

Efficient energy-based embedding models for link prediction in knowledge graphs

Pasquale Minervini¹ · Claudia d'Amato¹ · Nicola Fanizzi¹

Received: 8 October 2015 / Revised: 3 May 2016 / Accepted: 5 May 2016 /
Published online: 6 June 2016
© Springer Science+Business Media New York 2016

Abstract We focus on the problem of link prediction in Knowledge Graphs, with the goal of discovering new facts. To this purpose, Energy-Based Models for Knowledge Graphs that embed entities and relations in continuous vector spaces have been largely used. The main limitation in their applicability lies in the parameter learning phase, which may require a large amount of time for converging to optimal solutions. In this article, we first propose an unified view on different Energy-Based Embedding Models. Hence, for improving the model training phase, we propose the adoption of adaptive learning rates. We show that, by adopting adaptive learning rates during training, we can improve the efficiency of the parameter learning process by an order of magnitude, while leading to more accurate link prediction models in a significantly lower number of iterations. We extensively evaluate the proposed learning procedure on a variety of new models: our result show a significant improvement over state-of-the-art link prediction methods on two large Knowledge Graphs, namely WORDNET and FREEBASE.

Keywords Energy-based embedding models · Link predictions · RDF knowledge graphs

1 Introduction

Knowledge Graphs (KGs) are graph-structured knowledge bases, where factual knowledge is represented in the form of relationships between entities. We focus on KGs that adopt *Resource Description Framework* (RDF)¹ as their representation, since they constitute a

¹<http://www.w3.org/TR/rdf11-concepts/>

✉ Claudia d'Amato
claudia.damato@uniba.it

¹ Department of Computer Science, University of Bari, Bari, Italy

powerful instrument for search, analytics, recommendations, and data integration. Indeed, RDF is the Web standard for expressing information about resources.

A resource (hereafter also called *entity*) can be anything, including documents, people, physical objects, and abstract concepts. An RDF knowledge base (also called *RDF graph* as a KG) is a set of *RDF triples* of the form $\langle s, p, o \rangle$, where s , p and o denote the *subject*, the *predicate* (i.e. a *relation type*) and the *object* of the triple, respectively. Each triple $\langle s, p, o \rangle$ describes a statement, which can be interpreted as: *A relationship of type p holds between entities s and o .* The following example shows a set of RDF triples² describing the writer *William Shakespeare*:³

Example 1 (RDF Fragment)

```

(W. SHAKESPEARE, INFLUENCEDBY, G. CHAUCER)
(W. SHAKESPEARE, RELIGION, CHURCH OF ENGLAND)
(W. SHAKESPEARE, AUTHOR, HAMLET)
(HAMLET, GENRE, TRAGEDY)
(HAMLET, CHARACTER, OPHELIA)

```

Several RDF KGs are publicly available through the *Linked Open Data* (LOD) cloud, a collection of interlinked KGs such as Freebase (Bollacker et al. 2008), DBpedia (Bizer et al. 2009) and YAGO (Mahdisoltani et al. 2015). As of April 2014, the LOD cloud is composed of 1,091 interlinked KGs, globally describing more than 8×10^6 entities, and 188×10^6 relationships holding between them.⁴ However, KGs are often largely incomplete. For instance, 71 % of the persons described in Freebase⁵ have no known place of birth and 75 % of them have no known nationality (Dong et al. 2014).

For this reason, in this work, we focus on the problem of *predicting missing links* in large KGs, so as to discover new facts about a domain of interest. In the literature, this problem is referred to as *link prediction*, or *knowledge base completion*. The aim of this work is to provide an efficient and accurate model for predicting missing RDF triples in large RDF KGs (in a *link prediction* setting), without requiring extra background knowledge.

The link prediction task is well known in *Statistical Relational Learning* (SRL) (Getoor and Taskar 2007) which aims at modeling data from multi-relational domains, such as social networks, citation networks, protein interaction networks and knowledge graphs, and detecting missing links in such domains. Two main categories of models can be ascribed to SRL: *Probabilistic latent variable models* and *embedding models* (also frequently called *Energy-based models*). A detailed analysis of these two classes of models is reported in Section 4.

While appearing promising in terms of link prediction results, *Probabilistic latent variable models* showed limitations on scaling on large KG because of the complexity of the probabilistic inference and learning, which is intractable in general (Koller and Friedman 2009). Differently from them, *embedding models* have shown interesting ability

²This description is taken from the FREEBASE KG (Bollacker et al. 2008)

³For readability reasons, we describe entities and relations using an intuitive way of writing down triples as text rather than using the pure RDF syntax.

⁴State of the LOD Cloud 2014: <http://lod-cloud.net/>

⁵Available at <https://developers.google.com/freebase/data>

to scale on large KG while maintaining comparative performance in terms of predictive accuracy (Bordes and Gabrilovich 2015).

We focus specifically on a class of embedding models for KGs, named as *Energy-Based Embedding Models* (EBEMs), where entities and relations are embedded in continuous vector spaces, referred to as *embedding spaces*. In such models, the probability of an RDF triple to encode a true statement is expressed in terms of *energy* of the triple: this is an unnormalized score that is inversely proportional to such a probability value, and is computed as a function of the embedding vectors of the subject, the predicate and the object of the triple. The reason why we focus on this class of models, such as *Translating Embedding* (TransE) (Bordes et al. 2013) and other related ones (Bordes et al. 2011, 2014; Socher et al. 2013), is because it has been experimentally proved that they achieve state-of-the-art predictive accuracy results on link prediction tasks, while being able to scale to large and Web-scale KGs (Bordes et al. 2013; Dong et al. 2014; Bordes and Gabrilovich 2015). However, a major limiting factor for EBEMs lies in the parameter learning algorithm, which may require a long time (even days) to converge on large KGs (Chang et al. 2014).

In order to overcome such a limitation, we propose a method for reducing the learning time in EBEMs by an order of magnitude, while leading to more accurate link prediction models. Furthermore, we employ the proposed learning method for evaluating a family of novel EBEMs with useful properties. We experimentally tested our methods on two large and commonly used KGs: namely WORDNET and FREEBASE, following the same evaluation protocol used in (Bordes et al. 2013). Our results show a significant improvement over the state-of-the-art embedding models.

The rest of the paper is organized as follows. In Section 2, we introduce basics on Energy-Based Models. In Section 3 we propose: a) a framework for characterizing state-of-the-art EBEMs, b) a family of novel energy functions with useful properties, c) a method for improving the efficiency of the learning process in such models. In Section 4 the main related works falling in the *Probabilistic latent variable models* and *Embedding Models* categories are analyzed, while in Section 5 we empirically evaluate the proposed learning methods and energy functions. In Section 6 we summarize our work and outline future research directions.

2 Basics on energy-based models

Energy-Based Models (LeCun et al. 2006) are a versatile and flexible framework for modeling dependencies between variables. The key component is a scalar-valued *energy function* $E(\cdot)$, which associates a scalar *energy* with a configuration of variables. The energy of a configuration of variables is inversely proportional to the probability of such a configuration. Precisely, more likely configurations correspond to lower energy values, while less likely configurations correspond to higher energy values. Two main steps can be recognized in energy-based models: the *inference* step and the *learning* step.

The *inference* step consists in finding the most likely configuration of the variables of interest, that is the one that minimizes the energy function $E(\cdot)$. Given X and Y random variables, with values in \mathcal{X} and \mathcal{Y} , an example of the exploitation of the inference in energy-based models is given in the following.

Example 2 (Energy-Based Inference) Assuming that X describes the pixels of an image, while Y describes a discrete label associated with the image (such as “car” or “tree”), let $E : \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}$ be an energy function defined on the configurations of X and Y . The most

likely label $y^* \in \mathcal{Y}$ for an image $x \in \mathcal{X}$ can be inferred by finding the label in \mathcal{Y} that, given x , minimizes the energy function $E(\cdot)$:

$$y^* = \arg \min_{y \in \mathcal{Y}} E(x, y).$$

Learning in energy-based models consists in finding the most appropriate energy function within a family $\mathcal{F} = \{E_\theta \mid \theta \in \Theta\}$, indexed by parameters θ , that is in finding the function that is actually able to associate *lower energy states* with likely configurations of the variables of interests, and *higher energy states* to unlikely configurations of such variables. In practice, this corresponds to finding the energy function $E_\theta^* \in \mathcal{F}$ that minimizes a given *loss functional* \mathcal{L} , which measures the *quality* of the energy function on the data \mathcal{D} :

$$E_\theta^* = \arg \min_{E_\theta \in \mathcal{F}} \mathcal{L}(E_\theta, \mathcal{D}).$$

A normalized probability distribution can be derived from an energy-based model. Specifically, given an energy function $E : \mathcal{X} \mapsto \mathbb{R}$ defined on the possible configurations of a random variable X , it is possible to derive a corresponding probability distribution through the *Gibbs distribution*:

$$P(X = x) = \frac{1}{Z(\beta)} e^{-\beta E(x)}$$

where β is an arbitrary positive constant, and $Z(\beta) = \sum_{\tilde{x} \in \mathcal{X}} e^{-\beta E(\tilde{x})}$ is a normalizing factor⁶ referred to as the *partition function*.

3 A framework for energy-based embedding models

Energy-based models can be used for modeling the uncertainty in RDF KGs, in both statistical inference and learning tasks.

An RDF graph G can be viewed as a labeled directed multigraph, where entities are vertices, and each RDF triple is represented by a directed edge whose label is a predicate, and emanating from its source vertex to its object vertex. We denote with \mathcal{E}_G the set of all entities occurring as subjects or objects in G , formally:

$$\mathcal{E}_G = \{s \mid \exists (s, p, o) \in G\} \cup \{o \mid \exists (s, p, o) \in G\}$$

and we denote with \mathcal{R}_G the set of all relations appearing as predicates in G , formally:

$$\mathcal{R}_G = \{p \mid \exists (s, p, o) \in G\}.$$

Let $\mathcal{S}_G = \mathcal{E}_G \times \mathcal{R}_G \times \mathcal{E}_G$ be the space of *possible triples* of G , with $G \subseteq \mathcal{S}_G$, and let $E_\theta : \mathcal{S}_G \rightarrow \mathbb{R}$ (with parameters θ) be an energy function that defines an energy distribution over the set of possible triples \mathcal{S}_G .

The inference step consists in finding the most likely configuration of the variables of interest, i.e. the one that minimizes the energy function $E(\cdot)$. For instance, assume we need to know the most likely object o^* to appear in a triple with subject s (e.g. W. SHAKESPEARE) and predicate p (e.g. NATIONALITY). It can be inferred by finding the object o that minimizes the function $E(\cdot)$, as follows:

$$o^* = \arg \min_{o \in \mathcal{E}_G} E_\theta((s, p, o)).$$

⁶If X is a continuous random variable, then $Z(\beta) = \int_{\tilde{x} \in \mathcal{X}} e^{-\beta E(\tilde{x})}$.

Similar inference tasks can be performed with respect to the subject s , the predicate p , or a subset of such variables.

Learning consists in finding an energy function $E_\theta^* \in \mathcal{F}$, within a parametric family of energy functions $\mathcal{F} = \{E_\theta \mid \theta \in \Theta\}$ indexed by parameters θ , that minimizes a given loss functional \mathcal{L} defined on the RDF graph G , that is:

$$E_\theta^* = \arg \min_{E_\theta \in \mathcal{F}} \mathcal{L}(E_\theta, G).$$

Since the *energy* value for a triple expresses a quantity that is inversely proportional to the probability of the triple itself (see Section 2), in a *link prediction* setting, the energy function $E_\theta^*(\cdot)$ can be exploited for assessing a ranking of the so called *unobserved* triples, that are the triples in $\mathcal{S}_G \setminus G$. As such, triples associated with lower energy values (higher probabilities) will be more likely to be considered for a completion of the graph G , differently from triples associated with the higher energy values (lower probabilities).

On this point, it is important to note that *Open World Assumption* holds in RDF, which means that when a triple is missing in G , this does not have to be interpreted as that the corresponding statement is false (like for the case of the *Closed World Assumption* typically made in database settings), but rather that its truth value is *missing/unknown*, since it cannot be observed in the KG. We will refer to all triples in G as *visible triples*, and to all triples in $\mathcal{S}_G \setminus G$ as *unobserved triples*, which might encode true statements.

Within energy-based models, we particularly focus on *Energy-Based Embedding Models* (EBEMs) which are a specific class of models where each entity $x \in \mathcal{E}_G$ is mapped to a unique low-dimensional continuous vector $\mathbf{e}_x \in \mathbb{R}^k$, that is referred to as the *embedding vector* of x , and each predicate $p \in \mathcal{R}_G$ corresponds to an operation in the embedding vector space. As already pointed out in Section 1, the reason for such a choice is that EBEMs, such as *Translating Embedding* (TransE) (Bordes et al. 2013) and related models (Bordes et al. 2011, 2014; Socher et al. 2013), achieve state-of-the-art results in link prediction tasks, while being able to scale on very large (Web-scale) Knowledge Graphs (Bordes and Gabrilovich 2014).

In the following sections:

- (a) We present a unified general framework for formalizing EBEMs for KGs and we show that EBEMs proposed in the literature so far can be characterized with respect to their energy function. Additionally, we propose novel formulations of the energy functions with useful properties (see Section 3.1).
- (b) We then focus on the EBEM *learning* phase, by proposing a method for improving the efficiency of the parameters learning step (see Section 3.2).

3.1 Energy function characterization and new energy functions

The energy function $E_\theta : \mathcal{S}_G \rightarrow \mathbb{R}$ considered in the state-of-the art EBEMs for KG can be defined by using two types of parameters:

- **Shared Parameters:** used for computing the energy of all triples in the space of the possible triples \mathcal{S}_G of G .
- **Embedding Parameters:** used for computing the energy of triples containing a specific entity or relation $x \in \mathcal{E}_G \cup \mathcal{R}_G$. We denote such parameters by adding a subscript with the name of the entity or relation they are associated with. For instance, in the Translating Embeddings model, \mathbf{e}_s denotes the embedding vector representing a subject s , and \mathbf{e}_p denotes the translation vector representing a predicate p .

Both shared parameters and embedding parameters are learned from data. In particular, EBEMs for KGs associate each entity $x \in \mathcal{E}_G$ with a k -dimensional embedding vector $\mathbf{e}_x \in \mathbb{R}^k$, and each relation $p \in \mathcal{R}_G$ with a set of embedding parameters \mathbf{S}_p . Table 1 summarizes the energy functions that have been used by the state-of-the-art EBEMs for link prediction in KGs, and highlights the distinction between the two different kinds of parameters reported above. Please note that, in the table, the subscript for the parameters stand for the entity/predicate to which they refer to, e.g. the subscript p is added to the parameters associated with a particular predicate p .

The energy functions can be seen as sharing a common structure: given a RDF triple $\langle s, p, o \rangle$, its energy $E(\langle s, p, o \rangle)$ is computed by the following two-step process, also depicted in Fig. 1:

1. The embedding vectors $\mathbf{e}_s, \mathbf{e}_o \in \mathbb{R}^k$ respectively of the subject s and the object o of the triple, and the embedding parameters \mathbf{S}_p associated with the predicate p of the triple are used to obtain two new vectors $\mathbf{e}'_s, \mathbf{e}'_o \in \mathbb{R}^{k'}$ by means of two model-dependent functions $f_s(\cdot)$ and $f_o(\cdot)$:

$$\mathbf{e}'_s = f_s(\mathbf{e}_s, \mathbf{S}_p), \quad \mathbf{e}'_o = f_o(\mathbf{e}_o, \mathbf{S}_p).$$

2. The energy of a triple $\langle s, p, o \rangle$ is computed by a model-dependent function $g(\cdot)$, with $g : \mathbb{R}^{k'} \times \mathbb{R}^{k'} \mapsto \mathbb{R}$, applied to the vectors $\mathbf{e}'_s, \mathbf{e}'_o \in \mathbb{R}^{k'}$ resulting from the previous step:

$$E(\langle s, p, o \rangle) = g(\mathbf{e}'_s, \mathbf{e}'_o) = g(f_s(\mathbf{e}_s, \mathbf{S}_p), f_o(\mathbf{e}_o, \mathbf{S}_p)). \tag{1}$$

Please note that here, the proposed unifying framework is intended for *describing* EBEMs for KG: the choice for the functions $f_s(\cdot)$, $f_o(\cdot)$ and $g(\cdot)$ is model-dependent, and different models might correspond to different choices of such functions. As an example, in the following we show how the energy function adopted by the *Translating Embeddings* model (TransE) (Bordes et al. 2013), a state of the art EBEM for performing link prediction in KG, can be expressed by the use of the formalization presented above. TransE is particularly interesting: while its number of parameters grows *linearly* with the number of entities and relations in the KG, it yields state-of-the-art link prediction results on WORDNET and FREEBASE KGs (see the empirical comparison with other link prediction methods in Section 5.1).

Example 3 (Energy Function in TransE) In the formulation for the energy function of TransE (Bordes et al. 2013) (see also Table 1), each entity $x \in \mathcal{E}_G$ in an RDF graph G corresponds to a k -dimensional embedding vector $\mathbf{e}_x \in \mathbb{R}^k$, while each predicate $p \in \mathcal{R}_G$ corresponds to a *translation operation*, represented by a k -dimensional vector $\mathbf{e}_p \in \mathbb{R}^k$, in

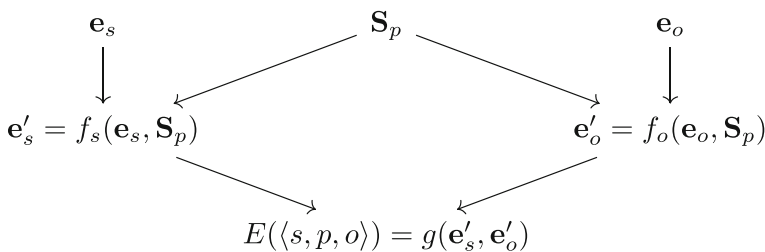


Fig. 1 Structure of the energy function in Energy-Based Embedding Models for KGs: $\mathbf{e}_s, \mathbf{S}_p$ and \mathbf{e}_o are the embedding parameters of s, p and o

Table 1 Energy-Based Embedding Models for knowledge graphs proposed in the literature, with their energy functions, shared and embedding parameters

Model	Energy function $E(s, p, o)$	Shared	Embedding
Unstructured (Bordes et al. 2014)	$\ \mathbf{e}_s - \mathbf{e}_o \ _1$		$\mathbf{e}_s, \mathbf{e}_o \in \mathbb{R}^k$
TransE (Bordes et al. 2013)	$\ (\mathbf{e}_s + \mathbf{e}_p) - \mathbf{e}_o \ _{1/2}$		$\mathbf{e}_s, \mathbf{e}_p, \mathbf{e}_o \in \mathbb{R}^k$
SE (Bordes et al. 2011)	$\ \mathbf{R}_{p,1} \mathbf{e}_s - \mathbf{R}_{p,2} \mathbf{e}_o \ _1$		$\mathbf{e}_s, \mathbf{e}_o \in \mathbb{R}^k, \mathbf{R}_{p, \cdot} \in \mathbb{R}^{n \times k}$
RESCAL (Nickel et al. 2011)	$\mathbf{e}_s^T \mathbf{R}_p \mathbf{e}_o$		$\mathbf{e}_s, \mathbf{e}_o \in \mathbb{R}^k, \mathbf{R}_p \in \mathbb{R}^{k \times k}$
SME lin. (Bordes et al. 2014)	$(\mathbf{R}_1 \mathbf{e}_s + \mathbf{R}_2 \mathbf{e}_p)^T (\mathbf{R}_3 \mathbf{e}_o + \mathbf{R}_4 \mathbf{e}_p)$	$\mathbf{R}_i \in \mathbb{R}^{n \times k}$	
SME bil. (Bordes et al. 2014)	$[(\mathbf{R}_1 \mathbf{e}_s) \times_3 (\mathbf{R}_2 \mathbf{e}_p)]^T [(\mathbf{R}_3 \mathbf{e}_o) \times_3 (\mathbf{R}_4 \mathbf{e}_p)]$	$\mathbf{R}_i \in \mathbb{R}^{n \times k}$	
NTN (Socher et al. 2013)	$\mathbf{u}_p^T \tanh(\mathbf{e}_s^T \mathbf{T}_p \mathbf{e}_o + \mathbf{R}_{p,1} \mathbf{e}_s + \mathbf{R}_{p,2} \mathbf{e}_o)$		$\mathbf{e}_s, \mathbf{e}_o \in \mathbb{R}^k, \mathbf{u}_p \in \mathbb{R}^n,$ $\mathbf{T}_p \in \mathbb{R}^{k \times k \times n}, \mathbf{R}_{p, \cdot} \in \mathbb{R}^{n \times k}$

the embedding vector space. As from Table 1, the energy function can be formulated by using the L_1 or the L_2 distance of the (translated) subject and object embedding vectors. In the case of L_1 formulation, the *energy* of an RDF triple $\langle s, p, o \rangle$ is given by the L_1 distance of $(\mathbf{e}_s + \mathbf{e}_p)$, corresponding to \mathbf{e}_s translated by \mathbf{e}_p , and \mathbf{e}_o :

$$E(\langle s, p, o \rangle) = \| (\mathbf{e}_s + \mathbf{e}_p) - \mathbf{e}_o \|_1.$$

This corresponds to the following choice of the functions $f_s(\cdot)$, $f_o(\cdot)$ and $g(\cdot)$:

$$f_s(\mathbf{e}_s, \{\mathbf{e}_p\}) = \mathbf{e}_s + \mathbf{e}_p, \quad f_o(\mathbf{e}_o, \{\mathbf{e}_p\}) = \mathbf{e}_o, \quad g(\mathbf{e}'_s, \mathbf{e}'_o) = \| \mathbf{e}'_s - \mathbf{e}'_o \|_1.$$

Besides proposing a general framework for expressing an energy function to be used by EBEMs, we also investigate whether the choice of other *affine transformations* for the functions $f_s(\cdot)$ and $f_o(\cdot)$, such as *scaling*, or *composition of translation and scaling*, leads to more accurate models than those generated by TransE (using the energy function reported in Table 1), while still having a number of parameters that scales linearly in the number of entities and relations in the KG. Specifically, we investigate the choices for the following $f_s(\cdot)$ and $f_o(\cdot)$ functions:

Translation: $f(\mathbf{e}_x, \{\mathbf{e}_p\}) = \mathbf{e}_x + \mathbf{e}_p,$
Scaling: $f(\mathbf{e}_x, \{\mathbf{e}_p\}) = \mathbf{e}_x \odot \mathbf{e}_p,$
Scaling \circ Translation: $f(\mathbf{e}_x, \{\mathbf{e}_{p,1}, \mathbf{e}_{p,2}\}) = (\mathbf{e}_x \odot \mathbf{e}_{p,1}) + \mathbf{e}_{p,2},$

where \circ denotes the composition operation between functions, and \odot denotes the Hadamard product, also referred to as element-wise product. The results of such a study are reported and discussed in Section 5.1.

In addition, we also evaluate the effect of enforcing the embedding vector of all entities to lie on the Euclidean unit $(k - 1)$ -sphere, that is $\mathbb{S}^{k-1} = \{\mathbf{x} \in \mathbb{R}^k \mid \|\mathbf{x}\|_2 = 1\}$. This is motivated by the fact that, in TransE (Bordes et al. 2013) and related-models, the L_2 norm of all entity embedding vectors is enforced to be 1: hence, we take into account the effect of normalizing the results of the functions $f_s(\cdot)$ and $f_o(\cdot)$, so that the resulting projections also lie on the Euclidean unit sphere (together with all entity embedding vectors).

In the next section, we discuss the *learning* step of EBEMs, which consists in finding the most appropriate energy function to be used during the *inference* step (see Section 2), and propose a method for improving both the efficiency of the learning process and the predictive accuracy of the learned model.

3.2 Learning the parameters of the energy function

As discussed in Section 2, *learning* in EBEMs for KGs corresponds to finding an energy function E_θ^* , within a family of functions $\mathcal{F} = \{E_\theta \mid \theta \in \Theta\}$ indexed by parameters θ , that minimizes a given *loss functional* \mathcal{L} measuring the compatibility of an energy function with respect to the RDF graph G :

$$E_\theta^* = \arg \min_{E_\theta \in \mathcal{F}} \mathcal{L}(E_\theta, G). \tag{2}$$

In the following, the definition for the *loss functional* \mathcal{L} is given. In agreement with the formalization presented in Section 3.1, a key point for learning the (best) energy function in EBEMs consists in learning the shared and embedding parameters to be used for computing the energy function. As for (Bordes et al. 2011, 2013, 2014), shared and embedding parameters are learned by using a *corruption process* $\mathcal{Q}(\tilde{x} \mid x)$ that, given a RDF triple $x \in G$,

Algorithm 1 Learning in EBEMs via *Stochastic Gradient Descent* (Bordes et al. 2013)

Input: Learning rate η , batch size n

Output: Optimal model parameters θ^*

- 1: Initialize model parameters θ_0
 - 2: **for** $t \in \langle 1, \dots, \tau \rangle$ **do**
 - 3: $\mathbf{e}_x \leftarrow \mathbf{e}_x / \|\mathbf{e}_x\|, \forall x \in \mathcal{E}_G$ {Normalize all entity embeddings}
 - 4: $T \leftarrow \text{SAMPLEBATCH}(G, n)$ {Sample observed and corrupted triples}
 - 5: $g_t \leftarrow \nabla \sum_{(x, \tilde{x}) \in T} [\gamma + E_\theta(x) - E_\theta(\tilde{x})]_+$ {Evaluate the gradient of \mathcal{L} w.r.t. θ }
 - 6: $\Delta_t \leftarrow -\eta g_t$ {Calculate the update to model parameters θ }
 - 7: $\theta_t \leftarrow \theta_{t-1} + \Delta_t$ { Update the model parameters θ }
 - 8: **end for**
 - 9: **return** θ_τ
-

produces a *corrupted* RDF triple \tilde{x} , uniformly sampled from the set of corrupted triples \mathcal{C}_x . Formally, given an RDF triple $\langle s, p, o \rangle$ from G , the set of corrupted triples for it is given by

$$\mathcal{C}_{\langle s, p, o \rangle} = \{ \langle \tilde{s}, p, o \rangle \mid \tilde{s} \in \mathcal{E}_G \} \cup \{ \langle s, p, \tilde{o} \rangle \mid \tilde{o} \in \mathcal{E}_G \}$$

that is the set obtained by replacing either the subject or the object of the triple with another entity from the set of entities \mathcal{E}_G .

The corruption process is applied to *positive* training RDF triples in order to generate *negative* examples that are actually missing in a KG. By corrupting the subject and the object of triples in the KG, the Local Closed World Assumption (LCWA) (Dong et al. 2014) is implicitly followed. In the LCWA, the idea is to consider the knowledge about a specific property p (e.g. NATIONALITY) of a resource s (e.g. W. SHAKESPEARE) to be *locally complete* if a value for p is already specified for the resource s . For instance, knowing that the triple $\langle \text{W. SHAKESPEARE, NATIONALITY, ENGLISH} \rangle$ is true (because observed in the KG), allows to assume that $\langle \text{W. SHAKESPEARE, NATIONALITY, AMERICAN} \rangle$ - i.e. a triple obtained by corrupting the object - is very likely to be false.

Since the final goal is to learn an energy function which associates lowest energy values (highest score) with observed triples, and highest energy values (lowest score) with unobserved triples, the corruption process $\mathcal{Q}(\tilde{x} \mid x)$ is used for defining the following margin-based stochastic ranking criterion over the triples in G :

$$\mathcal{L}(E_\theta, G) = \sum_{x \in G} \sum_{\tilde{x} \sim \mathcal{Q}(\tilde{x} \mid x)} [\gamma + E_\theta(x) - E_\theta(\tilde{x})]_+, \tag{3}$$

where $[x]_+ = \max\{0, x\}$, $\gamma > 0$ is a hyperparameter referred to as *margin*, and the embedding vector of each entity is enforced to have a unitary norm, i.e. $\forall x \in \mathcal{E}_G : \|\mathbf{e}_x\| = 1$. Actually, the loss functional in (3) enforces the energy of observed triples to be lower than the score of unobserved triples: the unitary norm constraints in the optimization problem prevent the training process to trivially solve it by increasing the entity embedding norms (Bordes et al. 2013).

The minimization problem in (2) can be solved by using projected Stochastic Gradient Descent (SGD) in mini-batch mode, as also proposed in (Bordes et al. 2011, 2013, 2014) and summarized in Alg. 1. The training algorithm works as follows: given an RDF graph G , at each iteration, it samples a batch of triples from G . Similarly to (Bordes et al. 2013), each batch is obtained by first randomly permuting all triples in G , then partitioning them into n_b batches of similar size, and iterating over them. A single pass over all triples in G is called an *epoch*. For each triple in the batch, the algorithm generates a *corrupted* triple by means

of the corruption process $\mathcal{Q}(\tilde{x} | x)$: this leads to a set T of observed and corrupted pairs of triples. Hence, the observed/corrupted triple pairs in T are used to evaluate the gradient of the loss functional \mathcal{L} in (3) with respect to the current model parameters θ . Finally, θ is updated in the steepest descent direction of the loss functional \mathcal{L} by a fixed learning rate η . This procedure is repeated until convergence (in (Bordes et al. 2013) the learning procedure was limited to 1000 epochs).

The main drawback of SGD is that it requires an initial, careful tuning of the learning rate η , that is also used across all parameters, without adapting to the characteristics of each parameter. However, if some entities are infrequent, the corresponding embedding vectors will tend to be updated less frequently during the learning process, and will require a longer time to be properly learned. For such a reason, the task of learning the model parameters in EBEMs by using SGD may require even days to terminate (Chang et al. 2014).

In order overcome such a limitation, we propose the adoption of *adaptive per-parameter learning rates* as a solution for reducing the learning time in EBEMs. The underlying idea consists in associating *smaller learning rates* to parameters updated more often (such as the embedding vectors of entities appearing more frequently) and *larger learning rates* to parameters updated less often. Specifically, while the SGD algorithm in Alg. 1 uses a global, fixed learning rate η , we propose relying on methods that estimate the optimal learning rate for each parameter while still being tractable for learning large models. In particular, we consider the following criteria for selecting the optimal learning rates: the Momentum method (Rumelhart et al. 1986), AdaGrad (Duchi et al. 2011) and AdaDelta (Zeiler 2012). Each of these methods can be implemented in Alg. 1, by replacing the update to model parameters on line 6 as specified in the following.

Momentum method The basic idea of this method is accelerating the progress along dimensions where the sign of the gradient does not change, while slowing the progress along dimensions where the sign of the gradient continues to change. This is done by keeping track of previous parameter updates with an exponential decay. The update step on line 6 of Alg. 1, in the Momentum method is given by:

$$\Delta_t \leftarrow \rho \Delta_{t-1} - \eta g_t,$$

where ρ is a hyperparameter controlling the decay of previous parameter updates.

AdaGrad The underlying idea in this method is that per parameter learning rates should grow with the inverse of gradient magnitudes: large gradients should have smaller learning rates, while small gradients should have larger learning rates, so that the progress along each dimension evens out over time. The update step on line 6 of Alg. 1, in AdaGrad is given by:

$$\Delta_t \leftarrow -\frac{\eta}{\sqrt{\sum_{j=1}^t g_j^2}} g_t,$$

where η is a global scaling hyperparameter. AdaGrad has been used on large scale learning tasks in a distributed environment (Dean et al. 2012).

AdaDelta This method uses an exponentially decaying average of squared gradients $E[g^2]$ and squared updates $E[\Delta^2]$, controlled by a decay term ρ , to give more importance to more recent gradients and updates. The update step on line 6 of Alg. 1, in AdaDelta is given by:

$$\Delta_t \leftarrow -\frac{\text{RMS}[\Delta]_{t-1}}{\text{RMS}[g]_t} g_t,$$

where $E[x]_t = \rho E[x]_{t-1} + (1 - \rho)x_t$ calculates the exponentially decaying average, $RMS[x]_t = \sqrt{E[x^2]_t + \epsilon}$, and ϵ is an offset hyperparameter.

All these methods leverage each parameter’s previous gradients for adaptively selecting the optimal learning rate. The additional space complexity provided by each of these methods is an additional accumulator for each parameter, containing its gradient history. We did not experience any sensible difference in runtimes in comparison with plain SGD.

4 Related works

In this section we survey the most representative related works in the categories of *Probabilistic latent variable models* and *Embedding Models*, by jointly highlighting their main peculiarities and drawbacks.

Probabilistic latent variable models Models in this class explain relations between entities by associating each entity to a set of intrinsic *latent attributes*. The term *latent* refers to the fact that the attributes are not directly observable in the data. Specifically, this class of models conditions the probability distribution of the relations between two entities on the latent attributes of such entities, and all relations are considered conditionally independent given the latent attributes. Similarly to Hidden Markov Models (Xu et al. 2006; Koller and Friedman 2009), this allows the information to *propagate* through the network of interconnected latent variables.

An early model in this family is the *Stochastic Block Model* (SB) (Wang and Wong 1987), which associates a *latent class* variable with each entity. In Fig. 2 (see left side), a simple SB for a social network is depicted. Here, each user $u \in U$ is associated with a latent class variable Z_u which conditions both its attributes A_u , and its relations R_i with other users. The *Infinite (Hidden) Relational Model* (Kemp et al. 2006; Xu et al. 2006) extends the SB by using Bayesian nonparametrics, so to automatically infer the optimal number of latent classes. The *Infinite Hidden Semantic Model* (Rettinger et al. 2009) further extends such model, to make use of constraints expressed in First Order Logic during the learning process, while the Mixed Membership Stochastic Block Model (Airoldi et al. 2008) extends the SB to allow entities to have mixed cluster-memberships. More recent works associate a set of *latent features* with each entity, instead of a single latent class. The *Nonparametric Latent Feature Relational Model* (Miller et al. 2009) is a latent feature model, which relies on Bayesian nonparametrics to automatically infer the optimal number of latent features

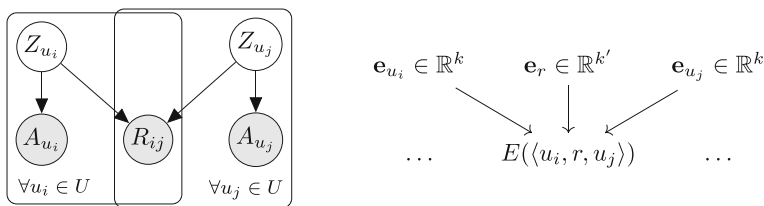


Fig. 2 *Left* – A simple SB for a social network: each user $u \in U$ is associated with a latent class variable Z_u which conditions both its attributes A_u , and its relations with other users. *Right* – An example of EBEM: the k -dimensional embedding vector $\mathbf{e}_u \in \mathbb{R}^k$ of an entity u (e.g. an user in a social network) is used for computing the *energy* of all RDF triples in which u appears in

during learning. Other approaches, such as (De Raedt et al. 2015), focus on the problem of inducing probabilistic logic programs. While showing interesting results in terms of predictive accuracy, a detailed analysis on the scalability issue, with particular reference to real-world large KGs is missing.

The main limitation of probabilistic latent variable models lies in the complexity of probabilistic inference and learning, which is intractable in general (Koller and Friedman 2009). As a consequence, these models do not result to be fully appropriate for modeling large KGs.

Embedding models Similarly to probabilistic latent feature models, in *Embedding Models* each entity in the KG is represented by means of a continuous k -dimensional *embedding vector* $\mathbf{e}_x \in \mathbb{R}^k$, encoding its intrinsic latent features within the KG. Nevertheless, models in this class do not necessarily rely on probabilistic inference for learning the optimal embedding vectors and this allows avoiding the issues related to the proper normalization of probability distributions, that may lead to intractable problems.

In RESCAL (Nickel et al. 2011), the problem of learning the embedding vector representations of all entities and predicates is cast as a *tensor factorization* problem: by relying on a bilinear model, and by using a squared reconstruction loss, an efficient learning algorithm, based on regularized *Alternating Least Squares*, is proposed. However, the number of parameters grows *super-linearly* with the number of predicates in the KG: for such a reason, RESCAL can hardly scale to highly-relational KGs (Jenatton et al. 2012). In EBEMs (depicted on the right side of Fig. 2), the *energy* of each RDF triple (s, p, o) is defined as a function of the embedding vectors \mathbf{e}_s and \mathbf{e}_o , associated with the subject s and the object o of the triple (as already detailed in Section 3). The major limitation in EBEMs is the *learning time*, i.e. the time required for learning the parameters of the energy function. Several options have been proposed for the choice of both the *energy function* and the *loss functional* for learning the embedding vectors representation (Bordes et al. 2011, 2013, 2014; Jenatton et al. 2012; Socher et al. 2013). These methods have been used to achieve state-of-the-art link prediction results while scaling on large KGs.

We outperform such methods both in terms of efficiency (by reducing the learning time by an order of magnitude) and effectiveness (by obtaining a more accurate model) - as shown by the empirical evaluations provided in Section 5.

5 Empirical evaluation

In this section, we present the empirical evaluation for our proposed solution. Particularly, we aim at answering the following questions:

- Q1:** Can adaptive learning rates, as proposed in Section 3.2, be used for improving the efficiency of parameters learning with respect to the current state-of-the-art EBEMs?
- Q2:** Do the energy functions proposed in Section 3.1 lead to more accurate link prediction models for KG completion?

In Section 5.1, we answer **Q1** by empirically evaluating the efficiency of the proposed learning procedure and the accuracy of the learned models. In Section 5.2, we answer **Q2** by evaluating the accuracy of models using the proposed energy functions in link prediction tasks.

In the following, we describe the KGs used for the evaluation, jointly with the adopted metrics.

Knowledge graphs As KGs, WORDNET (Miller 1995) and FREEBASE (FB15K) (Bollacker et al. 2008) have been adopted:

- WORDNET is a lexical ontology for the English language. It is composed of over 151×10^3 triples, describing 40943 entities and their relations by means of 18 predicate names.
- FREEBASE (FB15K) is a large collaborative knowledge base that is composed of over 592×10^3 triples, describing 14951 entities and their relations by means of 1345 predicate names.

As for the experiments, for comparison purpose, we use the very same training, validation and test sets adopted in Bordes et al. (2013). Specifically, as regards WORDNET, given the whole KG, 5000 triples were used for validation and 5000 were used for testing. As regards FB15K, 50000 triples were used for validation while 59071 were used for testing (the interested reader may refer to Bordes et al. (2013) for more informations about the creation of such datasets).

Evaluation metrics As for Bordes et al. (2013), the following metrics have been used:

- *averaged rank* (denoted as MEAN RANK)
- *proportion of ranks not larger than 10* (denoted as HITS@10).

They have been computed as follows. For each test triple $\langle s, p, o \rangle$, the object o is replaced by each entity $\tilde{o} \in \mathcal{E}_G$ in G thus generating a *corrupted* triple $\langle s, p, \tilde{o} \rangle$. The energy values of corrupted triples are computed by the model, and successively sorted in ascending order. The rank of the correct triple is finally stored. Similarly, this procedure is repeated by corrupting the subject s of each test triple $\langle s, p, o \rangle$. Aggregated over all test triples, this procedure leads to the two metrics: the *averaged rank* (denoted as MEAN RANK) that measures the average position of the true test triple in the ranking, and the *proportion of ranks not larger than 10* (denoted as HITS@10) that measures the number of times the true test triple is ranked among the most likely 10 triples. This setting is referred to as the RAW setting.

Please note that, if a generated corrupted triple already exists in the KG, ranking it before the original triple $\langle s, p, o \rangle$ is not wrong. For such a reason, an alternative setting, referred to as FILTERED setting (abbreviated with FILT.) is also considered. In this setting, corrupted triples that exist in either training, validation or test set are removed, before computing the rank of each triple.

In both RAW and FILTERED settings, it would be desirable to have lower MEAN RANK and higher HITS@10.

5.1 Evaluation of adaptive learning rates

In order to reply to question **Q1**, that is, for assessing whether Momentum, AdaGrad and AdaDelta are more efficient than SGD in minimizing the loss functional in (3), we empirically evaluated these methods on the task of learning the parameters in TransE on WORDNET and FREEBASE (FB15K) KGs, using the optimal settings described in Bordes et al. (2013) that is:

- $k = 20, \gamma = 2, d = L_1$ for WORDNET
- $k = 50, \gamma = 1, d = L_1$ for FB15K.

Following the empirical comparison of optimization methods in Schaul et al. (2014), we compared SGD, Momentum, AdaGrad and AdaDelta using an extensive grid of

hyperparameters. Specifically, given $\mathcal{G}_\eta = \{10^{-6}, 10^{-5}, \dots, 10^1\}$, $\mathcal{G}_\rho = \{1 - 10^{-4}, 1 - 10^{-3}, \dots, 1 - 10^{-1}, 0.5\}$ and $\mathcal{G}_\epsilon = \{10^{-6}, 10^{-3}\}$, the grids of hyperparameters for each of the optimization methods were defined as follows:

- **SGD and AdaGrad:** rate $\eta \in \mathcal{G}_\eta$.
- **Momentum:** rate $\eta \in \mathcal{G}_\eta$, decay rate $\rho \in \mathcal{G}_\rho$.
- **AdaDelta:** decay rate $\rho \in \mathcal{G}_\rho$, offset $\epsilon \in \mathcal{G}_\epsilon$.

For each possible combination of optimization method and hyperparameter values, we performed an evaluation consisting in 10 learning tasks, each time using a different random seed for initializing the model parameters in TransE. The same 10 random seeds were used for each of the evaluation tasks.

Figure 3 shows the behavior of the loss function for each of the optimization methods, for the best hyperparameter settings after 100 epochs, over the training set. It is immediate to see that, for both WORDNET and FB15K, AdaGrad (with $\eta = 0.1$) and AdaDelta (with $(1 - \rho) = 10^{-3}$ and $\epsilon = 10^6$) provide sensibly lower values of the loss functional \mathcal{L} than SGD and Momentum, even after a low number of iterations (< 10 epochs), and that AdaGrad and

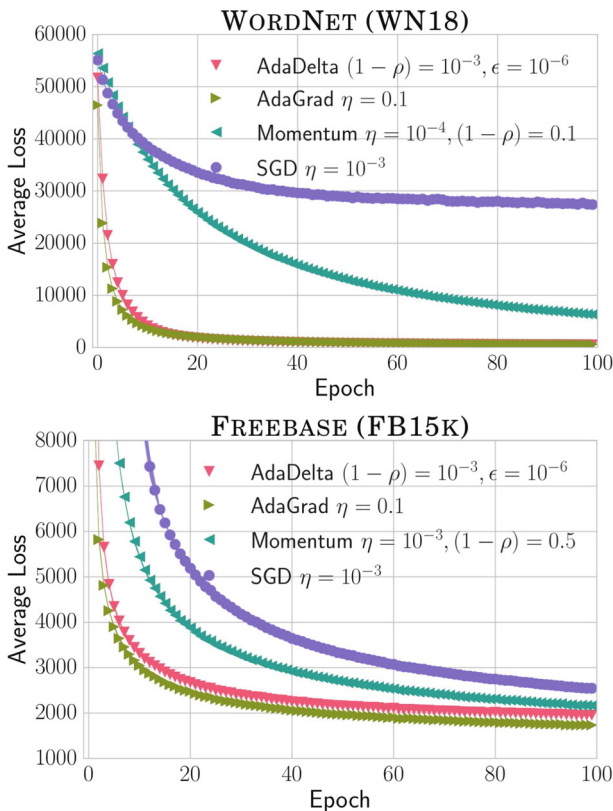


Fig. 3 Average loss across 10 TransE parameters learning tasks on WORDNET (*top*) and FREEBASE FB15K (*bottom*) knowledge graphs, using the optimal settings in Bordes et al. (2013). For each of the optimization methods, the hyperparameters settings that after 100 epochs achieve the lowest average loss are reported

AdaDelta, in their optimal hyperparameter settings, provide very similar loss values. Since AdaGrad has only one hyperparameter η and a lower complexity (it only requires one per parameter accumulator and a rescaling operation at each iteration) than AdaDelta, we select AdaGrad (with $\eta = 0.1$) as the optimization method of choice. Specifically, as a successive step, we needed to assess whether AdaGrad (with $\eta = 0.1$) leads to *more accurate models*, i.e. with lower MEAN RANK and higher HITS@10, than SGD. For the purpose, we trained TransE by using AdaGrad (with $\eta = 0.1$) for 100 epochs on a link prediction task on WORDNET and FB15K, under the same evaluation setting used in Bordes et al. (2013). Hyperparameters were selected according to the performance on the validation set using the same grid of hyperparameters adopted in Bordes et al. (2013). Specifically, we chose the margin $\gamma \in \{1, 2, 10\}$, the embedding vector dimension $k \in \{20, 50\}$, and the dissimilarity $d \in \{L_1, L_2\}$. Table 2 shows the results obtained by TransE trained using AdaGrad (with $\eta = 0.1$) for 100 epochs, in comparison with state-of-the-art results as reported in (Bordes et al. 2013).

From Table 2 can be noted that, despite of the sensibly lower number of training epochs (100, compared to 1000 used for training TransE with SGD, as reported by (Bordes et al. 2013)), TransE trained using AdaGrad provides more accurate link prediction models (i.e. lower MEAN RANK and higher HITS@10 values) than every other model in the comparison. A possible explanation for this phenomenon is the following. AdaGrad uses each parameter’s previous gradients for rescaling its learning rate: for such a reason, entities and predicates occurring less (resp. more) frequently will be associated with an higher (resp. lower) learning rate. As a result, the learning process for each parameter evens out over time, and all embedding parameters are learned at the same pace.

The results showed in this section largely prove that our solution is able to give a positive answer to Q1. Specifically, besides of experimentally proving that the adaptive learning rates proposed in Section 3.2 are able to improve the efficiency of parameters learning with respect to the current state-of-the-art EBEMs, we have also proved that the final learned model is able to outperform current state-of-the-art models in terms of MEAN RANK and HITS@10.

Table 2 Link Prediction Results: Test performance of several state-of-the-art Link Prediction methods on the WORDNET and FREEBASE (FB15K) KGs

Knowledge Graph	WORDNET				FREEBASE (FB15K)			
	MEAN RANK		HITS@10 (%)		MEAN RANK		HITS@10 (%)	
	RAW	FILT.	RAW	FILT.	RAW	FILT.	RAW	FILT.
Unstructured (Bordes et al. 2014)	315	304	35.3	38.2	1074	979	4.5	6.3
RESCAL (Nickel et al. 2011)	1180	1163	37.2	52.8	828	683	28.4	44.1
SE (Bordes et al. 2011)	1011	985	68.5	80.5	273	162	28.8	39.8
SME linear (Bordes et al. 2014)	545	533	65.1	74.1	274	154	30.7	40.8
SME bilinear (Bordes et al. 2014)	526	509	54.7	61.3	284	158	31.3	41.3
LFM (Jenatton et al. 2012)	469	456	71.4	81.6	283	164	26.0	33.1
TransE (Bordes et al. 2013)	263	251	75.4	89.2	243	125	34.9	47.1
TransE (AdaGrad)	169	158	80.5	93.5	189	73	44.0	60.1

Results show the MEAN RANK (the lower, the better) and HITS@10 (the higher, the better) for both the RAW and the FILTERED settings (Bordes et al. 2013)

5.2 Evaluation of the proposed energy functions

In this section, we evaluate the energy functions proposed in Section 3.1 in the definition of an EBEM, with the final goal of replying to question **Q2**, that is to assess whether the proposed energy functions lead to more accurate link prediction models for KGs completion than models at the state-of-the-art.

As from (1), the energy function of an EBEM can be rewritten as:

$$E((s, p, o)) = g(f_s(\mathbf{e}_s, \mathbf{S}_p), f_o(\mathbf{e}_o, \mathbf{S}_p))$$

where \mathbf{e}_s and \mathbf{e}_o denote the embedding vectors of the subject s and the object o of the triple, and \mathbf{S}_p denotes the set of embedding parameters associated with the predicate p . In Section 3.1 we proposed alternative choices for functions $f_s(\cdot)$ and $f_o(\cdot)$, that allow defining models whose number of parameters grows *linearly* with the number of entities and relations in the KG. Specifically, we proposed using *translation*, *scaling*, composition, and projection on the Euclidean unit sphere $n(\mathbf{x}) = \mathbf{x}/\|\mathbf{x}\|$. For each of the considered choices, we trained the corresponding EBEM on WORDNET and FB15K. Hyperparameters were selected on the basis of the model performance on the validation set: we selected the embedding vector dimension $k \in \{20, 50, 100\}$, the margin $\gamma \in \{2, 5, 10\}$, and the $g(\cdot)$ function $g(\mathbf{x}, \mathbf{y}) \in \{\|\mathbf{x} - \mathbf{y}\|_1, \|\mathbf{x} - \mathbf{y}\|_2, -\mathbf{x}^T \mathbf{y}\}$, corresponding to the L_1 and L_2 distances, and the negative dot product. Following the results from Section 5.1, model parameters were learned using AdaGrad (with $\eta = 0.1$) for 100 training epochs.

Table 3 shows the test results obtained with different choices of $f_s(\cdot)$ and $f_o(\cdot)$ function. Note that we used the notation $\mathbf{e}_{p,1}$ and $\mathbf{e}_{p,2}$ for referring to two distinct predicate embedding vectors, one used in the formulation of $f_s(\cdot)$ and the other in $f_o(\cdot)$, for avoiding name clashing. For the purposes of comparison, Table 3 also shows the results obtained, on the same link prediction tasks, by TransE (as reported in Bordes et al. (2013)) that is the best performing model in the literature.

From the table, it is interesting to note that, especially for highly multi-relational KGs such as FREEBASE (FB15K), *simpler models* for $f_s(\cdot)$ and $f_o(\cdot)$ provide better results than their more complex variants. A possible explanation is that many predicates in FB15K only occur in a limited number of triples (only 736 predicates out of 1345 occur in more than 20 triples) and in cases like this more expressive models are less likely to generalize correctly than simpler models. Given $f_o(\mathbf{e}_o, \{\mathbf{e}_p\}) = \mathbf{e}_o$, the best performing models, in terms of HITS@10, are:

- $f_s(\mathbf{e}_s, \{\mathbf{e}_p\}) = \mathbf{e}_s + \mathbf{e}_p$, representing the predicate-dependent *translation* of the subject's embedding vector
- $f_s(\mathbf{e}_s, \{\mathbf{e}_p\}) = \mathbf{e}_s \odot \mathbf{e}_p$, representing the predicate-dependent *scaling*.

This indicates that, despite the very different geometric interpretations, relying on simpler models improves link prediction results, especially in highly-relational KGs. This is probably due to the fact that, especially for the case of real-world KGs (which, by nature, tend to be very sparse), *simpler models tend to generalize better* and are less prone to over-fitting than more complex models. This characteristic can be very advantageous in real-world scenarios: relying on simpler models such as TransE (where the number of parameters scales *linearly* with the number of entities and predicates) can sensibly improve the training time, making learning from large and Web-scale KGs feasible.

Table 3 Link Prediction Results: Test performances of several EBEMs (on different choices of $f_s(\cdot)$ and $f_o(\cdot)$ functions) in comparison with TransE (Bordes et al. 2013) on WORDNET and FREEBASE (FB15K)

Knowledge Graph	WORDNET				FREEBASE (FB15K)			
	MEAN RANK		HITS@10 (%)		MEAN RANK		HITS@10 (%)	
	RAW	FILT.	RAW	FILT.	RAW	FILT.	RAW	FILT.
TransE (Bordes et al. 2013)	263	251	75.4	89.2	243	125	34.9	47.1
$f_s = e_s + e_p$ $f_o = e_o$	161	150	80.5	93.5	189	65	47.9	67.6
$f_s = e_s \odot e_p$ $f_o = e_o$	229	215	81.4	93.5	207	81	46.5	65.3
$f_s = (e_s \odot e_{p,1}) + e_{p,2}$ $f_o = e_o$	168	155	81.3	93.2	214	88	41.8	57.3
$f_s = e_s + e_{p,1}$ $f_o = e_o + e_{p,2}$	171	159	79.6	92.6	196	78	44.9	62.4
$f_s = e_s \odot e_{p,1}$ $f_o = e_o \odot e_{p,2}$	337	325	83.0	95.2	202	75	44.9	62.9
$f_s = (e_s \odot e_{p,1}) + e_{p,2}$ $f_o = e_o \odot e_{p,3}$	279	266	82.4	94.3	210	88	42.3	59.1
$f_s = (e_s \odot e_{p,1}) + e_{p,2}$ $f_o = (e_o \odot e_{p,3}) + e_{p,4}$	320	308	81.6	93.6	211	87	40.0	54.9
$f_s = n(e_s + e_p)$ $f_o = e_o$	211	200	75.7	88.7	237	115	39.5	55.4
$f_s = n(e_s \odot e_p)$ $f_o = e_o$	226	213	77.6	89.2	262	132	42.0	59.9
$f_s = n((e_s \odot e_{p,1}) + e_{p,2})$ $f_o = e_o$	160	148	77.7	88.7	239	103	42.8	59.1
$f_s = n(e_s + e_{p,1})$ $f_o = n(e_o + e_{p,2})$	262	251	79.3	91.6	206	86	47.5	66.5
$f_s = n(e_s \odot e_{p,1})$ $f_o = n(e_o \odot e_{p,2})$	761	750	73.4	83.5	249	120	42.0	61.0
$f_s = n(e_s \odot e_{p,1} + e_{p,2})$ $f_o = n(e_o \odot e_{p,3} + e_{p,4})$	624	613	74.7	83.6	238	114	42.7	60.4

Results show the MEAN RANK (the lower, the better) and HITS@10 (the higher, the better) in the RAW and FILTERED settings

We can conclude that, constraining the expressiveness of the models while using adaptive learning rates, yields a significant improvement over state-of-the-art methods discussed in (Bordes et al. 2013).

Source code and datasets for reproducing the experiments presented in this paper are available on-line.⁷

6 Conclusions and future works

We focused on Energy-Based Embedding Models, a novel class of link prediction models for knowledge graph completion where each entity in the graph is represented by a continuous embedding vector. Models in this class, like the *Translating Embedding* model (Bordes et al. 2013), have been used to achieve performance that is comparable with the main state-of-the-art methods while scaling on very large knowledge graphs.

In this work, we proposed: (i) a general framework for describing state-of-the-art Energy-Based Embedding Models, (ii) a family of novel energy functions, with useful properties,

⁷<https://github.com/pminervini/ebemkg/>

(iii) a method for improving the efficiency of the learning process by an order of magnitude, while leading to more accurate link prediction models.

We empirically evaluated the adoption of the proposed adaptive learning rates in the context of Energy-Based Embedding Models by showing that they provide more accurate link prediction models while reducing the learning time by an order of magnitude in comparison with state-of-the-art learning algorithms. We also empirically evaluated the newly proposed energy functions (with a number of parameters) that scales *linearly* with the number of entities and relations in the knowledge graph. Our results showed a significant improvement over state-of-the-art link prediction methods on the very same considered large KGs, which are WORDNET and FREEBASE.

For the future we plan to investigate on the formalization of Energy-Based Embedding Models that are able to take into account the available background knowledge. Other research directions include dynamically controlling the complexity of learned models, and further optimizing the learning process.

References

- Airoldi, E.M., Blei, D.M., Fienberg, S.E., & Xing, E.P. (2008). Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9, 1981–2014.
- Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, R., & Hellmann, S. (2009). DBPedia - A crystallization point for the web of data. *Journal of Web Seminars*, 7(3), 154–165.
- Bollacker, K.D., Evans, C., Paritosh, P., Sturge, T., & Taylor, J. (2008). Freebase: a collaboratively created graph database for structuring human knowledge, In Wang, J.T. (Ed.) *Proceedings of the ACM SIGMOD international conference on management of data, SIGMOD 2008* (pp. 1247–1250). Vancouver: ACM.
- Bordes, A., & Gabrilovich, E. (2015). Constructing and mining web-scale knowledge graphs. WWW 2015 Tutorial, In Gangemi, A., Leonardi, S., & Panconesi, A. (Eds.) *Proceedings of the 24th international conference on world wide web companion, WWW 2015 - companion volume.*: ACM.
- Bordes, A., & Gabrilovich, E. (2014). Constructing and mining web-scale knowledge graphs: KDD 2014 tutorial. In *In the 20th ACM SIGKDD international conference on knowledge discovery and data mining* (p. 1967): KDD '14.
- Bordes, A., Weston, J., Collobert, R., & Bengio, Y. (2011). Learning structured embeddings of knowledge bases, In Burgard, W. et al. (Eds.) *Proceedings of the twenty-fifth AAAI conference on artificial intelligence, AAAI 2011*. San Francisco: AAAI Press.
- Bordes, A., Usunier, N., García-Durán, A., Weston, J., & Yakhnenko, O. (2013). Translating embeddings for modeling multi-relational data, In Burges, C.J.C. et al. (Eds.) *Proceedings of the 27th Annual Conference on Neural Information Processing Systems* (pp. 2787–2795). Nevada: Lake Tahoe.
- Bordes, A., Glorot, X., Weston, J., & Bengio, Y. (2014). A semantic matching energy function for learning with multi-relational data - application to word-sense disambiguation. *Mach Learn*, 94(2), 233–259.
- Chang, K., Yih, W., Yang, B., & Meek, C. (2014). Typed tensor decomposition of knowledge bases for relation extraction, In Moschitti, A. et al. (Eds.) *Proceedings of the 2014 conference on empirical methods in natural language processing, EMNLP 2014, October 25-29, 2014. A meeting of SIGDAT, a Special Interest Group of the ACL* (pp. 1568–1579). Doha: ACL.
- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Le, Q.V., Mao, M.Z., Ranzato, M., Senior, A.W., Tucker, P.A., Yang, K., & Ng, A.Y. (2012). Large scale distributed deep networks, In Bartlett, P.L. et al. (Eds.) *Proceedings of the 26th Annual Conference on Neural Information Processing Systems* (pp. 1232–1240). Nevada: Lake Tahoe.
- De Raedt, L., Dries, A., Thon, I., Van den Broeck, G., & Verbeke, M. (2015). Inducing probabilistic relational rules from probabilistic examples, In Qiang Yang, Q., & Wooldridge, M. (Eds.) *Proceedings of 24th international joint conference on artificial intelligence (IJCAI), 2015* (pp. 1835–1843): AAAI press.
- Dong, X., Gabrilovich, E., Heitz, G., Horn, W., Lao, N., Murphy, K., Strohmman, T., Sun, S., & Zhang, W. (2014). Knowledge vault: a web-scale approach to probabilistic knowledge fusion, In Macskassy, S.A. et al. (Eds.) *The 20th ACM SIGKDD international conference on knowledge discovery and data mining, KDD '14* (pp. 601–610). New York: ACM.

- Duchi, J.C., Hazan, E., & Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12, 2121–2159.
- Getoor, L., & Taskar, B. (2007). Introduction to statistical relational learning. The MIT press.
- Jenatton, R., Roux, N.L., Bordes, A., & Obozinski, G. (2012). A latent factor model for highly multi-relational data. In Bartlett, P.L. et al. (Eds.) *Proceedings of the 26th Annual Conference on Neural Information Processing Systems*. (pp. 3176–3184). Nevada: Lake Tahoe.
- Kemp, C., Tenenbaum, J.B., Griffiths, T.L., Yamada, T., & Ueda, N. (2006). Learning systems of concepts with an infinite relational model. In *Proceedings, the twenty-first national conference on artificial intelligence and the eighteenth innovative applications of artificial intelligence conference* (pp. 381–388). Boston: AAAI press.
- Koller, D., & Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*: MIT Press.
- LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M., & Huang, F. (2006). A tutorial on energy-based learning. In Bakir, G. et al. (Eds.) *Predicting structured data*: MIT press.
- Mahdisoltani, F., Biega, J., & Suchanek, F.M. (2015). YAGO3: A Knowledge base from multilingual Wikipedias. In *CIDR, 2015, seventh biennial conference on innovative data systems research, Online Proceedings*.
- Miller, G.A. (1995). Wordnet: a lexical database for English. *Communications of the ACM*, 38(11), 39–41.
- Miller, K.T., Griffiths, T.L., & Jordan, M.I. Bengio, Y. et al. (Eds.) (2009). *Nonparametric latent feature models for link prediction*. Vancouver: Curran Associates, Inc.
- Nickel, M., Tresp, V., & Kriegel, H. (2011). A three-way model for collective learning on multi-relational data. In Getoor, L. et al. (Eds.) *Proceedings of the 28th international conference on machine learning, ICML 2011* (pp. 809–816). Bellevue: Omnipress.
- Rettinger, A., Nickles, M., & Tresp, V. (2009). Statistical relational learning with formal ontologies. In Buntine, W.L. et al. (Eds.) *Machine learning and knowledge discovery in databases, european conference, ECML PKDD 2009, Bled, Slovenia September 7-11, 2009, Proceedings, Part II. LNCS*, (Vol. 5782 pp. 286–301): Springer.
- Rumelhart, D.E., Hinton, G.E., & Wilson, R.J. (1986). Learning representations by back-propagating errors. *Nature*, 323, 533–536.
- Schaul, T., Antonoglou, I., & Silver, D. (2014). Unit tests for stochastic optimization. In *International conference on learning representations*. Banff.
- Socher, R., Chen, D., Manning, C.D., & Ng, A.Y. (2013). Reasoning with neural tensor networks for knowledge base completion. In Burges, C.J.C. et al. (Eds.) *Proceedings of the 27th Annual Conference on Neural Information Processing Systems*. (pp. 926–934). Nevada: Lake Tahoe.
- Wang, Y.J., & Wong, G.Y. (1987). Stochastic blockmodels for directed graphs. *Journal of the American Statistical Association*, 82(397), 8–19.
- Xu, Z., Tresp, V., Yu, K., & Kriegel, H. (2006). Infinite hidden relational models. In *UAI'06, Proceedings of the 22nd conference in uncertainty in artificial intelligence*. Cambridge: AUAI Press.
- Zeiler, M.D. (2012). ADADELTA: An adaptive learning rate method. arXiv:[1212.5703](https://arxiv.org/abs/1212.5703).