**Big Data Analytics Summer School, June 2020** 

# Knowledge Graph Embedding

#### Dr. Sahar Vahdati







## The Challenges of Knowledge Graphs (KGs)

Methods



The represented knowledge can be **incorrect**.

The represented knowledge can be (incomplete)

#### **Knowledge Graph Refinement**

ML-based Approaches for **Completion** of Knowledge Graphs

Error detection Approaches for **Correction** of Knowledge Graphs

Knowledge Graph Refinement: A Survey of Approaches and Evaluation Methods

Editor(s): Philipp Cimiano, Universität Bielefeld, Germany Solicited review(s): Natasha Noy, Google Inc., USA; Philipp Cimiano, Universität Bielefeld, Germany; two anonymous reviewers

Heiko Paulheim,

Data and Web Science Group, University of Mannheim, B6 26, 68159 Mannheim, Germany E-mail: heiko@informatik.uni-mannheim.de

## **Graph Completion with Machine Learning**

Hypothesis: ML can help to learn new knowledge from already existing ones!

Machine learning builds a <u>mathematical model</u> based on sample data ( set of attributes as training data D\_n with n samples).

**Supervised Learning** 

**Un-Supervised Learning** 

Semi-Supervised Learning

$$D_n = \{x_i, y_i\}_{i=1}^n \qquad D_n = \{x_i\}_{i=1}^n$$

$$D_n = \{x_i, y_i\}_{i=1}^{n_1} \cup \{x_i\}_{i=n_1+1}^n$$

#### **MI-based Mathematical Model for Learning**



## Requirement

The input and output of such a Mathematical Model should be numeric!



## What is an embedding?

**Wikipedia:** In mathematics, an embedding is one instance of some mathematical structure contained within another instance, such as a group that is a subgroup.

**Google**: An embedding is a relatively low-dimensional space into which you can translate high-dimensional vectors.





#### **Graph Embedding**

	•				• •	•	• • •			
/		•			•					-
000	0001		1000	110	0101010	0011100	1101101	0110111	1011011	1110000
000	0010		1000	101	0101001	0011111	1101110	0110100	1011000	1110011
000	0100		1000	011	0101111	0011001	1101000	0110010	1011110	1110101
000	1000		1001	111	0100011	0010101	1100100	0111110	1010010	1111001
001	0000		1010	111	0111011	0001101	1111100	0111110	1001010	1100001
010	0000		1100	111	0001011	0111101	1001100	0010110	1111010	1010001
100	0000		0000	111	1101011	1011101	0101100	1110110	0011010	0110001
000	0011		1000	100	0101000	0011110	1101111	0110101	1011001	1110010
000	0110		1000	001	0101101	0011011	1101010	0110000	1011100	1110111
000	1100		1001	011	0100111	0010001	1100000	0111010	1010110	1111101
001	1000		1011	111	0110011	0000101	1110100	0101110	1000010	1101001
000	1010		1001	101	0100001	0010111	1100110	0111100	1010000	1111011
001	0100		1010	011	0111111	0001001	1111000	0100010	1001110	1100101
001	0010		1010	101	0111001	0001111	1111110	0100100	1001000	1100011
000	1110		1001	001	0100101	0010011	1100010	0111000	1010100	1111111





#### **Node Embedding**





#### **Knowledge Graph Embedding**



#### Symbolic Representation

**Vector Representation** 

## **Knowledge Graph Embedding Models**

Goal: Predicting missing links between nodes!



## Learning Steps of a KGE

Assign random vectors to each entity instance or relation!

Give these vectors to the KGE model

Let the model learn the embeddings

Return degree of plausibility

Sahar Vahdati - Knowledge Graph Embeddings - BDA School 2020

Optimize the learned embeddings

## **Designing Knowledge Graph Embeddings**



## Datasets/KGs

Dataset: Knowledge Graph

Triple facts are shown in the form of (h,r,t)

Vectors are shown in the form of (h, r, t)

(Barack Obama, WasBornIn, Honolulu)

(Barack Obama, WasBornin, Honolulu)\*

x i=  $(\vec{h}, \vec{r}, \vec{t})^*$ 

y i= 1 (True), O(False) with a degree of plausibility.

Obama



#### **Learning Datasets**

Train: Use to learn embeddings.

**Test**: Check the correctness of the results on one gold standard.

Validation: Check the behavior of the model in different situations.

**Negative Sampling**: let the model also learn from incorrect samples.

#### How to arrange the vectors?



Symbolic Representation

**Vector Representation** 

#### **Score function**

Provides the degree of correctness for a triple:

**TransE Score Function** 

 $f_r^{TransE}(h,t) = \|\overrightarrow{h+r} - \overrightarrow{t}\|$ 

 $f_{wasBornIn}(\overrightarrow{BarackObama}, \overrightarrow{Honolulu}) = 0, \rightarrow (BarackObama, wasBornIn, Honolulu) is True$  $f_{wasBornIn}(\overrightarrow{BarackObama}, \overrightarrow{Berlin}) = 9 \rightarrow (BarackObama, wasBornIn, Berlin) is False$ 

#### **Purpose: Vector Arrangements**

TransE Embedding Model: subject + property = object



#### **Purpose: Vector Arrangements**

**RotatE Score Function** 

 $subject \bigcirc property = object$ 



## Optimization

To optimize the random vectors in a way that the scores are closer to what we have in the KG and the model knows about them!

In order to design an optimizer, we need to know what do we want to optimize with a criteria.

Margin: The difference between correctness of positive triples and negative one

#### **Loss Function**

Adjust the embedding vectors for entities and relations i.e., (h, r, t) to enforce the criteria.



## **Designing Knowledge Graph Embeddings**



#### Example use-case Co-author Recommendation

## **Heterogeneous Scholarly Metadata and Enquiries!**

- Enormous resources / metadata providers
- Different types
- Different formats
- Large-scale metadata
- Diverse
- Broad scientific domains
- Value and importance
- Increasing publishing rate



## **Machine Learning Support for Scholarly Domain**

Many metaresearch enquiries of scholars remain unrevealed.

Link prediction for recommendation-based services:

- Who can be the best candidate for collaboration?
- Which group/university can be the best candidate for collaboration?
- Which groups in different fields can work on an effective research together?
- ...

ML-based approaches for link prediction: Knowledge Graphs Embedding

## **Ontology of an Scholarly Knowledge Graph**



#### **Knowledge Graphs Embedding Models**

TransE:  $f_r(h,t) = \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|$ 

RotatE: 
$$f_r(h,t) = \|\mathbf{h} \circ \mathbf{r} - \mathbf{t}\|$$



$$\mathsf{Loss Functions} \left\{ \begin{aligned} \mathcal{L} &= \sum_{(h,r,t)\in S^+} \sum_{(h',r',t')\in S^-} [f_r(h,t) + \gamma - f_r(h',t')]_+ \\ \mathcal{L}_{RS} &= \sum \sum [f_r(h,t) + \gamma - f_r(h',t')]_+ + \lambda [f_r(h,t) - \gamma_1]_+ \end{aligned} \right.$$

#### **Loss function**

Loss function helps to obtain embedding vectors for entities and relations i.e., (h, r, t).



#### **Loss function**

Problems of Margin Ranking loss:

• Margin sliding: there are infinite solutions for the score of positive and negative samples

$$f_r(h,t) = 0 \text{ and } f_r(h',t') = \gamma, \text{ or}$$

$$f_r(h,t) = \gamma \text{ and } f_r(h',t') = 2\gamma, \text{ or}$$

$$\vdots$$

$$f_r(h,t) = (n-1)\gamma \text{ and } f_r(h',t') = n\gamma$$

• Embeddings are adversely affected by false negative samples.

#### **Loss function**

Limit-based scoring loss:



## **Challenge of Co-author Recommendation**

Problems:

- Many to Many relation
  - Coauthor relation
  - o Citations
- In the existence of many-to many relations, the rate of false negative samples increases.

- Generating negative samples are based on a random corruption.
- Assuming that N= 1000 is the number of all authors in a SKG, the probability of generating false negatives for an author with 100 true or sensible but unknown collaborations becomes 100/1000= 10%.

## **Optimization of Margin ranking loss**



#### **Soft Marginal Loss**



#### **Experiments and Results**

		FB15	ők	WN18		
	FMR	FMRR	FHits@10	FMR	FMRR	FHits@10
TransE [4]	125	_	47.1	251	-	89.2
ComplEx [19]	106	67.5	82.6	543	94.1	94.7
ConvE [7]	51	68.9	85.1	504	94.2	95.5
RotatE[15]	49	$\overline{68.8}$	85.9	388	94.6	$\overline{95.5}$
TransE-RS [28]	<u>38</u>	57.2	82.8	<u>189</u>	$\overline{47.9}$	$\overline{95.1}$
TransE-SM	46	64.8	87.2	201	47.8	95.2
RotatE-SM	<b>40</b>	<b>70.4</b>	<b>87.2</b>	213	94.7	96.1

#### **Collaboration Recommendations**



## **More Application of KGEs**

**Question Answering Systems** 

**Recommendation Systems** 

**Prediction Systems** 

٠

.

•

## Workshop of Knowledge Representation & Representation Learning

ECAI 2020 in Santiago de Compostela, June 2020

Thank you!