# Box Embeddings: An open-source library for representation learning using geometric structures

**Tejas Chheda**[*†], **Purujit Goyal**[*†], **Trang Tran**[*†‡], **Dhruvesh Patel** [†],
**Michael Boratko**[†], **Shib Sankar Dasgupta**[†], and **Andrew McCallum**[†]

[†] College of Information and Computer Sciences
University of Massachusetts Amherst, MA 01003, USA
[‡] MassMutual Data Science, MA 01002, USA
{tchheda,purujitgoyal,ttrang,dhruveshpate}@cs.umass.edu
{mboratko,ssdasgupta,mccallum}@cs.umass.edu

## Abstract

A major factor contributing to the success of modern representation learning is the ease of performing various vector operations. Recently, objects with geometric structures (eg. distributions, complex or hyperbolic vectors, or regions such as cones, disks, or boxes) have been explored for their alternative inductive biases and additional representational capacities. In this work, we introduce Box Embeddings, a Python library that enables researchers to easily apply and extend probabilistic box embeddings. [1] Fundamental geometric operations on boxes are implemented in a numerically stable way, as are modern approaches to training boxes which mitigate gradient sparsity. The library is fully open-source, and compatible with both PyTorch and TensorFlow, which allows existing neural network layers to be replaced with or transformed into boxes effortlessly. In this work, we present the implementation details of the fundamental components of the library, and the concepts required to use box representations alongside existing neural network architectures.

## 1 Introduction

Much of the success of modern deep learning rests on the ability to learn representations of data compatible with the structure of deep architectures used for training and inference (Hinton, 2007; LeCun et al., 2015). Vectors are the most common choice of representation, as linear transformations are well understood and element-wise non-linearities offer increased representational capacity while being straightforward to implement. Recently, various alternatives to vector representations have been explored, each with different inductive biases or capabilities. Vilnis and McCallum (2015) represent words using Gaussian distributions, which can be thought of as a vector representation with an explicit parameterization of variance. This variance was demonstrated to be capable of capturing the generality of concepts, and KL-divergence provides a natural asymmetric operation between distributions, ideas which were expanded upon in Athiwaratkun and Wilson (2018). Nickel and Kiela (2017), on the other hand, change the embedding space itself from Euclidean to hyperbolic space, where the negative curvature has been shown to provide a natural inductive bias toward modeling tree-like graphs (Nickel and Kiela, 2018; Weber, 2020; Weber and Nickel, 2018).

A subset of these alternative approaches explores *region-based* representations, where entities are not represented by a single point in space but rather explicitly parameterized regions whose volumes and intersections are easily calculated. Order embeddings (Vendrov et al., 2016) represent elements using infinite cones in $\mathbb{R}^n_+$ and demonstrate their efficacy of modeling *partial orders*. Lai and Hockenmaier (2017) endow order embeddings with probabilistic semantics by integrating the space under a negative exponential measure, allowing the calculation of arbitrary marginal, joint, and conditional probabilities. Cone representations are not particularly flexible, however - for instance, the resulting probability model cannot represent negative correlation - motivating the development of *probabilistic box embeddings* (Vilnis et al., 2018), where entities are represented by $n$-dimensional rectangles (i.e. Cartesian products of intervals) in Euclidean space.

Probabilistic box embeddings have undergone several rounds of methodological improvements.

---

[*] Equal Contributions.

[1] The source code and the usage and API documentation for the library is available at https://github.com/iesl/box-embeddings and https://www.iesl.cs.umass.edu/box-embeddings/main/index.html, respectively. A quick video tutorial is available at https://youtu.be/MEPDw8sIwUY.

The original model used a surrogate function to pull disjoint boxes together, which was improved upon in Li et al. (2018) via Gaussian convolution of box indicator functions, resulting in a smoother loss landscape and better performance as a result. Dasgupta et al. (2020) improved box training further by using a latent random variable approach, where the corners of boxes are modeled using Gumbel random variables. These latter models lacked valid probabilistic semantics, however, a fact rectified in Boratko et al. (2021).

While each methodological improvement demonstrated better performance on various modeling tasks, the implementations grew more complex, bringing with it various challenges related to performance and numerical stability. Various applications of probabilistic box embeddings (eg. modeling joint-hierarchies (Patel et al., 2020), uncertain knowledge graph representation (Chen et al., 2021), or fine-grained entity typing (Onoe et al., 2021)) have relied on bespoke implementations, adding unnecessary difficulty and differences in implementation when applying box embeddings to new tasks. To mitigate this issue and make applying and extending box embeddings easier, we saw the need to introduce a reusable, unified, stable library that provides the basic functionalities needed in studying box embeddings. To this end, we introduce "Box Embeddings", a fully open-source Python library hosted on PyPI. The contributions of this work are as follows:

- Provide a modular and reusable library that aids the researchers in studying probabilistic box embeddings. The library is compatible with both of the most popular Machine Learning libraries: PyTorch and TensorFlow.
- Create extensive documentation and example code, demonstrating the use of the library to make it easy to adapt to existing code-bases.
- Rigorously unit-test the codebase with high coverage, ensuring an additional layer of reliability.

## 2   Box Embeddings

Formally, a "box" is defined as a Cartesian product of closed intervals,

$$B(\theta) = \prod_{i=1}^{n} [z_i(\theta), Z_i(\theta)]$$
$$= [z_1(\theta), Z_1(\theta)] \times \cdots \times [z_n(\theta), Z_n(\theta)],$$

| | |
|---|---|
| $z$ | the lower-left coordinate of the boxes |
| $Z$ | the top-right coordinate of the boxes |
| centre | the center coordinate of the boxes, $\frac{z+Z}{2}$ |
| box_shape | shape of the center coordinates (or $z$, $Z$) |
| box_reshape | if possible, reshapes the box_shape into the target_shape |
| broadcast | if possible, adds new dimensions to the box_shape to make it compatible with the target_shape |

Table 1: BoxTensor Properties

where $\theta$ represent some latent parameters. In the simplest case, $\theta \in \mathbb{R}^{2n}$ are free parameters, and $z_i, Z_i$ are projections onto the $i$ and $n + i$ components, respectively. In general, however, the parameterization may be more complicated, eg. $\theta$ may be the output from a neural network. For brevity, we omit the explicit dependency on $\theta$. The different operations (such as volume and intersection) commonly used when calculating probabilities from box embeddings can all be defined in terms of $z_i, Z_i$ - the *min* and *max* coordinates of the interval in each dimension.

### 2.1   Parameterizations

The fundamental component of the library is the `BoxTensor` class, a wrapper around the `torch.Tensor` and `tensorflow.Tensor` class that represents a tensor/array of boxes. `BoxTensor` is an opaque wrapper, in that it exposes the operations and properties necessary to use the box representations (see table 1) irrespective of the specific way in which the parameters $\theta$ are related to $z_i, Z_i$. The main two properties of the `BoxTensor` are z and Z, which represent the *min* and *max* coordinates of an instance of `BoxTensor`. Listing 1 shows how to create an instance of `BoxTensor` consisting of two 2-dimensional boxes in Figure 1.

```python
import torch
from box_embeddings.parameterizations import
    BoxTensor
theta = torch.tensor(
    [[[-2, -2], [-1, -1]], [[1, 0], [3, 4]]]
)
box = BoxTensor(theta)
A = box[0]
B = box[1]
```

Listing 1: Manually initializing a `BoxTensor` consisiting for the 2-D boxes depicted in Figure 1.

Given a `torch.Tensor` corresponding to the parameters $\theta$ of a `BoxTensor`, one can obtain a box representation in multiple ways depending on the constraints on the *min* and *max* coordinates of the box representations as well as the the range of values in $\theta$. The `BoxTensor` class itself simply splits $\theta$ in half on the last dimension, using $\theta[\ldots, 1 : n]$ as $z$ and $\theta[\ldots, n + 1 : 2n]$ as $Z$.
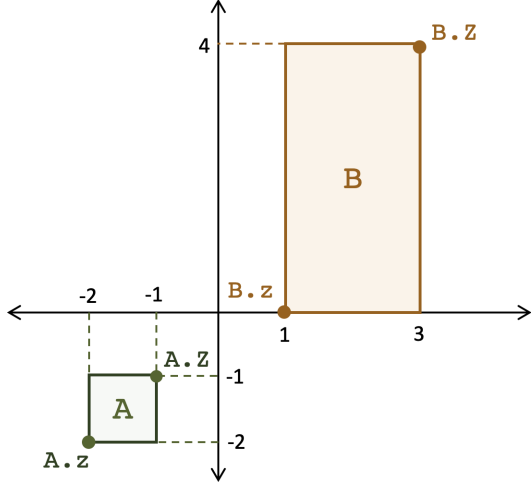
Figure 1: Box Parameterization

Here, the `Ellipsis` "…" denotes any number of leading dimensions, for instance, batch, sequence-length, etc. For the sake of simplifying the notations, from here on, the presence of the leading dimensions will not be explicitly denoted using the `Ellipsis`. Moreover, all the indexing operations can be assumed to be operating only on the last dimension, unless stated otherwise.

```
from box_embeddings.parameterizations import
    BoxTensor, MinDeltaBoxTensor, SigmoidBoxTensor
box_tensor = BoxTensor(theta)
box_tensor_pos_sides = MinDeltaBoxTensor(theta)
box_tensor_in_unit_cube = SigmoidBoxTensor(theta)
```

Listing 2: Converting latent vectors to boxes, for various choices of box parameterizations.

Any box can be represented in this fashion, however some settings of $\theta$ may lead to situations where $z_i > Z_i$. This scenario is invalid under conventional box models (Vilnis et al., 2018; Li et al., 2018), and although valid for models which interpret these coordiantes as parameters of a latent random variable (Dasgupta et al., 2020; Boratko et al., 2021) it is often still desirable to constrain side-lengths to be non-negative. `MinDeltaBoxTensor` represents boxes that are unbounded and have non-negative side-length in each dimension. That is, it outputs boxes with $z, Z \in \mathbb{R}^n$ and $z_i \leq Z_i$, and furthermore any such box has a corresponding $\theta$ under this parameterization. A valid probabilistic interpretation of box embeddings requires that their embedding space has finite measure, however. One trivial way to accomplish this is to parameterize boxes to remain within the unit hypercube, which can be accomplished via the `SigmoidBoxTensor` or `TanhBoxTensor` classes. The specific mathematical operations re-

| Parameterization | $z$ | $Z$ |
|---|---|---|
| BoxTensor | $\theta[1:n]$ | $\theta[n+1:2n]$ |
| MinDeltaBoxTensor | $\theta[1:n]$ | $z + \text{softplus}(\theta[n+1:2n])$ |
| SigmoidBoxTensor | $\sigma(\theta[1:n])$ | $z + (1-z)\sigma(\theta[n+1:2n])$ |
| TanhBoxTensor | $\frac{\tanh(\theta[1:n])+1}{2}$ | $z + \frac{(1-z)\tanh(\theta[n+1:2n])}{2}$ |

Table 2: The different subclasses of `BoxTensor` and how they represent boxes using the learnable parameters $\theta \in \mathbb{R}^{2n}$ taken as input.

lating the $\theta$ variables to their $z, Z$ coordinates are found in Table 2, and example usage can be found in Listing 2.[2]

### 2.2 Operations on BoxTensor

We provide a variety of modules that implement different operations on the box-tensors, such as `Intersection`, `Volume`, `Pooling` and `Regularization`. We also implemented a `BoxEmbedding` layer that, just like a vector embedding layer, provides index lookup. However, unlike a vector embedding layer, this returns boxes instead of vectors. We discuss these layers in detail below.

#### 2.2.1 Intersection

Given two instances of `BoxTensor` with compatible shapes, this operation performs the intersection between the two box-tensors and returns an instance of `BoxTensor` as the result. For two instances of `BoxTensor` $A$ and $B$ with coordinates $(z_A, Z_A)$ and $(z_B, Z_B)$ respectively, the $(z, Z)$ coordinates of the resulting intersection box for the two types of intersection operations, `HardIntersection` (Vilnis et al., 2018; Li et al., 2018) and `GumbelIntersection` (Dasgupta et al., 2020), are shown in Table 3, and corresponding codes are provided in Listing 3.

```
from box_embeddings.parameterizations import
    BoxTensor
from box_embeddings.modules.intersection import
    HardIntersection, GumbelIntersection

boxA = BoxTensor(theta_a)
boxB = BoxTensor(theta_b)

hard_intersection = HardIntersection()
gumbel_intersection =
    GumbelIntersection(intersection_temperature=0.8)

hard_ab = hard_intersection(boxA, boxB)
gumbel_ab = gumbel_intersection(boxA, boxB)
```

Listing 3: Various approaches to computing the intersection of two box tensors.

---

[2]The TensorFlow version for all the code snippets is provided in Appendix.

| Intersection type | $z$ | $Z$ |
|---|---|---|
| `HardIntersection` | $\max(z_A, z_B)$ | $\min(Z_A, Z_B)$ |
| `GumbelIntersection` | $\beta \operatorname{LSE}(\frac{z_A}{\beta}, \frac{z_B}{\beta})$ | $-\beta \operatorname{LSE}(-\frac{Z_A}{\beta}, -\frac{Z_B}{\beta})$ |

Table 3: Expressions for the two kinds of intersection layers. Here, LSE denotes logsumexp, i.e., $\operatorname{LSE}(x, y) := \log(\exp(x) + \exp(y))$

### 2.2.2 Volume

Boxes (or intersections of boxes) are typically queried for their *volumes*. Our `HardVolume` layer implements the volume calculation as originally introduced in Vilnis et al. (2018), which is simply a direct multiplication of side-lengths. It is in this setting where bounded parameterizations such as `SigmoidBoxTensor` and `TanhBoxTensor` are particularly useful, as the resulting volumes can be interpreted as yielding a valid marginal or joint probability. Note, however, that the guarantees of positive side-lengths do not apply when taking the intersection of two disjoint boxes, in which case the resulting box should have zero volume.

Our `SoftVolume` layer implements the volume function proposed by Li et al. (2018), which mitigates the training difficulties that arise when disjoint boxes should overlap. Finally, our `BesselApproxVolume` layer implements the volume function proposed in Dasgupta et al. (2020), which approximates the expected volume of a box where the coordinates are interpreted as location parameters of Gumbel random variables. The expressions and the code snippets for the various volume operations are given in Table 4 and 4, respectively.

**Remark 1.** Note that due to the presence of the product, the naive implementation of volume computations as shown in Table 4 will often result in numerical overflow or underflow for dimensions greater than 5. Hence, we provide an option to compute the volume in log-space, which is *on* by default.

```
from box_embeddings.modules.volume import
    HardVolume, SoftVolume, BesselApproxVolume

hard_volume = HardVolume()
log_volA = hard_volume(boxA)

soft_volume = SoftVolume(volume_temperature=5.0)
log_vol_ab = soft_volume(hard_ab)

bessel_volume =
    BesselApproxVolume(volume_temperature=5.0,
    intersection_temperature=0.8)
log_vol_ab = bessel_volume(gumbel_ab)
```

Listing 4: Different proposed methods for computing box volume, of increasing "smoothness".

| Intersection type | Volume |
|---|---|
| `HardVolume` | $\prod_{i=1}^{n} \max(Z_i - z_i, 0)$ |
| `SoftVolume` | $\prod_{i=1}^{n} T * \operatorname{softplus}(\frac{Z_i - z_i}{T})$ |
| `BesselApproxVolume` | $\prod_{i=1}^{n} T * \operatorname{softplus}(\frac{Z_i - z_i - 2\gamma\beta}{T})$ |

Table 4: The expressions for different volume implementations. Here, $(z, Z)$ are the min-max coordinates of the input `BoxTensor`, $T$ is the volume temperature hyperparameter, $\gamma$ is the Euler-Mascheroni constant, $\beta$ is the gumbel intersection parameter, and $\operatorname{softplus}(x) = \log(1 + \exp x)$.

### 2.2.3 Pooling

The library also provides pooling operations that take as input an instance of `BoxTensor` and reduce one of the leading dimensions by pooling across it. Currently, there are two types of pooling operations implemented – intersection based, which takes intersection across all the boxes in a particular dimension, and mean based, which takes the arithmetic mean of the min and max coordinates of the boxes across a dimension.

### 2.2.4 Regularization

There is an excessive slackness in the learning objective defined using containment conditions on boxes, which leads to large flat regions of local minima resulting in poor training. In order to mitigate this problem, Patel et al. (2020) introduces volume based regularization for boxes, which augments the loss with a penalty if the box volume exceeds a certain threshold. This penalty reduces the size of the flat local minima facilitating better training of boxes.

```
from box_embeddings.modules.pooling import
    HardIntersectionBoxPooler
from box_embeddings.modules.regularization import
    L2SideBoxRegularizer

pooler = HardIntersectionBoxPooler()
pooled_box = pooler(box)

box_regularizer =
    L2SideBoxRegularizer(log_scale=True)

vol_box = soft_volume(pooled_box)
loss = loss_fn(vol_box) +
    box_regularizer(pooled_box)
```

Listing 5: Box pooling and regularization operations.

### 2.3 Embedding

`BoxTensor` and its children classes, do not store learnable parameters directly, they simply wrap the input tensor and provide an interface which interprets the wrapped tensor as box representation. However, when working with a shallow model

206

(embedding only model), one needs an embedding layer that owns its parameters and outputs boxes corresponding to the input indices. The library provides `BoxEmbedding` layer that works like a native embedding layer in PyTorch or TensorFlow, i.e., it performs index lookup, but instead of returning an instance of the native tensor, it returns instance of `BoxTensor`.

### 2.3.1 Initializers

We also provide an abstract interface `BoxInitializer` to implement various methods for initializing the learnable parameters of the `BoxEmbedding` layer. As a concrete example we implement `UniformBoxInitializer`, which initializes boxes with uniformly random min coordinates and side lengths. This is used as the default initializer for the `BoxEmbedding` layer unless specified otherwise.

## 3 Applications

In this section, we demonstrate the Box Embeddings library by using it to implement models for two real-world tasks: a representation learning task of hierarchical graph modeling (Nickel and Kiela, 2017; Vilnis et al., 2018), and the NLP task of natural language inference (Dagan et al., 2005; Bowman et al., 2015). We first demonstrate the intuition behind the containment-based loss function used to train these models using a toy example involving two 2-dimensional boxes.

### 3.1 Toy example

For the purpose of demonstration, we set up a toy example which embeds a simple graph with just two nodes, $X, Y$ and one edge $(X, Y)$. We start with two non-overlapping boxes at initialization: box$_X$ and box$_Y$, and use SGD to train the parameters that minimize the following loss function

$$\mathcal{L}(\theta) = -\log \frac{\text{Vol}\left(\text{B}(\theta_X) \cap \text{B}(\theta_Y)\right)}{\text{Vol}\left(\text{B}(\theta_Y)\right)}.$$

Geometrically, this encourages box$_Y \subseteq$ box$_X$. If using a box embedding with valid probabilistic semantics, this loss function can be interpreted as binary cross-entropy with $P(X|Y) = 1$.[3] The code for this example can be found in Appendix A.2. We visualize the containment training process in Figure 3. Each line represents the edge of the

box in one dimension, with the left endpoint of a blue or orange line to be the minimum coordinate of a box, and the right endpoint of a line to be the maximum coordinate of a box.

### 3.2 Representing hierarchical graph

Representing relations between the nodes of a hierarchy is useful for various NLP and Machine Learning tasks such as natural language inference (Wang et al., 2019; Sharma et al., 2019), entity typing (Onoe et al., 2021), multi-label classification (Chatterjee et al., 2021), and question answering (Jin et al., 2019; Fang et al., 2020). For example, in Figure 2, knowing the hypernym relationship between the pairs *(herb, basil)*, *(herb, thyme)*, and *(herb, rosemary)* can help paraphrase the sentence "This dish requires basil, thyme and rosemary" into "This dish requires several herbs.". Additionally, knowing the relationship between *(herb, banana)*, and *(fruit, banana)* can help answer questions such as "What is both a herb and a fruit?" Note that this latter example maps directly onto the notion of box intersection, as we are seeking an element contained in both "herb" and "fruit".

For demonstration, we train box embeddings to represent the hypernym graph of WordNet (Miller et al., 1990). Hypernym or IS-A is a transitive relation between a pair of words, where one word (hypernym) represents a general/broader concept, and the other word (hyponym) is a more specific sub-concept (Yu et al., 2015). The transitive reduction of the WordNet noun hierarchy contains 82,114 entities and 84,363 edges. The learning task is framed as an edge classification task where, given a pair of nodes $(h, t)$, the model outputs the probability of existence of an edge from $h$ to $t$. Following Patel et al. (2020), we train an edge classification model using the transitive reduction edges augmented with varying percentages of the transitive closure edges (10%, 25%, 50%) as positive examples and randomly sampled negative examples with positive to negative ratio of 1:10. The `BoxEmbedding` layer is initialized with random boxes representing the nodes of the hypernym graph. For each input pair $x = (h_i, t_i)$, the probability of existence of the edge $h_i \rightarrow t_i$ is computed as

$$P(h_i \rightarrow t_i) = \frac{\text{Vol}\left(\text{B}(\theta_{h_i}) \cap \text{B}(\theta_{t_i})\right)}{\text{Vol}\left(\text{B}(\theta_{t_i})\right)}.$$

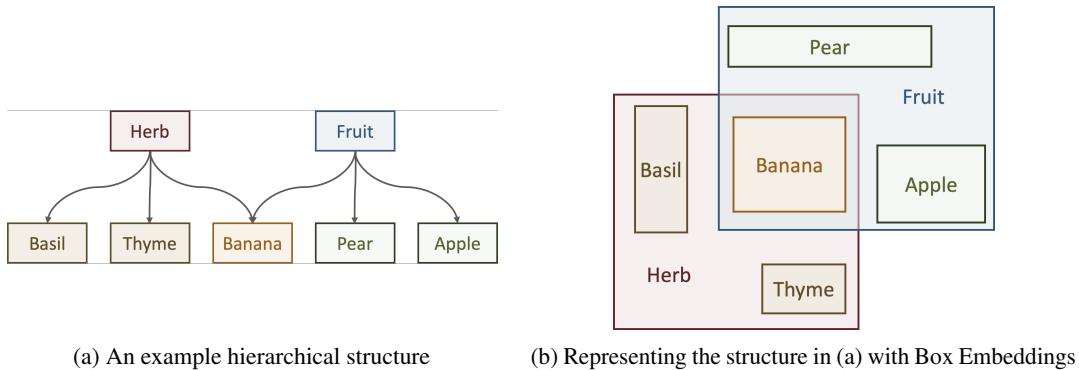In our case, we use `MinDeltaBoxTensor` parameterization, `HardIntersection` and

---

[3]To understand further the motivation for this choice of graph embedding, see Vilnis et al. (2018).

(a) An example hierarchical structure     (b) Representing the structure in (a) with Box Embeddings

Figure 2: Box Embeddings can capture hierarchical structures commonly observed in natural language



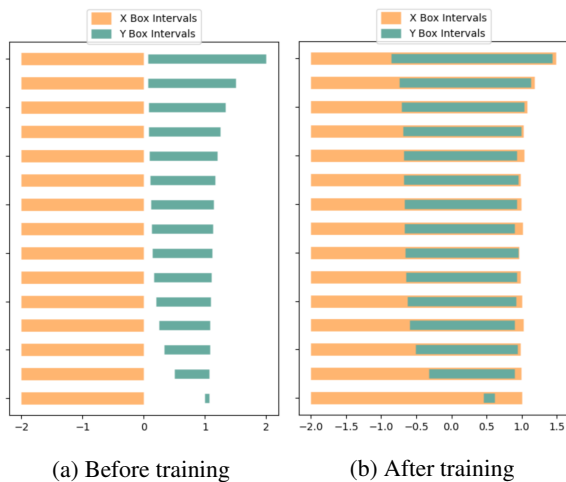(a) Before training     (b) After training

Figure 3: Visualization of two 15-dimensional boxes before and after containment training described in Section 3.1. The green box $B(\theta_Y)$ has been trained to be entirely contained in the orange box $B(\theta_X)$.

| TC Edges | 0% | 10% | 25% | 50% |
|---|---|---|---|---|
| w/o Regularization | 44.2% | 71.3% | 81.1% | 89.1% |
| w Regularization | 59.4% | 90.3% | 91.9% | 94.2% |

Table 5: Test F1 scores for predicting the transitive closure of WordNet's hypernym relations when training on increasing amounts of edges from the transitive closure

`SoftVolume`. Binary cross-entropy loss is used to train the model for edge classification. The test set consists of positive edges sampled from the rest of the transitive closure (not seen during training) and a fixed set of random negatives with the same positive to negative ratio as training. As seen in Table 5, we are able to replicate the result from Patel et al. (2020).

### 3.3 Natural Language Inference (NLI)

Natural language inference (Dagan et al., 2005; Bowman et al., 2015) is a task where, given two sentences, premise and hypothesis, the model is required to pick whether the premise entails the hypothesis, contradicts the hypothesis, or whether neither relationship holds. The task of NLI is setup as multi-class classification, and in the two-class version, the model is only required to decide whether the premise entails the hypothesis or not (Mishra et al., 2021). Although NLI deals with a pair of sentences at a time, in the space of all possible sentences the transitive relation of entailment establishes a partial order. If the sentences are encoded as boxes then we can train box containment to capture the transitive entailment relation. To demonstrate this, we choose the MNLI corpus (Williams et al., 2018) from the GLUE benchmark (Wang et al., 2018). Since the MNLI dataset presents the NLI task as a three-class problem, we collapse *contradiction* and *neutral* labels into a single label called *not-entails* to obtain a two-class problem with class labels *entails* and *not-entails*.

In order to obtain box representation for the premise and hypothesis sentences, we use a neural network $E$ to first get vector representations $v_p$ and $v_h$ for the premise and the hypothesis, respectively. Both these vectors are then interpreted as the parameters $\theta_p := v_p$ and $\theta_h := v_h$ of a box tensor. Finally, the probability of the *entails* class is computed as

$$P(\text{entails}) = \frac{\text{Vol}\left(\mathbf{B}(\theta_p) \cap \mathbf{B}(\theta_h)\right)}{\text{Vol}\left(\mathbf{B}(\theta_h)\right)}.$$

The parameters of the encoder are trained using the ADAM optimizer (Kingma and Ba, 2014) with binary cross-entropy as the loss. Table 6 shows the test accuracy with two different encoders. As seen, the performance is much higher than random or majority class baselines.

| Neural Network Encoder (E) | Accuracy |
|---|---|
| RoBERTa | 78% |
| LSTM | 73% |
| Random Baseline | 50% |
| Majority Baseline | 66% |

Table 6: Test accuracy on MNLI task using box embeddings

## 4 Conclusion

In this paper, we have introduced Box Embeddings, the first Python library focused on allowing region-based representations to be used with deep learning libraries. Our library implements proposed training methods and geometric operations on probabilistic box embeddings in a well-tested and numerically-stable fashion. We described the concepts needed to understand and apply this library to novel tasks, and applied the library to graph modeling and natural language inference, demonstrating both shallow and deep contextualized box representations. We hope the release of this package will aid researchers in using region-based representations in their work, and that the well-documented codebase will facilitate additional methodological extensions to probabilistic box embedding models.

## Acknowledgements

## References

Ben Athiwaratkun and Andrew Gordon Wilson. 2018. Hierarchical density order embeddings. In *International Conference on Learning Representations*.

Michael Boratko, Javier Burroni, Shib Sankar Dasgupta, and Andrew McCallum. 2021. Min/max stability and box distributions. In *Conference on Uncertainty in Artificial Intelligence*. PMLR.

Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.

Soumya Chatterjee, Ayush Maheshwari, Ganesh Ramakrishnan, and Saketha Nath Jagaralpudi. 2021. Joint learning of hyperbolic label embeddings for hierarchical multi-label classification. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2829–2841, Online. Association for Computational Linguistics.

Xuelu Chen, Michael Boratko, Muhao Chen, Shib Sankar Dasgupta, Xiang Lorraine Li, and Andrew McCallum. 2021. Probabilistic box embeddings for uncertain knowledge graph reasoning. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 882–893.

Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Machine Learning Challenges Workshop*, pages 177–190. Springer.

Shib Dasgupta, Michael Boratko, Dongxu Zhang, Luke Vilnis, Xiang Li, and Andrew McCallum. 2020. Improving local identifiability in probabilistic box embeddings. *Advances in Neural Information Processing Systems*, 33.

Yuwei Fang, Siqi Sun, Zhe Gan, Rohit Pillai, Shuohang Wang, and Jingjing Liu. 2020. Hierarchical graph network for multi-hop question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8823–8838.

Geoffrey E. Hinton. 2007. Learning multiple layers of representation. *Trends in Cognitive Sciences*, 11(10).

Hai Jin, Yi Luo, Chenjing Gao, Xunzhu Tang, and Pingpeng Yuan. 2019. Comqa: Question answering over knowledge base via semantic matching. *IEEE Access*, 7:75235–75246.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Alice Lai and Julia Hockenmaier. 2017. Learning to predict denotational probabilities for modeling entailment. In *EACL*.

Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature*, 521.

Xiang Li, Luke Vilnis, Dongxu Zhang, Michael Boratko, and Andrew McCallum. 2018. Smoothing the geometry of probabilistic box embeddings. In *International Conference on Learning Representations*.

George A Miller, Richard Beckwith, Christiane Fellbaum, Derek Gross, and Katherine J Miller. 1990. Introduction to wordnet: An on-line lexical database. *International journal of lexicography*, 3(4):235–244.

Anshuman Mishra, Dhruvesh Patel, Aparna Vijayakumar, Xiang Lorraine Li, Pavan Kapanipathi, and Kartik Talamadupula. 2021. Looking beyond sentence-level natural language inference for question answering and text summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1322–1336, Online. Association for Computational Linguistics.

Maximillian Nickel and Douwe Kiela. 2017. Poincaré embeddings for learning hierarchical representations. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Maximillian Nickel and Douwe Kiela. 2018. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *International Conference on Machine Learning*, pages 3779–3788. PMLR.

Yasumasa Onoe, Michael Boratko, Andrew McCallum, and Greg Durrett. 2021. Modeling fine-grained entity types with box embeddings. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2051–2064, Online. Association for Computational Linguistics.

Dhruvesh Patel, Shib Sankar Dasgupta, Michael Boratko, Xiang Li, Luke Vilnis, and Andrew McCallum. 2020. Representing joint hierarchies with box embeddings. In *Automated Knowledge Base Construction*.

Soumya Sharma, Bishal Santra, Abhik Jana, Santosh Tokala, Niloy Ganguly, and Pawan Goyal. 2019. Incorporating domain knowledge into medical NLI using knowledge graphs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6092–6097, Hong Kong, China. Association for Computational Linguistics.

Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 2016. Order-embeddings of images and language. In *ICLR*.

Luke Vilnis, Xiang Li, Shikhar Murty, and Andrew McCallum. 2018. Probabilistic embedding of knowledge graphs with box lattice measures. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 263–272.

Luke Vilnis and Andrew McCallum. 2015. Word representations via gaussian embedding. In *ICLR*.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, Brussels, Belgium. Association for Computational Linguistics.

Xiaoyan Wang, Pavan Kapanipathi, Ryan Musa, Mo Yu, Kartik Talamadupula, Ibrahim Abdelaziz, Maria Chang, Achille Fokoue, Bassem Makni, Nicholas Mattei, and Michael Witbrock. 2019. Improving natural language inference using external knowledge in the science questions domain. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):7208–7215.

Melanie Weber. 2020. Neighborhood growth determines geometric priors for relational representation learning. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 266–276. PMLR.

Melanie Weber and Maximilian Nickel. 2018. Curvature and representation learning: Identifying embedding spaces for relational data. *NeurIPS Relational Representation Learning*.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.

Zheng Yu, Haixun Wang, Xuemin Lin, and Min Wang. 2015. Learning term embeddings for hypernymy identification. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.

# A Appendix

## A.1 TensorFlow (TF) version

```python
import tensorflow as tf
from box_embeddings.parameterizations import
    TFBoxTensor
theta = tf.Variable(
    [[[0, 0], [2, 2]], [[4, 0], [8, 4]]]
)
box = BoxTensor(theta)
boxA = box[0]
boxB = box[1]
```

Listing 6: TF code for initializing a `BoxTensor`.

```python
from box_embeddings.parameterizations import
    TFMinDeltaBoxTensor, TFSigmoidBoxTensor,
    TFTanhBoxTensor
box_tensor = TFMinDeltaBoxTensor(theta)
box_tensor_pos_sides = TFSigmoidBoxTensor(theta)
box_tensor_in_unit_cube = TFTanhBoxTensor(theta)
```

Listing 7: TF code for converting theta vectors to boxes.

```python
from box_embeddings.parameterizations import
    TFBoxTensor
from box_embeddings.modules.intersection import
    TFHardIntersection
from box_embeddings.modules.intersection import
    TFGumbelIntersection

boxA = TFBoxTensor(theta_a)
boxB = TFBoxTensor(theta_b)

hard_intersection = TFHardIntersection()
gumbel_intersection = TFGumbelIntersection()

hard_ab = hard_intersection(boxA, boxB)
gumbel_ab = gumbel_intersection(boxA, boxB)
```

Listing 8: TF code for computing the intersection of two box tensors.

```python
from box_embeddings.modules.volume import
    TFHardVolume
from box_embeddings.modules.volume import
    TFSoftVolume
from box_embeddings.modules.volume import
    TFBesselApproxVolume

hard_volume = TFHardVolume()
volA = hard_volume(boxA)

soft_volume = TFSoftVolume()
vol_ab = soft_volume(hard_ab)

bessel_volume = TFBesselApproxVolume()
vol_ab = bessel_volume(gumbel_ab)
```

Listing 9: TF code for computing the volume of a box.

```python
from box_embeddings.modules.pooling import
    TFHardIntersectionBoxPooler
from box_embeddings.modules.regularization import
    TFL2SideBoxRegularizer

pooler = TFHardIntersectionBoxPooler()
pooled_box = pooler(box)

box_regularizer =
    TFL2SideBoxRegularizer(log_scale=True)

vol_box = soft_volume(pooled_box)
loss = loss_fn(vol_box) +
    box_regularizer(pooled_box)
```

Listing 10: TF code for performing pooling and regularization operations over a box.

## A.2 Toy Example

```python
import torch
import numpy
from box_embeddings.parameterizations.box_tensor
    import BoxTensor
from box_embeddings.modules.volume.volume import
    Volume
from box_embeddings.modules.intersection import
    Intersection

# Initialization
x_z = numpy.array([-2.0 for n in range(1, 16)])
x_Z = numpy.array([0.0 for k in (x_z)])
data_x = torch.tensor([x_z, x_Z],
    requires_grad=True)
box_H = BoxTensor(data_x)

y_z = numpy.array([1/n for n in range(1, 16)])
y_Z = numpy.array([1 + k for k in reversed(y_z)])
data_y = torch.tensor([y_z, y_Z],
    requires_grad=True)
box_T = BoxTensor(data_y)

# Training function
learning_rate = 0.1
def train(box_1, box_2, optimizer, epochs=1):
    best_loss = int()
    best_box_1 = None
    best_box_2 = None
    box_vol = Volume(volume_temperature=0.1,
        intersection_temperature=0.0001)
    box_int =
        Intersection(intersection_temperature=0.0001)
    for e in range(epochs):
        loss = box_vol(box_2) -
            box_vol(box_int(box_1, box_2))
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
        if best_loss < loss.item():
            best_loss = loss.item()
            best_box_2 = box_2
            best_box_1 = box_1
        print('Iteration %d, loss = %.4f' % (e,
            loss.item()))
    return best_box_1, best_box_2

# Train
optimizer = torch.optim.SGD([data_x, data_y],
    lr=learning_rate)
best_box_H, best_box_T = train(box_H, box_T,
    optimizer, epochs=50)
```

Listing 11: Training Pipeline for the Toy Example (3.1)