

Evaluating Text Representations Techniques for Hypernymy Detection: The Case of Arabic Language

Randah Alharbi and Husni Al-Muhtaseb*

Department of Information and Computer Science, * IRC for Intelligent Secure Systems
King Fahd University of Petroleum and Minerals (KFUPM)
Dhahran, Saudi Arabia
raharbi@uqu.edu.sa and muhtaseb@kfupm.edu.sa

Abstract. Text representation is a key performance component for any hypernymy-related task. In this study, we investigate representation techniques to understand which features best represent the hypernymy relation, focusing on three factors of representation: word embedding, embedding combination techniques, and using features. The results indicate that different embeddings have different effects on performance; concatenation, 'addition and subtraction' have led to better performance, and using unsupervised measures has a negative effect on performance.

1 Introduction

Text representation is a fundamental step in all NLP and IE tasks. Numerous types of word representation exist, from basic sparse and dense representations to complex representations such as neural embeddings. Many tasks have adopted the use of traditional neural word embedding, such as GloVe [1]. Contextual embeddings, such as Bidirectional Encoder Representations from Transformers (BERT), give different representations for a term based on its context [2]. These general word embeddings can model semantic similarity and relatedness between terms that encode various lexico-semantic and topical relations such as synonymy, antonymy, and hypernymy [3]. Some studies have proposed hypernymy-specific representations to better model hypernymy-relation for hypernymy-related tasks. In this work, we experiment with three representation factors for the hypernymy detection task, focusing on Arabic. The objectives of our experimentation are: (1) Evaluate hypernymy-specific embeddings against traditional embedding and contextual embedding. To the best of our knowledge, no other study compares the performance of these types of embedding. (2) Test the effectiveness of vectors' combination techniques for Arabic. (3) Evaluate the effect of adding unsupervised measure values to the representation. We select several unsupervised measures, Weeds Precision (WeedsPrec) measure [4], Clarck Degree of Entailment (ClarckDE) [5], as unidirectional similarity measures, InvCL [6] as directional similarity measure, SLQS_cos [7] as an entropy-based distributional measure and cosine similarity as baseline.

2 Methodology

The main goal of our study is to find the best text representation that models hypernymy relations. We aim to test if hypernymy-specific embedding is better at modeling the hypernymy relation for the hypernymy detection task. To

conduct the evaluation experiments, we have selected GloVe embedding as the traditional embedding baseline and BERT as the contextual embedding. For hypernymy-specific embedding, we have selected two retrofitted embeddings, Lexical Entailment Attract-Repel (LEAR) and Generalized Lexical ENtailmnet (GLEN), and two geometrical-based embeddings, Poincare for hierarchical data and Poincare GloVe. To mitigate external effects on the performance, we have tried to control most of the models' hyperparameters and the experimental setups. We also evaluate several mathematical operations for combining terms embedding. We have applied mathematical operations on our baseline embedding and three hypernymy-specific embeddings: LEAR, Poincare GloVe, and GLEN. Finally, we have experimented with enhancing term representations with the values of unsupervised measures as input features. We have used features embedding for one feature and for a combination of two, three, four, and five features. The feature embedding is concatenated with the terms embeddings before being used as input to the model. In the following subsections, we highlight the details of embedding training corpus, datasets, classification models, and experimental setup.

2.1 Corpus and Datasets

AraBERT corpus: We have trained all word embedding on half the corpus used to train an Arabic version of BERT called AraBERT [8]. It is a collection of Arabic text of 77GB in size and with a vocabulary of 12+ million words. We have used half of the AraBERT corpus for training for resource utilization.

Arabic Semantic Relation Dataset (ASRD): We have created our in-house dataset for Arabic lexical semantic relationships extracted from multiple Arabic semantic resources. It contains one-word examples for hypernym, hyponym, has_instance, is_instance, entailment, synonym, meronym, holonym, attribute, antonym, cause, similar, and verb_group. The number of examples is 958341. ARSD datasets will be publicly available.

Evaluation Datasets: We utilized lexical-semantic constraints extracted from the ASRD to train both the embedding models and the classification models. This suggests that a shared vocabulary might influence the performance of the embeddings. To mitigate this effect, we have selected seven English benchmark datasets containing hypernymy relations from HypEval and translated them into Arabic using Google Translate, namely, BLESS, ENTAILMENT, Lenci/Benotto, Weeds, BIBLESS, and two versions of Root9. We have filtered the terms in these datasets to include only single-word entries that were present in the AraBERT training corpus.

2.2 Representations Training

All representations are trained on half AraBERT corpus except Poincare, which is trained on ASRD hypernymy pairs, and BERT, which is pre-trained on the full AraBERT corpus. For **GloVe** we preprocess the corpus using AraBERT preprocessor and we have used the original GloVe code to train our version. **LEAR** is retrofitting-based embedding that takes pre-trained embeddings as input and modifies the embedding according to lexical-semantic relations con-

straints. We used GloVe embedding and ASRD constraints as input. LEAR needs synonyms and hypernyms for its Attract objective and antonym for its repel objective. We have used the official Python implementation of LEAR with slight modifications to adapt it to our data and trained for 5, 20, and 100 iterations. **GLEN** also takes pre-trained embeddings and lexical-semantic constraints, but it generates generalized modified embeddings for all vocabulary, even those with no constraints. We have used the official implementation of GLEN with default hyperparameters except for the number of iterations to stop training if there is no improvement. **Poincare GloVe** uses a modified GloVe objective to generate new word embeddings in hyperbolic space. It does not necessarily use pre-trained word embedding or lexical-semantic constraints; rather, we have used the co-occurrence calculation file generated by GloVe training as a basis for its calculation. We have trained two versions **100D Poincare GloVe** trained using 100D Poincare ball and all vocab, and $50 \times 2D$ Poincare GloVe trained in the cartesian product of 50 2D Poincare balls and most frequent words of the vocabulary. **Poincare Embedding** is trained using lexical-semantic constraints with a tree-like structure and we used hypernym and has_instance from ASRD to train it with negative examples set to 5. We have used its Gensim implementation. We have used pre-trained **BERT** for Arabic (AraBERT V2), and for each term, we have extracted features of the final layer output.

2.3 Classification Models and tasks

To Assess the effectiveness of the chosen embeddings in modeling hypernymy relations, we have used the resulting embeddings from each model as input to the hypernymy detection task. The detection model will classify input examples as hypernymy or not, with two classes as output. The positive examples are hypernyms, entailment, and has_instance; other relations are considered negative examples. The goal of our evaluation was not to achieve the highest performance but rather to fairly evaluate representation models by keeping experiment variables consistent among different experiments. Thus, we have used a simple feed-forward neural classification model for each task with an embedding layer, one hidden layer, and an output layer. We have trained a model per embedding. To evaluate the classification models, we test the trained model on several datasets, including the test set of ASRD.

We have evaluated vector combination techniques and the incorporation of unsupervised measures using hypernymy detection models trained on the ASRD dataset and tests on the testing datasets. The representation combination techniques experiments use the same model of hypernymy detection, but we varied the mathematical operations used to combine the pair vectors and the size of the resulting combined embedding. In the unsupervised measures experiments, we have trained models to incorporate every single feature and combination of two, three, four, and five features with term embeddings from GloVe, LEAR5, and 100D Poincare Embeddings. We have created an input features sub-network to enable the learning of feature embeddings. We use cross-entropy loss, Stochastic Gradient Descent (SGD) optimizer, 150 dimensions hidden layer, and ReLU activation function on the output layer and trained for 50 epochs except for models

that use BERT representation due to time and computing power limitations.

3 Results and Discussion

Experiment 1: Evaluating Word representations:

Table 1 shows the result of hypernymy detection using different representations. On ASRD, the best-performing model is Poincare embedding, followed by 100D Poincare GloVe. This is reasonable since Poincare embedding is trained solely on hypernymy examples of ASRD. Moreover, all hypernymy-specific embeddings outperform the GloVe baseline. On two datasets, Poincare GloVe models outperform others and score similarly to the best embedding on the other two datasets. GLEN outperforms other embeddings on two other datasets and performs similarly to the best embedding model on the three datasets. LEAR5 outperforms other embeddings on two other datasets and performs similarly to the best-performing embedding on four other datasets. These results suggest that the performance of the hypernymy detection task is highly affected by the datasets. For example, on both BIBLESS and ENTAILMENT datasets, which have the same type of positive and negative examples, the best-performing embeddings are GLEN and LEAR5. Our findings are similar to the findings of [9] for unsupervised hypernymy detection for English, which shows that no unsupervised measure outperforms others on all of their testing datasets because of how the negative samples in a dataset are constructed. Surprisingly, BERT is the least-performing model on the ASRD dataset, which might indicate the difficulty of the hypernymy detection task.

Data set	#Examples	GloVe	Lear 5	LEAR 20	LEAR 100	Poincare Embedding	100D Poincare GloVe	50x2E Poincare GloVe	GLEN	BERT
ASRD	191668	0.7631	0.8069	0.8016	0.7970	0.8394	0.8297	0.8037	0.8093	0.6186
BLESS	9486	0.5546	0.5282	0.5245	0.5282	0.5509	0.5495	0.5326	0.4914	0.4873
BIBLE	1167	0.6715	0.7124	0.6648	0.6531	0.4804	0.6258	0.6446	0.7187	0.5966
ENTAIL	1837	0.6333	0.6483	0.6267	0.6158	0.5079	0.6309	0.6372	0.6701	0.5549
LB	3253	0.5487	0.5723	0.5600	0.5479	0.4421	0.5447	0.5634	0.5443	0.5228
Weeds	3253	0.5839	0.5850	0.5693	0.5613	0.4606	0.5758	0.5833	0.5820	0.4933
Root9d	6181	0.5934	0.5538	0.5389	0.5455	0.5012	0.5968	0.5745	0.5638	0.5262
Root9h	9233	0.5436	0.5237	0.5201	0.5270	0.5397	0.5537	0.5380	0.5098	0.5049

Table 1: F1-score results for the hypernymy detection task (BIBLE. is BIBLESS, ENTAIL. is ENTAILMENT, and LB is LenciBenotto)

Experiment 2: Embedding Combination Techniques:

Table 2 shows the result of arithmetic operations to combine terms' vectors for the hypernymy detection task for GloVe, LEAR5, 100D Poincare GloVe, and GLEN. The results indicate the least effective combination techniques are addition and multiplication. Concatenation, subtraction, 'addition and subtraction', and 'concatenation and subtraction' have similar effects on performance. Meanwhile, concatenation, 'addition and subtraction', and 'concatenation and subtraction' perform slightly better than other operations. Our results for the Arabic language are similar to [10] results for Vietnamese, unlike [3] results, which found that vector difference and concatenation are the best operations for

English. This might indicate that vector operations have different effects on different target languages. Nevertheless, in this study, we have used concatenation to preserve all the information of both terms of the relation

Operations	GloVe	LEAR 5	100D Poincare GloVe	GLEN
concat	0.7631	0.8069	0.8297	0.8093
add	0.3815	0.3824	0.3955	0.3850
sub	0.7532	0.8096	0.8223	0.8068
mul	0.3817	0.4353	0.3815	0.3815
add_sub	0.7578	0.8110	0.8315	0.8110
concat_sub	0.7541	0.8080	0.8320	0.8076

Table 2: F1-score for various embedding techniques and operations.

Experiment 3: Using Unsupervised Measures as Input Features

Unlike our assumption, the results show that no single feature or combination of features is better than the baselines of using no features. Moreover, there is no difference among feature types in their effect on performance. This might be attributed to the new size of the representation, which makes it harder for the model to learn the patterns that indicate hypernymy from the representations. Also, the quality of the learned feature embeddings might be affected. Features embedding might need more epochs to be learned. Moreover, testing on ASRD, the used features might add little information beyond terms embedding. Testing on the other dataset might have led to different results, as the features might add more information about the terms. Table 3 shows the effect of incorporating a combination of four and five features as input.

Features	GloVe	LEAR5	100D Poincare
No features	0.76	0.81	0.83
invCL, clarkeDE, weeds_prec, cosine	0.75	0.76	0.77
invCL, clarkeDE, weeds_prec, SLQS_Cos	0.75	0.77	0.77
invCL, clarkeDE, cosine, SLQS_Cos	0.76	0.76	0.77
invCL, weeds_prec, cosine, SLQS_Cos	0.75	0.77	0.77
clarkeDE, weeds_prec, cosine, SLQS_Cos	0.75	0.77	0.77
invCL, clarkeDE, weeds_prec, cosine, SLQS_Cos	0.74	0.77	0.77

Table 3: F1-score results for using combinations of four and five features as input

3.1 Discussion

From the experiments’ results, we observe that, despite being trained without lexical-semantics constraints, Poincare GloVe performs well in the hypernymy detection task. Additionally, 100D Poincare GloVe slightly outperforms its counterpart, 50x2D Poincare GloVe. This highlights the effectiveness of modeling the hypernymy relation in hyperbolic space. On the other hand, GLEN outperforms other representations on some of the hypernymy detection datasets, and it falls short with ASRD because GLEN is known to have more impact when used with datasets with fewer known constraints [11]. Generally, the results reveal that no single representation consistently outperforms the others across all evaluation datasets. This suggests that the way training and evaluation datasets are constructed plays a crucial role in performance. This encourages future researchers

to create a dataset that uses multiple strategies for negative hypernymy examples, leading to more robust performance. The results for vector combination techniques on hypernymy detection for Arabic indicate that when computing resources are limited, one might choose the mathematical operation that results in a vector with the least dimensions size (subtraction). Enhancing the representation of terms' embeddings with more information extracted from calculating the unsupervised measures does not enhance the performance of the hypernymy detection task. Further incorporation techniques other than feature embedding should be evaluated.

4 Conclusion

In this work, we investigated the impact of various types of embeddings on the performance of hypernymy detection tasks. Our findings suggest that the choice of dataset used in the training and evaluation has a significant effect on model performance. Moreover, all vector operations perform similarly except addition and multiplication.

Acknowledgment: This work is supported by the Interdisciplinary Research Center for Intelligent Secure Systems IRC-ISS at KFUPM through the grant INSS2405.

References

- [1] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019.
- [3] Julie Weeds, Daoud Clarke, Jeremy Reffin, David Weir, and Bill Keller. Learning to distinguish hypernyms and co-hyponyms. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, 2014.
- [4] Julie Weeds and David Weir. A general framework for distributional similarity. In *Proceedings of the 2003 conference on Empirical methods in natural language processing*, 2003.
- [5] Daoud Clarke. Context-theoretic semantics for natural language: an overview. In *Proceedings of the workshop on geometrical models of natural language semantics*, 2009.
- [6] Alessandro Lenci and Giulia Benotto. Identifying hypernyms in distributional semantic spaces. In Eneko Agirre, Johan Bos, Mona Diab, Suresh Manandhar, Yuval Marton, and Deniz Yuret, editors, **SEM 2012: The First Joint Conference on Lexical and Computational Semantics*, 2012.
- [7] Enrico Santus, Alessandro Lenci, Qin Lu, and Sabine Schulte im Walde. Chasing hypernyms in vector spaces with entropy. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, volume 2: Short Papers*, 2014.
- [8] Wissam Antoun, Fady Baly, and Hazem Hajj. Arabert: Transformer-based model for arabic language understanding. In *LREC 2020 Workshop Language Resources and Evaluation Conference*, 2020.
- [9] Haw-Shiuan Chang, ZiYun Wang, Luke Vilnis, and Andrew McCallum. Distributional inclusion vector embedding for unsupervised hypernymy detection. In *North American Chapter of the Association for Computational Linguistics*, 2017.
- [10] Bui Van Tan, Nguyen Phuong Thai, and Nguyen Minh Thuan. Enhancing performance of lexical entailment recognition for vietnamese based on exploiting lexical structure features. In *2018 10th International Conference on Knowledge and Systems Engineering (KSE)*, 2018.
- [11] Goran Glavaš and Ivan Vulic. Generalized tuning of distributional word vectors for monolingual and cross-lingual lexical entailment. In *the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.