

Statistical Inference on KGs

Knowledge bases on the Web



Motivation

- There are some serious challenges related to Knowledge Graphs
 - **Growth:** knowledge graphs are highly incomplete!
 - *Knowledge extraction:* extract new entities and relations
 - *Link prediction:* add (or guess) new relations between entities
 - **Validation:** knowledge graphs are not always correct!
 - *Entity resolution:* merge duplicate entities, split wrongly merged ones
 - *Error detection:* remove false assertions
 - **Interface:** how to make it easier to access knowledge?
 - *Semantic parsing:* interpret the meaning of queries
 - *Question answering:* computing answers using knowledge graphs
 - **Intelligence:** can AI emerge from knowledge graphs?
 - *Automatic reasoning* and planning
 - Generalization and abstraction

Motivation

- There are some serious challenges related to Knowledge Graphs
 - **Growth:** knowledge graphs are highly incomplete!
 - Knowledge extraction: extract new entities and relations
 - *Link prediction:* add (or guess) new relations between entities
 - **Validation:** knowledge graphs are not always correct!
 - *Entity resolution:* merge duplicate entities, split wrongly merged ones
 - *Error detection:* remove false assertions
 - **Interface:** how to make it easier to access knowledge?
 - *Semantic parsing:* interpret the meaning of queries
 - *Question answering:* computing answers using knowledge graphs
 - **Intelligence:** can AI emerge from knowledge graphs?
 - *Automatic reasoning* and planning
 - Generalization and abstraction

We talked about it in
the previous lectures

Motivation

- There are some serious challenges related to Knowledge Graphs
 - **Growth:** knowledge graphs are highly incomplete!
 - *Knowledge extraction:* extract new entities and relations
 - *Link prediction:* add (or guess) new relations between entities
 - **Validation:** knowledge graphs are not always correct!
 - *Entity resolution:* merge duplicate entities, split wrongly merged ones
 - *Error detection:* remove false assertions
 - **Interface:** how to make it easier to access knowledge?
 - *Semantic parsing:* interpret the meaning of queries
 - *Question answering:* computing answers using knowledge graphs
 - **Intelligence:** can AI emerge from knowledge graphs?
 - *Automatic reasoning* and planning
 - Generalization and abstraction

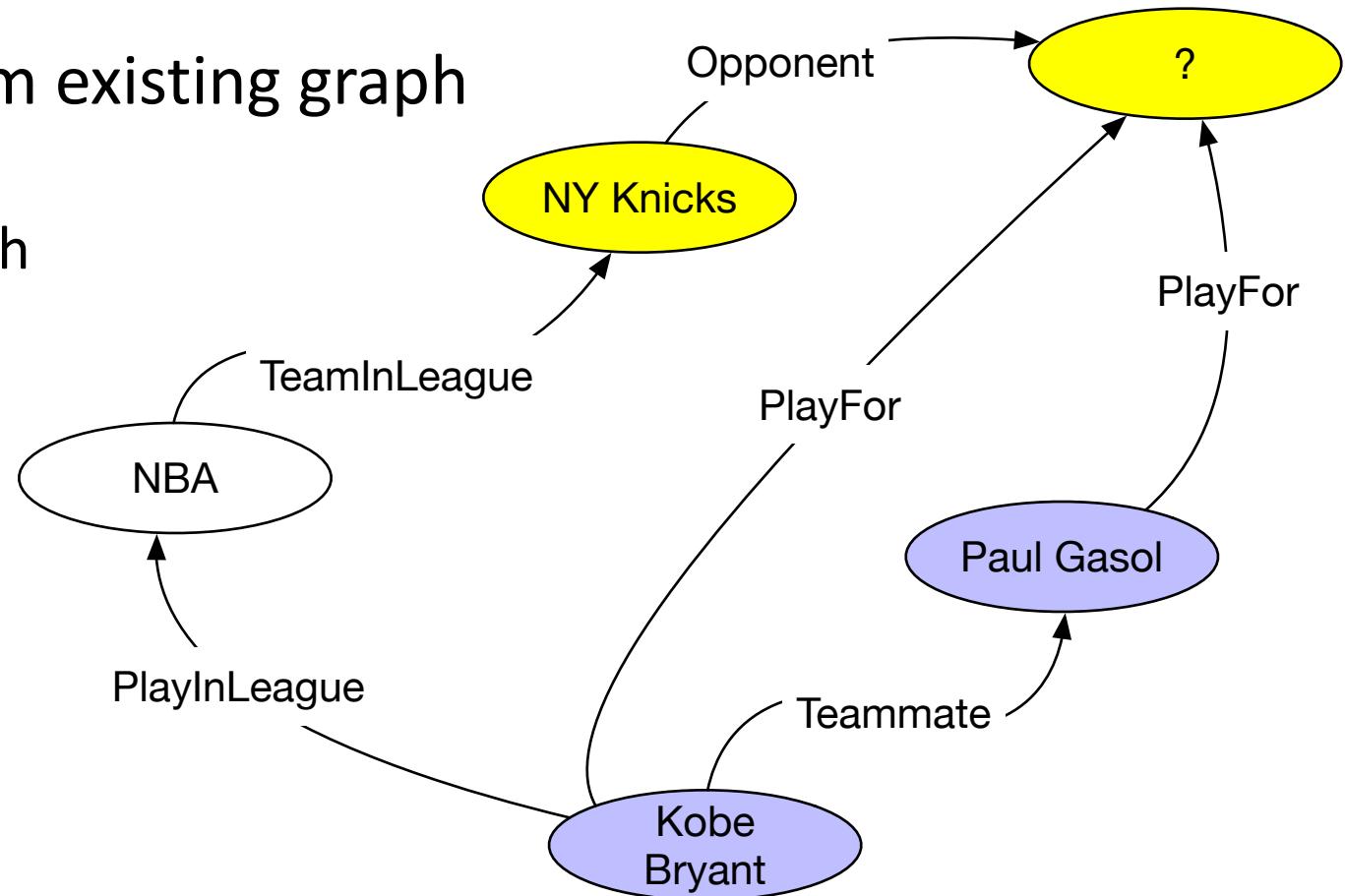
Link Prediction

- **Goal:** Add knowledge from existing graph

- No external source
- Reasoning within the graph

- **Some solutions**

- Rule-based methods
- Probabilistic models
- Factorization models
- Embedding models



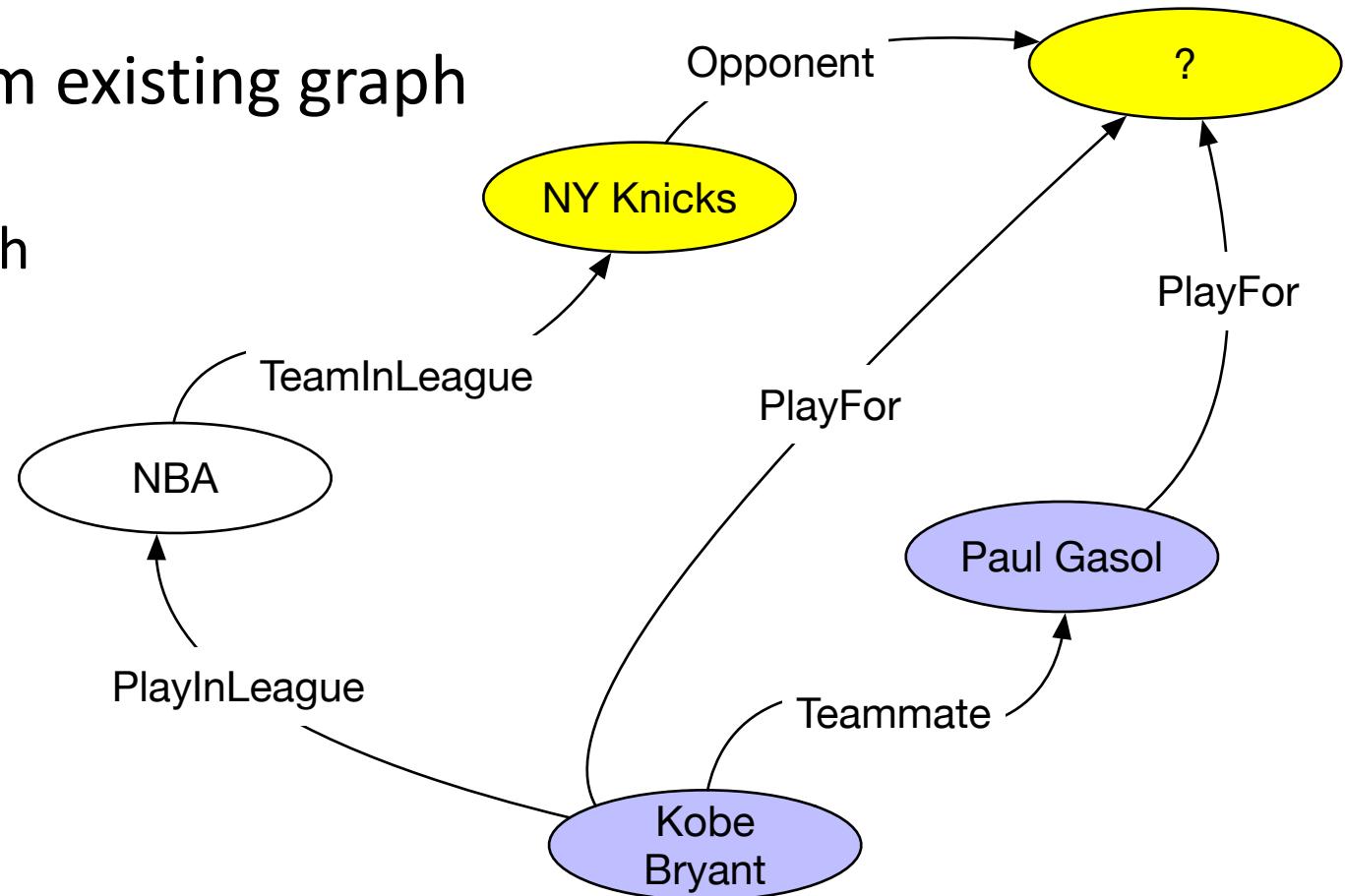
Link Prediction

- **Goal:** Add knowledge from existing graph

- No external source
- Reasoning within the graph

- **Some solutions**

- *Rule-based methods*
- Probabilistic models
- Factorization models
- Embedding models



Link Prediction with Rules

- Goal: We want to make the predictions by applying rules that we extract from the graph => *rule mining*
- Rule mining is an important research topic in several communities
- **Association rule mining:** process of discovering rules that associate items together

$I = \{i_1, i_2, \dots, i_m\}$ is a set of **items**

Transaction $t \subseteq I$ is a subset of items

Transaction database $T = \{t_1, t_2, t_3, \dots, t_n\}$. We assume T is large

An **association rule** is defined as $X \rightarrow Y$ where $X, Y \subseteq I$ and $X \cap Y = \emptyset$

- A famous algorithm for association mining is the *Apriori* algorithm [1]

Link Prediction with Rules

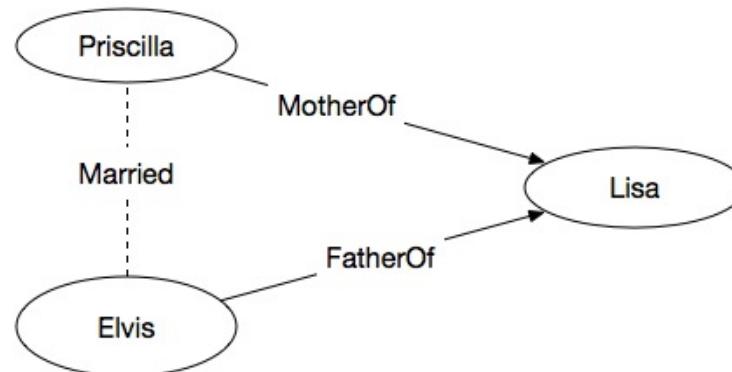
- **Association rule mining:** process of discovering rules that associate items together

In our case, we can assume that a transaction is a connected subgraph

$$t_i = \{mother(Priscilla, Lisa), father(Elvis, Lisa), married(Priscilla, Elvis)\}$$

If we observe that this transaction occurs *often*, then we can infer the rule

$$mother(x, y), father(z, y) \rightarrow married(x, z)$$



Link Prediction with Rules

- **Inductive Logic Programming (ILP)**: field of machine learning that deduces logical rules from ground facts

Input:

- A background knowledge B (mostly a set of atoms, e.g. *married(obama,michelle)*)
- A set of positive examples E^+ (some facts that hold)
- A set of negative examples E^- (some facts that do not hold)

Goal:

Calculate an hypothesis h s.t. $B \wedge h$ entails E^+ and does not entail E^-

- Several ILP systems have been designed to extract rules from the Web, e.g. SHERLOCK [1] and DL-Learner [2]

- [1] S. Schoenmackers, O. Etzioni, D. S. Weld, and J. Davis, “Learning First-order Horn Clauses from Web Text,” in EMNLP 2010, pp. 1088–1098.
- [2] J. Lehmann, “DL-Learner: learning concepts in description logics,” *Journal of Machine Learning Research*, vol. 10, no. Nov, pp. 2639–2642, 2009.

Link Prediction with Rules: AMIE

One big problem of rule association mining and ILP methods is that they either assume CWA or require negative samples

Knowledge Graphs operate under OWA: what is not known might still be true!

Link Prediction with Rules: AMIE

- **AMIE [1]** is a system where rules are extracted with a “rule-association” mining mindset
- Horn rules are constructed incrementally and in parallel
 - They start considering a rule with a single atom, e.g. *fatherOf(X,Y)*
 - Then add another atom which shares at least one variable, e.g. *motherOf(Z,Y)*
 - When the rule is closed (e.g. each variable appears in at least two atoms), it is evaluated using confidence: Only rules with high confidence are kept

$$conf(rule) = \frac{count(body \wedge head)}{count(body)}$$

[1] L. A. Galárraga, C. Teflioudi, K. Hose, and F. Suchanek, “AMIE: Association Rule Mining Under Incomplete Evidence in Ontological Knowledge Bases,” in *WWW*, 2013, pp. 413–422.

Link Prediction with Rules: AMIE

- An important novelty of the AMIE system is the introduction of the notion of *Partial Completeness Assumption*

Definition: Let KG' be the graph which contain all true facts which are not known in KG . The assumption consists of assuming that if $r(x, y) \in KG$ then $\nexists y'. r(x, y') \in KG'$

Standard confidence

$$conf(\vec{B} \Rightarrow r(x, y)) := \frac{supp(\vec{B} \Rightarrow r(x, y))}{\#(x, y) : \exists z_1, \dots, z_m : \vec{B}}$$

PCA-based confidence

$$conf_{pca}(\vec{B} \Rightarrow r(x, y)) := \frac{supp(\vec{B} \Rightarrow r(x, y))}{\#(x, y) : \exists z_1, \dots, z_m, y' : \vec{B} \wedge r(x, y')}$$

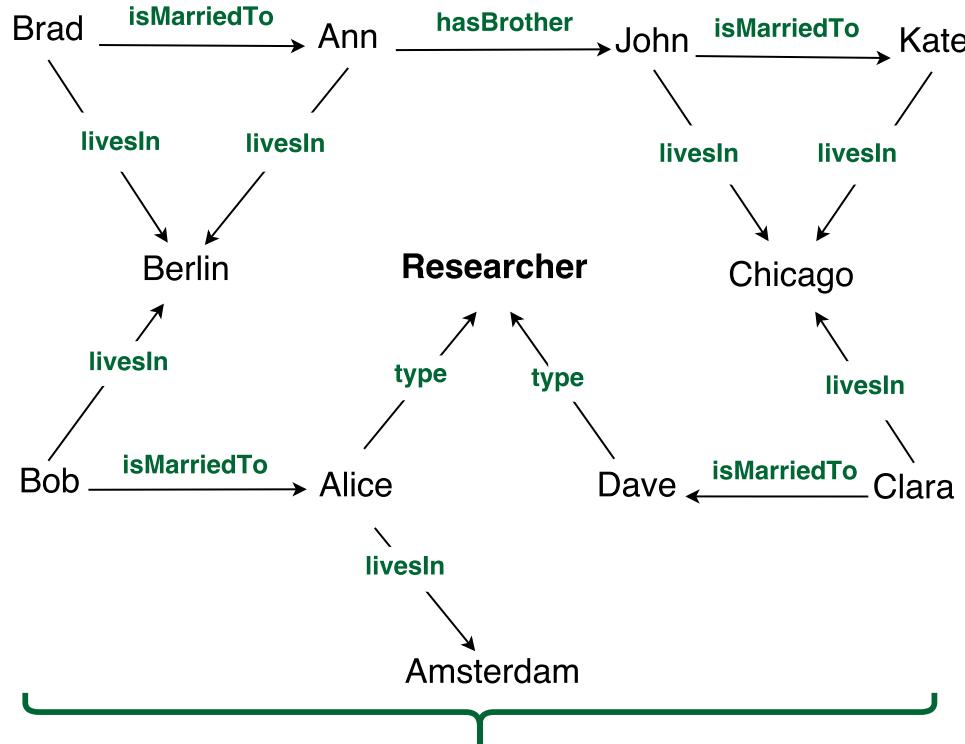
Evaluation

System	Top n	Predictions	Precision
ALEPH	7	2997	27%
AMIE	13	3180	66%
ALEPH	9	5031	26%
AMIE	29	5003	47%
ALEPH	17	8457	30%
AMIE	52	8686	45%

	Top 20 rules	Top 30 rules	All rules
Confidence	0.76	0.63	0.33
PCA Confidence	0.32	0.29	0.29

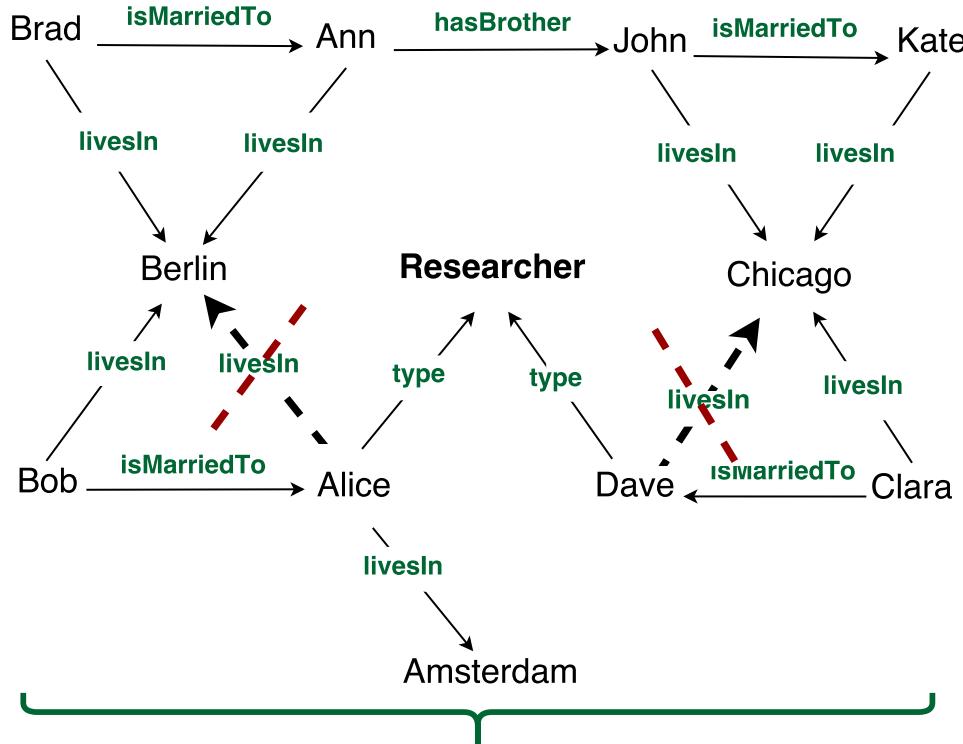
$$isCitizenOf(x, y) \Rightarrow livesIn(x, y)$$
$$hasAdvisor(x, y) \wedge graduatedFrom(x, z) \Rightarrow worksAt(y, z)$$
$$wasBornIn(x, y) \wedge isLocatedIn(y, z) \Rightarrow isCitizenOf(x, z)$$
$$hasWonPrize(x, G. W. Leibniz) \Rightarrow livesIn(x, Germany)$$

Link Prediction with Rules: Non-monotonic Learning



$r : \text{livesIn}(X, Z) \leftarrow \text{isMarriedTo}(Y, X), \text{livesIn}(Y, Z)$

Link Prediction with Rules: Non-monotonic Learning

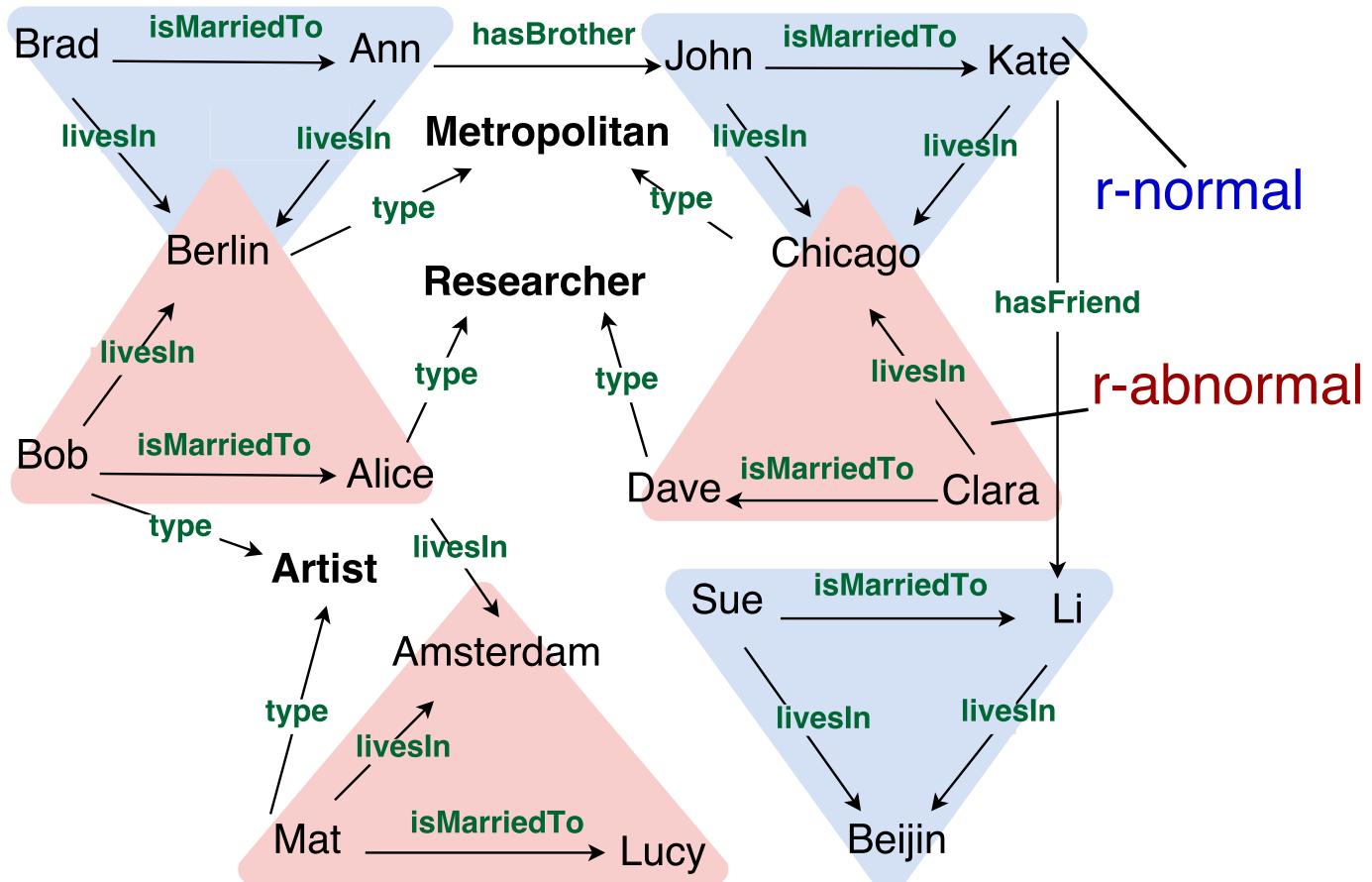


$r : \text{livesIn}(X, Z) \leftarrow \text{isMarriedTo}(Y, X), \text{livesIn}(Y, Z), \text{not researcher}(X)$

Link Prediction with Rules: Non-monotonic Learning

- We developed a technique to mine non-monotonic rules [1].
- Focus on **unary** rules (binary rules are converted with *propositionalization*)
- **Step 1.** Mine association Horn clauses from the KG. We can use FP-Growth, an off-the-shelf popular algorithm
- **Step 2.** Determine normal and abnormal substitutions for every rule (see next slide)
- **Step 3.** Find all exception candidates (=predicates for unary rules) for every rule
- **Step 4.** Rank exception candidates and select the best ones

Abnormal substitutions



$r : \text{livesIn}(X, Z) \leftarrow \text{isMarriedTo}(Y, X), \text{livesIn}(Y, Z)$

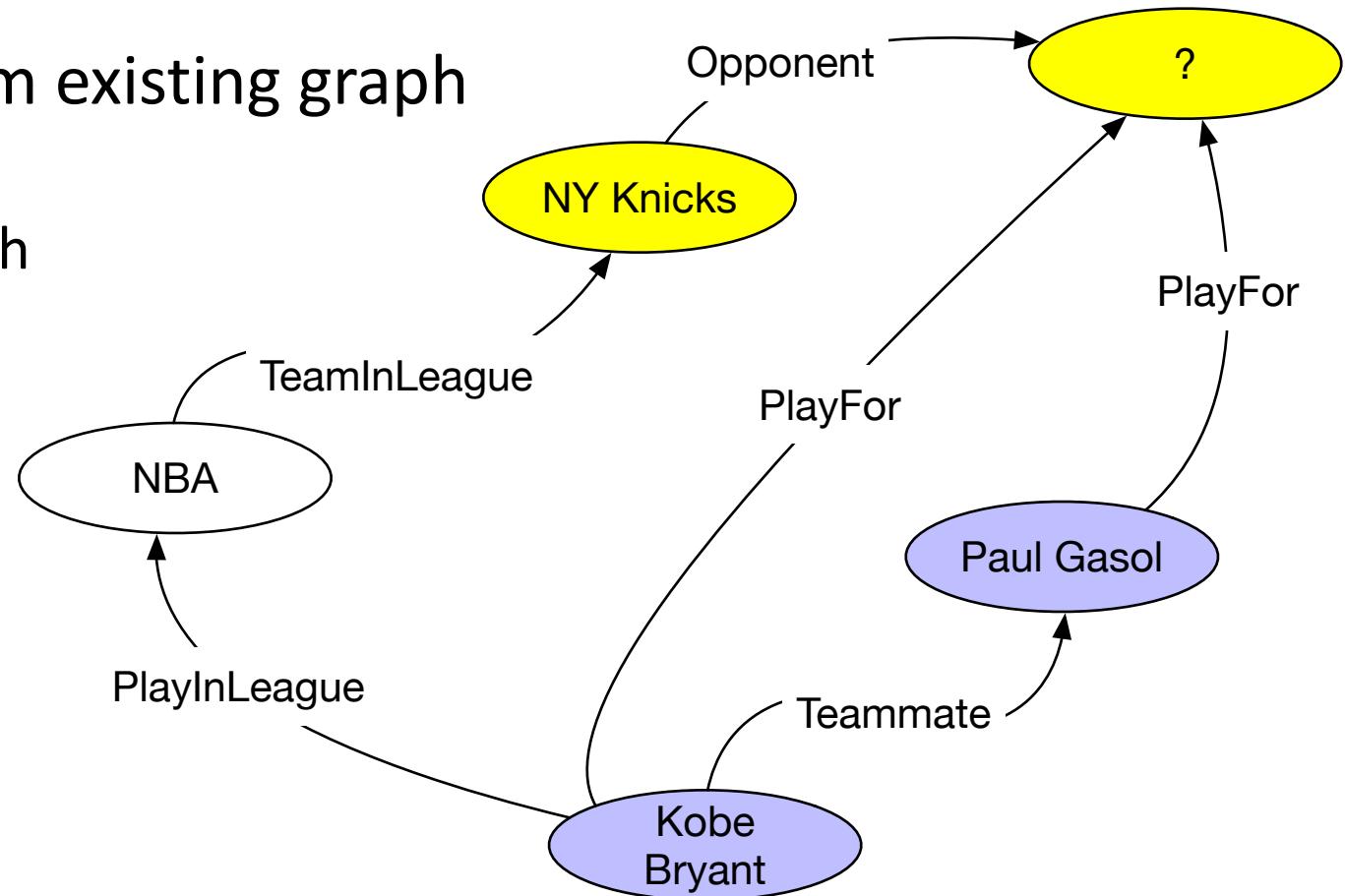
Link Prediction

- **Goal:** Add knowledge from existing graph

- No external source
- Reasoning within the graph

- **Some solutions**

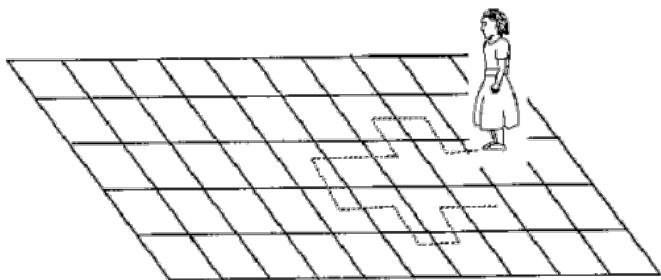
- Rule-based methods
- *Probabilistic models*
- Factorization models
- Embedding models



Link Prediction with Random Walks

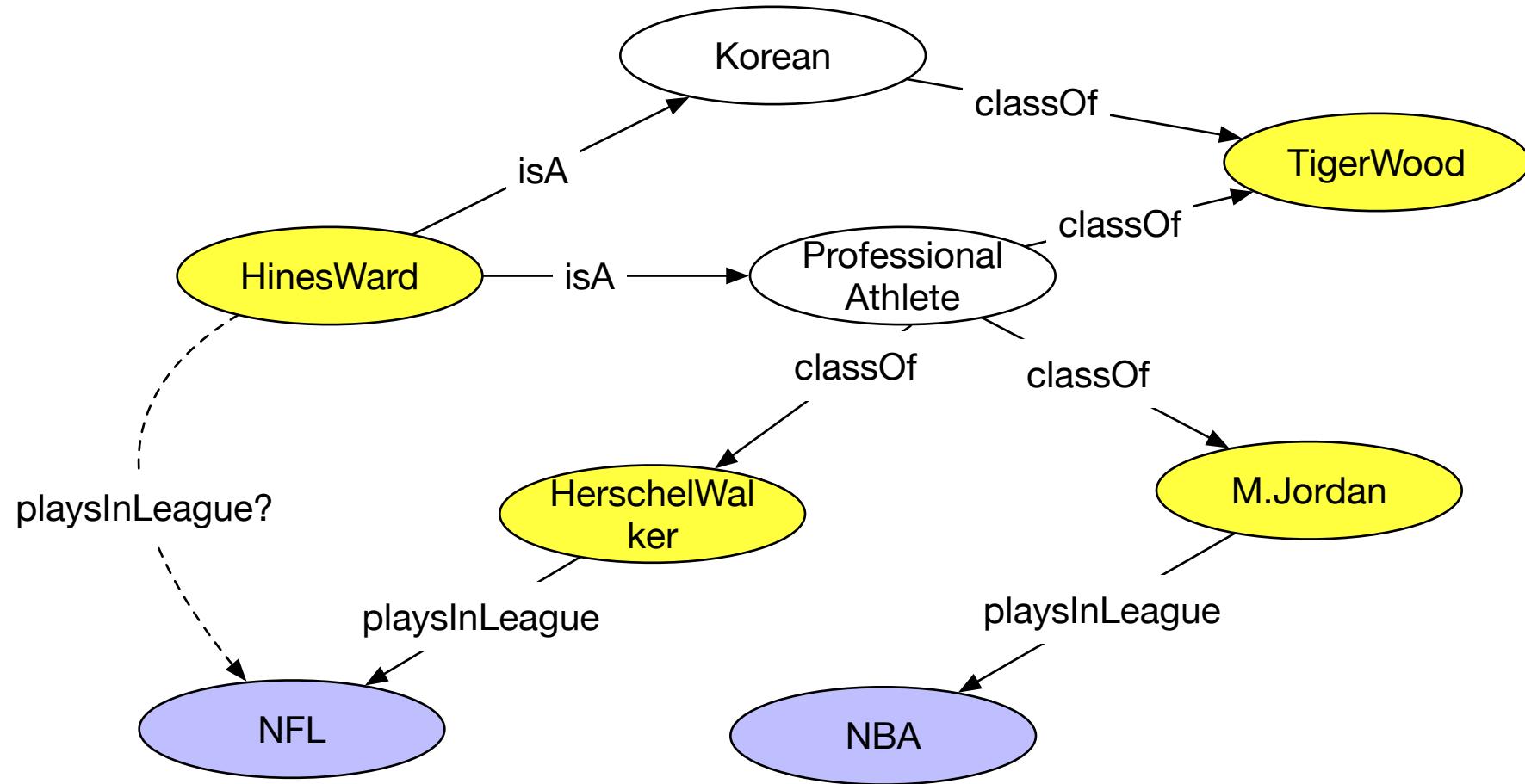
- Rules are expensive to calculate and are “brittle”. [1] proposes a different approach, called Path Ranking Algorithm (PRA), which is a process based on *random walks*

Definition: Given a graph and a starting point, we select a neighbor of it at random, and move to this neighbor; then we select a neighbor of this point at random, and move to it etc. The (random) sequence of points selected this way is a random walk on the graph. [2]

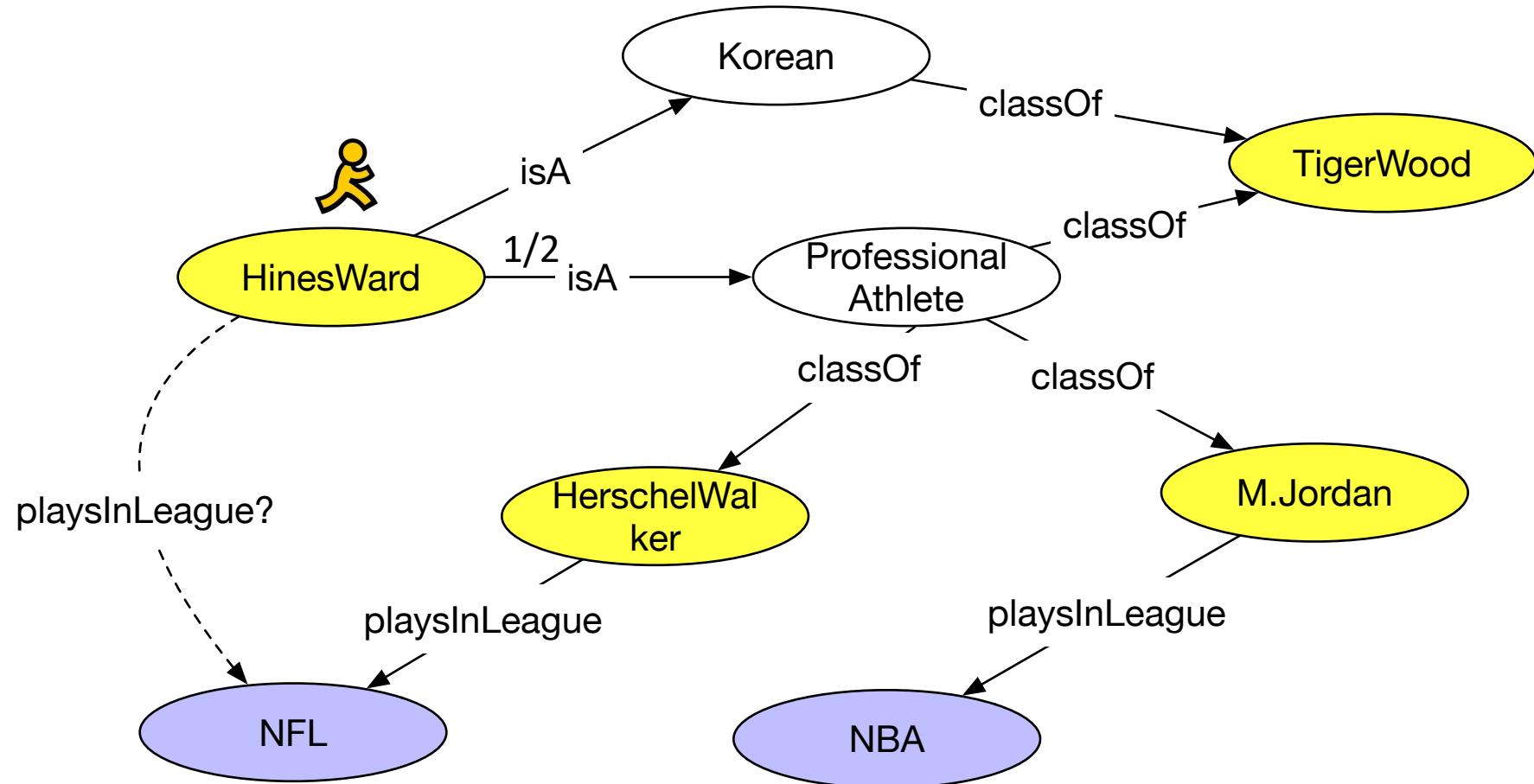


- [1] N. Lao, T. Mitchell, and W. W. Cohen, “Random walk inference and learning in a large scale knowledge base,” in *EMNLP*, 2011, pp. 529–539.
[2] L. Lovász, “Random walks on graphs: A survey,” *Combinatorics, Paul erdos is eighty*, vol. 2, no. 1, pp. 1–46, 1993.

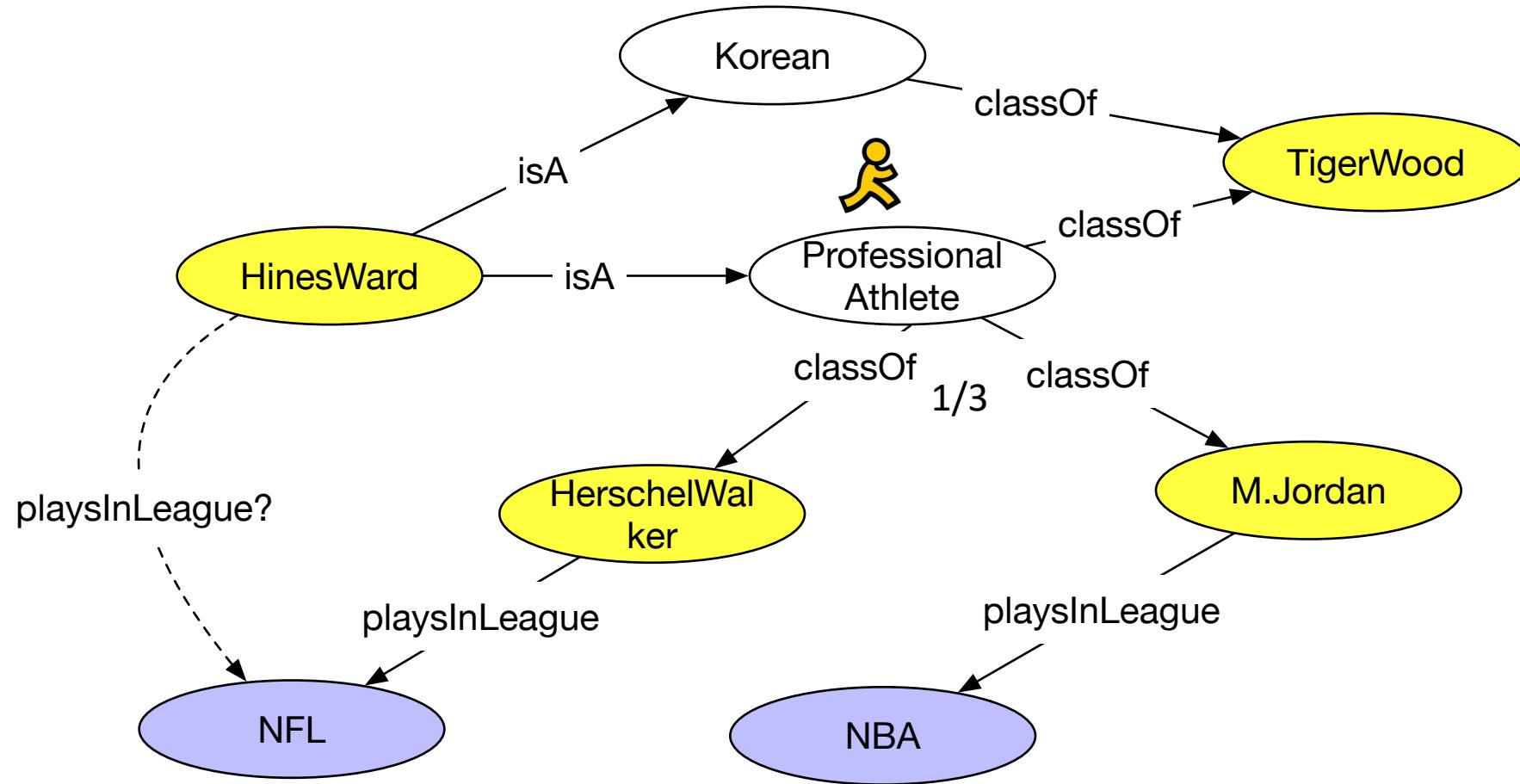
Link Prediction with Random Walks



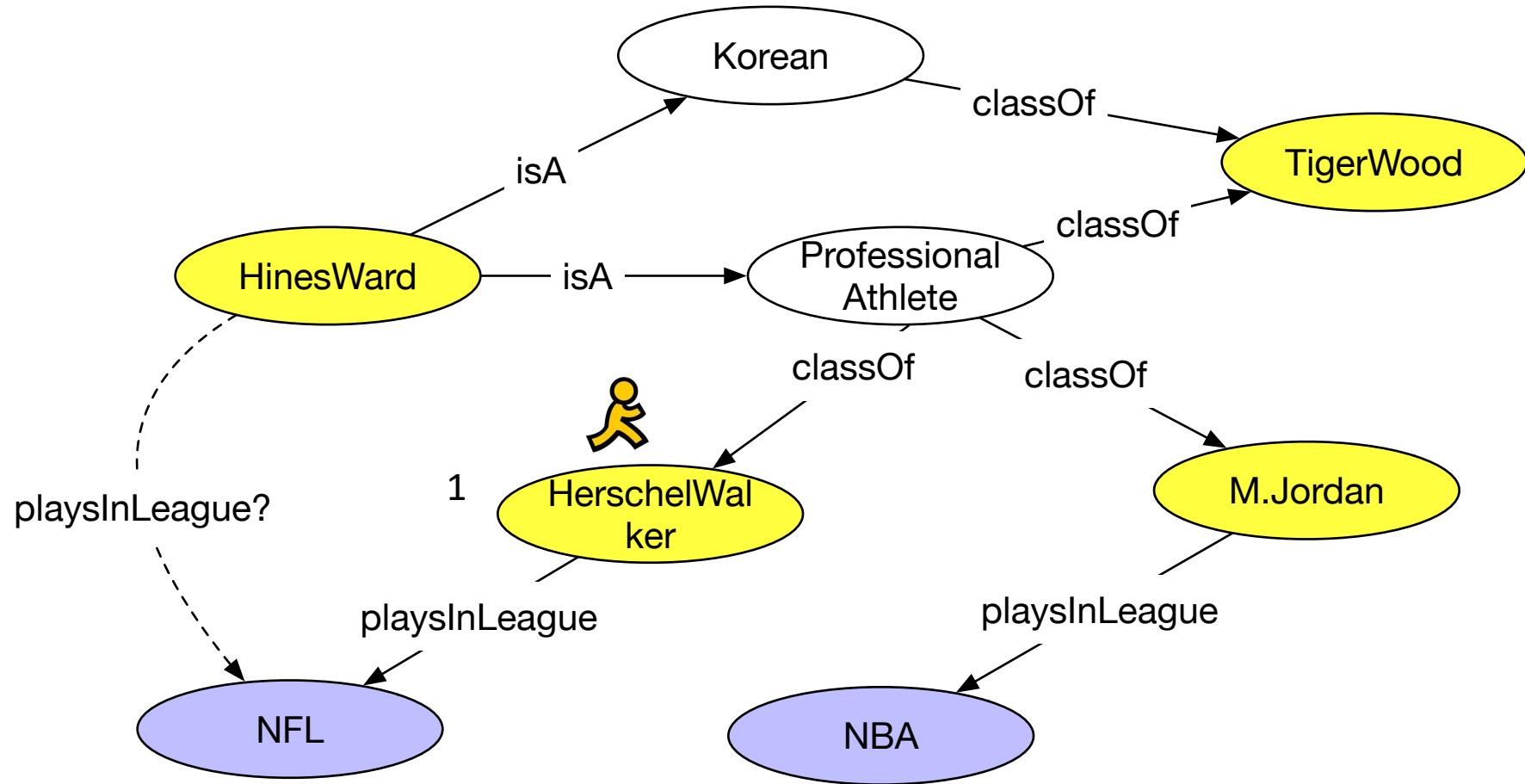
Link Prediction with Random Walks



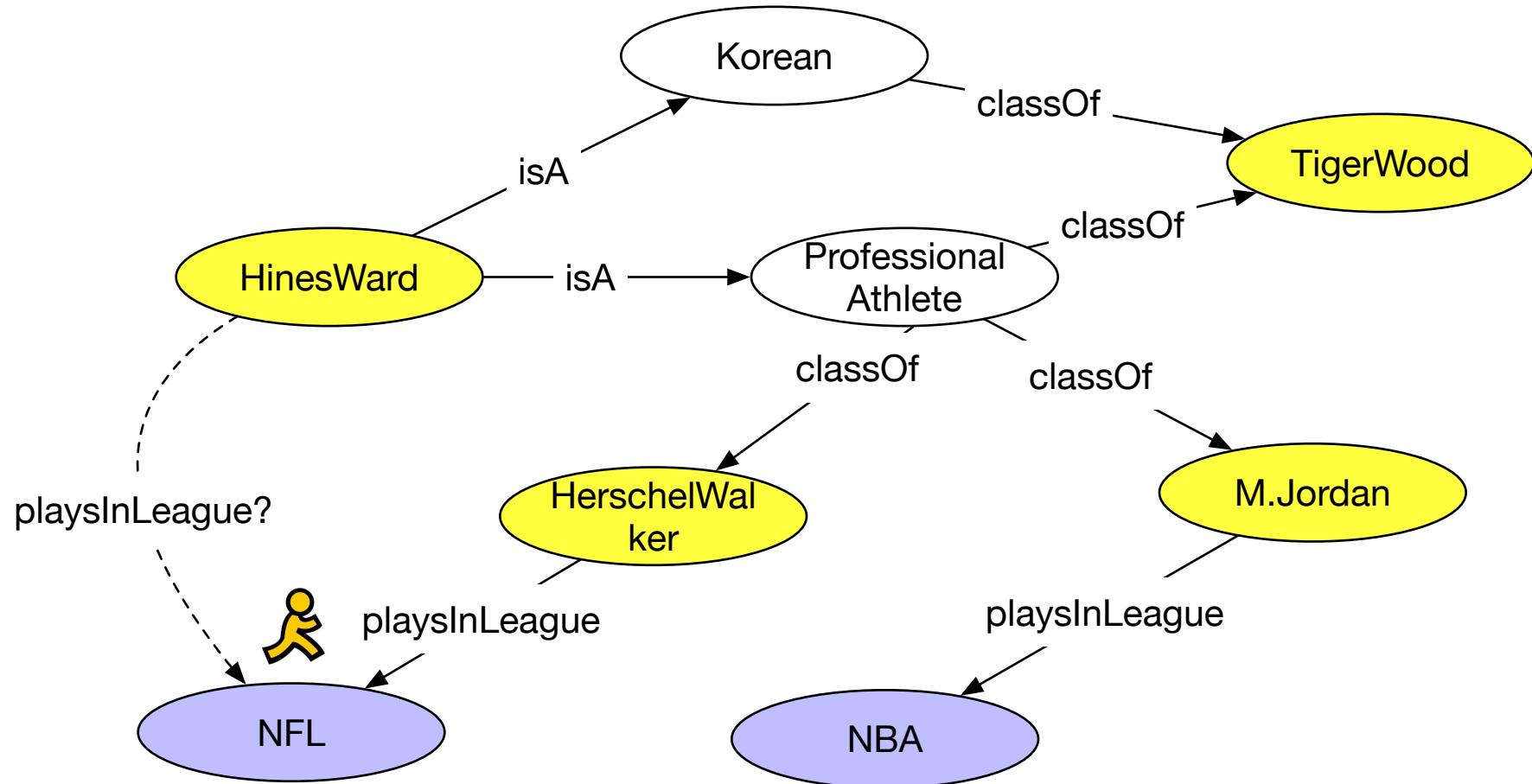
Link Prediction with Random Walks



Link Prediction with Random Walks



Link Prediction with Random Walks



What's the probability that Hines Ward plays in NFL? $1/2 * 1/3$

Link Prediction with Random Walks

- Instead of doing random walks, we can constrain them following the links established by some rules we can mine from the graph

$$isA(x_1, c) \wedge classOf(c, x_2) \wedge playsInLeague(x_2, l) \rightarrow playsInLeague(x_1, l)$$

- We can group the rules by the predicate in the head (e.g. *playInLeague*), and perform many random walks to estimate the probabilities with some training triples

Link Prediction with Random Walks

- Suppose we have n rules with *playsInLeague* has head. We can construct a feature vector of n elements and train a linear regression model with the objective function

$$score(e; s) = \sum_{P \in P} h_{s,P}(e) \Theta_P \text{ where } h_{s,P}(e) \text{ represents the probability of reaching } e \text{ from } s \text{ following the path } p$$

- The algorithm constructs one model per relation

Link Prediction with Random Walks

Evaluation

- Used NELL KG
- Relied on human assessment (5 humans, 75% agreement) on sample queries of the type $R(x,?)$
- If we look at the top-10 results, the method performs as good as a N-FOIL -- an inductive baseline
- If we move to the top-100 the results get significantly better

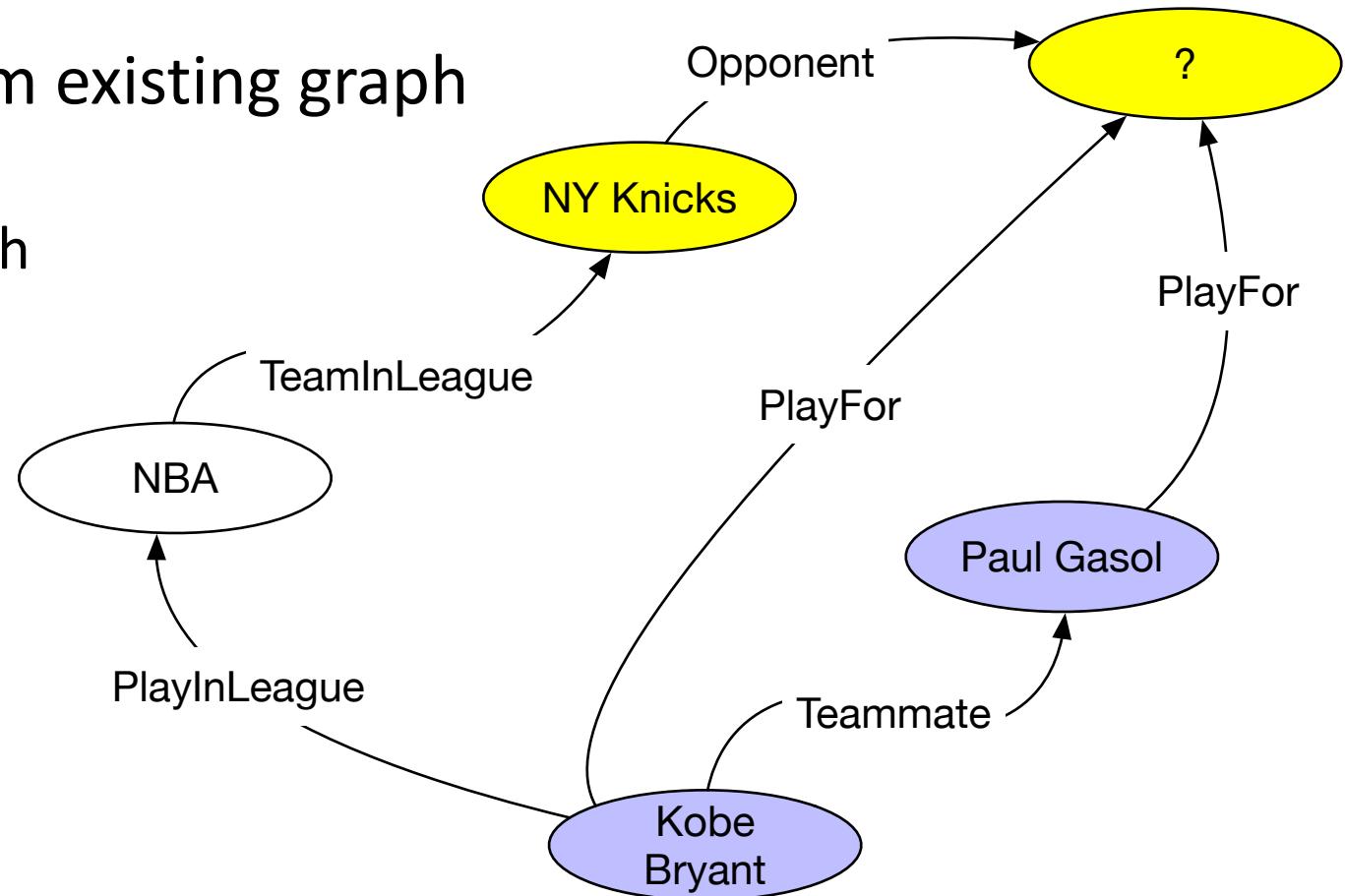
Link Prediction

- **Goal:** Add knowledge from existing graph

- No external source
- Reasoning within the graph

- **Some solutions**

- Rule-based methods
- Probabilistic models
- *Factorization models*
- Embedding models



Link Prediction with Tensors

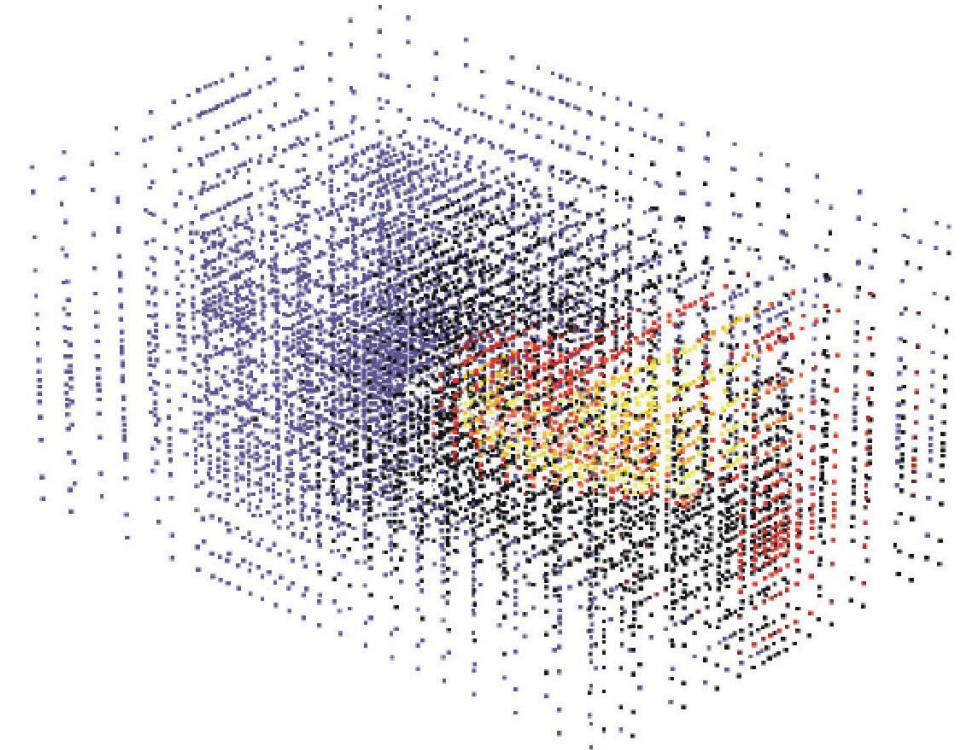
- **Definition:**

- 1-dimension array:** vector

- 2-dimension array:** matrix

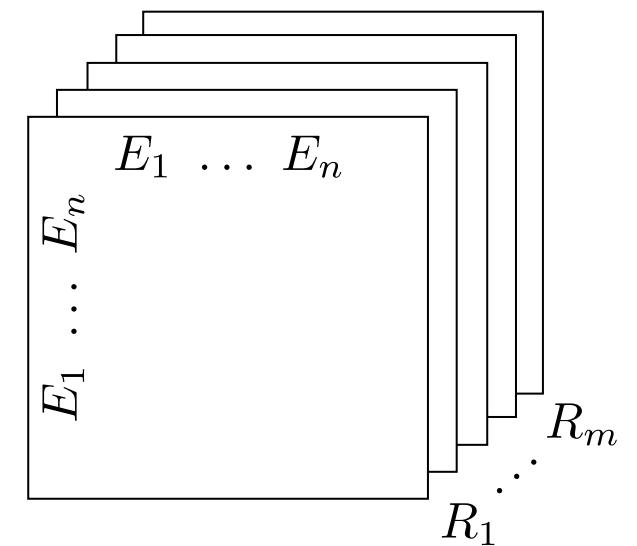
- 3 or more-dimensions array:** tensor

Multiple binary relations can be expressed naturally as a three-way tensor



Link Prediction with Tensors

- Tensors are well-studied data structures in Machine Learning [1]
- Suppose our knowledge graph has m relations and n entities
- We can assign two *modes* (=ways) to the entities and one *mode* to the relations
- The entries of the tensor are 1 when a relation between two entities exists. Otherwise, the value is 0
- E.g. $(2,3,1) = 1$ if $R_1(e_2,e_3)$ exists in the KG



Link Prediction with Tensors

- **RESCAL [1]** is a method that performs link predictions by *factorizing* the tensor

$$X_k \Rightarrow A \times R_k \times A^T$$
$$X_k \approx A R_k A^T$$

The diagram shows the decomposition of a tensor X_k into three components. On the left, a large rectangular block labeled X_k is shown with several smaller rectangles stacked behind it, representing its multi-dimensional nature. An arrow points from this to the right side of the equation. On the right, the tensor is represented as a product of three matrices: A , R_k , and A^T . The matrix A is a tall rectangle, R_k is a shorter rectangle with multiple horizontal layers, and A^T is a wide rectangle. The multiplication symbol (\times) is placed between A and R_k , and between R_k and A^T .

- Let n be the number of entities and r the number of dimensions
- A is a $n \times r$ matrix that represents the *global* entity-latent-components
- R_k is a asymmetric $r \times r$ matrix that specifies the interaction of the latent components per predicate

Link Prediction with Tensors

The factorization is translated as an optimization problem

$$\min_{A, R_k} \text{loss}(A, R_k) + \text{reg}(A, R_k)$$

where *loss* is the loss function

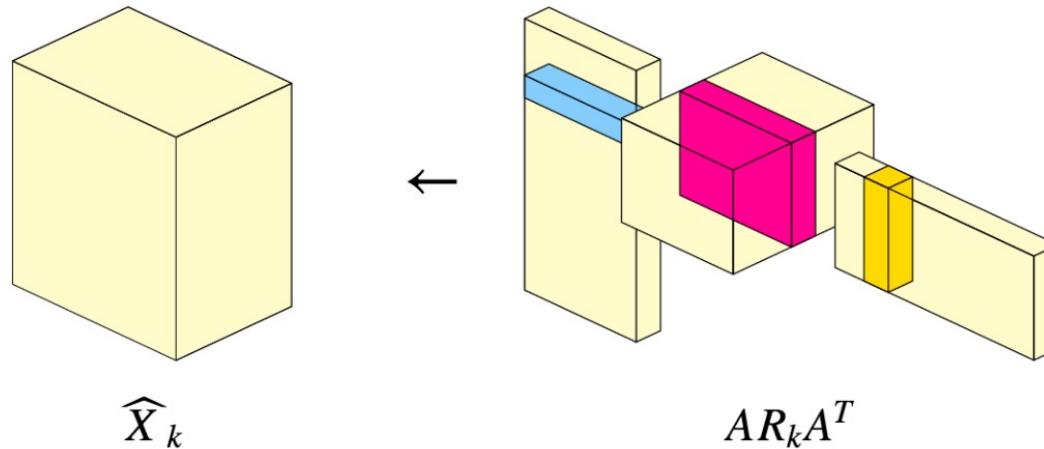
$$\text{loss}(A, R_k) = \frac{1}{2} \sum_k \| \mathcal{X}_k - A R_k A^T \|_F^2$$

and *reg* is the regularization function

$$\text{reg}(A, R_k) = \frac{1}{2} \lambda \left(\|A\|_F^2 + \sum_k \|R_k\|_F^2 \right)$$

The calculation is performed using the Alternating-least Square algorithm (ALS)

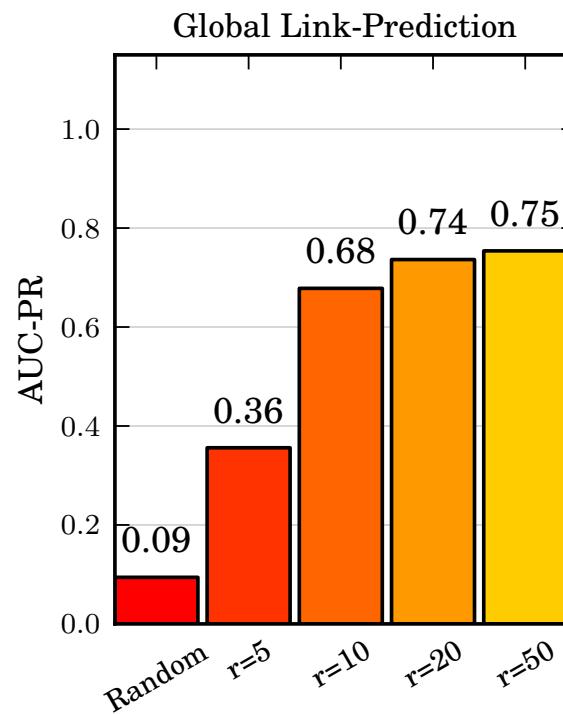
Link Prediction with Tensors



- **Basic procedure** of link prediction with RESCAL
 - First, we compute the “rank-reduced” factorization
 - Second, entries after reconstruction (i.e. $\widehat{x}_{ijk} = \vec{a}_i^T \textcolor{magenta}{R}_k \vec{a}_j$) can be regarded as *confidence values* of the model that the corresponding triple exists

Factorizing YAGO

- As an experiment, the authors of [1] took a subset of YAGO of 2.6M entities and 64M facts
- Type of queries: $p(?,o)$ over 38 relations



[1] M. Nickel, V. Tresp, and H.-P. Kriegel, “Factorizing YAGO: Scalable Machine Learning for Linked Data,” in *WWW*, 2012, pp. 271–280.

Link Prediction with Tensors

- Collective learning is performed via the entities' latent-component representation
- A is a matrix that contains **embeddings** of the entities. We can use it in a similar way as for word embeddings (word2vec)
- A major problem of RESCAL is that it has to learn a very large number of parameters (esp. R). Not so cheap to compute

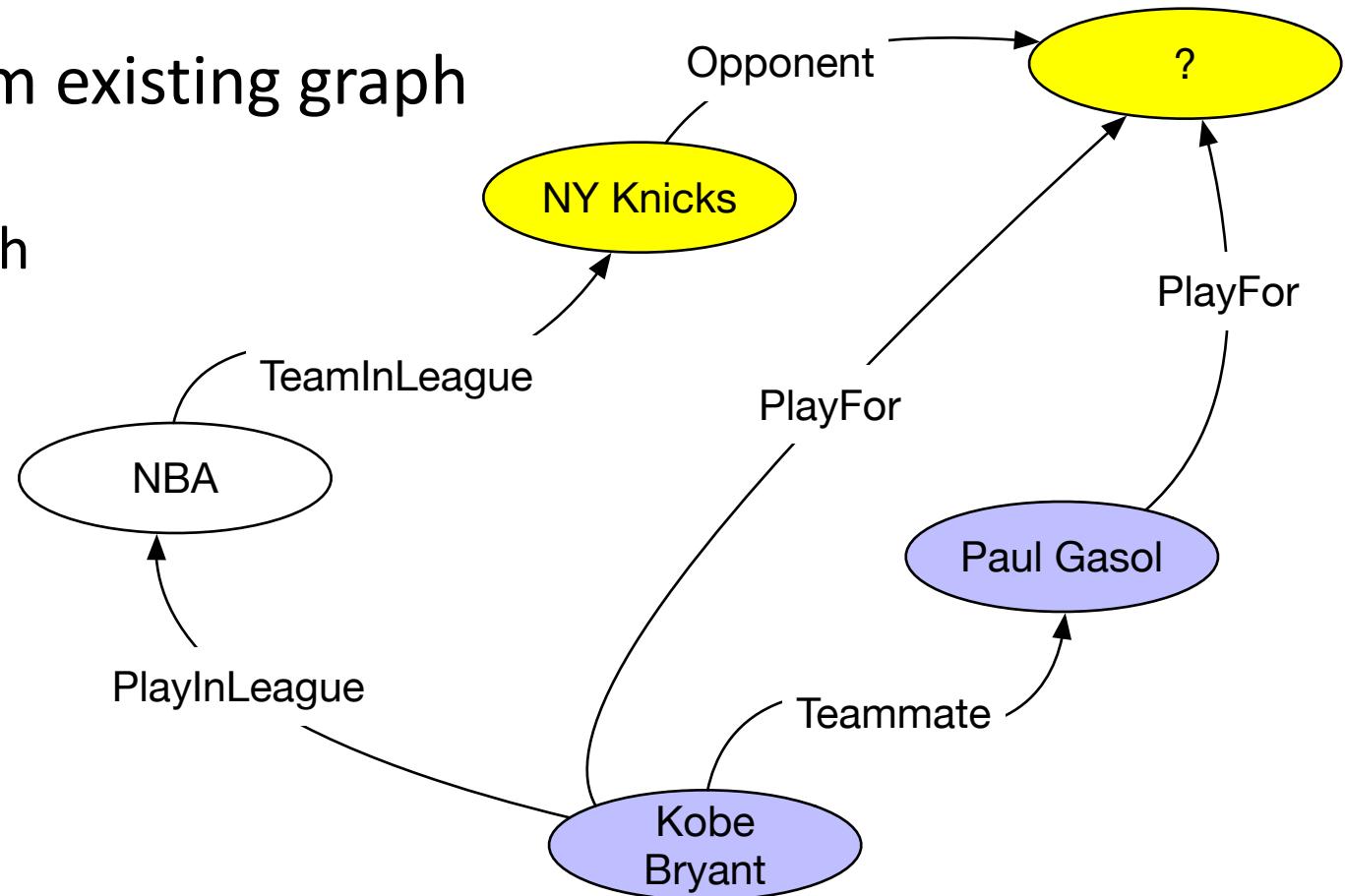
Link Prediction

- **Goal:** Add knowledge from existing graph

- No external source
- Reasoning within the graph

- **Some solutions**

- Rule-based methods
- Probabilistic models
- Factorization models
- *Embedding models*



Latent-features model

- *RESCAL* is a latent-features model. Predictions are conditionally independent w.r.t. a given set of global latent features [1]
- Number of parameters in the RESCAL model = $rd^2 + nd$
where d = # dimensions, n = # entities, r = # relations
- Can we lower the number of parameters?

Latent distance models

- Latent distance models are a subclass of latent feature models where the likelihood of links depends on the distance between the latent representation of the entities
- Consider one relation (e.g. *friendsOf*) and two entities (e.g. Mark and John)
 - $\text{friendsOf}(\text{Mark}, \text{John}) = d(\mathbf{e}_{\text{mark}}, \mathbf{e}_{\text{john}})$ where d is a arbitrary distance measure (e.g. Euclidean distance) [1]
- See the connection with word2vec?

TransE

- TransE [1] is a latent distance model with less parameters than RESCAL
- Every entity and relation get a vector of d elements.
- The distance in TransE is: $d(\mathbf{e}_i + \mathbf{r}_k, \mathbf{e}_j)$ which reads as: “given an entity and a relation, what is the closest vertex to their sum?
- The model has $rd + nd$ parameters, much less than RESCAL

[1] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, “Translating Embeddings for Modeling Multi-relational Data,” in *NIPS* 2013.

TransE

Overview training

1. Create a vector of d elements for each entity $e \in E$ and $r \in R$
2. Take as input a set of triples $(h, l, t) \in S$ (S is our knowledge graph)
3. Minimize the function with stochastic gradient descent)

$$\mathcal{L} = \sum_{(h, \ell, t) \in S} \sum_{(h', \ell, t') \in S'_{(h, \ell, t)}} [\gamma + d(\mathbf{h} + \boldsymbol{\ell}, \mathbf{t}) - d(\mathbf{h}' + \boldsymbol{\ell}, \mathbf{t}')]_+$$

where

$$S'_{(h, \ell, t)} = \{(h', \ell, t) | h' \in E\} \cup \{(h, \ell, t') | t' \in E\}$$

TransE

Overview testing

1. Given an “unseen” triple $(h, l, t) \in S'$, we calculate $d(\mathbf{h} + \mathbf{l}, \mathbf{e}_v)$ for all $\forall v \in E$ and rank the results
2. Then, we check the average position of t (mean rank) or whether t occurs in the top k positions (hits@ k).
3. We can look at the position of t w.r.t. all other entities (*raw position*) or remove all the “answers” that are already explicit in the KB (*filtered position*)

TransE

- Evaluation

DATASET <i>METRIC Eval. setting</i>	WN				FB15K				FB1M	
	MEAN RANK		HITS@10 (%)		MEAN RANK		HITS@10 (%)		MEAN RANK	HITS@10 (%)
	<i>Raw</i>	<i>Filt.</i>	<i>Raw</i>	<i>Filt.</i>	<i>Raw</i>	<i>Filt.</i>	<i>Raw</i>	<i>Filt.</i>	<i>Raw</i>	<i>Raw</i>
Unstructured [2]	315	304	35.3	38.2	1,074	979	4.5	6.3	15,139	2.9
RESCAL [11]	1,180	1,163	37.2	52.8	828	683	28.4	44.1	-	-
SE [3]	1,011	985	68.5	80.5	273	162	28.8	39.8	22,044	17.5
SME(LINEAR) [2]	545	533	65.1	74.1	274	154	30.7	40.8	-	-
SME(BILINEAR) [2]	526	509	54.7	61.3	284	158	31.3	41.3	-	-
LFM [6]	469	456	71.4	81.6	283	164	26.0	33.1	-	-
TransE	263	251	75.4	89.2	243	125	34.9	47.1	14,615	34.0

TransE

Open problems: Negation

- In TransE, we create random triples to use as evidence
- However, the OWA does not allow to make such assumption
- Some approaches proposed to leverage the schema to create negative examples:
 - *Functional properties* => if $\langle a \rangle \text{ } \langle \text{bornIn} \rangle \text{ } \langle b \rangle$, then $\langle a \rangle \text{ } \langle \text{bornIn} \rangle \text{ } \langle c \rangle$ where $b \neq c$ is surely false
 - Class hierarchies => if $\langle p \rangle$ has a range C, then we can pick an instance of a class $D \neq C$ to construct negative samples
- These approaches introduce biases and do not work well because some classes are too popular and will get only few negative examples

TransE

Open problems: Data sparsity

- Some relations appear only a few times to be properly trained.
- This is an open problem. There is not yet an effective solution for that.

TransE

Open problems: Lack of established evaluation methodology

- There is not a single metric that can be representative for the overall performance of the method
- Datasets are too different from each others. TransE returns very different results depending on the used KG
 - Freebase 74% (hit@10)
 - Wordnet 94%
 - YAGO 22%
- The latest research is focusing on trying to replicate reasoning (logical rules) using these models!

TransE

Open problems: Scalability

- Hogwild! [1] is a well-known technique to parallelize the learning of *sparse* models

Algorithm 1 HOGWILD! update for individual processors

```
1: loop
2:   Sample  $e$  uniformly at random from  $E$ 
3:   Read current state  $x_e$  and evaluate  $G_e(x_e)$ 
4:   for  $v \in e$  do  $x_v \leftarrow x_v - \gamma G_{ev}(x_e)$ 
5: end loop
```

E : training set

x_e : parameters that relate to e

G : gradient of the loss function on e

TransE

Open problems: Scalability

- Hogwild! [1] is a well-known technique to parallelize the learning of *sparse* models.
- It can be trivially expanded to an arbitrary number of threads. It works only if the problem is *sparse*
- Are KG embeddings an example of a sparse problem? Preliminary evaluation indicate it is