

---

# Bridging Weighted Rules and Graph Random Walks for Statistical Relational Models

Seyed Mehran Kazemi<sup>1,\*</sup> and David Poole<sup>1</sup>

<sup>1</sup>Laboratory of Computational Intelligence, Computer Science Department,  
University of British Columbia, Vancouver, BC, Canada

Correspondence\*:  
Seyed Mehran Kazemi  
smkazemi@cs.ubc.ca

## 2 ABSTRACT

3 The aim of statistical relational learning is to learn statistical models from relational or graph-  
4 structured data. Three main statistical relational learning paradigms include weighted rule  
5 learning, random walks on graphs, and tensor factorization. These paradigms have been mostly  
6 developed and studied in isolation for many years, with few works attempting at understanding  
7 the relationship among them or combining them. In this paper, we study the relationship between  
8 the path ranking algorithm (PRA), one of the most well-known relational learning methods in  
9 the graph random walk paradigm, and relational logistic regression (RLR), one of the recent  
10 developments in weighted rule learning. We provide a simple way to normalize relations and  
11 prove that relational logistic regression using normalized relations generalizes the path ranking  
12 algorithm. This result provides a better understanding of relational learning, especially for the  
13 weighted rule learning and graph random walk paradigms. It opens up the possibility of using the  
14 more flexible RLR rules within PRA models and even generalizing both by including normalized  
15 and unnormalized relations in the same model.

16 **Keywords:** Statistical Relational Artificial Intelligence, Relational Learning, Weighted Rule Learning, Graph Random Walk, Relational  
17 Logistic Regression, Path Ranking Algorithm

## 1 INTRODUCTION

18 Traditional machine learning algorithms learn mappings from a feature vector indicating categorical and  
19 numerical features to an output prediction of some form. Statistical relational learning (Getoor and Taskar,  
20 2007), or statistical relational AI (StarAI) (De Raedt et al., 2016), aims at probabilistic reasoning and  
21 learning when there are (possibly various types of) relationships among the objects. The relational models  
22 developed in StarAI community have been successfully applied to several applications such as knowledge  
23 graph completion (Lao et al., 2011; Nickel et al., 2012; Bordes et al., 2013; Pujara et al., 2013; Trouillon  
24 et al., 2016), entity resolution (Singla and Domingos, 2006; Bhattacharya and Getoor, 2007; Pujara and  
25 Getoor, 2016; Fatemi, 2017), tasks in scientific literature (Lao and Cohen, 2010b), stance classification  
26 (Sridhar et al., 2015; Ebrahimi et al., 2016), question answering (Khot et al., 2015; Dries et al., 2017), etc.

27 During the past decade and more, three paradigms of statistical relational models have appeared. The  
28 first paradigm is the weighted rule learning where first-order rules are learned from data and a weight is  
29 assigned to each rule indicating a score for the rule. The main difference among these models is in the types

30 of rules they allow and their interpretation of the weights. The models in this paradigm include Problog  
31 (De Raedt et al., 2007), Markov logic (Domingos et al., 2008), probabilistic interaction logic (Hommersom  
32 and Lucas, 2011), probabilistic soft logic (Kimmig et al., 2012), and relational logistic regression (Kazemi  
33 et al., 2014). Recent

34 The second paradigm is the random walk on graphs, where several random walks are performed on  
35 a graph each starting at a random node and probabilistically transitioning to neighbouring nodes. The  
36 probability of each node being the answer to a query is proportional to the probability of the random walks  
37 ending up at that node. The main difference among these models is in the way they walk on the graph  
38 and how they interpret obtained results from the walks. Examples of relational learning algorithms based  
39 on random walk on graphs include PageRank (Page et al., 1999), FactRank (Jain and Pantel, 2010), path  
40 ranking algorithm (Lao and Cohen, 2010b; Lao et al., 2011), and HeteRec (Yu et al., 2014).

41 The third paradigm is the tensor factorization paradigm, where for each object and relation an embedding  
42 is learned. The probability of two objects participating in a relation is a simple function of the objects'  
43 and relation's embeddings (e.g., the sum of the element-wise product of the three embeddings). The main  
44 difference among these models is in the type of embeddings and the function they use. Examples of models  
45 in this paradigm include YAGO (Nickel et al., 2012), TransE (Bordes et al., 2013), and ComplEx (Trouillon  
46 et al., 2016).

47 The models in each paradigm have their own advantages and disadvantages. Kimmig et al. (2015)  
48 survey the models based on weighted rule learning. Nickel et al. (2016) survey models in all paradigms  
49 for knowledge graph completion. Kazemi et al. (2017) compare several models in these paradigms for  
50 relational aggregation. None of these surveys, however, aims at understanding the relationship among these  
51 paradigms. In fact, these paradigms have been mostly developed and studied in isolation with few works  
52 aiming at understanding the relationship among them or combining them (Riedel et al., 2013; Nickel et al.,  
53 2014; Lin et al., 2015).

54 With several relational paradigms/models developed during the past decade and more, understanding  
55 the relationship among them and pruning the ones that either do not work well or are subsets of the other  
56 models is crucial. In this paper, we study the relationship between two relational learning paradigms: graph  
57 random walk and weighted rule learning. In particular, we study the relationship among path ranking  
58 algorithm (PRA) (Lao and Cohen, 2010b) and relational logistic regression (RLR) (Kazemi et al., 2014).  
59 The former is one of the most well-known relational learning tools in graph random walk paradigm, and  
60 the latter is one of the recent developments in weighted rule learning paradigm. By imposing restrictions  
61 on the rules that can be included in models, we identify a subset of RLR models that we call RC-RLR.  
62 Then we provide a simple way to normalize relations and prove that PRA models correspond to RC-RLR  
63 models using normalized relations. Other strategies for walking randomly on the graph (e.g., data-driven  
64 path finding (Lao et al., 2011)) can then be viewed as structure learning methods for RC-RLR. Our result  
65 can be extended to several other weighted rule learning and graph random walk models.

66 The relationship between weighted rules and graph random walks has not been discovered before. For  
67 instance, Nickel et al. (2016) describe them as two separate classes of models for learning from relational  
68 data in their survey. Lao et al. (2011) compare their instance of PRA to a model based on weighted rules  
69 empirically, reporting their PRA model outperforms the weighted rule model, but not realizing that their  
70 PRA model could be a subset of the weighted rule model if they had normalized the relations.

71 Our result is beneficial for both graph random walk and weighted rule learning paradigms, as well as for  
 72 researchers working on theory and applications of statistical relational learning. Below is a list of potential  
 73 benefits our result provides:

- 74 • It provides a clearer intuition and understanding on two relational learning paradigms thus facilitating  
 75 further improvements of both.
- 76 • It opens up the possibility of using the more flexible RLR rules within PRA models.
- 77 • It opens up the possibility of generalizing both PRA and RLR models by using normalized and  
 78 unnormalized relations in the same model.
- 79 • It sheds light on the shortcomings of graph random walk algorithms and points out potential ways to  
 80 improve them.
- 81 • One of the claimed advantages of models based on weighted rule learning compared to other relational  
 82 models is that they can be easily explained to a broad range of people (Nickel et al., 2016). Our result  
 83 improves the explainability of models learned through graph random walk, by providing a weighted  
 84 rule interpretation for them.
- 85 • It identifies a sub-class of weighted rules that can be evaluated efficiently and have a high modelling  
 86 power as they have been successfully applied to several applications. The evaluation of these weighted  
 87 rules can be even further improved using sampling techniques developed within graph random walk  
 88 community (e.g., see Fogaras et al. (2005); Lao and Cohen (2010a); Lao et al. (2011)). Several structure  
 89 learning algorithms (corresponding to random walk strategies) have been already developed for this  
 90 sub-class.
- 91 • It facilitates leveraging new insights and techniques developed within each paradigm (e.g., weighted  
 92 rule models that leverage deep learning techniques (Šourek et al., 2015; Kazemi and Poole, 2018), or  
 93 reinforcement learning based approaches to graph walk (Das et al., 2017)) to the other paradigm.
- 94 • For those interested in the applications of relation learning, our result facilitates decision making on  
 95 selecting the paradigm or the relational model to be used in their application.

## 2 BACKGROUND AND NOTATIONS

96 In this section, first we define some basic terminology. Then we introduce a running example which will  
 97 be used throughout the paper. Then we describe relational logistic regression and path ranking algorithm  
 98 for relational learning. While semantically identical, our descriptions of these two models may be slightly  
 99 different from the descriptions in the original articles as we aim at describing the two algorithms in a way  
 100 that simplifies our proofs.

### 101 2.1 Terminologies

102 Throughout the paper, we assume True is represented by 1 and False is represented by 0.

103 A **population** is a finite set of objects (or individuals). A **logical variable (logvar)** is typed with a  
 104 population. We represent logvars with lower-case letters. The population associated with a logvar  $x$  is  $\Delta_x$ .  
 105 The cardinality of  $\Delta_x$  is  $|\Delta_x|$ . For every object, we assume there exists a unique *constant* denoting that  
 106 object. A lower-case letter in bold represents a tuple of logvars and an upper-case letter in bold represents  
 107 a tuple of constants. An **atom** is of the form  $V(t_1, \dots, t_k)$  where  $V$  is a functor and each  $t_i$  is a logvar or  
 108 a constant. When  $range(V) \in \{0, 1\}$ ,  $V$  is a predicate. A **unary** atom contains exactly one logvar and a  
 109 **binary** atom contains exactly two logvars. We write a **substitution** as  $\theta = \{x_1, \dots, x_k\} / \{t_1, \dots, t_k\}$

	Paper <sub>1</sub>	Paper <sub>2</sub>	Paper <sub>3</sub>	Paper <sub>4</sub>	Paper <sub>5</sub>	Paper <sub>6</sub>
Paper <sub>1</sub>			1			
Paper <sub>2</sub>						
Paper <sub>3</sub>				1	1	
Paper <sub>4</sub>			1		1	1
Paper <sub>5</sub>		1	1	1		1
Paper <sub>6</sub>	1	1				

(a)

	Paper <sub>1</sub>	Paper <sub>2</sub>	Paper <sub>3</sub>	Paper <sub>4</sub>	Paper <sub>5</sub>	Paper <sub>6</sub>
Paper <sub>1</sub>			1			
Paper <sub>2</sub>						
Paper <sub>3</sub>				$\frac{1}{2}$	$\frac{1}{2}$	
Paper <sub>4</sub>			$\frac{1}{3}$		$\frac{1}{3}$	$\frac{1}{3}$
Paper <sub>5</sub>		$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$		$\frac{1}{4}$
Paper <sub>6</sub>	$\frac{1}{2}$	$\frac{1}{2}$				

(b)

**Figure 1.** (a) A relation showing citations among papers (papers on the  $Y$  axis cite papers on the  $X$  axis), (b) The relation in part (a) after row-wise count normalization.

110 where each  $x_i$  is a different logvar and each  $t_i$  is a logvar or a constant in  $\Delta_{x_i}$ . A **grounding** of an atom  
 111  $V(x_1, \dots, x_k)$  is a substitution  $\theta = \{\langle x_1, \dots, x_k \rangle / \langle X_1, \dots, X_k \rangle\}$  mapping each of its logvars  $x_i$  to an  
 112 object in  $\Delta_{x_i}$ . Given a set  $\mathcal{A}$  of atoms, we denote by  $\mathcal{G}(\mathcal{A})$  the set of all possible groundings for the atoms  
 113 in  $\mathcal{A}$ . A **value assignment** for a set of groundings  $\mathcal{G}(\mathcal{A})$  maps each grounding  $V(\mathbf{X}) \in \mathcal{G}(\mathcal{A})$  to a value in  
 114  $range(V)$ .

115 A **literal** is an atom or its negation. A **formula**  $\varphi$  is a literal, a disjunction  $\varphi_1 \vee \varphi_2$  of formulae or a  
 116 conjunction  $\varphi_1 \wedge \varphi_2$  of formulae. Our formulae correspond to open formulae in negation normal form in  
 117 logic. An **instance** of a formula  $\varphi$  is obtained by replacing each logvar  $x$  in  $\varphi$  by one of the objects in  $\Delta_x$ .  
 118 Applying a **substitution**  $\theta = \{\langle x_1, \dots, x_k \rangle / \langle t_1, \dots, t_k \rangle\}$  on a formula  $\varphi$  (written as  $\varphi\theta$ ) replaces each  $x_i$   
 119 in  $\varphi$  with  $t_i$ . A *weighted formula (WF)* is a pair  $\langle w, \varphi \rangle$  where  $w$  is a weight and  $\varphi$  is a formula.

120 A binary predicate  $S(x, y)$  can be viewed as a function whose domain is  $\Delta_x$  and whose range is  $2^{\Delta_y}$ :  
 121 each  $X \in \Delta_x$  is mapped to  $\{Y : S(X, Y)\}$ . Following Lao and Cohen (2010b), we consider  $S^{-1}$  as the  
 122 inverse of  $S$  whose domain is  $\Delta_y$  and whose range is  $2^{\Delta_x}$ , such that  $S^{-1}(x, y)$  holds iff  $S(y, x)$  holds. A  
 123 **path relation**  $\mathcal{PR}$  is of the form  $x_0 \xrightarrow{R_1} x_1 \xrightarrow{R_2} \dots \xrightarrow{R_l} x_l$  where  $R_1, R_2, \dots, R_l$  are predicates,  $x_0, \dots, x_l$   
 124 are different logvars,  $domain(R_i) = \Delta_{x_{i-1}}$  and  $range(R_i) = \Delta_{x_i}$ . We define  $domain(\mathcal{PR}) = \Delta_{x_0}$   
 125 and  $range(\mathcal{PR}) = \Delta_{x_l}$ . Applying a substitution  $\theta = \{\langle x_1, \dots, x_k \rangle / \langle t_1, \dots, t_k \rangle\}$  on a path relation  $\mathcal{PR}$   
 126 (written as  $\mathcal{PR}\theta$ ) replaces each  $x_i$  in  $\mathcal{PR}$  with  $t_i$ . A **weighted path relation (WPR)** is a pair  $\langle w, \mathcal{PR} \rangle$   
 127 where  $w$  is a weight and  $\mathcal{PR}$  is a path relation.

## 128 2.2 Running Example

129 As a running example, we use the *reference recommendation* problem: finding relevant citations for a new  
 130 paper. We consider three populations: the population of new papers for which relevant citations are to be  
 131 found, the population of existing papers whose citations are known, and the population of publication years.  
 132 The atoms that will be used for this problem throughout the paper are the following.  $WillCite(q, p)$  is the  
 133 atom to be predicted and indicates whether a query/new paper  $q$  will cite an existing paper  $p$ .  $Cited(p_1, p_2)$   
 134 shows whether or not an existing paper  $p_1$  has cited another existing paper  $p_2$ .  $PubIn(p, y)$  shows that  $p$  has  
 135 been published in year  $y$ .  $ImBef(y_1, y_2)$  indicates that  $y_2$  is the year immediately before  $y_1$ . The reference  
 136 recommendation problem can be viewed as follows: given a query paper  $Q$ , find a subset of existing papers  
 137 that  $Q$  will cite (i.e. find any paper  $P$  such that  $WillCite(Q, P)$  holds).

### 138 2.3 Relational Logistic Regression

139 Relational logistic regression (Kazemi et al., 2014) defines conditional probabilities based on weighted  
 140 rules. It can be viewed as the directed analogue of logistic regression, and as the directed analogue of  
 141 Markov logic (Domingos et al., 2008).

142 Let  $V(\mathbf{x})$  be an atom whose probability depends on a set  $\mathcal{A}$  of atoms,  $\psi$  be a set of WFs containing only  
 143 atoms from  $\mathcal{A}$ ,  $\hat{I}$  be a value assignment for the groundings in  $\mathcal{G}(\mathcal{A})$ ,  $\mathbf{X}$  be an assignment of objects to  $\mathbf{x}$ ,  
 144 and  $\{\mathbf{x}/\mathbf{X}\}$  be a substitution mapping logvars  $\mathbf{x}$  to objects  $\mathbf{X}$ .

**Relational logistic regression (RLR)** defines the probability of  $V(\mathbf{X})$  given  $\hat{I}$  as follows:

$$Prob_{\psi}(V(\mathbf{X}) = True \mid \hat{I}) = \sigma\left(\sum_{\langle w, \varphi \rangle \in \psi} w * \eta(\varphi\{\mathbf{x}/\mathbf{X}\}, \hat{I})\right) \quad (1)$$

145 where  $\eta(\varphi\{\mathbf{x}/\mathbf{X}\}, \hat{I})$  is the number of instances of  $\varphi\{\mathbf{x}/\mathbf{X}\}$  that are True with respect to  $\hat{I}$  and  $\sigma$  is the  
 146 sigmoid function. RLR makes the closed-world assumption: any ground atom that has not been observed  
 147 to be True is False. Note that  $\eta(True, \hat{I}) = 1$ .

148 Following Kazemi et al. (2014) and Fatemi et al. (2016), we assume that formulae in WFs have no  
 149 disjunction and replace conjunction with multiplication. Then atoms whose functors have a continuous  
 150 range can be also allowed in formulae. For instance if a value assignment maps  $R(X)$  to 1,  $S(X)$  to 0.9  
 151 and  $T(X)$  to 0.3, then the formula  $R(X) * S(X) * T(X)$  evaluates to  $1 * 0.9 * 0.3 = 0.27$ .

152 **EXAMPLE 1.** An RLR model may use the following WFs to define the conditional probability of  
 153  $WillCite(q, p)$  in our running example:

$$WF_0 : \langle w_0, True \rangle$$

$$WF_1 : \langle w_1, Publn(q, y) * ImBef(y, y') * Publn(p, y') \rangle$$

$$WF_2 : \langle w_2, Publn(q, y) * Publn(p', y) * Cited(p', p) \rangle$$

$$WF_3 : \langle w_3, Cited(p_1, p_2) * Cited(p_2, p) \rangle$$

157  $WF_0$  is a bias.  $WF_1$  considers existing papers that have been published a year before the query paper. A  
 158 positive weight for this WF indicates that papers published a year before the query paper are more likely to  
 159 be cited.  $WF_2$  considers existing papers cited by the other papers published in the same year as the query  
 160 paper. A positive weight for this WF indicates that as the number of times a paper has been cited by the  
 161 other papers published in the same year as the query paper grows, the chances of the query paper citing  
 162 that paper increases.  $WF_3$  considers existing papers that have been cited by other papers that have been  
 163 themselves cited by other papers. Note that the score of the last WF only depends on the paper being cited,  
 164 not the paper citing.

165 Consider the citations among existing papers in Fig. 1(a) and let the publication year for all the six  
 166 papers be 2017. Suppose we have a query paper  $Q$  which is to be published in 2017 and we want to find  
 167 the probability of  $WillCite(Q, Paper_2)$  according to the WFs above. Applying the substitution  $\{\langle q, p \rangle /$   
 168  $\langle Q, Paper_2 \rangle\}$  to the above four WFs gives the following four WFs respectively:

$$WF_0 : \langle w_0, True \rangle$$

$$WF_1 : \langle w_1, Publn(Q, y) * ImBef(y, y') * Publn(Paper_2, y') \rangle$$

$$WF_2 : \langle w_2, \text{Publn}(Q, y) * \text{Publn}(p', y) * \text{Cited}(p', \text{Paper}_2) \rangle$$

$$WF_3 : \langle w_3, \text{Cited}(p_1, p_2) * \text{Cited}(p_2, \text{Paper}_2) \rangle$$

170

171 Then we evaluate each WF. The first one evaluates to  $w_0$ . The second evaluates to 0 as  $Q$  is being  
 172 published in 2017 and  $\text{Paper}_2$  has also been published in 2017. The third WF evaluates to  $w_2 * 2$  as there  
 173 are 2 papers that have been published in the same year as  $Q$  and cite  $\text{Paper}_2$ . And the last WF evaluates  
 174 to  $w_3 * 4$  as  $\text{Paper}_5$  and  $\text{Paper}_6$  (that cite  $\text{Paper}_2$ ) are each cited by two other papers. Therefore, the  
 175 conditional probability of  $\text{WillCite}(Q, \text{Paper}_2)$  is as follows:

$$\sigma(w_0 + w_2 * 2 + w_3 * 4)$$

## 176 2.4 Path Ranking Algorithm

177 Let  $V(s, e)$  be a target binary predicate, i.e. for a query object  $S \in \Delta_s$ , we would like to find the  
 178 probability of any  $E \in e$  having the relation  $V$  with  $S$ . **Path ranking algorithm (PRA)** (Lao and Cohen,  
 179 2010b) defines this probability using a set of WPRs  $\Psi$ . The first logvar of each path relation in  $\Psi$  is either  $s$   
 180 or a logvar other than  $s$  and  $e$ , the last logvar is always  $e$ , and the middle logvars are neither  $s$  nor  $e$ .

181 In PRA, each path relation  $\mathcal{PR} = x_0 \xrightarrow{R_1} x_1 \xrightarrow{R_2} \dots \xrightarrow{R_l} e$  defines a distribution over the objects in  
 182  $\Delta_e$ . This distribution corresponds to the probability of following  $\mathcal{PR}$  and landing at each of the objects  
 183 in  $\Delta_e$ , and is computed as follows. Firstly, a uniform distribution  $D_0$  is considered on the objects in  $\Delta_{x_0}$ ,  
 184 corresponding to the probability of landing at each of these objects if the object is selected randomly. For  
 185 instance if there are  $\alpha$  objects in  $\Delta_{x_0}$ ,  $D_0$  for all objects is  $\frac{1}{\alpha}$ . Then, the distribution  $D_1$  over the objects  
 186 in  $\Delta_{x_1}$  is calculated by marginalizing over the variables in  $D_0$  and following a random step on  $R_1$ . For  
 187 instance for an object  $X_1 \in \Delta_{x_1}$ , assume  $R_1(x_0, X_1)$  holds only for two objects  $X_0$  and  $X'_0$  in  $\Delta_{x_0}$ . Also  
 188 assume  $X_0$  and  $X'_0$  have the  $R_1$  relation with  $\beta$  and  $\gamma$  objects in  $x_1$  respectively. Then the probability of  
 189 landing at  $X_1$  is  $\frac{1}{\alpha} * \frac{1}{\beta} + \frac{1}{\alpha} * \frac{1}{\gamma}$ . The following distributions  $D_2, \dots, D_l$  can be computed similarly.  $D_l$   
 190 gives the probability of landing at any object in  $\Delta_e$ .

191 Let  $\theta = \{\langle s, e \rangle / \langle S, E \rangle\}$ . In order to find  $\text{Prob}(V(S, E))$ , for each path relation  $\mathcal{PR} \in \Psi$ , PRA calculates  
 192 the probability of landing at  $E$  according to  $\mathcal{PR}\theta$  (denoted by  $h(\mathcal{PR}\theta)$ ), and calculates  $\text{Prob}(V(S, E))$  by  
 193 taking the sigmoid of the weighted sum of these probabilities as follows:

$$\text{Prob}(V(S, E)) = \sigma\left(\sum_{\langle w, \mathcal{PR} \rangle \in \Psi} w \cdot h(\mathcal{PR}\theta)\right) \quad (2)$$

194 Algorithm 1 shows a recursive algorithm for calculating  $h(\mathcal{PR})$  for a path relation  $\mathcal{PR}$ . The first  
 195 if statement specifies that the walk starts randomly at any object in  $\Delta_{x_0}$ .  $\text{uniform}(\Delta_{x_0})$  indicates a  
 196 uniform probability over the objects in  $\Delta_{x_0}$ . This is the termination criterion of the recursion. When  
 197  $\mathcal{PR} = x_0 \xrightarrow{R_1} x_1 \xrightarrow{R_2} \dots \xrightarrow{R_l} x_l$  is not empty ( $l \neq 0$ ), first the probability of landing at any object  $E'$   
 198 in the range of  $\mathcal{PR}' = x_0 \xrightarrow{R_1} x_1 \xrightarrow{R_2} \dots \xrightarrow{R_{l-1}} x_{l-1}$  is calculated using a recursive call to  $h(\mathcal{PR}')$  and  
 199 stored in  $p\text{Land}_{l-1}$ . The probability of landing at any object  $E$  in range of  $\mathcal{PR}$  by randomly walking  
 200 on  $\mathcal{PR}$  can then be calculated as the sum of the probabilities of landing at each object  $E'$  by randomly  
 201 walking on  $\mathcal{PR}'$  multiplied by the probability of reaching  $E$  from  $E'$  by a random walk according to the  
 202 predicate  $R_l$ . The two nested for loops calculate the probability of landing at any object  $E \in \text{range}(\mathcal{PR})$   
 203 according to  $R_l$ .  $R_l(E', E)$  indicates whether there is a link from  $E'$  to  $E$  (otherwise the probability of

**Algorithm 1**  $h(\mathcal{PR})$ 

**Input:** Relation path  $\mathcal{PR} = x_0 \xrightarrow{R_1} x_1 \xrightarrow{R_2} \dots \xrightarrow{R_l} x_l$

**Output:** Probability of landing at any object in  $\Delta_{x_l}$  when starting randomly at any object in  $\Delta_{x_0}$  and walking on  $\mathcal{PR}$ .

```

1: if  $l = 0$  then
2:   return  $uniform(\Delta_{x_0})$ 
3:  $\mathcal{PR}' = x_0 \xrightarrow{R_1} x_1 \xrightarrow{R_2} \dots \xrightarrow{R_{l-1}} x_{l-1}$ 
4:  $pLand_{l-1} = h(\mathcal{PR}')$ 
5: for  $E \in range(\mathcal{PR})$  do
6:    $pLand_l(E) = 0$ 
7: for  $E' \in range(\mathcal{PR}')$  do
8:    $C_{R_l}(E') = \#E \in range(\mathcal{PR})$  s.t.  $R_l(E', E) = 1$ 
9:   for  $E \in range(\mathcal{PR})$  do
10:     $pWalk(E', E) = \frac{R_l(E', E)}{C_{R_l}(E')}$ 
11:     $pLand_l(E) += pLand_{l-1}(E') * pWalk(E', E)$ 
12: return  $pLand_l$ 

```

204 transitioning from  $E'$  to  $E$  according to  $R_l$  is 0) and  $C_{R_l}$  is a normalization constant indicating the number  
 205 of possible transitions from  $E'$  according to  $R_l$ .  $pWalk(E', E)$  indicates the probability of walking from  
 206  $E'$  to  $E$  if one of the objects connected to  $E'$  through  $R_l$  is selected uniformly at random, which equals  
 207  $\frac{R_l(E', E)}{C_{R_l}}$ .  $pLand_l$  stores the probability of landing at any object  $E$  in the range of  $(\mathcal{PR})$  following  $\mathcal{PR}$ , and  
 208 is returned as the output of the function.

209 **EXAMPLE 2.** A PRA model may use the following WPRs to define the conditional probability of  
 210  $WillCite(q, p)$  in our running example:

$$\begin{aligned}
 &WPR_0 : \langle w_0, p \rangle \\
 &WPR_1 : \left\langle w_1, q \xrightarrow{PubIn} y \xrightarrow{ImBef} y' \xrightarrow{PubIn^{-1}} p \right\rangle \\
 &WPR_2 : \left\langle w_2, q \xrightarrow{PubIn} y \xrightarrow{PubIn^{-1}} p' \xrightarrow{Cited} p \right\rangle \\
 &WPR_3 : \left\langle w_3, p_1 \xrightarrow{Cited} p_2 \xrightarrow{Cited} p \right\rangle
 \end{aligned}$$

214  $WPR_0$  is a bias,  $WPR_1$  considers the papers published a year before the query paper,  $WPR_2$  considers  
 215 papers cited by other papers published in the same year as the query paper, and  $WPR_3$  mimics PageRank  
 216 algorithm for finding important papers in terms of citations (cf. Lao and Cohen (2010b) for more detail).  
 217 Consider the citations among existing papers in Fig. 1(a) and let the publication year for all the six papers  
 218 be 2017. Suppose we have a query paper  $Q$  which is to be published in 2017 and we want to find the  
 219 probability of  $WillCite(Q, Paper_2)$  according to the PRA model above. Applying the substitution  $\{\langle q, p \rangle /$   
 220  $\langle Q, Paper_2 \rangle\}$  to the above WPRs gives the following WPRs respectively:

$$\begin{aligned}
 &WPR_0 : \langle w_0, Paper_2 \rangle \\
 &WPR_1 : \left\langle w_1, Q \xrightarrow{PubIn} y \xrightarrow{ImBef} y' \xrightarrow{PubIn^{-1}} Paper_2 \right\rangle
 \end{aligned}$$

$$WPR_2 : \left\langle w_2, Q \xrightarrow{\text{Publn}} y \xrightarrow{\text{Publn}^{-1}} p' \xrightarrow{\text{Cited}} Paper_2 \right\rangle$$

$$WPR_3 : \left\langle w_3, p_1 \xrightarrow{\text{Cited}} p_2 \xrightarrow{\text{Cited}} Paper_2 \right\rangle$$

222

223  $WPR_0$  evaluates to  $w_0$ .  $WPR_1$  evaluates to 0.  $WPR_2$  evaluates to  $w_2 * (\frac{1}{6} * \frac{1}{4} + \frac{1}{6} * \frac{1}{2}) = w_2 * 0.125$  as for  
 224 the path  $y \xrightarrow{\text{Publn}^{-1}} p'$  there is  $\frac{1}{6}$  probability for randomly walking to either  $Paper_5$  or  $Paper_6$  and then there  
 225 is  $\frac{1}{4}$  probability to walk randomly from  $Paper_5$  to  $Paper_2$  and  $\frac{1}{2}$  probability to walk randomly from  $Paper_6$   
 226 to  $Paper_2$  according to Cited relation.  $WPR_3$  evaluates to  $w_3 * \frac{1}{6} * (\frac{1}{2} * \frac{1}{4} + \frac{1}{3} * (\frac{1}{4} + \frac{1}{2}) + \frac{1}{4} * \frac{1}{2}) \approx w_3 * 0.083$ .  
 227 The  $\frac{1}{6}$  outside parenthesis is the probability of randomly starting at any paper,  $\frac{1}{2} * \frac{1}{4}$  is the probability  
 228 of transitioning from  $Paper_3$  to  $Paper_5$  and then to  $Paper_2$ , and so forth. Therefore, the conditional  
 229 probability of WillCite( $Q, Paper_2$ ) is as follows:

$$\sigma(w_0 + w_2 * 0.125 + w_3 * 0.083)$$

### 3 RLR WITH NORMALIZED RELATIONS GENERALIZES PRA

230 In order to prove that RLR with normalized relations generalizes PRA, we first define relation chains and  
 231 describe some of their properties.

#### 232 3.1 Relations Chain

233 **DEFINITION 1.** We define a **relations chain** as a list of binary atoms  $V_1(x_0, x_1), \dots, V_m(x_{m-1}, x_m)$   
 234 such that for each  $V_i$  and  $V_{i+1}$ , the second logvar of  $V_i$  is the same as the first logvar of  $V_{i+1}$ ,  $x_0, \dots, x_m$   
 235 are different logvars, and  $V_i$  and  $V_j$  can be the same or different predicates.

236 **EXAMPLE 3.**  $V_1(x, y), V_2(y, z)$  is a relations chain, and  $V_1(x, y), V_2(z, y)$  and  $V_1(x, y), V_2(y, z), V_3(z, x)$   
 237 are not relations chains.

238 **DEFINITION 2.** A first-order formula corresponds to a relations chain if all its literals are binary  
 239 predicates and non-negated, and there exists an ordering of the literals that is a relations chain.

240 **EXAMPLE 4.** The first-order formula  $V_1(x_1, x_2) * V_2(x_3, x_1)$  corresponds to a relations chain as the  
 241 order  $V_2(x_3, x_1), V_1(x_1, x_2)$  is a relations chain.

242 It follows from RLR definition that re-ordering the literals in each of its WFs does not change the  
 243 distribution. For any WF whose formula corresponds to a relations chain, we assume hereafter that its  
 244 literals have been re-ordered to match the order of the corresponding relations chain.

245 **DEFINITION 3.** Let  $V(x, y)$  be a target atom. *Relations chain RLR (RC-RLR)* is a subset of RLR for  
 246 defining a conditional probability distribution for  $V(x, y)$  where:

- 247 • formulae of WFs correspond to relations chains,
- 248 • for each WF, the second logvar of the last atom is  $y$ ,
- 249 •  $x$  may only appear as the first logvar of the first atom,
- 250 •  $y$  may only appear as the second logvar of the last atom.

251 For RLR models, in order to evaluate a formula, one may have nested loops over logvars of the formula  
 252 that do not appear in the target atom, or conjoin all literals one by one and then count. WFs of RC-RLR,

**Algorithm 2**  $Eval(\varphi)$ **Input:** Formula  $\varphi = R_1(x_0, x_1) * R_2(x_1, x_2) * \dots * R_l(x_{l-1}, x_l)$ .**Output:** Evaluation of  $\varphi$ .

```

1: if  $l = 0$  then
2:   return  $ones(|\Delta_{x_0}|)$ 
3:  $\varphi' = R_1(x_0, x_1) * R_2(x_1, x_2) * \dots * R_{l-1}(x_{l-2}, x_{l-1})$ 
4:  $eval_{l-1} = Eval(\varphi')$ 
5: for  $E \in \Delta_{x_l}$  do
6:    $eval_l(E) = 0$ 
7:   for  $E' \in \Delta_{x_{l-1}}$  do
8:     for  $E \in \Delta_{x_l}$  do
9:        $canWalk(E', E) = R_l(E', E)$ 
10:       $eval_l(E) += eval_{l-1}(E') * canWalk(E', E)$ 
11: return  $eval_l$ 

```

253 however, can be evaluated in a special way. In order to evaluate a formula in RC-RLR, starting from the  
 254 end (or beginning), the effect of each literal can be calculated and then the literal can be removed from the  
 255 formula. Algorithm 2 indicates how a formula corresponding to a relations chain can be evaluated. This  
 256 evaluation grows with the product of the number of literals in the formula and the number of observed data  
 257 which makes it highly scalable.

258 When  $l = 0$ , the formula corresponds to True and evaluates to 1 for any  $X_0 \in x_0$ . Therefore, in  
 259 this case the algorithm returns a vector of ones of size  $|\Delta_{x_0}|$ . Otherwise, the algorithm first evaluates  
 260  $\varphi' = R_1(x_0, x_1) * R_2(x_1, x_2) * \dots * R_{l-1}(x_{l-2}, x_{l-1})$  using a recursive call to the  $Eval$  function. The  
 261 resulting vector is stored in  $eval_{l-1}$  such that for a  $E' \in \Delta_{x_{l-1}}$ ,  $eval_{l-1}[E']$  indicates the result of evaluating  
 262  $\varphi' = R_1(x_0, x_1) * R_2(x_1, x_2) * \dots * R_{l-1}(x_{l-2}, E')$ . Then in order to evaluate  $\varphi$  for some  $E \in \Delta_{x_l}$ , we sum  
 263  $eval_{l-1}[E']$ s for any  $E' \in \Delta_{x_{l-1}}$  such that  $R_l(E', E)$  holds.  $canWalk$  in the algorithm is 1 if  $R_l(E', E)$   
 264 holds and 0 otherwise, and  $eval_l(E) += eval_{l-1}(E') * canWalk(E', E)$  adds  $eval_{l-1}[E']$  to  $eval_l[E]$  if  
 265  $canWalk$  is 1.

266 PROPOSITION 1. *Algorithm 2 is correct.*

267 PROOF. Let  $\varphi = R_1(x_0, x_1) * R_2(x_1, x_2) * \dots * R_l(x_{l-1}, x_l) * eval_l(x_l)$  ( $eval_l(x_l)$  can be initialized to a  
 268 vector of ones at the beginning of the algorithm). Since by definition of relations chain  $x_l$  only appears in  $R_l$   
 269 and  $eval_l(x_l)$ , for any  $X_{l-1} \in \Delta_{x_{l-1}}$  we can evaluate  $eval_{l-1}(X_{l-1}) = \sum_{X_l \in \Delta_{x_l}} R_l(X_{l-1}, X_l) * eval_l(X_l)$   
 270 separately and replace  $R_l(x_{l-1}, x_l) * eval_l(x_l)$  with  $eval_{l-1}(x_{l-1})$  thus getting  $\varphi' = R_1(x_0, x_1) * R_2(x_1, x_2) * \dots * R_{l-1}(x_{l-2}, x_{l-1}) * eval_{l-1}(x_{l-1})$ . The same procedure can compute  $\varphi'$ .

**272 3.2 From PRA to Relation Chains**

273 PROPOSITION 2. *A path relation corresponds to a relations chain.*

274 PROOF. Let  $\mathcal{PR} = x_0 \xrightarrow{R_1} x_1 \xrightarrow{R_2} \dots \xrightarrow{R_l} x_l$  be a path relation. We create a relation atom  $R_i(x_{i-1}, x_i)$   
 275 for any sub-path  $x_{i-1} \xrightarrow{R_i} x_i$  resulting in relations  $R_1(x_0, x_1), R_2(x_1, x_2), \dots, R_l(x_{l-1}, x_l)$ . By definition  
 276 of path relations, the second logvar of any relation  $R_i$  is the same as the first logvar of the next relation.  
 277 Since by definition the logvars in a path relation are different, the second logvar of any relation  $R_i$  is only  
 278 equivalent to the first logvar of the next relation.

279 EXAMPLE 5. Consider the path relation  $q \xrightarrow{\text{Publn}} y \xrightarrow{\text{Publn}^{-1}} p' \xrightarrow{\text{Cited}} p$  from Example 2. This path  
 280 relation corresponds to a relations chain with atoms  $\text{Publn}(q, y)$ ,  $\text{Publn}^{-1}(y, p')$  and  $\text{Cited}(p', p)$ .

### 281 3.3 Row-Wise Count Normalization

282 Having a binary predicate  $V(x, y)$  and a set of pairs of objects for which  $V$  holds, one may consider the  
 283 importance of these pairs to be different. For instance, if a paper has cited only 20 papers, the importance of  
 284 these citations may be more than the importance of citations for a paper citing 100 papers. One way to take  
 285 the importance of the pairs into account is to normalize the relations. A simple way to normalize a relation  
 286 is to normalize it by row-wise counts. For some  $X \in \Delta_x$ , let  $\alpha$  represent the number of  $Y' \in \Delta_y$  such  
 287 that  $V(X, Y')$  holds. When  $\alpha \neq 0$ , instead of considering  $V(X, Y) = 1$  for a pair  $\langle X, Y \rangle$ , we normalize  
 288 it to  $V(X, Y) = \frac{1}{\alpha}$ . After this normalization, the citations of a paper with 20 citations are 5 times more  
 289 important than the citations of a paper with 100 citations overall. Note that when  $\alpha = 0$ , we do not change  
 290 any values. We refer to this normalization method as *row-wise count (RWC) normalization*. Fig. 1(b) show  
 291 the result of applying RWC normalization to the relation in Fig. 1(a). Note that there may be several other  
 292 ways to normalize a relation; here we introduced RWC because, as we will see in the upcoming sections, it  
 293 is the normalization method used in PRA.

### 294 3.4 Main Theorem

295 THEOREM 1. Any PRA model is equivalent to an RC-RLR model with RWC normalization.

296 PROOF. Let  $\Psi = \{\langle w_0, \mathcal{PR}_0 \rangle, \dots, \langle w_k, \mathcal{PR}_k \rangle\}$  represent a set of WPRs used by a PRA model. We  
 297 proved in Proposition 2 that any path relation  $\mathcal{PR}_i$  in  $\Psi$  corresponds to a relations chain. By multiplying  
 298 the relations in the relation chain, one gets a formula  $\varphi_i$  for each  $\mathcal{PR}_i$  and this formula is by construction  
 299 guaranteed to correspond to a relations chain. We construct an RC-RLR model whose WFs are  $\psi =$   
 300  $\{\langle v_0, \varphi_0 \rangle, \dots, \langle v_k, \varphi_k \rangle\}$ . Given that the relations (and their order) used in  $\mathcal{PR}_i$  and  $\varphi_i$  are the same for any  
 301  $i$ , the only differences between the evaluation of  $\mathcal{PR}_i$  and  $\varphi_i$  according to Algorithm 1 and Algorithm 2  
 302 are: 1- Algorithm 1 divides  $R_l(E', E)$  by  $C_{R_l}(E')$  while Algorithm 2 does not, and 2- in the termination  
 303 condition, Algorithm 1 returns a uniform distribution over objects in  $\Delta_{x_0}$  while Algorithm 2 returns a  
 304 vector of ones of size  $|\Delta_{x_0}|$ . Dividing  $R_l(E', E)$  by  $C_{R_l}(E')$  is equivalent to RWC normalization and the  
 305 difference in the constant value of the function in the termination condition gets absorbed in the weights  
 306 that are multiplied to each path relation or formula. Therefore, the RC-RLR model with WFs  $\psi$  is identical  
 307 to the PRA model with WPRs  $\Psi$  after normalizing the relations using RWC.

308 EXAMPLE 6. Consider the PRA model in Example 2. For the four WPRs in that model, we create the  
 309 following corresponding WFs for an RC-RLR model by multiplying the relations in the path relations:

$$\begin{aligned}
 & \langle v_0, \text{True} \rangle \\
 310 & \langle v_1, \text{Publn}(q, y_1) * \text{ImBef}(y_1, y_2) * \text{Publn}^{-1}(y_2, p) \rangle \\
 311 & \langle v_2, \text{Publn}(q, y_1) * \text{Publn}^{-1}(y_1, p') * \text{Cited}(y_1, p) \rangle \\
 312 & \langle v_3, \text{Cited}(p_1, p_2) * \text{Cited}(p_2, p) \rangle
 \end{aligned}$$

313 Consider computing  $\text{WillCite}(Q, \text{Paper}_2)$  according to an RC-RLR model with the above WFs, where  
 314 all existing papers and  $Q$  have been published in 2017 and the relations have been normalized using RWC  
 315 normalization (e.g., as in Fig. 1(b) for relation Cited). Then the first formula evaluates to  $v_0$ . The second

316 WF evaluates to 0. The third WF evaluates to  $v_2 * \frac{1}{6} * (\frac{1}{4} + \frac{1}{2})$  as the values in relation  $\text{Publn}^{-1}$  have been  
 317 normalized to  $\frac{1}{6}$  for year 2017 and the values in relation Cited have been normalized to  $\frac{1}{4}$  and  $\frac{1}{2}$  for  $\text{Paper}_5$   
 318 and  $\text{Paper}_6$  as in Fig. 1(b). The last WF evaluates to  $v_3 * (\frac{1}{2} * \frac{1}{4} + \frac{1}{3} * (\frac{1}{4} + \frac{1}{2}) + \frac{1}{4} * \frac{1}{2})$ . The  $\frac{1}{2} * \frac{1}{4}$  comes  
 319 from  $\text{Cited}(\text{Paper}_3, \text{Paper}_5) * \text{Cited}(\text{Paper}_5, \text{Paper}_2)$ ,  $\frac{1}{3} * (\frac{1}{4} + \frac{1}{2})$  comes from  $\text{Cited}(\text{Paper}_4, \text{Paper}_5) * \text{Cited}(\text{Paper}_5, \text{Paper}_2)$   
 320 and  $\text{Cited}(\text{Paper}_4, \text{Paper}_6) * \text{Cited}(\text{Paper}_6, \text{Paper}_2)$  and  $\frac{1}{4} * \frac{1}{2}$  comes from  
 321  $\text{Cited}(\text{Paper}_5, \text{Paper}_6) * \text{Cited}(\text{Paper}_6, \text{Paper}_2)$ . As it can be viewed from Example 2, after creating the  
 322 equivalent RC-RLR model and normalizing the relations using RWC normalization, all WPRs evaluate to  
 323 the same value as their corresponding WF, except the last WF. The  $\frac{1}{6}$  before the parenthesis in Example 2  
 324 is missing when evaluating the last WF. This  $\frac{1}{6}$ , however, is a constant independent of the query (it is the  
 325 constant value of the uniform distribution in the if statement corresponding to the termination criteria in  
 326 Algorithm 1). Assuming  $v_3 = w_3 * \frac{1}{6}$  and all other  $v_i$ s are the same as  $w_i$ s, the conditional probability of  
 327  $\text{Cited}(Q, \text{Paper}_2)$  according to the RC-RLR model above will be the same as the PRA model in Example 2.

### 328 3.5 From Random Walk Strategies to Structure Learning

329 The restrictions imposed on the formulae by path relations in PRA reduces the number of possible  
 330 formulae to be considered in a model compared to RLR models. However, there may still be many possible  
 331 path relations and considering all possible path relations for a PRA model may not be practical.

332 Lao and Cohen (2010b) allow the random walk to follow any path, but restrict the maximum number of  
 333 steps. In particular, they only allow for path relations whose length is less than some  $l$ . The value of  $l$  can  
 334 be selected based on the number of objects, relations, available hardware, and the amount of time one can  
 335 afford for learning/inference. This strategy automatically gives a (very simple) structure learning algorithm  
 336 for RC-RLR by considering only formulae whose number of relations are less than  $l$ .

337 Lao et al. (2011) follow a more sophisticated approach for limiting the number of path relations. Besides  
 338 limiting the maximum length of the path relations to  $l$ , Lao et al. (2011) impose two more restrictions:  
 339 for any path relation to be included, 1- the probability of reaching the target objects must be non-zero  
 340 for at least a fraction  $\alpha$  of the training query objects, and 2- it should at least retrieve one target object  
 341 in the training set. During parameter learning, they impose a Laplacian prior on their weights to further  
 342 reduce the number of path relations. In an experiment on knowledge completion for NELL (Carlson et al.,  
 343 2010), they show that these two restrictions plus the Laplacian prior reduce the number of possible path  
 344 relations by almost 99.6 and 99.99 percents when  $l = 3$  and  $l = 4$  respectively. Therefore, their random  
 345 walk strategy is capable of taking more steps (i.e. selecting a larger value for  $l$ ) and capture features that  
 346 require longer chains of relations. This random walk strategy is called *data-driven path finding*.

347 Both restrictions in data-driven path finding can be easily verified for RC-RLR formulae and the set of  
 348 possible formulae can be restricted accordingly. Furthermore, during parameter learning, a Laplacian prior  
 349 can be imposed on the weights of the weighted formulae. RC-RLR models learned in this way corresponds  
 350 to PRA models learned using data-driven path finding. Therefore, data-driven path finding can be also  
 351 considered as a structure learning algorithm for RC-RLR. With the same reasoning, several other random  
 352 walk strategies can be considered as structure learning algorithms for RC-RLR, and vice versa. This allows  
 353 for faster development of the two paradigms by leveraging ideas developed in each community in the other.

## 4 PRA VS. RLR

354 An advantage of PRA models over RLR models is their efficiency: there is a smaller search space for WFs  
 355 and all WFs can be evaluated efficiently. Such efficiency makes PRA scale to larger domains where models

356 based on weighted rule learning such as RLR often have scalability issues. It also allows PRA models to  
 357 scale to and capture features that require longer chains of relations. However, the efficiency comes at the  
 358 cost of losing modelling power. In the following subsections, we discuss such costs.

#### 359 4.1 Shortcomings of Relations Chains

360 Since PRA models restrict themselves to relations chains of a certain type, they lose the chance to leverage  
 361 many other WFs. As an example, in order to predict  $\text{Cites}(p_1, p_2)$  for the reference recommendation task,  
 362 suppose we would like to recommend papers published a year before the target paper that have been cited  
 363 by the papers published in the same year as the target paper. Such a feature requires the following formula:  
 364  $\text{Publn}(p_1, y) * \text{Before}(y, y') * \text{Publn}(p_2, y') * \text{Cites}(p', p_2) * \text{Publn}(p', y)$ . It is straightforward to verify that  
 365 this formula cannot be included in RC-RLR (and consequently in PRA) as  $p_2$  (the second logvar of the  
 366 target atom) is appearing twice in the formula, thus violating the last condition in Definition 3. While  
 367 restricting the formulae to the ones that correspond to relations chain may speed up learning and reasoning,  
 368 it reduces the space of features that can be included in a relational learning model, thus potentially  
 369 decreasing accuracy.

#### 370 4.2 Non-binary Atoms

371 One issue with PRA models is the difficulty in including unary atoms in such models. As an example,  
 372 suppose in Example 2 we would like to treat conference papers and journal papers differently. For an  
 373 RLR model, this can be easily done by including  $\text{Conference}(p)$  or  $\text{Journal}(p)$  as an extra atom in the  
 374 formulae. For PRA, however, this cannot be done. The way unary atoms are currently handled in PRA  
 375 models is through  $\text{isA}$  and  $\text{isA}^{-1}$  relations (Lao et al., 2011). For instance, a path relation may contain  
 376  $\text{paper} \xrightarrow{\text{isA}} \text{type}$ , but the only next predicate that can be applied to this path is  $\text{isA}^{-1}$  giving the other papers  
 377 with the same type as the paper in the left of the arrow. This is, however, limiting and does not allow for,  
 378 e.g., treating conference and journal papers differently.

379 Atoms with more than two logvars are another issue for PRA models since they restrict their models to  
 380 binary atoms. While any relation with more than two arguments can be converted into several binary atoms,  
 381 the random walk strategies used for PRA models (and the probabilities for making these random steps)  
 382 make it unclear how atoms with more than two logvars can be leveraged in PRA models.

#### 383 4.3 Continuous Atoms

384 For any sub-path  $x \xrightarrow{R} y$  in a path relation of a PRA model,  $R$  typically has a range  $\{0, 1\}$ : for any  
 385 object  $X \in \Delta_x$ , this sub-path gives the objects in  $\Delta_y$  participating in relation  $R$  with  $X$ . PRA models can  
 386 be extended to handle some forms of continuous atoms. For instance for the reference recommendation  
 387 problem, suppose we have an atom  $\text{Sim}(p, p')$  indicating a measure of similarity between the titles of two  
 388 papers. The higher the  $\text{Sim}(p, p')$ , the more similar the titles of the two papers. A sensible WF for an RLR  
 389 model predicting  $\text{Cites}(p_1, p_2)$  may be  $\text{Sim}(p_1, p') * \text{Cites}(p', p_2)$ . In order to extend PRA models to be able  
 390 to leverage such continuous atoms, one has to change line 8 in Algorithm 1 to sum the values of  $R_l(E', E)$   
 391 instead of counting how many times the relation holds.

392 For many types of continuous atoms, however, it is not straightforward to extend PRA models to leverage  
 393 them. As an example, suppose we have an atom  $\text{Temperature}(r, d)$  showing the temperature of a region in  
 394 a specific date. It is not clear how a random walk step can be made based on this atom as the temperature  
 395 can, e.g., be positive or negative.

#### 396 4.4 Relational Normalization

397 Normalizing the relations is often ignored in models based on weighted rule learning. For the most part,  
398 this ignorance may be because several of these models cannot handle continuous atoms. Given that PRA is  
399 a special form of weighted rule learning models such as RLR with RWC normalization, not normalizing  
400 the relations may be the reason why in Lao et al. (2011)'s experiments, PRA outperforms the weighted rule  
401 learning method FOIL (Quinlan, 1990) for link prediction in NELL (Carlson et al., 2010).

402 The type of normalization used in PRA (RWC) may not be the best option in many applications. As  
403 an example, suppose for the reference recommendation task we want to find papers similar to the query  
404 paper in terms of the words they use. Let  $\text{Contains}^{-1}(w, p)$  show the relation for words in each paper. It  
405 is well-known in information retrieval that words do not have equal importances and a normalization of  
406  $\text{Contains}^{-1}(w, p)$  is necessary to take such importance into account. PRA models consider the importance  
407 of each word  $W$  as  $\text{Score}_1(W) = \frac{1}{f(W)}$ , where  $f(W)$  is the number of papers containing the word  $W$   
408 (see e.g., (Lao and Cohen, 2010b)). However, it has been well-known in information retrieval community  
409 for several decades, and information theoretically justified more than a decade ago (Robertson, 2004),  
410 that  $\text{Score}_2(W) = \log\left(\frac{\#papers}{f(W)}\right)$  provides a better importance score. Most TF-IDF (Salton and Buckley,  
411 1988) based information retrieval algorithms currently rely on  $\text{Score}_2$ . It is straightforward to include  
412 the latter score in an RLR model: one only has to multiply the formulae using word information by  
413  $\text{Score}_2(W)$ , without normalizing the  $\text{Contains}^{-1}(w, p)$  relation (see, e.g., (Fatemi, 2017)). However, it is  
414 not straightforward how such a score can be incorporated into PRA models as they do not include unary or  
415 continuous atoms.

#### 416 4.5 Evaluating Formulae

417 Evaluating the formulae in models based on weighted rule learning is known to be expensive, especially  
418 for relations with lower sparsities and for longer formulae. In practice, approximations are typically used  
419 for scaling the evaluations. Since formulae in RC-RLR correspond to path relations, these formulae can be  
420 approximated efficiently using sampling techniques developed within graph random walk community such  
421 as fingerprinting (Fogaras et al., 2005; Lao and Cohen, 2010a), weighted particle filtering (Lao and Cohen,  
422 2010a), and low-variance sampling (Lao et al., 2011), without noticeably affecting the accuracy. Extending  
423 sampling ideas to other formulae is an interesting future direction.

## 5 CONCLUSION

424 With abundance of relational and graph data, statistical relational learning has gained great amounts of  
425 attention. Three main relational learning paradigms have been developed during the past decade and more:  
426 weighted rule learning, graph random walk, and tensor factorization. These paradigms have been mostly  
427 developed and studied in isolation with few works aiming at understanding the relationship among them  
428 or combining them. In this paper, we studied the relationship between two relational learning paradigms:  
429 weighted rule learning and graph random walk. In particular, we studied the relationship between relational  
430 logistic regression (RLR), one of the recent developments in weighted rule learning paradigm, and path  
431 ranking algorithm (PRA), one of the most well-known algorithms in graph random walk paradigm. Our  
432 main contribution was to prove that PRA models correspond to a subset of RLR models after row-wise  
433 count normalization. We discussed the advantages this proof provides for both paradigms as well as for  
434 statistical relational AI community in general. Our result sheds light on several issues with both paradigms  
435 and possible ways to improve them.

## REFERENCES

- 436 Bhattacharya, I. and Getoor, L. (2007). Collective entity resolution in relational data. *ACM Transactions*  
437 *on Knowledge Discovery from Data (TKDD)* 1, 5
- 438 Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., and Yakhnenko, O. (2013). Translating embeddings  
439 for modeling multi-relational data. In *NIPS*. 2787–2795
- 440 Carlson, A., Betteridge, J., Kisiel, B., Settles, B., Jr, E. R. H., and Mitchell, T. M. (2010). Toward an  
441 architecture for never-ending language learning. In *AAAI*. vol. 5, 3
- 442 Das, R., Dhuliawala, S., Zaheer, M., Vilnis, L., Durugkar, I., Krishnamurthy, A., et al. (2017). Go for a  
443 walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning.  
444 *arXiv preprint arXiv:1711.05851*
- 445 De Raedt, L., Kersting, K., Natarajan, S., and Poole, D. (2016). Statistical relational artificial intelligence:  
446 Logic, probability, and computation. *Synthesis Lectures on Artificial Intelligence and Machine Learning*  
447 10, 1–189
- 448 De Raedt, L., Kimmig, A., and Toivonen, H. (2007). Problog: A probabilistic prolog and its application in  
449 link discovery. In *IJCAI*. vol. 7
- 450 Domingos, P., Kok, S., Lowd, D., Poon, H., Richardson, M., and Singla, P. (2008). Markov logic. In  
451 *Probabilistic Inductive Logic Programming*, eds. L. D. Raedt, P. Frasconi, K. Kersting, and S. Muggleton  
452 (New York: Springer). 92–117
- 453 Dries, A., Kimmig, A., Davis, J., Belle, V., and De Raedt, L. (2017). Solving probability problems in  
454 natural language. In *IJCAI*
- 455 Ebrahimi, J., Dou, D., and Lowd, D. (2016). Weakly supervised tweet stance classification by relational  
456 bootstrapping. In *EMNLP*. 1012–1017
- 457 Fatemi, B. (2017). *Finding a Record in a Database*. Master’s thesis, University of British Columbia
- 458 Fatemi, B., Kazemi, S. M., and Poole, D. (2016). A learning algorithm for relational logistic regression:  
459 Preliminary results. *arXiv preprint arXiv:1606.08531*
- 460 Fogaras, D., Rácz, B., Csalogány, K., and Sarlós, T. (2005). Towards scaling fully personalized pagerank:  
461 Algorithms, lower bounds, and experiments. *Internet Mathematics* 2, 333–358
- 462 Getoor, L. and Taskar, B. (2007). *Introduction to statistical relational learning* (MIT press)
- 463 Hommersom, A. and Lucas, P. J. (2011). Generalising the interaction rules in probabilistic logic. In *IJCAI*
- 464 Jain, A. and Pantel, P. (2010). Factrank: Random walks on a web of facts. In *Proceedings of the 23rd*  
465 *International Conference on Computational Linguistics*. 501–509
- 466 Kazemi, S. M., Buchman, D., Kersting, K., Natarajan, S., and Poole, D. (2014). Relational logistic  
467 regression. In *KR*
- 468 Kazemi, S. M., Fatemi, B., Kim, A., Peng, Z., Tora, M. R., Zeng, X., et al. (2017). Comparing aggregators  
469 for relational probabilistic models. *arXiv preprint arXiv:1707.07785*
- 470 Kazemi, S. M. and Poole, D. (2018). Relnn: A deep neural model for relational learning. In *AAAI*
- 471 Khot, T., Balasubramanian, N., Gribkoff, E., Sabharwal, A., Clark, P., and Etzioni, O. (2015). Markov  
472 logic networks for natural language question answering. *arXiv preprint arXiv:1507.03045*
- 473 Kimmig, A., Bach, S., Broecheler, M., Huang, B., and Getoor, L. (2012). A short introduction  
474 to probabilistic soft logic. In *Proceedings of the NIPS Workshop on Probabilistic Programming:*  
475 *Foundations and Applications*. 1–4
- 476 Kimmig, A., Mihalkova, L., and Getoor, L. (2015). Lifted graphical models: a survey. *Machine Learning*  
477 99, 1–45

- 478 Lao, N. and Cohen, W. W. (2010a). Fast query execution for retrieval models based on path-constrained  
479 random walks. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge*  
480 *discovery and data mining* (ACM), 881–888
- 481 Lao, N. and Cohen, W. W. (2010b). Relational retrieval using a combination of path-constrained random  
482 walks. *Machine learning* 81, 53–67
- 483 Lao, N., Mitchell, T., and Cohen, W. W. (2011). Random walk inference and learning in a large scale  
484 knowledge base. In *EMNLP*. 529–539
- 485 Lin, Y., Liu, Z., Luan, H., Sun, M., Rao, S., and Liu, S. (2015). Modeling relation paths for representation  
486 learning of knowledge bases. *arXiv preprint arXiv:1506.00379*
- 487 Nickel, M., Jiang, X., and Tresp, V. (2014). Reducing the rank in relational factorization models by  
488 including observable patterns. In *NIPS*. 1179–1187
- 489 Nickel, M., Murphy, K., Tresp, V., and Gabrilovich, E. (2016). A review of relational machine learning for  
490 knowledge graphs. *Proceedings of the IEEE* 104, 11–33
- 491 Nickel, M., Tresp, V., and Kriegel, H.-P. (2012). Factorizing yago: scalable machine learning for linked  
492 data. In *Proceedings of the 21st international conference on World Wide Web* (ACM), 271–280
- 493 Page, L., Brin, S., Motwani, R., and Winograd, T. (1999). *The PageRank citation ranking: Bringing order*  
494 *to the web*. Tech. rep., Stanford InfoLab
- 495 Pujara, J. and Getoor, L. (2016). Generic statistical relational entity resolution in knowledge graphs. *arXiv*  
496 *preprint arXiv:1607.00992*
- 497 Pujara, J., Miao, H., Getoor, L., and Cohen, W. W. (2013). Knowledge graph identification. In *International*  
498 *Semantic Web Conference (ISWC)* (Springer)
- 499 Quinlan, J. R. (1990). Learning logical definitions from relations. *Machine learning* 5, 239–266
- 500 Riedel, S., Yao, L., McCallum, A., and Marlin, B. M. (2013). Relation extraction with matrix factorization  
501 and universal schemas. In *HLT-NAACL*. 74–84
- 502 Robertson, S. (2004). Understanding inverse document frequency: on theoretical arguments for idf. *Journal*  
503 *of documentation* 60, 503–520
- 504 Salton, G. and Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information*  
505 *processing & management* 24, 513–523
- 506 Singla, P. and Domingos, P. (2006). Entity resolution with markov logic. In *Data Mining, 2006. ICDM'06.*  
507 *Sixth International Conference on* (IEEE), 572–582
- 508 Šourek, G., Aschenbrenner, V., Železny, F., and Kuželka, O. (2015). Lifted relational neural networks. In  
509 *Proceedings of the 2015th International Conference on Cognitive Computation: Integrating Neural and*  
510 *Symbolic Approaches-Volume 1583* (CEUR-WS. org), 52–60
- 511 Sridhar, D., Foulds, J. R., Huang, B., Getoor, L., and Walker, M. A. (2015). Joint models of disagreement  
512 and stance in online debate. In *ACL (1)*. 116–125
- 513 Trouillon, T., Welbl, J., Riedel, S., Gaussier, É., and Bouchard, G. (2016). Complex embeddings for simple  
514 link prediction. In *ICML*. 2071–2080
- 515 Yu, X., Ren, X., Sun, Y., Gu, Q., Sturt, B., Khandelwal, U., et al. (2014). Personalized entity  
516 recommendation: A heterogeneous information network approach. In *Proceedings of the 7th ACM*  
517 *international conference on Web search and data mining* (ACM), 283–292