels of the top- k nodes for which we are most certain (min uncertainty). Notice that Line 19-25 in Alg. 1 leverages unlabeled nodes directly. Future work will investigate different update equations that leverage these weights differently. For instance, instead of combining these weights into \mathbf{w}_i^R and \mathbf{w}_i^I , we can define different weight vectors to maintain these weights independently of those from Line 12-18.

3.5 Classification

We now reinterpret the problem for classification, but it may also be used for regression in a straightforward fashion. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix where the rows represent training instances and the columns represent features. Further, let \mathbf{Z} be a matrix of test instances, then the class of a single test instance \mathbf{z}_j is predicted as follows. First, the similarity of \mathbf{z}_j with respect to each training example in \mathbf{X} is computed. For instance, suppose \mathbf{x}_i belongs to class $k \in \mathcal{C}$, then $\mathbf{S}(\mathbf{x}_i, \mathbf{z}_j)$ is added to the k th element w_k of the weight vector \mathbf{w} . The similarity of the instances in \mathbf{X} of class k with respect to the test object \mathbf{z}_j is formalized as,

$$w_k = \sum_{\mathbf{x}_i \in \mathcal{X}_k} \mathbf{S}(\mathbf{x}_i, \mathbf{z}_j) \quad (2)$$

where \mathcal{X}_k is the set of training objects from \mathbf{X} of class k . Thus

Table 1: Experiments comparing RSM and the adapted iRML-based WVRN. These results demonstrate the effectiveness of RSM. In all cases, RSM outperforms WVRN, and the difference in accuracy is significant. Results that are significant are bolded. Moreover, RSM naturally supports interactive RML, while also significantly more flexible than WVRN (*e.g.*, naturally supports graph-based semi-supervised learning).

Relational data	Hyperparameters				Accuracy	
	$ C $	σ	α	ω	RSM	WVRN
aff-polbooks	3	0.3	0.7	0.6	84.73	74.41
bio-Gene	2	0.2	0.5	0.5	79.65	72.41
bio-Enzymes349	2	0.1	0.25	0	77.81	62.50
aff-musicGenre	8	0.1	0.75	0.5	74.90	60.59

supported by RSM. In particular, RSM is shown to be fast, parallel, space-efficient, amenable to streaming and dynamic queries/updates, and most importantly, naturally supports real-time interactive learning and inference, and provides rapid immediate (and visual) feedback to the user at real-time interactive response times. This provides context and allows for quick and intuitive understanding of the intermediate results. Upon each change (an additional node/edge is selected visually by the user, ...), a dynamic graph query is issued, and all models, graph properties and statistics are immediately updated and displayed to the user.

We also investigated relational learning for sparsely labeled graph data. For instance, we used *ca-cora* with 10% labels,

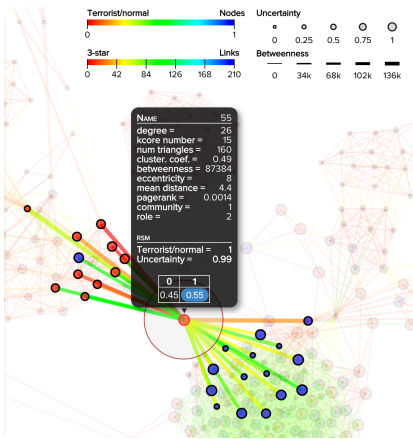


Figure 3: Understanding model uncertainty via interactive visual exploration. Nodes are colored by class label and weighted (sized) by the uncertainty. Above RSM enables the user to visually explore the prediction confidence of local nodes, as well as the uncertainty, and other model and graph features via simple and intuitive local dynamic queries (initiated by mousing over a node, or selecting a group of nodes).

and only used graph features (that is, no initial attributes were used). The experiments suggest that other RSM variants do not work as well due to the high homophily present in *ca-cora*. On the other hand, WVRN is known to work extremely well in data with precisely these characteristics. Despite this fact, the graph-based RSM variant that leverages the k -hop neighborhood and the neighbors (graph) features outperforms it significantly. In particular, it achieves 84.44% accuracy using only the similarity between the k -hop neighbors, compared to 78.73% accuracy given by WVRN. This demonstrates the generality and flexibility of the proposed family of RSM methods. While existing relational learning methods have been designed to work for relational data with high homophily, the experiments demonstrated that RSM is able to model data with varying levels of homophily and still significantly outperform others that are more specialized and less flexible to handle such data characteristics.

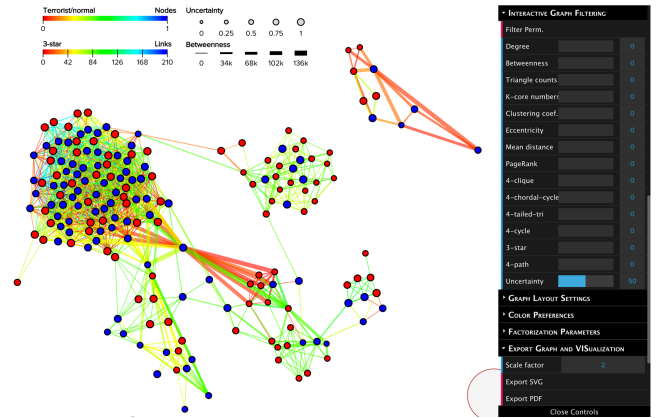


Figure 4: Interactive filtering via uncertainty. All nodes and edges with high confidence are filtered (≤ 0.50), leaving only those that are most uncertain. Relative entropy is used to measure the uncertainty of the learned probability distribution. Thus uncertainty is bounded between 0 and 1.

As expected, we also find that the weight or contribution assigned to the various types of relational dependencies may significantly impact the predictive quality of the model. In particular, the classification accuracy of the model is significantly different as the amount of relational information incorporated into the model is varied. Many results were removed for brevity.

Recent work by McDowell *et al.* found that neighbor attributes can improve classification accuracy. In addition to these types of features, we also investigate using graph features such as the frequency of 4-cliques centered on a vertex (edge) for improving relational learning tasks. For instance, motifs of size $k=4$ and greater are computed using the recent graphlet decomposition algorithm proposed in [3]. We also investigated other graph features including community label, betweenness, eccentricity, mean distance, k -core numbers, color class, PageRank, triangle counts, clustering coef., degree, roles. A key finding of these experiments is that the structural features alone may significantly improve predictive performance. This result implies that the connectivity of the graph gives rise to graph motif patterns that are correlated with the different node labels. For instance, suppose a web page representing a

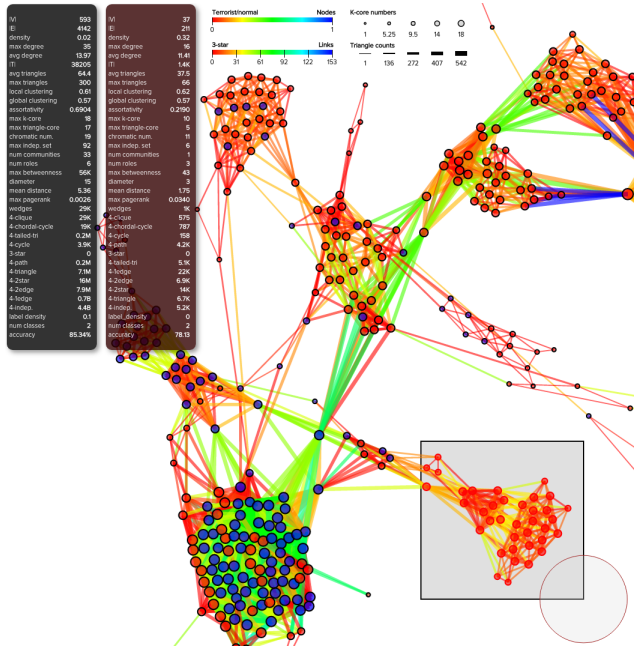


Figure 5: The proposed relational learning framework is not only fast and accurate, but also flexible. This gives the user the ability to interactively learn models in real-time and adjust them accordingly, as well as tune parameters, explore/construct features, among many other possibilities. In the screenshot above (from soc-terror), the user interactively learns a global model (see the leftmost black side panel for stats and accuracy), and then selects a subgraph H by visually selecting nodes and edges with a simple and intuitive drag-of-the-mouse/gesture. From this selected subgraph (in real-time), a local RSM model is learned for H and stats and accuracy (using same k-fold CV) are reported in the red side panel on the left side. Node color represents the class label (terrorist/normal), whereas the link color encodes the number of 3-star graphlet patterns centered at each edge (both can be adapted in real-time by the user). The weight (or strength) of the nodes and links represent the local max k-core number and triangle counts, respectively.

node in a web graph forms a large number of 3-star motifs, then this node is likely to be a malicious page (spammer).

An overview of the *iRML* system for RSM is shown in Figure 2. In that example, we first interactively learn a model, then select the misclassified nodes for further analysis. The global statistics of the selected subgraph are shown in the rightmost panel. Node color represents the model’s uncertainty using an entropy-based measure, whereas the size of the node indicates whether it was correctly classified or not. In Figure 2, misclassified nodes are given a larger size so that they can easily be identified for further exploration. Uncertainty (and the estimated class prob. distribution, statistics, etc.) of a node or set of nodes may also be displayed by selecting or mousing over those nodes of interest (Fig. 3). RSM also supports interactive real-time visual graph filtering (e.g., remove all uncertain nodes above a user-specified threshold, e.g., see

Fig. 4 and Fig. 5). In addition, all visualizations are interactive and support brushing, linking, zooming, panning, tooltips, etc (Fig. 5). Efficient update rules are also derived to avoid relearning the model (after each user interaction/visual query). For example, after the deletion (or insertion) of a node, we can update the global relational model via a fast localized update. These local updates enable real-time exploration capabilities by leveraging fast exact or approximate solutions (e.g., to support real-time interactive queries seamlessly, see Fig. 5).

Many of the components in RSM may be explored using interactive visualization and analytic techniques, including the attribute to predict, initial features to use (non-relational and graph-based features), local model for estimation, kernel function (RBF, linear, polynomial, etc.), hyper-parameters (for selected kernel), node- and feature-wise normalization scheme (L1, min-max, etc.), as well as whether to use semi-supervised learning (SSL), and meta-features (based on current estimates). For example, see Fig. 2 (right panel). Interactive link prediction methods and many other important learning components are also included in our *iRML* system and can be leveraged directly by RSM for improving learning and inference.

Complexity: Given n training instances, m test instances, and d features. For learning and inference, RSM takes $\mathcal{O}(nd)$ time on a single test instance (in the stream), and therefore $\mathcal{O}(nmd)$ for all m test examples. The complexity for both sparse and dense training feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ is given below. Consider the case where \mathbf{X} is sparse and stored as a sparse matrix using compressed sparse column/row format. Let $\Omega_{\mathbf{X}}$ denote the number of nonzeros in \mathbf{X} , then the cost of a single test example is $\mathcal{O}(|\Omega_{\mathbf{X}}|)$ linear in the number of nonzeros in \mathbf{X} . Further, assuming p processing units, then $\mathcal{O}(|\Omega_{\mathbf{X}}|/p)$, and hence is very scalable for real-time systems. Now suppose \mathbf{X} is a dense matrix. Given a dense training set $\mathbf{X} \in \mathbb{R}^{n \times d}$ (few zeros), it takes $\mathcal{O}(md)$ time for each test object, and $\mathcal{O}(\frac{md}{p})$ assuming p processing units. Note that the cost may also be significantly reduced by selecting a representative set of training objects using one of the previously proposed methods.

Finally, we discuss the space complexity of the proposed relational learning framework. Given a single test node to predict, RSM takes only $\mathcal{O}(k)$ space where k is the number of classes. This is of course in addition to the graph and features. For parallel learning and inference using p workers, the lock-free version of RSM takes $\mathcal{O}(pk)$ space. To avoid locks, each worker maintains a k -dimensional vector of similarity scores, which are then combined upon completion. However, if locks are used then the initial cost of $\mathcal{O}(k)$ holds, but at the expense of runtime. Now, suppose collective inference is used, then RSM takes only $\mathcal{O}(nk)$ space where n is the number of test nodes with unknown class labels.

6. CONCLUSION

This paper proposes a graph-based learning framework called *relational similarity machines* (RSM) for (semi-supervised) learning in large and noisy networks with arbitrary relational autocorrelation. Most importantly, RSM is a fast, accurate, and flexible *relational learning framework* based on the notion of maximizing similarity. Our approach naturally handles both binary and multi-class classification problems in large attributed and possibly streaming networks. Despite the importance of relational learning, existing methods are unable

to handle relational data with varying levels of relational autocorrelation, and therefore limited in their ability (and utility) in many situations and data characteristics. In contrast, RSM naturally handles network data with arbitrary levels of autocorrelation by adjusting a simple hyperparameter. Furthermore, RSM has many other advantages over existing methods, including: (a) it is space- and time-efficient, (b) easily parallelizable and thus scalable for large relational data, (c) naturally allows for varying levels of relational autocorrelation, (d) flexible as many components are interchangeable (and can be learned from the data directly or fine-tuned by an expert), (e) amenable to (attributed) graph streams, and (f) has excellent accuracy. In addition, both node and link classification problems are easily handled in RSM. Moreover, RSM is able to leverage node and link attributes simultaneously as well as heterogeneous networks with multiple node and edge types.

7. REFERENCES

- [1] C. Ahlberg, C. Williamson, and B. Shneiderman. Dynamic queries for information exploration: An implementation and evaluation. In *Proc. of SIGCHI*, pages 619–626, 1992.
- [2] N. K. Ahmed, J. Neville, and R. Kompella. Space-efficient sampling from social activity streams. In *SIGKDD BigMine*, pages 53–60, 2012.
- [3] N. K. Ahmed, J. Neville, R. A. Rossi, and N. Duffield. Efficient graphlet counting for large networks. In *ICDM*, pages 1–10, 2015.
- [4] N. K. Ahmed, J. Neville, R. A. Rossi, N. Duffield, and T. L. Willke. Graphlet decomposition: Framework, algorithms, and applications. *KAIS*, pages 1–32, 2016.
- [5] N. K. Ahmed and R. A. Rossi. Interactive visual graph analytics on the web. In *ICWSM*, pages 566–569, 2015.
- [6] N. Andrienko and G. Andrienko. A visual analytics framework for spatio-temporal analysis and modelling. *DMKD*, 27(1):55–83, 2013.
- [7] C. R. Aragon, S. S. Poon, G. S. Aldering, R. C. Thomas, and R. Quimby. Using visual analytics to develop situation awareness in astrophysics. *Info. Vis.*, 8(1):30–41, 2009.
- [8] R. Arias-Hernández, J. Dill, B. Fisher, and T. M. Green. Visual analytics and human-computer interaction. *Interact.*, 18(1):51–55, 2011.
- [9] D. S. Bassett and E. Bullmore. Small-world brain networks. *neurosci.*, page 512, 2006.
- [10] M. Bilgic, L. Mihalkova, and L. Getoor. Active learning for networked data. In *ICML*, 2010.
- [11] M. Bilgic, G. M. Namata, and L. Getoor. Combining collective classification and link prediction. In *ICDM Workshops*, 2007.
- [12] J. Choo, H. Lee, J. Kihm, and H. Park. ivisclassifier: An interactive visual analytics system for classification based on supervised dimension reduction. In *VAST*, 2010.
- [13] L. De Raedt and K. Kersting. *Probabilistic Inductive Logic Programming*. Springer, 2008.
- [14] M. C. Du Plessis and M. Sugiyama. Semi-supervised learning of class balance under class-prior change by distribution matching. *Neural Networks*, 50:110–119, 2014.
- [15] J. A. Fails and D. R. Olsen Jr. Interactive machine learning. In *IUI*, pages 39–45, 2003.
- [16] N. Friedman, L. Getoor, D. Koller, and A. Pfeffer. Learning probabilistic relational models. In *IJCAI*, pages 1300–1309. Springer-Verlag, 1999.
- [17] L. Getoor and B. Taskar, editors. *Intro. to SRL*. MIT Press, 2007.
- [18] D. H. Jeong, C. Ziemkiewicz, B. Fisher, W. Ribarsky, and R. Chang. iPCA: An Interactive System for PCA-based Visual Analytics. In *C. Grap.*, volume 28, pages 767–774, 2009.
- [19] A. Kapoor, B. Lee, D. S. Tan, and E. Horvitz. Performance and preferences: Interactive refinement of machine learning procedures. In *AAAI*, 2012.
- [20] D. A. Keim, J. Kohlhammer, G. Ellis, and F. Mansmann. *Mastering the information age-solving problems with visual analytics*. 2010.
- [21] J. Kielman, J. Thomas, and R. May. Foundations and frontiers in visual analytics. *Info. Vis.*, 8(4):239, 2009.
- [22] S. Kim, Y. Jang, A. Mellema, D. S. Ebert, and T. Collins. Visual analytics on mobile devices for emergency response. In *VAST*, pages 35–42, 2007.
- [23] A. Kimmig, S. Bach, M. Broecheler, B. Huang, and L. Getoor. A short introduction to probabilistic soft logic. In *NIPS Prob. Prog. Workshop*, 2012.
- [24] S. Macskassy and F. Provost. A simple relational classifier. In *SIGKDD MRDM*, pages 64–76, 2003.
- [25] S. A. Macskassy and F. Provost. Classification in networked data: A toolkit and a univariate case study. *JMLR*, 8:935–983, 2007.
- [26] P. Mahadevan, D. Krioukov, M. Fomenkov, X. Dimitropoulos, A. Vahdat, et al. The Internet AS-level Topology: three data sources and one definitive metric. *SIGCOMM*, 36(1):17–26, 2006.
- [27] A. Malik, R. Maciejewski, T. F. Collins, and D. S. Ebert. Visual analytics law enforcement toolkit. In *HST*, pages 222–228, 2010.
- [28] A. Malik, R. Maciejewski, B. Maule, and D. S. Ebert. A visual analytics process for maritime resource allocation and risk assessment. In *VAST*, 2011.
- [29] L. McDowell, K. Gupta, and D. Aha. Meta-prediction for coll. classif. In *FLAIRS*, 2010.
- [30] L. K. McDowell and D. W. Aha. Labels or attributes? Rethinking the neighbors for collective classification in sparsely-labeled networks. In *CIKM*, pages 847–852, 2013.
- [31] L. K. McDowell, K. M. Gupta, and D. W. Aha. Cautious collective classification. *JMLR*, 10:2777–2836, 2009.
- [32] A. McGovern, N. Collier, I. Matthew Gagne, D. Brown, and A. Rodger. Spatiotemporal relational probability trees: An introduction. In *ICDM*, pages 935–940, 2008.
- [33] M. McPherson, L. Smith-Lovin, and J. M. Cook. Birds of a feather: Homophily in social networks. *Annual Review of Sociology*, 27(1):415–444, 2001.
- [34] J. Neville and D. Jensen. Relational dependency networks. *Journal of Machine Learning Research*, under review.
- [35] J. Neville, D. Jensen, L. Friedland, and M. Hay. Learning relational probability trees. In *SIGKDD*, 2003.
- [36] J. Neville, D. Jensen, and B. Gallagher. Simple estimators for relational Bayesian classifiers. In *ICDM*, 2003.
- [37] D. Oglic, D. Paurat, and T. Gärtner. Interactive knowledge-based kernel pca. In *PKDD*. 2014.
- [38] G. Robertson, D. Ebert, S. Eick, D. Keim, and K. Joy. Scale and complexity in visual analytics. *Info. Vis.*, 8(4):247–253, 2009.
- [39] R. A. Rossi and N. K. Ahmed. An interactive data repository with visual analytics. *SIGKDD Explor.*, 17(2):37–41, 2016.
- [40] R. A. Rossi, L. K. McDowell, D. W. Aha, and J. Neville. Transforming graph data for statistical relational learning. *JAIR*, 45:363–441, 2012.
- [41] R. A. Rossi and J. Neville. Modeling the evolution of discussion topics and communication to improve relational classification. In *SOMA*, pages 89–97, 2010.
- [42] R. A. Rossi and R. Zhou. Toward interactive relational learning. In *AAAI*, pages 4383–4384, 2016.
- [43] A. Savikhin, H. C. Lam, B. Fisher, and D. S. Ebert. An experimental study of financial portfolio selection with visual analytics for decision support. In *HICSS*, 2011.
- [44] M. Sedlmair, P. Isenberg, D. Baur, M. Mauerer, C. Pigorsch, and A. Butz. Cardiogram: visual analytics for automotive engineers. In *SIGCHI*, pages 1727–1736, 2011.
- [45] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad. Collective classification in network data. *AI Magazine*, 29(3):93, 2008.
- [46] J. J. Thomas and K. A. Cook. *Illuminating the Path: the research and dev. agenda for visual analytics*. IEEE, 2005.
- [47] X. Wang, R. Garnett, and J. Schneider. Active search on graphs. In *SIGKDD*, pages 731–738, 2013.
- [48] X. Zhu and A. B. Goldberg. Intro. to Semi-Supervised Learning. *Synthesis lectures on AI and ML*, 3(1), 2009.