

Word2Box: Capturing Set-Theoretic Semantics of Words using Box Embeddings

Shib Sankar Dasgupta*, Michael Boratko*, Siddhartha Mishra ,
Shriya Atmakuri , Dhruvesh Patel , Xiang Lorraine Li , Andrew McCallum

Manning College of Information & Computer Sciences

University of Massachusetts Amherst

{ssdasgupta,mboratkro,siddharthami,satmakuri,
dhruveshpate,xiangl,mccallum}@cs.umass.edu

Abstract

Learning representations of words in a continuous space is perhaps the most fundamental task in NLP, however words interact in ways much richer than vector dot product similarity can provide. Many relationships between words can be expressed set-theoretically, for example adjective-noun compounds (eg. “red cars” \subseteq “cars”) and homographs (eg. “tongue” \cap “body” should be similar to “mouth”, while “tongue” \cap “language” should be similar to “dialect”) have natural set-theoretic interpretations. Box embeddings are a novel region-based representation which provide the capability to perform these set-theoretic operations. In this work, we provide a fuzzy-set interpretation of box embeddings, and learn box representations of words using a set-theoretic training objective. We demonstrate improved performance on various word similarity tasks, particularly on less common words, and perform a quantitative and qualitative analysis exploring the additional unique expressivity provided by WORD2BOX.

1 Introduction

The concept of learning a distributed representation for a word has fundamentally changed the field of natural language processing. The introduction of efficient methods for training vector representations of words in Word2Vec (Mikolov et al., 2013), and later GloVe (Pennington et al.) as well as FastText (Bojanowski et al., 2017) revolutionized the field, paving the way for the recent wave of deep architectures for language modeling, all of which implicitly rely on this fundamental notion that a word can be effectively represented by a vector.

While now ubiquitous, the concept of representing a word as a single point in space is not particularly natural. All senses and contexts, levels of abstraction, variants and modifications which the word may represent are forced to be captured by

*Equal Contributions.

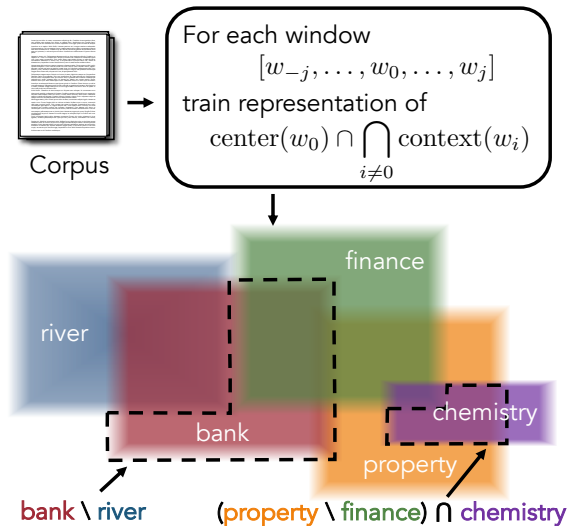


Figure 1: Given a corpus, Gumbel Boxes are trained as a fuzzy sets representing sets of windows with given center or context words. The representations can then be queried using multiple set-theoretic operations. In the graphic, $\text{bank} \setminus \text{river}$ overlaps highly with finance , and we would also expect high overlap with other boxes (not depicted) such as firm or brokerage (see Table 6). Similarly, we would expect boxes for chemical properties such as hardness or solubility to overlap with the dotted region indicating $\text{property} \setminus \text{finance} \cap \text{chemistry}$, and indeed we do observe such overlaps in the WORD2BOX model (see Table 7).

the specification of a single location in Euclidean space. It is thus unsurprising that a number of alternatives have been proposed.

Gaussian embeddings (Vilnis and McCallum, 2015) propose modeling words using densities in latent space as a way to explicitly capture uncertainty. Poincaré embeddings (Tifrea et al., 2019) attempt to capture a latent hierarchical graph between words by embedding words as vectors in hyperbolic space. Trained over large corpora via similar unsupervised objectives as vector baselines, these models demonstrate an improvement on word

similarity tasks, giving evidence to the notion that vectors are not capturing all relevant structure from their unsupervised training objective.

A more recent line of work explores region-based embeddings, which use geometric objects such as disks (Suzuki et al., 2019), cones (Vendrov et al., 2016; Lai and Hockenmaier, 2017; Ganea et al., 2018), and boxes (Vilnis et al., 2018) to represent entities. These models are often motivated by the need to express asymmetry, benefit from particular inductive biases, or benefit from calibrated probabilistic semantics. In the context of word representation, their ability to represent words using geometric objects with well-defined intersection, union, and difference operations is of interest, as we may expect these operations to translate to the words being represented in a meaningful way.

In this work, we introduce WORD2BOX, a region-based embedding for words where each word is represented by an n -dimensional hyperrectangle or “box”. Of the region-based embeddings, boxes were chosen as the operations of intersection, union, and difference are easily calculable. Specifically, we use a variant of box embeddings known as Gumbel boxes, introduced in (Dasgupta et al., 2020). Our objective (both for training and inference) is inherently set-theoretic, not probabilistic, and as such we first provide a fuzzy-set interpretation of Gumbel boxes yielding rigorously defined mathematical operations for intersection, union, and difference of Gumbel boxes.

We train boxes on a large corpus in an unsupervised manner with a continuous bag of words (CBOW) training objective, using the intersection of boxes representing the context words as the representation for the context. The resulting model demonstrates improved performance compared to vector baselines on a large number of word similarity benchmarks. We also compare the models’ abilities to handle set-theoretic queries, and find that the box model outperforms the vector model 90% of the time. Inspecting the model outputs qualitatively also demonstrates that WORD2BOX can provide sensible answers to a wide range of set-theoretic queries.

2 Background

Notation Let $V = \{v_i\}_{i=1}^N$ denote the vocabulary, indexed in a fixed but arbitrary order. A sentence $\mathbf{s} = (s_1, \dots, s_j)$ is simply a (variable-length) sequence of elements in our vocab $s_i \in V$. We

view our corpus $C = \{\mathbf{s}_i\}$ as a multiset¹ of all sentences in our corpus. Given some fixed “window size” ℓ , for each word s_i in a sentence \mathbf{s} we can consider the window centered at i ,

$$\mathbf{w}_i = [s_{i-\ell}, \dots, s_i, \dots, s_{i+\ell}],$$

where we omit any indices exceeding the bounds of the sentence. Given a window \mathbf{w}_i we denote the center word using $\text{cen}(\mathbf{w}_i) = s_i$, and denote all remaining words as the context $\text{con}(\mathbf{w}_i)$. We let C_W be the multiset of all windows in the corpus.

2.1 Fuzzy sets

Given any ambient space U a set $S \subseteq U$ can be represented by its characteristic function $\mathbb{1}_S : U \rightarrow \{0, 1\}$ such that $\mathbb{1}_S(u) = 1 \iff u \in S$. This definition can be generalized to consider functions $m : U \rightarrow [0, 1]$, in which case we call the pair $A = (U, m)$ a *fuzzy set* and $m = m_A$ is known as the *membership function* (Zadeh, 1965; Klir and Yuan, 1996). There is historical precedent for the use of fuzzy sets in computational linguistics (Zhelezniak et al., 2019a; Lee and Zadeh, 1969). More generally, fuzzy sets are naturally required any time we would like to learn a set representation in a gradient-based model, as hard membership assignments would not allow for gradient flow.

In order to extend the notion of set intersection to fuzzy sets, it is necessary to define a *t-norm*, which is a binary operation $\top : [0, 1] \times [0, 1] \rightarrow [0, 1]$ which is commutative, monotonic, associative, and equal to the identity when either input is 1. The min and product operations are common examples of t-norms. Given any t-norm, the intersection of fuzzy sets A and B has membership function $m_{A \cap B}(x) = \top(m_A(x), m_B(x))$. Any t-norm has a corresponding t-conorm which is given by $\perp(a, b) = 1 - \top(1 - a, 1 - b)$; for min the t-conorm is max, and for product the t-conorm is the probabilistic sum, $\perp_{\text{sum}}(a, b) = a + b - ab$. This defines the union between fuzzy sets, where $m_{A \cup B}(x) = \perp(m_A(x), m_B(x))$. Finally, the complement of a fuzzy set simply has member function $m_{A^c}(x) = 1 - m_A(x)$.

2.2 Box embeddings

Box embeddings, introduced in (Vilnis et al., 2018), represent elements \mathbf{x} of some set X as a Cartesian

¹A *multiset* is a set which allows for repetition, or equivalently a sequence where order is ignored.

product of intervals,

$$\begin{aligned} \text{Box}(\mathbf{x}) &:= \prod_{i=1}^d [x_i^-, x_i^+] \\ &= [x_1^-, x_1^+] \times \cdots \times [x_d^-, x_d^+] \subseteq \mathbb{R}^d. \end{aligned} \quad (1)$$

The volume of a box is simply the multiplication of the side-lengths,

$$|\text{Box}(\mathbf{x})| = \prod_{i=1}^d \max(0, x_i^+ - x_i^-),$$

and when two boxes intersect, their intersection is

$$\begin{aligned} \text{Box}(\mathbf{x}) \cap \text{Box}(\mathbf{y}) &= \prod_{i=1}^d [\max(x_i^-, y_i^-), \min(x_i^+, y_i^+)]. \end{aligned}$$

Boxes are trained via gradient descent, and these hard min and max operations result in large areas of the parameter space with no gradient signal. Dasgupta et al. (2020) address this problem by modeling the corners of the boxes $\{x_i^\pm\}$ with Gumbel random variables, $\{X_i^\pm\}$, where the probability of any point $\mathbf{z} \in \mathbb{R}^d$ being inside the box $\text{Box}_G(\mathbf{x})$ is given by

$$P(\mathbf{z} \in \text{Box}_G(\mathbf{x})) = \prod_{i=1}^d P(z_i > X_i^-) P(z_i < X_i^+).$$

For clarity, we will denote the original (“hard”) boxes as Box , and the Gumbel boxes as Box_G . The Gumbel distribution was chosen as it was min/max stable, thus the intersection $\text{Box}_G(\mathbf{x}) \cap \text{Box}_G(\mathbf{y})$ which was defined as a new box with corners modeled by the random variables $\{Z_i^\pm\}$ where

$$Z_i^- := \max(X_i^-, Y_i^-) \text{ and } Z_i^+ := \min(X_i^+, Y_i^+)$$

is actually a Gumbel box as well. Boratko et al. (2021) observed that

$$\begin{aligned} P(\mathbf{z} \in \text{Box}_G(\mathbf{x}) \cap \text{Box}_G(\mathbf{y})) &= \\ P(\mathbf{z} \in \text{Box}_G(\mathbf{x})) P(\mathbf{z} \in \text{Box}_G(\mathbf{y})), \end{aligned} \quad (2)$$

and also provided a rigorous probabilistic interpretation for Gumbel boxes when embedded in a space of finite measure, leading to natural notions of “union” and “intersection” based on these operations of the random variables (Boratko et al., 2021).

In this work, we do not embed the boxes in a space of finite measure, but instead interpret them as *fuzzy sets*, where the above probability (of a point \mathbf{z} being inside the Gumbel box) acts as a soft membership function.

3 Fuzzy Sets of Windows

In this section, we describe the motivation for using fuzzy sets to represent words, starting with an approach using traditional sets.

First, given a word $v \in V$, we can consider the windows centered at v ,

$$\text{cen}_W(v) := \{w \in C_W : \text{cen}(w) = v\},$$

and the set of windows whose context contains v ,

$$\text{con}_W(v) := \{w \in C_W : \text{con}(w) \ni v\}.$$

Note that cen_W is a function which takes in a word and returns a set of windows, whereas cen is a function which takes in a window and returns the center word, and a similar distinction holds for con_W and con .

A given window is thus contained inside the intersection of the sets described above, namely

$$\begin{aligned} [w_{-j}, \dots, w_0, \dots, w_j] \\ \in \text{cen}_W(w_0) \cap \bigcap_{i \neq 0} \text{con}_W(w_i). \end{aligned}$$

As an example, the window

$$\mathbf{w} = \text{“quick brown fox jumps over”},$$

is contained inside the cen_W (“fox”) set, as well as con_W (“quick”), con_W (“brown”), con_W (“jumps”), con_W (“over”). With this formulation, the intersection of the con_W sets provide a natural choice of representation for the context. We might hope that $\text{cen}_W(v)$ provides a reasonable representation for the word v itself, however by our set theoretic definition for any $u \neq v$ we have $\text{cen}_W(u) \cap \text{cen}_W(v) = \emptyset$.

We would like the representation of u to overlap with v if u has “similar meaning” to v , i.e. we would like to consider

$$\widetilde{\text{cen}}_W(v) := \{w \in W : \text{cen}(w) \text{ similar to } v\}.$$

A crisp definition of *meaning* or *similarity* is not possible (Hill et al., 2015; Finkelstein et al., 2001) due to individual subjectivity. Inter-annotator agreement for Hill et al. (2015) is only 0.67, for example, which makes it clear that $\widetilde{\text{cen}}_W(v)$ could not possibly be represented as a traditional set. Instead, it seems natural to consider $\widetilde{\text{cen}}_W(v)$ as represented by a fuzzy set (W, m) , where $m(w) \in$

$[0, 1]$ can be thought of as capturing graded similarity between v and $\text{cen}(w)$.² In the same way, we can define

$$\widetilde{\text{con}}_W(v) := \{w \in W : \text{con}(v) \ni w \text{ similar to } v\},$$

which would also be represented as a fuzzy set.³

As we wish to capture these similarities with a machine learning model, we now must find trainable representations of fuzzy sets.

Remark 1. Our objective of learning trainable representations for these sets provides an additional practical motivation for using fuzzy sets - namely, the hard assignment of elements to a set is not differentiable. Any gradient-descent based learning algorithm which seeks to represent sets will have to consider a smoothed variant of the characteristic function, which thus leads to fuzzy sets.

4 Gumbel Boxes as Fuzzy Sets

In this section we will describe how we model fuzzy sets using Gumbel boxes (Dasgupta et al., 2020). As noted in Section 2.2, the Gumbel Box model represents entities $\mathbf{x} \in X$ by $\text{Box}_G(\mathbf{x})$ with corners modeled by Gumbel random variables $\{X_i^\pm\}$. The probability of a point $\mathbf{z} \in \mathbb{R}^d$ being inside this box is

$$P(\mathbf{z} \in \text{Box}_G(\mathbf{x})) = \prod_{i=1}^d P(z_i > X_i^-)P(z_i < X_i^+).$$

Since this is contained in $[0, 1]$, we have that $(\mathbb{R}^d, P(\mathbf{z} \in \text{Box}_G(\mathbf{x})))$ is a fuzzy set. For clarity, we will refer to this fuzzy set as $\text{Box}_F(\mathbf{x})$.

The set complement operation has a very natural interpretation in this setting, as $\text{Box}_F(\mathbf{x})^c$ has membership function $1 - P(\mathbf{z} \in \text{Box}_G(\mathbf{x}))$, that is, the probability of \mathbf{z} not being inside the Gumbel box. The product t-norm is a very natural choice as well, as the intersection $\text{Box}_F(\mathbf{x}) \cap \text{Box}_F(\mathbf{y})$ will have membership function $P(\mathbf{z} \in \text{Box}_G(\mathbf{x}))P(\mathbf{z} \in \text{Box}_G(\mathbf{y}))$, which is precisely the membership function associated with $\text{Box}_G(\mathbf{x}) \cap \text{Box}_G(\mathbf{y})$, where here the intersection is between Gumbel boxes as defined in Dasgupta et al. (2020). Finally, we find that the membership function for

²For an even more tangible definition, we can consider $m(w)$ the percentage of people who consider u to be similar to $\text{cen}(w)$ when used in context $\text{con}(w)$.

³Note that this gives a principled reason to use different representation for $\widetilde{\text{con}}_W(v)$ and $\widetilde{\text{con}}_W(v)$, as they fundamentally represent different sets.

the union $\text{Box}_F(\mathbf{x}) \cup \text{Box}_F(\mathbf{y})$ is given (via the t-conorm) by

$$P(\mathbf{z} \in \text{Box}_G(\mathbf{x})) + P(\mathbf{z} \in \text{Box}_G(\mathbf{y})) - P(\mathbf{z} \in \text{Box}_G(\mathbf{x})P(\mathbf{z} \in \text{Box}_G(\mathbf{y}))). \quad (3)$$

Remark 2. Prior work on Gumbel boxes had not defined a union operation on Gumbel boxes, however (3) has several pleasing properties apart from being a natural consequence of using the product t-norm. First, it can be directly interpreted as the probability of \mathbf{z} being inside $\text{Box}_G(\mathbf{x})$ or $\text{Box}_G(\mathbf{y})$. Second, if the Gumbel boxes were embedded in a space of finite measure, as in Boratko et al. (2021), integrating (3) would yield the probability corresponding to $P(\text{Box}(\mathbf{x}) \cup \text{Box}(\mathbf{y}))$.

To calculate the size of the fuzzy set $\text{Box}_F(\mathbf{x})$ we integrate the membership function over \mathbb{R}^d ,

$$|\text{Box}_F(\mathbf{x})| = \int_{\mathbb{R}^d} P(\mathbf{z} \in \text{Box}_G(\mathbf{x})) d\mathbf{z}.$$

The connection between this integral and that which was approximated in Dasgupta et al. (2020) is provided by Lemma 3 of Boratko et al. (2021), and thus we have

$$|\text{Box}_F(\mathbf{x})| \approx \prod_{i=1}^d \beta \log \left(1 + \exp \left(\frac{\mu_i^+ - \mu_i^-}{\beta} - 2\gamma \right) \right)$$

where μ_i^-, μ_i^+ are the location parameters for the Gumbel random variables X_i^-, X_i^+ , respectively. As mentioned in Section 2.2, Gumbel boxes are closed under intersection, i.e. $\text{Box}_G(\mathbf{x}) \cap \text{Box}_G(\mathbf{y})$ is also a Gumbel box, which implies that the size of the fuzzy intersection

$$\begin{aligned} & |\text{Box}_F(\mathbf{x}) \cap \text{Box}_F(\mathbf{y})| \\ &= \int_{\mathbb{R}^d} P(\mathbf{z} \in \text{Box}_G(\mathbf{x}))P(\mathbf{z} \in \text{Box}_G(\mathbf{y})) d\mathbf{z} \\ &= \int_{\mathbb{R}^d} P(\mathbf{z} \in \text{Box}_G(\mathbf{x}) \cap \text{Box}_G(\mathbf{y})) d\mathbf{z} \end{aligned}$$

can be approximated as well. As both of these are tractable, integrating (3) is also possible via linearity. Similarly, we can calculate the size of fuzzy set differences, such as

$$\begin{aligned} & |\text{Box}_F(\mathbf{x}) \setminus \text{Box}_F(\mathbf{y})| = \\ & \int_{\mathbb{R}^d} P(\mathbf{z} \in \text{Box}_G(\mathbf{x}))[1 - P(\mathbf{z} \in \text{Box}_G(\mathbf{y}))] d\mathbf{z}. \end{aligned}$$

By exploiting linearity and closure under intersection, it is possible to calculate the size of arbitrary fuzzy intersections, unions, and set differences, as well as any combination of such operations.

Remark 3. If our boxes were embedded in a space of finite measure, as in [Boratko et al. \(2021\)](#), the sizes of these fuzzy sets would correspond to the intersection, union, and negation of the binary random variables they represent.

5 Training

In this section we describe our method of training fuzzy box representations of words, which we refer to as WORD2BOX.

In [Section 3](#) we defined the fuzzy sets $\widetilde{\text{cen}}_W(v)$ and $\widetilde{\text{con}}_W(v)$, and in [Section 4](#) we established that Gumbel boxes can be interpreted as fuzzy sets, thus for WORD2BOX we propose to learn center and context box representations

$$\begin{aligned} \text{cen}_B(v) &:= \text{Box}_F(\widetilde{\text{cen}}_W(v)) \\ \text{con}_B(v) &:= \text{Box}_F(\widetilde{\text{con}}_W(v)). \end{aligned}$$

Given a window, $\mathbf{w} = [w_{-j}, \dots, w_0, \dots, w_j]$, we noted that \mathbf{w} must exist in the intersection,

$$\widetilde{\text{cen}}_W(w_0) \cap \bigcap_{i \neq 0} \widetilde{\text{con}}_W(w_i) \quad (4)$$

and thus we consider a max-margin training objective where the score for a given window is given as

$$f(\mathbf{w}) := \left| \text{cen}_B(w_0) \cap \bigcap_{i \neq 0} \text{con}_B(w_i) \right|. \quad (5)$$

To create a negative example \mathbf{w}' we follow the same procedure as CBOW from [Mikolov et al. \(2013\)](#), replacing center words with a word sampled from the unigram distribution raised to the $3/4$. We also subsample the context words as in [Mikolov et al. \(2013\)](#). As a vector baseline, we compare with a WORD2VEC model trained in CBOW-style. We attach the source code with supplementary material.

6 Experiments and Results

We evaluate both WORD2VEC and WORD2BOX on several quantitative and qualitative tasks that cover the aspects of semantic similarity, relatedness, lexical ambiguity, and uncertainty. Following the previous relevant works ([Athiwaratkun and Wilson, 2018](#); [Meyer and Lewis, 2020](#); [Baroni et al., 2012](#)), we train on the lemmatized WaCkypedia corpora ([Baroni et al., 2009](#)), specifically ukWaC which is an English language corpus created by web crawling. After additional

pre-processing (details in [appendix A](#)) the corpus contains around 0.9 billion tokens, with just more than 112k unique tokens in the vocabulary. Noting that an n -dimensional box actually has $2n$ parameters (for min and max coordinates), we compare 128-dimensional WORD2VEC embeddings and 64-dimensional WORD2BOX embeddings for all our experiments. We train over 60 different models for both the methods for 10 epochs using random sampling on a wide range of hyperparameters (please refer to [appendix C](#) for details including learning rate, batch size, negative sampling, sub-sampling threshold, etc.). In order to ensure that the only difference between the models was the representation itself, we implemented a version of WORD2VEC in PyTorch, including the negative sampling and sub-sampling procedures recommended in ([Mikolov et al., 2013](#)), using the original implementation as a reference. As we intended to train on GPU, however, our implementation differs from the original in that we use Stochastic Gradient Descent with varying batch sizes. We provide our source code at <https://github.com/iesl/word2box>.

6.1 Word Similarity Benchmarks

We primarily evaluate our method on several word similarity benchmarks: SimLex-999 ([Hill et al., 2015](#)), WS-353 ([Finkelstein et al., 2001](#)), YP-130 ([Yang and Powers, 2006](#)), MEN ([Bruni et al., 2014](#)), MC-30 ([Miller and Charles, 1991](#)), RG-65 ([Rubenstein and Goodenough, 1965](#)), VERB-143 ([Baker et al., 2014](#)), Stanford RW ([Luong et al., 2013](#)), Mturk-287 ([Radinsky et al., 2011](#)) and Mturk-771 ([Halawi et al., 2012](#)). These datasets consist of pairs of words (both noun and verb pairs) that are annotated by human evaluators for semantic similarity and relatedness.

In [table 1](#) we compare the WORD2BOX and WORD2VEC models which perform best on the similarity benchmarks. We observe that WORD2BOX outperforms WORD2VEC (as well as the results reported by other baselines) in the majority of the word similarity tasks. We outperform WORD2VEC by a large margin in Stanford RW and YP-130, which are the rare-word datasets for noun and verb respectively. Noticing this effect, we enumerated the frequency distribution of each dataset. The datasets fall in different sections of the frequency spectrum, e.g., Stanford RW ([Luong et al., 2013](#)) only contains rare words which make its median frequency to be 5,683, whereas

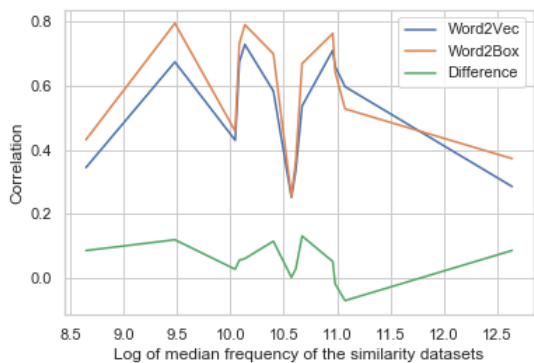


Figure 2: This plot depicts the gain in correlation score for WORD2BOX against WORD2VEC is much higher for the low and mid frequency range.

WS-353 (Rel) (Finkelstein et al., 2001) contains many more common words, with a median frequency of 64,490. We also observe a larger relative performance improvement over WORD2VEC on other datasets which have low to median frequency words, e.g. MC-30, MEN-Tr-3K, and RG-65, all with median frequency less than 25k. The order they appear in the table and the subsequent plots is lowest to highest frequency, left to right. Please refer to Appendix B for details.

In figure 2, we see that WORD2BOX outperforms WORD2VEC more significantly with less common words. In order to investigate further, we selected four datasets (RW-Stanford (rare words), Simelex-999, SimVerb-3500, WS-353 (Rel)), truncated them at a frequency threshold, and calculated the correlation for different levels of this threshold. In figure 3, we demonstrate how the performance gap between WORD2BOX and WORD2VEC changes as increasing amount frequent words are added to these similarity datasets. We posit that the geometry of box embeddings is more flexible in the way it handles sets of mutually disjoint words (such as rare words) which all co-occur with a more common word. Boxes have exponentially many corners, relative to their dimension, allowing extreme flexibility in the possible arrangements of intersection to represent complicated co-occurrences.

6.2 Set Theoretic Operations

All the senses, contexts and abstractions of a word can not be captured accurately using a point vector, and must be captured with sets. In this section, we evaluate our models capability of representing sets by performing set operations with the trained models.

6.2.1 Dataset

Homographs, words with identical spelling but distinct meanings, and polysemous words are ideal choice of probe for this purpose, as demonstrated by the “bank”, “river” and “finance” example of Figure 1. We constructed set-theoretic logical operations on words based on common polysemous words and homographs (Nelson et al., 1980). For example, the word “property” will have association with words related both “asset” and “attribute”, and thus the union of the later two should be close to the original word “property”. Likewise, the intersection set of “property” and “math” should contain many words related to mathematical properties or concepts.

To this end, we created a dataset consisting of triples (A, B, C) where $A \circ B$ should yield a set similar to C , for various set-theoretic operations \circ . In this task, given two words A and B and a set theoretic operation \circ , we try to find the rank of word C in the sorted list based on the set similarity (vector similarity scores for the vectors) score between $A \circ B$ and all words in the vocab. The dataset consists of 52 examples for both Union and Negation, and 20 examples for Intersection. The details of the dataset can be found in Appendix D.

6.2.2 Quantitative Results

In Table 2, we report the percentage of times WORD2BOX outperforms WORD2VEC, i.e. the model yields better rank for the word C . Note that it is not clear how to design the union, difference or the intersection operations with vectors. We consider several relevant choices, including component-wise operations (addition, subtraction, min and max) which yield a representation for $A \circ B$, as well as operations which operate on the scores - eg. score max pooling ranks each word X using $\max(A \cdot X, B \cdot X)$, and similarly for score min pooling. The purpose of these operations is to mimic the essence of union and intersection in the vector space, however, it is evident that the trained vector geometry is not harmonious to this construction as well.

We observe that almost of all the values are more than 0.9, meaning that WORD2BOX yields a higher rank for the target C than WORD2VEC over 90% of the time. This empirically validates that our model is indeed capturing the underlying set theoretic aspects of the words in the corpus.

	Stanford RW	RG-65	YP-130	MEN	MC-30	Mturk-287	SimVerb-3500	SimLex-999	Mturk-771	WS-353 (Sim)	WS-353 (All)	WS-353 (Rel)	VERB-143
*Poincaré	—	75.97	—	—	80.46	—	18.90	31.81	—	—	62.34	—	—
*Gaussian	—	71.00	41.50	71.31	70.41	—	—	32.23	—	76.15	65.49	58.96	—
WORD2VEC	40.25	66.80	43.77	68.45	75.57	61.83	23.58	37.30	59.90	75.81	69.01	61.29	31.97
WORD2BOX	45.08	81.45	51.6	73.68	87.12	70.62	29.71	38.19	68.51	78.60	68.68	60.34	48.03

Table 1: Similarity: We evaluate our box embedding model WORD2BOX against a standard vector baseline WORD2VEC. For comparison, we also include the reported results for Gaussian and Poincaré embeddings, however we note that these may not be directly comparable as many other aspects (eg. corpus, vocab size, sampling method, training process, etc.) may be different between these models.

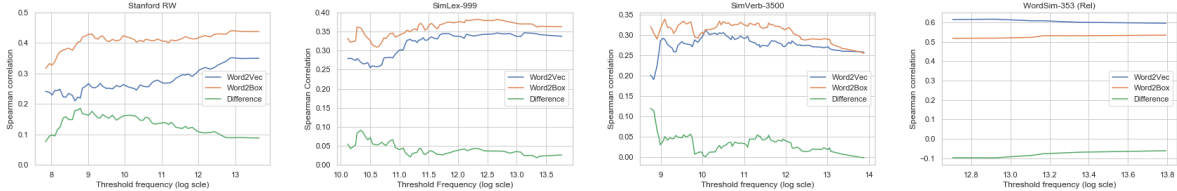


Figure 3: We plot the Spearman’s correlation score vs Threshold frequency in log scale for Stanford RW, Simelex-999 SimVerb-3500, WS-353 (Rel). The correlation value is calculated on the word pairs where both of them have frequency less than the threshold frequency.

WORD2BOX \ WORD2VEC	$A \cap B$	$A \setminus B$	$A \cup B$
$(A + B) \cdot X$	0.90	0.92	0.98
$(A - B) \cdot X$	0.90	0.65	0.80
$\max(A, B) \cdot X$	0.95	0.86	0.86
$\min(A, B) \cdot X$	0.90	0.75	0.92
$\max(A \cdot X, B \cdot X)$	0.95	0.84	0.94
$\min(A \cdot X, B \cdot X)$	1.0	0.80	0.84

Table 2: Percentage of queries for which WORD2BOX set operations return the target word with higher rank than the given vector operation for WORD2VEC. Scores higher than 0.5 means that WORD2BOX outperformed WORD2VEC. For subsequent qualitative comparisons we take the vector operation which performs most favorably for WORD2VEC.

6.2.3 Qualitative Analysis

In this section, we present some interesting examples of set theoretic queries on words, with different degrees of complexities. For all the tables in this section, we perform the set-operations on the query words and present the ranked list of most similar words to the output query. Many of these queries are based on the aforementioned homographs, for which there are natural expectations of what various set-theoretic operations should capture. Our results are presented in Table 3-7.

The results in Table 4 look reasonable for both models, as is to be expected since this is simply the similarity function for each model. Even increasing to a single intersection, as in Table 5, starts to demonstrate that WORD2VEC may often return very low-frequency words. In Table 6, we observe that set difference of “property” and “land” yields a

set of words that are related to attributes of science subjects, eg. algebra or chemistry. We wanted to examine how the model would handle more complicated queries, for example if we first perform “property” \ “finance” and then further intersect with “algebra” or “chemistry”, does the introduction of the relatively high-frequency “finance” term cause the model to struggle to recapture these items? In Table 7 we observe that the outputs for WORD2BOX do indeed correspond to properties of those sub-fields of science, whereas the results in WORD2VEC focus strongly on “finance”. In general, we observe better consistency of WORD2BOX with all the example logical queries.

7 Related Work

Learning distributional vector representations from a raw corpus was introduced in Mikolov et al. (2013), quickly followed by various improvements (Pennington et al.; Bojanowski et al., 2017). More recently, vector representations which incorporate contextual information have shown significant improvements (Peters et al., 2018; Devlin et al., 2019; Radford et al., 2019; Brown et al., 2020). As these models require context, however, Word2Vec-style approaches are still relevant in settings where such context is unavailable.

Hyperbolic representations (Nickel and Kiela, 2017; Ganea et al., 2018; Chamberlain et al., 2017) have become popular in recent years. Most related to our setting, Tifrea et al. (2019) propose a hyperbolic analog to GloVe, with the motivation that the hyperbolic embeddings will discover a la-

WORD2BOX			WORD2VEC					
$(\text{bank} \cap \text{river}) \cap X$	$(\text{bank} \cup \text{river}) \cap X$	$(\text{bank} \setminus \text{river}) \cap X$	$(\text{bank} + \text{river}) \cdot X$	$(\text{bank} - \text{river}) \cdot X$	$\max(\text{bank}, \text{river}) \cdot X$	$\max(\text{bank}, \text{river}) \cdot X$	$\max(\text{bank} - X, \text{river} - X)$	$\min(\text{bank} - X, \text{river} - X)$
headwaters	tributary	barclays	tributaries	cheques	tributary	vipava	tributaries	gauley
tributary	valley	hsbc	tributary	tymoshenko	tributaries	quabbin	headwaters	pymatuning
lake	headwaters	banking	headwaters	receivables	prut	irwell	tributary	'utricularia
basin	reservoir	citigroup	nakdong	citibank	chambal	trabajadores	headwater	luangwa
estuary	gorge	citibank	vipava	eurozone	headwaters	chattahoochee	distributaries	vipava
creek	lake	firm	estuary	brinks	larrys	tributaries	larrys	guadalquivir
valley	dam	ipo	larrys	defrauded	nakdong	belait	kobuk	suir
reservoir	headwater	brokerage	headwater	courtaulds	waterway	bougouriba	estuary	meenachil
canal	junction	interbank	distributary	refinance	loyalsock	canal	ijssel	tributary
floodplain	creek	kpmg	luangwa	mortgage	'hyperolius	glomma	distributary	battuta

Table 3: Output of WORD2BOX and WORD2VEC for various set operations

A	WORD2BOX $A \cap X$	WORD2VEC $A \cdot X$
bank	capital settlement airline hotel gateway treasury firm government loan casino	debit depositors securities kaupthing interbank subprime counterparty citibank fdic nasdaq
economics	education architecture politics economy literature faculty agriculture phd journalism	microeconomic keynesian microeconomics minored macroeconomics econometrics sociology thermodynamics evolutionism structuralist
microeconomics	economics mathematics physics philosophy theory technology economist principle research analysis	microeconomic initio germline instantiation zachman macroeconomics oxoglutarate glycemic noncommutative pubmed
property	land register status manor purpose locality premise landmark site residence	easement infringes burgage krajobrazowy chattels policyholder leasehold intestate liabilities ceteris
rock	music pop mountain cave band blues dance groove hot disco	shoegaze rhyolitic punk britpop mafic outcrops metalcore bluesy sedimentary quartzite

Table 4: Similarity outputs for WORD2BOX and WORD2VEC

A	B	WORD2BOX $(A \cap B) \cap X$	WORD2VEC $(A + B) \cdot X$
girl	boy	kid girls schoolgirl teenager woman boys child baby teenage orphan	shoeshine nanoha soulja schoolgirl yeller beastie jeezy crudup 'girl rahne
property	burial	cemetery bury estate grave interment tomb dwelling site gravesite sarcophagus	interment moated interred dunams ceteris burials catafalque easement deeded inhumation
	historical	historic estate artifact archaeological preserve ownership patrimony heritage landmark site	krajobrazowy burgage easement kravis dilapidation tohono intangible domesday moated laertius
	house	estate mansion manor residence houses tenement building premise buildings site	leasehold mansion tenements outbuildings estate burgage bedrooms moated burgesses manor
tongue	body	eye mouth ear limb lip forehead anus neck finger penis	tubercle ribcage meatus diverticulum forelegs radula tuberosity elastin foramen nostrils
	language	dialect idiom pronunciation meaning cognate word accent colloquial speaking speak	fluently dialects vowels patois languages loanwords phonology lingala tigrinya fluent

Table 5: Comparison of set intersection operation

A	B	WORD2BOX $(A \setminus B) \cap X$	WORD2VEC $(A - B) \cdot X$
algebra	finance	homomorphism isomorphism automorphism abelian algebraic bilinear topological morphism spinor homeomorphism	homeomorphic unital homomorphisms nilpotent algebraically projective holomorphic propositional nondegenerate endomorphism
bank	finance	wensum junction neman mouth tributary downstream corner embankment forks sandwich	shaddai takla thrombus gauley paria epenthetic chibchan urubamba foremast bolshaya
	river	barclays hsbc banking citigroup citibank firm ipo brokerage interbank kpmg	cheques tymoshenko receivables citibank eurozone brinks defrauded courtaulds refinance mortgage
chemistry	finance	biochemistry superconductor physics physic eutectic heat isotope fluorescence yttrium spectroscopy	augite alkyne desorption phosphorylating dimorphism fumarate hypertrophic empedocles hydratase enantiomer
property	land	homotopy isomorphism involution register bijection symplectic eigenvalue idempotent compactification lattice	brst stieltjes l'p repressor absurdum doesn conjugates nonempty didn wouldn

Table 6: Comparison of set difference operation

A	B	C	WORD2BOX $((A \setminus B) \cap C) \cap X$	WORD2VEC $(A - B + C) \cdot X$
property	finance	algebra	laplacian nilpotent antiderivative lattice surjective automorphism invertible homotopy integer integrand	expropriate extort refco underwrite reimburse refinance parmalat refinancing brokerage privatizing
		chemistry	eutectic desiccant allotrope phenocryst hardness solubility monoclinic hygroscopic nepheline trehalose	refinance brokerage burgage stockbroking refinancing warranties reimburse madoff privatizing valorem

Table 7: Comparison of set difference followed by intersection operation

tent hierarchical structure between words.⁴ Vilnis and McCallum (2015) use Gaussian distributions to represent each word, and KL Divergence as a score function.⁵ Athiwaratkun and Wilson (2018) extended such representations by adding certain thresholds for each distribution. For a different purpose, Ren and Leskovec (2020) use Beta Distributions to model logical operations between words. Our work can be seen as a region-based analog to these models.

Of the region-based embeddings, Suzuki et al. (2019) uses hyperbolic disks, and Ganea et al. (2018) uses hyperbolic cones, however these are not closed under intersection nor are their intersections easily computable. Vendrov et al. (2016) and Lai and Hockenmaier (2017) use an axis-aligned cone to represent a specific relation between words/sentences, for example an entailment relation. Vilnis et al. (2018) extends Lai and Hockenmaier (2017) by adding an upper-bound, provably increasing the representational capacity of the model. Li et al. (2019) and Dasgupta et al. (2020) are improved training methods to handle the difficulties inherent in gradient-descent based region learning. Ren et al. (2020) and Abboud et al. (2020) use a box-based adjustment of their loss functions, which suggest learning per-entity thresholds are beneficial. Chen et al. (2021) use box embeddings to model uncertain knowledge graphs, Onoe et al. (2021) use boxes for fined grained entity typing, and Patel et al. (2022) use boxes for multi-label classification.

Fuzzy sets, a generalization of sets, have been widely studied in the context of clustering (Bezdek and Harris, 1978), decision theory (Zimmermann, 1987) and linguistics (De Cock et al., 2000). However, the use of fuzzy sets in NLP has been fairly limited. Bhat et al. (2020) normalized each dimension of a word vector against all the word vectors

in the vocabulary and interpret them as probability features that enabled them to perform fuzzy set theoretic operations with the words. Zhao and Mao (2018) and Zhelezniak et al. (2019b) build fuzzy set representations of sentences using pre-trained vector embeddings for words and show the usefulness such representations on semantic textual similarity (STS) tasks. Jimenez et al. (2013, 2014) use the soft-cardinality features for a fuzzy set representation of a sentence to perform the task of entailment and textual relatedness. All these works, use pre-trained vector embeddings for the words to form fuzzy sets representing sentences. However, in this work we learn fuzzy set representations for words from corpus.

8 Conclusion

In this work we have demonstrated that box embeddings can not only effectively train to represent pairwise similarity but also can capture the rich set-theoretic structure of words via unsupervised training. This is a consequence of the fact that Gumbel boxes are an efficient parameterization of fuzzy sets, with sufficient representational capacity to model complicated co-occurrence interactions while, at the same time, allowing for tractable computation and gradient-based training of set-theoretic queries. The set-theoretic representation capabilities of box models allow them to generalize in a calibrated manner, leading to a more coherent and self-consistent model of sets.

Acknowledgments

The authors would like to thank the members of the Information and Extraction Synthesis Laboratory (IESL) at UMass Amherst for helpful discussions. This work was partially supported by IBM Research AI through the AI Horizons Network and the Chan Zuckerberg Initiative under the project Scientific Knowledge Base Construction. Additional support was provided by the National

⁴Reported results are included in table 1 as ‘‘Poincaré’’

⁵Reported results are included in table 1 as ‘‘Gaussian’’

Science Foundation (NSF) under Grant Numbers IIS-1514053 and IIS-2106391, the Defense Advanced Research Projects Agency (DARPA) via Contract No. FA8750-17-C-0106 under Subaward No. 89341790 from the University of Southern California, and the Office of Naval Research (ONR) via Contract No. N660011924032 under Subaward No. 123875727 from the University of Southern California. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IBM, CZI, NSF, DARPA, ONR, or the U.S. Government.

References

- Ralph Abboud, İsmail İlkan Ceylan, Thomas Lukasiewicz, and Tommaso Salvatori. 2020. *Boxe: A box embedding model for knowledge base completion*. In *Proceedings of the 34th Annual Conference on Neural Information Processing Systems, NeurIPS*.
- Ben Athiwaratkun and Andrew Gordon Wilson. 2018. *Hierarchical density order embeddings*. In *International Conference on Learning Representations*.
- Simon Baker, Roi Reichart, and Anna Korhonen. 2014. *An unsupervised model for instance level subcategorization acquisition*. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 278–289, Doha, Qatar. Association for Computational Linguistics.
- Marco Baroni, Raffaella Bernardi, Ngoc-Quynh Do, and Chung-chieh Shan. 2012. *Entailment above the word level in distributional semantics*. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 23–32, Avignon, France. Association for Computational Linguistics.
- Marco Baroni, Silvia Bernardini, A. Ferraresi, and E. Zanchetta. 2009. *The wacky wide web: a collection of very large linguistically processed web-crawled corpora*. *Language Resources and Evaluation*, 43:209–226.
- James C Bezdek and J Douglas Harris. 1978. *Fuzzy partitions and relations; an axiomatic basis for clustering*. *Fuzzy sets and systems*, 1(2):111–127.
- Siddharth Bhat, Alok Debnath, Souvik Banerjee, and Manish Shrivastava. 2020. *Word embeddings as tuples of feature probabilities*. In *Proceedings of the 5th Workshop on Representation Learning for NLP*, pages 24–33, Online. Association for Computational Linguistics.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. *Enriching word vectors with subword information*. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Michael Boratko, Javier Burrone, Shib Sankar Dasgupta, and Andrew McCallum. 2021. *Min/max stability and box distributions*. In *Uncertainty in Artificial Intelligence*, pages 2146–2155. PMLR.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. In *Advances in Neural Information Processing Systems*.
- Elia Bruni, Nam Khanh Tran, and Marco Baroni. 2014. *Multimodal distributional semantics*. *J. Artif. Int. Res.*, 49(1):1–47.
- Benjamin Paul Chamberlain, James R. Clough, and Marc Peter Deisenroth. 2017. *Neural embeddings of graphs in hyperbolic space*. *ArXiv*, abs/1705.10359.
- Xuelu Chen, Michael Boratko, Muhao Chen, Shib Sankar Dasgupta, Xiang Lorraine Li, and Andrew McCallum. 2021. *Probabilistic box embeddings for uncertain knowledge graph reasoning*. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*.
- Shib Sankar Dasgupta, Michael Boratko, Dongxu Zhang, Luke Vilnis, Xiang Lorraine Li, and Andrew McCallum. 2020. *Improving local identifiability in probabilistic box embeddings*. In *Advances in Neural Information Processing Systems*.
- Martine De Cock, Ulrich Bodenhofer, and Etienne E Kerre. 2000. *Modelling linguistic expressions using fuzzy relations*. In *Proc. 6th Int. Conf. on Soft Computing (IIZUKA2000)*, pages 353–360. Citeseer.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. *BERT: Pre-training of deep bidirectional transformers for language understanding*. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Minneapolis, Minnesota. Association for Computational Linguistics.
- Lev Finkelstein, Evgeniy Gabilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppin. 2001. *Placing search in context: The*

- concept revisited. In *Proceedings of the 10th international conference on World Wide Web*, pages 406–414.
- Octavian-Eugen Ganea, Gary Bécigneul, and Thomas Hofmann. 2018. Hyperbolic entailment cones for learning hierarchical embeddings. In *International Conference on Machine Learning*.
- Guy Halawi, Gideon Dror, Evgeniy Gabrilovich, and Yehuda Koren. 2012. [Large-scale learning of word relatedness with constraints](#). KDD '12, page 1406–1414, New York, NY, USA. Association for Computing Machinery.
- Felix Hill, Roi Reichart, and Anna Korhonen. 2015. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695.
- Sergio Jimenez, Claudia Bécerra, and Alexander Gelbukh. 2013. [SOFTCARDINALITY-CORE: Improving Text Overlap with Distributional Measures for Semantic Textual Similarity](#). In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, pages 194–201, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Sergio Jimenez, George Dueñas, Julia Baquero, and Alexander Gelbukh. 2014. [UNAL-NLP: Combining Soft Cardinality Features for Semantic Textual Similarity, Relatedness and Entailment](#). In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 732–742, Dublin, Ireland. Association for Computational Linguistics.
- George J Klir and Bo Yuan. 1996. Fuzzy sets, fuzzy logic, and fuzzy systems: selected papers by lotfi a. zadeh.
- Alice Lai and Julia Hockenmaier. 2017. Learning to predict denotational probabilities for modeling entailment. In *EACL*.
- E.T. Lee and L.A. Zadeh. 1969. [Note on fuzzy languages](#). *Information Sciences*, 1(4):421–434.
- Xiang Li, Luke Vilnis, Dongxu Zhang, Michael Boratko, and Andrew McCallum. 2019. Smoothing the geometry of probabilistic box embeddings. *ICLR*.
- Thang Luong, Richard Socher, and Christopher Manning. 2013. [Better word representations with recursive neural networks for morphology](#). In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 104–113, Sofia, Bulgaria. Association for Computational Linguistics.
- Francois Meyer and Martha Lewis. 2020. [Modelling lexical ambiguity with density matrices](#). pages 276–290.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *NIPS*.
- George A Miller and Walter G Charles. 1991. [Contextual correlates of semantic similarity](#). *Language & Cognitive Processes*, 6(1):1–28.
- Douglas Nelson, Cathy Mcevoy, John Walling, and Joseph Wheeler. 1980. [The university of south florida homograph norms](#). *Behavior research methods, instruments, computers: a journal of the Psychonomic Society, Inc*, 12:16–37.
- Maximilian Nickel and Douwe Kiela. 2017. Poincaré embeddings for learning hierarchical representations. In *Neural Information Processing Systems*.
- Yasumasa Onoe, Michael Boratko, and Greg Durrett. 2021. Modeling fine-grained entity types with box embeddings. *Association for Computational Linguistics*.
- Dhruvesh Patel, Pavitra Dangati, Jay-Yoon Lee, Michael Boratko, and Andrew McCallum. 2022. [Modeling label space interactions in multi-label classification using box embeddings](#). In *International Conference on Learning Representations*.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. New Orleans, Louisiana. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, R. Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Kira Radinsky, Eugene Agichtein, Evgeniy Gabrilovich, and Shaul Markovitch. 2011. [A word at a time: Computing word relatedness using temporal semantic analysis](#). WWW '11, page 337–346, New York, NY, USA. Association for Computing Machinery.
- Hongyu Ren, Weihua Hu, and Jure Leskovec. 2020. Query2box: Reasoning over knowledge graphs in vector space using box embeddings. In *8th International Conference on Learning Representations*. OpenReview.net.
- Hongyu Ren and Jure Leskovec. 2020. Beta embeddings for multi-hop logical reasoning in knowledge graphs. *arXiv preprint arXiv:2010.11465*.
- Herbert Rubenstein and John B. Goodenough. 1965. [Contextual correlates of synonymy](#). *Commun. ACM*, 8(10):627–633.

- Ryota Suzuki, Ryusuke Takahama, and Shun Onoda. 2019. Hyperbolic disk embeddings for directed acyclic graphs. In *International Conference on Machine Learning*, pages 6066–6075. PMLR.
- Alexandru Tifrea, Gary Becigneul, and Octavian-Eugen Ganea. 2019. **Poincare glove: Hyperbolic word embeddings**. In *International Conference on Learning Representations*.
- Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 2016. Order-embeddings of images and language. In *International Conference on Learning Representations*.
- Luke Vilnis, Xiang Li, Shikhar Murty, and Andrew McCallum. 2018. Probabilistic embedding of knowledge graphs with box lattice measures. In *Association for Computational Linguistics*.
- Luke Vilnis and Andrew McCallum. 2015. Word representations via gaussian embedding. In *ICLR*.
- Wikipedia contributors. 2022. List of english homographs — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=List_of_English_homographs&oldid=1074954944. [Online; accessed 9-March-2022].
- Dongqiang Yang and David M. W. Powers. 2006. Verb similarity on the taxonomy of wordnet. In *In the 3rd International WordNet Conference (GWC-06), Jeju Island, Korea*.
- L.A. Zadeh. 1965. **Fuzzy sets**. *Information and Control*, 8(3):338–353.
- Rui Zhao and Kezhi Mao. 2018. **Fuzzy Bag-of-Words Model for Document Representation**. *IEEE Transactions on Fuzzy Systems*, 26(2):794–804. Conference Name: IEEE Transactions on Fuzzy Systems.
- Vitalii Zhelezniak, Aleksandar Savkov, April Shen, Francesco Moramarco, Jack Flann, and Nils Y. Hammerla. 2019a. **Don’t settle for average, go for the max: Fuzzy sets and max-pooled word vectors**. In *International Conference on Learning Representations*.
- Vitalii Zhelezniak, Aleksandar Savkov, April Shen, Francesco Moramarco, Jack Flann, and Nils Y Hammerla. 2019b. **DON’T SETTLE FOR AVERAGE, GO FOR THE MAX: FUZZY SETS AND MAX-POOLED WORD VECTORS**. In *International Conference on Learning Representations*.
- Hans-Jürgen Zimmermann. 1987. *Fuzzy sets, decision making, and expert systems*, volume 10. Springer Science & Business Media.

A Preprocessing

The WaCKyedia corpus has been tokenized and lemmatized. We used the lemmatized version of the corpus, however it was observed that various tokens were not split as they should have been (eg. “1.5billion” -> “1.5 billion”). We split tokens using regex criteria to identify words and numbers. All punctuation was removed from the corpus, all numbers were replaced with a “<num>” token, and all words were made lowercase. We also removed any words which included non-ascii symbols. After this step, the entire corpus was tokenized once more, and any token occurring less than 100 times was dropped.

B Dataset Analysis

Dataset	Median Frequency
Men-Tr-3K	23942
Mc-30	25216
Mturk-771	43128
Simlex-999	40653
Verb-143	309192
Yp-130	23044
Rw-Stanford	5683
Rg-65	13088
Ws-353-All	58803
Ws-353-Sim.	57514
Ws-353-Rel	64490
Mturk-287	32952
Simverb-3500	39020

Table 8: Median Frequency of each similarity dataset.

C Hyperparameters

As discussed in Section 6, we train on 128 dimensional WORD2VEC and 64 dimensional WORD2BOX models for 10 epochs. We ran at least 60 runs for each of the models with random seed and randomly chose hyperparameter from the following range - batch_size:[2048, 4096, 8192, 16384, 32768], learning rate log_uniform[exp(-1), exp(-10)], Window_size: [5, 6, 7, 8, 9, 10], negative_samples: [2, 5, 10, 20], sub_sampling threshold: [0.001, 0.0001]. The best working hyperparameter sets and the corresponding checkpoints can be found here:

A	B	$A \cap B$
<i>girl</i>	<i>boy</i>	<i>child</i>
<i>pet</i>	<i>wolf</i>	<i>dog</i>
<i>winner</i>	<i>medal</i>	<i>gold</i>
<i>video</i>	<i>entertainment</i>	<i>movie</i>
<i>ocean</i>	<i>sound</i>	<i>wave</i>
<i>finance</i>	<i>river</i>	<i>bank</i>
<i>parent</i>	<i>woman</i>	<i>mother</i>
<i>bird</i>	<i>America</i>	<i>eagle</i>
<i>car</i>	<i>sea</i>	<i>boat</i>
<i>farm</i>	<i>animal</i>	<i>cow</i>
<i>fruit</i>	<i>yellow</i>	<i>banana</i>
<i>house</i>	<i>royal</i>	<i>palace</i>
<i>property</i>	<i>chemistry</i>	<i>solubility</i>
<i>bank</i>	<i>river</i>	<i>basin</i>
<i>policy</i>	<i>government</i>	<i>legislation</i>
<i>incense</i>	<i>odor</i>	<i>candle</i>
<i>spirit</i>	<i>drink</i>	<i>beer</i>
<i>dance</i>	<i>song</i>	<i>ballad</i>
<i>work</i>	<i>art</i>	<i>painting</i>
<i>instrument</i>	<i>wind</i>	<i>flute</i>

Table 9: Set theoretic queries for *Intersection*. In this task, given A and B , the model need to predict the word for $A \cap B$.

D Dataset for Set Theoretic Queries

In this section, we describe the dataset for set theoretic evaluation. We evaluate on *set-intersection*, *set-difference*, *set-union* queries. For each of these tasks, we create queries of the form $\langle A, B, A \circ B \rangle$, where, \circ is any of the mentioned set operation. In case of *set-union*, we find homographs to be an excellent choice as they are words describing multiple different choices of words. We choose commonly used homographs from list of homographs available in wikipedia (Wikipedia contributors, 2022) to construct this dataset. We manually eliminated many of the words which are rare or when the homographs are referring to concepts which are semantically similar. We provide some examples of the dataset for union queries in table 10. Also, note that we can perform the task for *set-difference* by just swapping the B and $A \cup B$, since $B = (A \cup B) \setminus A$, i.e., if we subtract one concept from the homographs then we must get back a set containing the other concept. So the same table is being used for *set-difference* task. We manually create a small evaluation set for the *set-intersection* task, listed in Table 9.

<i>A</i>	<i>B</i>	$A \cup B$
<i>table</i>	<i>chair</i>	<i>furniture</i>
<i>car</i>	<i>plane</i>	<i>transportation</i>
<i>city</i>	<i>village</i>	<i>location</i>
<i>wolf</i>	<i>bear</i>	<i>animal</i>
<i>shirt</i>	<i>pant</i>	<i>clothes</i>
<i>computer</i>	<i>phone</i>	<i>electronics</i>
<i>red</i>	<i>blue</i>	<i>color</i>
<i>movie</i>	<i>book</i>	<i>entertainment</i>
<i>school</i>	<i>college</i>	<i>education</i>
<i>doctor</i>	<i>engineer</i>	<i>profession</i>
<i>box</i>	<i>circle</i>	<i>shape</i>
<i>big</i>	<i>small</i>	<i>size</i>
<i>dog</i>	<i>tree</i>	<i>bark</i>
<i>fish</i>	<i>tone</i>	<i>bass</i>
<i>sports</i>	<i>wing</i>	<i>bat</i>
<i>carry</i>	<i>animal</i>	<i>bear</i>
<i>sadness</i>	<i>color</i>	<i>blue</i>
<i>bend</i>	<i>weapon</i>	<i>bow</i>
<i>hit</i>	<i>food</i>	<i>buffet</i>
<i>combine</i>	<i>building</i>	<i>compound</i>
<i>happy</i>	<i>list</i>	<i>content</i>
<i>acquire</i>	<i>agreement</i>	<i>contract</i>

Table 10: Examples of set theoretic queries for *Union*. In this task, given *A* and *B*, the model need to predict the word for $A \cap B$. Also note that, we use the same table for the *set difference* queries by treating swapping the *B* and $A \cup B$ columns.