

Optimizing Class Subsumption through Controlled Dynamics of n-Balls in Vector Space

Aniket Mitra¹[0009-0007-6343-3830] and Vinu E. Venugopal²[0000-0003-4429-9932]

International Institute of Information Technology, Bangalore, India
{aniket.mitra,vinu.ev}@iiitb.ac.in

Abstract. Representing entities from an ontology as geometric shapes (such as balls, boxes, etc.) in a low-dimensional vector space, known as Region-based Geometric Knowledge Graph Embedding (RKGE), has demonstrated the ability to outperform traditional knowledge graph embedding methods in reasoning tasks while preserving the structural properties and syntactic characteristics of ontological axioms. In this study, we introduce a novel approach to enhance the subsumption capability of geometric embeddings based on *n-balls*. Additionally, we propose techniques to enhance the quality of such embeddings by extracting meta-information from the information-rich lexicons or annotations within the domain ontology.

Keywords: $\mathcal{EL}++$ DL · n-Ball Embedding · Knowledge Graph Embeddings

1 Introduction

Model theoretic languages like Description Logic (DL) are used to represent the semantics of OWL ontology axioms. The \mathcal{EL} , a sub-language of DL, is widely used to represent large biomedical ontologies such as GO and SNOMED due to their fast reasoning (tractable) property and support of major symbolic logic constructs including concept intersection, existential relations between concepts, role chains, etc. Notably, $\mathcal{EL}++$ TBox axioms (an extension of \mathcal{EL}) can be reduced to one of the below normal forms (NF) in linear time maintaining resultant normalized TBox size linear to the original TBox thereby still guaranteeing tractable reasoning property [1].

– **NF 1-4 (Concept axioms):** $C \sqsubseteq D, C \sqcap D \sqsubseteq E, \exists R.C \sqsubseteq D, C \sqsubseteq \exists R.D$

The bottom concept axioms and role inclusion axioms can be represented as shown below. We numbered them NF 5-7 for ease of addressing.

– **NF 5 (Bottom concept axioms):** $C \sqcap D \sqsubseteq \perp, \exists R.C \sqsubseteq \perp, C \sqsubseteq \perp$

– **NF 6-7 (Role axioms):** $R \sqsubseteq S, R_1 \circ R_2 \sqsubseteq S$

where $\{C, D, E, \perp\} \in N_c$ & $\{R, R_1, R_2, S\} \in N_r$. Here N_c and N_r denote the set of classes and roles in the ontology respectively.

In EmEL++ model [4] (a.k.a. *n-ball* approach), the authors attempted geometric construction of each concept and role as *balls* and translation vectors

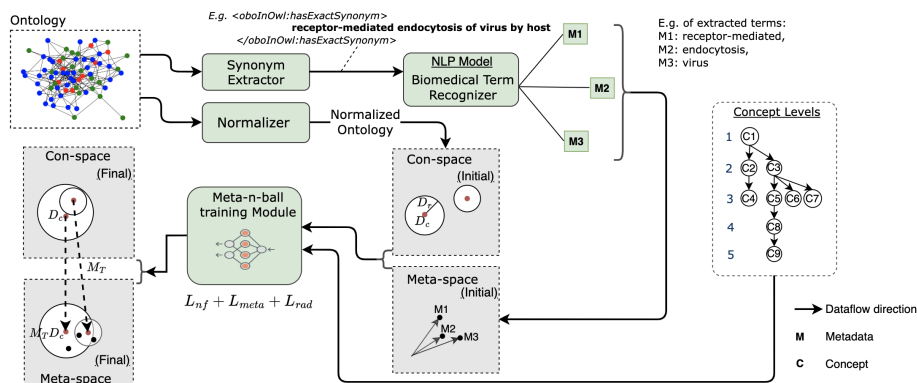


Fig. 1. Dataflow showing the generation of Meta-n-ball from domain ontology

respectively in vector space \mathbb{R}^n . They formulated specific loss functions preserving the semantic meaning of each of these $\mathcal{E}\mathcal{L} + \text{NFs}$ and optimized the balls via a Machine Learning (ML) model trained on these loss functions. Proper training will assign close vector representations to subsumption balls and bigger radius to the super-concept ball such that the sub-concept ball is totally engulfed by its super-concept ball. An alternative technique for embedding ontological data utilizes the graph-walk approach, albeit it ignores the structural nuances and characteristics inherent in the underlying ontology. In this method, each concept is represented as a node, while the relationships are illustrated as edges within the graph. Multiple rounds of randomized walks are executed on this graph structure to generate embeddings. Among these methodologies, OWL2Vec* [2] stands out by harnessing the substantial *meta-information* inherent in ontologies, including labels, synonyms, definitions, and more. By integrating this information into the graph structure, OWL2Vec* achieves a more comprehensive representation of the ontology, thereby enhancing its reasoning capabilities. The ball method conducts reasoning operations in linear time utilizing a straightforward ML model with minimal hyper-parameters, as opposed to complex graph-walk models that entail multiple path explorations, extensive pre-trained language models, and numerous parameters. Nevertheless, the accuracy of RKGE is significantly contingent upon the design of its loss functions. However, the meta-information present in the ontology is largely ignored while fine-tuning the model. We hypothesize that since a sub-ball is nothing but the more specific version of its super-ball, they must have common metadata terms that can be utilized to push the sub-balls properly inside their correct super-balls by tailoring accurate loss functions.

2 Proposed Approach: Meta-n-ball Model

Figure 1 outlines the general workflow of our proposed approach, ¹Meta-n-ball Model. We gather meta-information from the ontology and feed it through a

¹ <https://github.com/Aniket-Mitra/Meta-N-Ball>

pre-trained NLP model to extract biomedical terms. These terms are then initialized as n-dimensional (n-D) vectors in the metadata vector space, also known as *meta-space*. Simultaneously, concepts are initialized in the Concept Vector Space, or *con-space*, as n-D balls. Our Meta-n-ball module subsequently refines these balls to generate the final embeddings. To preserve the distinctiveness of the explicit concepts from the ontology and the extracted meta-information, we decided to represent them in two separate vector spaces, con-space and meta-space respectively, inspired by [3]. The combined loss function for the con-space and the meta-space vectors contains three components as shown in Figure 1. The loss functions of the seven NFs, as denoted by L_{oth} , are designed for training the con-space entities based on [4]. Additionally, L_{rad} incorporates concept-level information. To enhance the quality of the ball’s radius, the shortest distance to a specific concept from the root concept node in the concept hierarchy (a DAG) also called level is taken into consideration. This approach emphasizes the loss function more prominently for smaller balls. In Eqn. 1, Le_i (0 if unavailable) denotes level of a concept with radius R_i and γ denotes margin loss parameter.

$$L_{rad} = \begin{cases} [R_i]_- + \gamma & \text{if } Le_i = 0 \\ \left[\sqrt{Le_i} * R_i \right]_- & \text{if } Le_i \geq 1 \end{cases} \quad (1)$$

The loss function L_{meta} , representing the third type of loss function, is dedicated to learning meta-information. In Eq. 2, we use transformation matrix M_T which maps the concepts to meta-space and ensures that the metadata m resides inside its correct concept C (center C_c & radius R_c) in meta-space.

$$L_{meta} = \left[\|M_T C_c - m\| - R_c - \gamma \right]_+ + \left| \|C_c\| - 1 \right| + \left| \|m\| - 1 \right| \quad (2)$$

During training, the concepts, relations, metadata embeddings and M_T are optimized at every iteration in mini-batches based on the aggregation of these loss functions to reach their final values.

3 Results & Conclusion

Evaluation Metrics. The existing evaluation metrics primarily focus on calculating the distance between the centers of subsumption balls, often overlooking the quality aspects of the generated balls. This includes whether the radii are positive and whether the super-ball has a greater radius, the quality of subsumption (total-*ideal*, partial), etc. Keeping these in mind we design the following evaluation criteria: (1) *Valid Radius Proportion (VRP)*. The proportion of test cases where both radii are positive and the radius of the super-concept ball is larger than that of sub-concept. (2) *Overall Distance between Centers (ODC)*. To assess whether the distance between the centers of sub and super balls has decreased across all test cases in our new model, we conducted a one-sided t-test to determine which model exhibits greater distance values. The reporting format is as follows: if there is a significant reduction (p-val<0.05) in distance for the meta-n-ball model, it is denoted as (T-stat, emel>meta-nball), and vice versa. (3) *Perfect Overlapping Proportion (TOP)*. The proportion of valid VRP test

Table 1. Comparing n-ball (EmEL++) and Meta-n-ball approaches.

Dataset	Evaluation Metrics	EmEL++	Meta-NBall	EmEL++	Meta-NBall	EmEL++	Meta-NBall
Test Splits		Split1		Split2		Split3	
GO	VRP	0.707	0.742	0.596	0.640	0.587	0.627
	ODC(T-test)	3.26, meta-nball>emel		3.06, meta-nball>emel		No significant change	
	TOP	0.241	0.287	0.242	0.270	0.234	0.256
HPO	VRP	0.621	0.688	0.380	0.450	0.402	0.466
	ODC(T-test)	No significant change		5.73,emel>meta-nball		2.62, meta-nball>emel	
	TOP	0.203	0.219	0.091	0.182	0.150	0.158

cases where the sub-concept ball is completely inside it’s super-concept ball.

Experiments & Results. We utilized the Python package Scispacy with the *en_core_sci_md* NLP model, pretrained on 50k relevant biomedical entities. Metadata was extracted from synonym data found in the Gene Ontology² (GO) and Human Phenotype Ontology³ (HPO). We set hyper-parameters for both models as follows: n=100, =-0.1, epoch=1000. Concept levels were extracted from the respective .obo files available on official websites using the Python goatools package. Our model was tested on three separate valid-test splits, as shown in Table-1, with each test sample containing the true subsumption relation in the format (<sub-ball> <super-ball>). Our meta-n-ball approach consistently enhances the quality of ball radii and rectifies numerous imperfect subsumptions, as evidenced by improved VRP and TOP values across all test splits in Table 1. This observation supports our hypothesis and serves as motivation for further exploration. However, the decline in ODC performance in several test splits may be attributed to the inclusion of biomedical terms that are not suitable for our final model. Hence, implementing filtering criteria to include only relevant terms is essential. Moreover, upon closer examination of extracted biomedical terms, it is apparent that some metadata are overly generic (e.g., “activity”, “positive”) or near-duplicates (e.g., “ureteric reflux” & “ureteral reflux”). Filtering out such imperfect metadata terms and exploring additional meta-information beyond synonyms, such as definitions and labels, while maintaining performance, poses an intriguing challenge for future research.

References

1. Baader, F., Brandt, S., Lutz, C.: Pushing the EL envelope. In: Proceedings of the IJCAI. pp. 364–369. Professional Book Center (2005)
2. Chen, J., Hu, P., Jiménez-Ruiz, E., Holter, O.M., Antonyrajah, D., Horrocks, I.: OWL2Vec*: embedding of OWL ontologies. *Mach. Learn.* **110**(7), 1813–1845 (2021)
3. Hao, J., Chen, M., Yu, W., Sun, Y., Wang, W.: Universal representation learning of knowledge bases by jointly embedding instances and ontological concepts. In: Proceedings of the 25th ACM SIGKDD. pp. 1709–1719. ACM (2019). <https://doi.org/10.1145/3292500.3330838>
4. Mondal, S., Bhatia, S., Mutharaju, R.: Emel++: Embeddings for EL++ description logic. In: Proceedings of the AAI-MAKE. vol. 2846. CEUR-WS.org (2021)

² <https://geneontology.org/docs/download-ontology/>

³ <https://hpo.jax.org/app/data/annotations>