

Joint Word and Entity Embeddings for Entity Retrieval from a Knowledge Graph

Fedor Nikolaev^{1,2} and Alexander Kotov¹

¹ Textual Data Analytics (TEANA) Lab, Wayne State University, USA

² Kazan Federal University, Russia



Knowledge graphs

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

- *Knowledge graphs* are a way to represent knowledge as a set of subject-predicate-object (SPO) triples
- An *entity* is an abstract or material object designated by an identifier (e.g. URI http://dbpedia.org/resource/Barack_Obama, in the case of DBpedia)
- *Entities* are always *subjects* in SPO triples
- Entities are connected with other entities, literals or scalars by relations or *predicates* (e.g. *dbo:genre*, *dbo:knownFor*, *dbo:spouse*, *dbp:memberOf*, etc.)
- Each SPO triple represents a simple fact (e.g. $dbr:Barack_Obama \xrightarrow{dbo:spouse} dbr:Michelle_Obama$)



Existing knowledge graphs

Knowledge graphs

Related Work

Problem

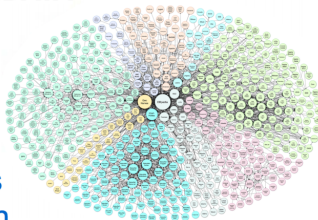
Method

Experiments

Conclusions



Facebook's
Entity Graph



Microsoft's
Satori



OpenIE
(Reverb, OLLIE)

Google's
Knowledge Graph

DBpedia entity page (rendered)

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

dbpedia.org/page/Barack_Obama

DBpedia

Browse using

Formats

Faceted Browser

Sparql Endpoint

About: Barack Obama

An Entity of Type : agent, from Named Graph : <http://dbpedia.org>, within Data Space : dbpedia.org

Barack Hussein Obama II (US /bəˈrɑːk huːˈseɪn əˈbɑːmə/; born August 4, 1961) is the 44th and current President of the United States, and the first African American to hold the office. Born in Honolulu, Hawaii, Obama is a graduate of Columbia University and Harvard Law School, where he served as president of the Harvard Law Review. He was a community organizer in Chicago before earning his law degree.

Property	Value
db:abstract	<ul style="list-style-type: none">Barack Hussein Obama II (US /bəˈrɑːk huːˈseɪn əˈbɑːmə/; born August 4, 1961) is the 44th and current President of the United States, and the first African American to hold the office. Born in Honolulu, Hawaii, Obama is a graduate of Columbia University and Harvard Law School, where he served as president of the Harvard Law Review. He was a community organizer in Chicago before earning his law degree. He worked as a civil rights attorney and taught constitutional law at University of Chicago Law School between 1992 and 2004. He served three terms representing the 13th District in the Illinois Senate from 1997 to 2004, running unsuccessfully for the United States House of Representatives in 2000.In 2004, Obama received national attention during his campaign to represent Illinois in the United States Senate with his victory in the March Democratic Party primary, his keynote address at the Democratic National Convention in July, and his election to the Senate in November. He began his presidential campaign in 2007 and, after a close primary campaign against Hillary Rodham Clinton in 2008, he won sufficient delegates in the Democratic Party primaries to receive the presidential nomination. He then defeated Republican nominee John McCain in the general election, and was inaugurated as president on January 20, 2009. Nine months after his inauguration, Obama was named the 2009 Nobel Peace Prize laureate.During his first two

- **Entities**

- *dbr:Barack_Obama*
- *dbr:Michelle_Obama*

- **Categories**

- *dbc:Presidents_of_the_United_States*
- *dbc:Critics_of_Islamophobia*

- **Literals**

- *dbr:Barack_Obama* *dbo:birthDate* "1961-08-04"
- *dbr:Barack_Obama* *foaf:gender* "male"

- **Predicates**

- *dbo:birthDate*
- *dbo:spouse*

Entity retrieval from a knowledge graph

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

- **Entity Search:** finding an entity based on its description
 - *"Ben Franklin"*
 - *"Einstein Relativity theory"*
- **List Search:** finding a set of entities based on their description
 - *"Formula 1 drivers who won the Monaco Grand Prix"*
 - *"animals lay eggs mammals"*
- **Attribute Search:** find a property of an entity
 - *"When was Intel founded?"*
 - *"What is the elevation of Karakoram?"*

Term-based KG entity retrieval

Knowledge
graphs

Related Work

Problem

Method

Experiments

Conclusions

Traditionally, entities are represented as multi-field documents and retrieved using structured document retrieval models:

- Fielded Sequential Dependence Model (FSDM) [Zhitsov et al., SIGIR 2015]
- Parametrized Fielded Sequential Dependence Model (PFSDM) [Nikolaev et al., SIGIR 2016]
- BM25F [Robertson and Zaragoza, Foundations and Trends in IR, 2009]

Key limitation: matching of queries to entities is performed at the word level

Network embedding methods

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

- Aim to embed network nodes into a low-dimensional vector space
- **Main idea:** apply of word embedding methods to sequences obtained using *random walks* on a given network
- Popular methods:
 - DeepWalk [Perozzi et al., KDD 2014]
 - LINE [Tang et al., WWW 2015]
 - node2vec [Grover and Leskovec, KDD 2016]
 - struc2vec [Ribeiro et al., KDD 2017]

Problems with network embeddings

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

- 1 We can apply network embeddings to knowledge graphs, but can't utilize entity embedding obtained this way directly in word-based retrieval models
- 2 We can use only word embeddings, but they utilize no information from a given knowledge graph

We propose **K**nowledge graph **E**ntity and **W**ord **E**MBEDDINGS for **R**etrieval (KEWER), a method that given a KG G :

- learns distributed representations of words (in predicates, literals, entity and category names) as well as entities and categories in G in the same embedding space
- utilizes the local structure of G when learning these embeddings

KEWER consists of three steps:

1 Random Walks from Knowledge Graph Entities

Starting from each KG entity, generate γ random walks of length $\leq t$.

Example:

dbr: Pierre_Curie $\xrightarrow{\text{dbp:spouse}}$ *dbr: Marie_Curie* $\xrightarrow{\text{dbp:knownFor}}$ *dbr: Radioactivity*

KEWER consists of three steps:

1 Random Walks from Knowledge Graph Entities

Starting from each KG entity, generate γ random walks of length $\leq t$.

Example:

dbr: Pierre_Curie $\xrightarrow{dbp:spouse}$ *dbr: Marie_Curie* $\xrightarrow{dbp:knownFor}$ *dbr: Radioactivity*

2 Replacement with Surface Forms

Randomly replace entity and category URIs with their surface forms (i.e. word tokens) in sequences of entity and category URIs, predicates and literals generated by random walks on G . The surface form of an entity or category for URI replacement is chosen uniformly at random from a set of available surface forms.

3 Learn Embeddings

Learn embeddings of words, entities and categories by maximizing the log-likelihood of observing other KG elements (word, entity or category) ξ_{i+j} in the context of each KG element ξ_i :

$$\frac{1}{T} \sum_{i=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(\xi_{i+j} | \xi_i), \quad \xi_{1 \dots T} \in \Xi,$$

$$\Xi = E \cup N \begin{cases} \cup K, & \text{if } \mathbf{categories} \text{ are used} \\ \cup V, & \text{if } \mathbf{literals} \text{ are used} \\ \cup P, & \text{if } \mathbf{predicates} \text{ are used.} \end{cases}$$

where $p(\xi_o | \xi_l)$ is defined using softmax:

$$p(\xi_o | \xi_l) = \frac{\exp(\mathbf{v}_{\xi_o}^{\top} \mathbf{v}_{\xi_l})}{\sum_{k=1}^{|\Xi|} \exp(\mathbf{v}_{\xi_k}^{\top} \mathbf{v}_{\xi_l})}$$

Entity retrieval using KEWER embeddings

Knowledge
graphs

Related Work

Problem

Method

Experiments

Conclusions

Embedding of a query q is a weighted sum of the embeddings of individual query words \mathbf{v}_{q_i} [Arora et al., ICLR 2017]:

$$\mathbf{q} = \sum_{i=1}^k \frac{a}{p(q_i) + a} \mathbf{v}_{q_i}$$

Entity retrieval using KEWER embeddings

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

Embedding of a query q is a weighted sum of the embeddings of individual query words \mathbf{v}_{q_i} [Arora et al., ICLR 2017]:

$$\mathbf{q} = \sum_{i=1}^k \frac{a}{p(q_i) + a} \mathbf{v}_{q_i}$$

Entities are scored according to the cosine similarity between entity embedding and query embedding:

$$KEWER(q, e) = \cos(\mathbf{q}, \mathbf{v}_e)$$

Entity retrieval using KEWER embeddings

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

Embedding of a query q is a weighted sum of the embeddings of individual query words \mathbf{v}_{q_i} [Arora et al., ICLR 2017]:

$$\mathbf{q} = \sum_{i=1}^k \frac{a}{p(q_i) + a} \mathbf{v}_{q_i}$$

Entities are scored according to the cosine similarity between entity embedding and query embedding:

$$KEWER(q, e) = \cos(\mathbf{q}, \mathbf{v}_e)$$

These scores can be interpolated with BM25F scores:

$$MM(q, e) = \beta KEWER(q, e) + (1 - \beta) BM25F(q, e), \quad 0 \leq \beta \leq 1$$

Utilizing entity linking

Knowledge
graphs

Related Work

Problem

Method

Experiments

Conclusions

To fine-tune query's vector representation, we can perform entity linking on a query and add embeddings of the linked entities to the query embedding:

$$\mathbf{q}_{el} = \sum_{i=1}^k \frac{a}{p(q_i) + a} \mathbf{v}_{q_i} + \sum_{i=1}^m s(e_i) \mathbf{v}_{e_i},$$

where $s(e_i)$ is the entity linker's annotation score for the entity e_i .

As a baseline, we used our implementation of the Jointly word and entity embedding method [Wang et al., EMNLP 2014]:

$$\mathcal{L}_J = \mathcal{L}_K + \mathcal{L}_T + \mathcal{L}_A$$

- Knowledge component loss \mathcal{L}_K is a translation-based loss for triples (similar to TransE [Bordes et al., NIPS 2013]).
- Text component loss \mathcal{L}_T corresponds to CBOW word embeddings trained on entity abstracts.
- Alignment loss \mathcal{L}_A aligns embeddings for words and entities based on entity abstracts.

Several similar models [Xie et al., AAAI 2016; Zhong et al., EMNLP 2015] were proposed for KG link prediction and triplet classification tasks.

Usefulness of KG structural components

Knowledge graphs

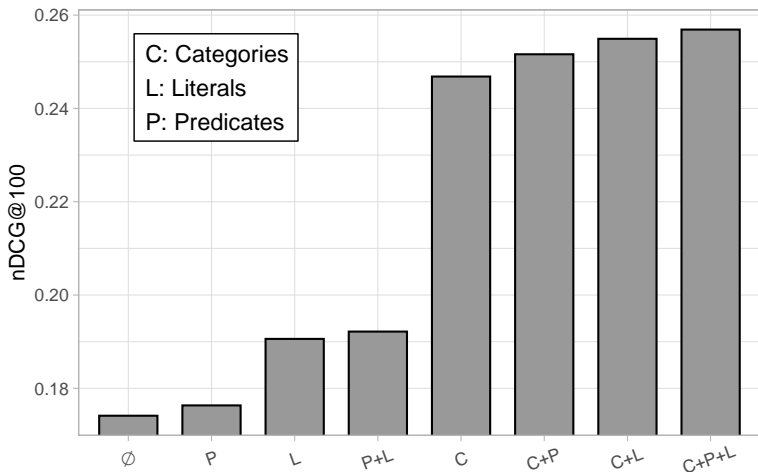
Related Work

Problem

Method

Experiments

Conclusions



nDCG₁₀₀ when using different combinations of categories, literals and predicates to train KEWER embeddings

Retrieval performance with different entity linkers

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

Sp stands for DBpedia Spotlight [Daiber et al., I-SEMANTICS 2013], *SM* for SMAPH [Cornolti et al., WWW 2016], *N* for Nordlys [Hasibi et al., SIGIR 2017].

Model	nDCG ₁₀	nDCG ₁₀₀	MAP
KEWER	0.2102	0.2569	0.1449
KEWER _{<i>el-Sp</i>}	0.2417	0.2803	0.1579
KEWER _{<i>el-SM</i>}	0.2704	0.3098	0.1780
KEWER _{<i>el-N</i>}	0.2660	0.3083	0.1775
Jointly (desp)	0.0486	0.0547	0.0211
Jointly _{<i>el-Sp</i>} (desp)	0.1603	0.1587	0.0838
Jointly _{<i>el-SM</i>} (desp)	0.1981	0.1924	0.1014
Jointly _{<i>el-N</i>} (desp)	0.1870	0.1814	0.0981
Jointly (sf)	0.0291	0.0393	0.0137
Jointly _{<i>el-Sp</i>} (sf)	0.1365	0.1357	0.0684
Jointly _{<i>el-SM</i>} (sf)	0.1685	0.1627	0.0795
Jointly _{<i>el-N</i>} (sf)	0.1624	0.1598	0.0836

Re-ranking performance

Knowledge graphs

Related Work

Problem

Method

Experiments

Conclusions

Statistically significant improvements (determined by a randomized test with $\alpha = 0.05$) over BM25F and BM25F+word2vec are indicated by “*” and “†”, respectively.

SemSearch ES

Model	nDCG ₁₀	nDCG ₁₀₀	MAP
BM25F	0.6606	0.7391	0.5693
BM25F+word2vec	0.6798*	0.7445	0.5712
BM25F+KEWER	0.6606	0.7333	0.5627
BM25F+KEWER _{el-SM}	0.6619	0.7409	0.5690

INEX-LD

Model	nDCG ₁₀	nDCG ₁₀₀	MAP
BM25F	0.4456	0.5127	0.3271
BM25F+word2vec	0.4591	0.5227	0.3406*
BM25F+KEWER	0.4676*	0.5298*	0.3417*
BM25F+KEWER _{el-SM}	0.4577*	0.5215*	0.3363*

ListSearch

Model	nDCG ₁₀	nDCG ₁₀₀	MAP
BM25F	0.4287	0.4989	0.3506
BM25F+word2vec	0.4235	0.5055*	0.3551
BM25F+KEWER	0.4402†	0.5210*†	0.3752*†
BM25F+KEWER _{el-SM}	0.4451*†	0.5251*†	0.3777*†

QALD-2

Model	nDCG ₁₀	nDCG ₁₀₀	MAP
BM25F	0.3442	0.4375	0.2861
BM25F+word2vec	0.3567*	0.4504*	0.2986*
BM25F+KEWER	0.3859*†	0.4743*†	0.3154*†
BM25F+KEWER _{el-SM}	0.3800*†	0.4700*†	0.3081*†

All queries

Model	nDCG ₁₀	nDCG ₁₀₀	MAP
BM25F	0.4631	0.5416	0.3792
BM25F+word2vec	0.4730*	0.5504*	0.3874*
BM25F+KEWER	0.4831*†	0.5602*†	0.3955*†
BM25F+KEWER _{el-SM}	0.4807*†	0.5601*†	0.3944*†

Example query

Knowledge
graphs

Related Work

Problem

Method

Experiments

Conclusions

Top 10 entities for the query “*wonders of the ancient world*” when using term-based retrieval with BM25F and cosine similarity based on query and entity embeddings. Relevant results are *italicized* and highly relevant results are **boldfaced**.

BM25F	KEWER
Seven Wonders of the Ancient World	Colossus of Rhodes
<i>7 Wonders of the Ancient World (video game)</i>	Statue of Zeus at Olympia
<i>Wonders of the World</i>	Temple of Artemis
<i>Seven Ancient Wonders</i>	List of archaeoastronomical sites by country
The Seven Fabulous Wonders	Hanging Gardens of Babylon
The Seven Wonders of the World (album)	Antikythera mechanism
Times of India's list of seven wonders of India	Timeline of ancient history
<i>Lighthouse of Alexandria</i>	<i>Wonders of the World</i>
7 Wonders (board game)	<i>Lighthouse of Alexandria</i>
Colossus of Rhodes	Great Pyramid of Giza

Conclusions

Knowledge
graphs

Related Work

Problem

Method

Experiments

Conclusions

- 1 Using all KG structural components (entities, categories, literals, and predicates) to learn KEWER embeddings results in the highest retrieval accuracy on DBpedia-Entity v2.
- 2 KEWER is particularly suitable for improving the ranking of results of complex entity search queries, such as question answering, list search, and keyword queries, where it can provide semantic relevance signal not captured by the retrieval models based on term matching.

Code, runs, and embeddings are available at
<https://github.com/teanalab/kewer>

Thank you! Questions?