

Learning to Predict Denotational Probabilities For Modeling Entailment

Alice Lai and Julia Hockenmaier

Department of Computer Science
University of Illinois at Urbana-Champaign
{aylai2, juliahmr}@illinois.edu

Abstract

We propose a framework that captures the denotational probabilities of words and phrases by embedding them in a vector space, and present a method to induce such an embedding from a dataset of denotational probabilities. We show that our model successfully predicts denotational probabilities for unseen phrases, and that its predictions are useful for textual entailment datasets such as SICK and SNLI.

1 Introduction

In order to bridge the gap between vector-based distributional approaches to lexical semantics that are intended to capture which words occur in similar contexts on the one hand, and logic-based approaches to compositional semantics that are intended to capture the truth conditions under which statements hold on the other hand, Young et al. (2014) introduced the concept of "denotational similarities" that are intended to measure the similarity of simple, declarative statements in terms of the similarity of their truth conditions. Young et al. borrowed the notion of the denotation of a declarative sentence s , $\llbracket s \rrbracket$, as the set of possible worlds in which the sentence is true from classical truth-conditional semantics, and applied it to the domain of image descriptions by defining the *visual denotation* of a sentence (or phrase) as the (sub)set of images it can describe. This allows them to estimate the denotational probabilities of phrases from Flickr30K, a corpus of 30,000 images, each paired with five descriptive captions. Young et al. (2014) and Lai and Hockenmaier (2014) showed that these similarities are complementary to, and potentially more useful for semantic tasks that involve entailment than, standard distributional similarities. However, the systems presented in these

papers were restricted to looking up the denotational similarities of frequent phrases in the training data. In this paper, we go beyond this prior work and define a model that can *predict* the denotational probabilities of novel phrases and sentences. Our experimental results indicate that these predicted denotational similarities are useful for the SICK (Marelli et al., 2014) and SNLI (Bowman et al., 2015) textual entailment datasets.

2 Textual Entailment in SICK and SNLI

The goal of textual entailment is to predict whether a hypothesis sentence is true, false, or neither based on the premise text (Dagan et al., 2013). Due in part to the Recognizing Textual Entailment (RTE) challenges (Dagan et al., 2006), the task of textual entailment recognition has received a lot of attention in recent years. Although full entailment recognition systems typically require a complete NLP pipeline, including coreference resolution etc., this paper considers a simplified variant of this task in which premises and hypotheses are both individual sentences, allowing us to ignore the complexities arising from having to understand longer texts, and instead focus on the purely semantic problem of how to represent the meaning of sentences. This simplified task has been popularized by two datasets, the Sentences Involving Compositional Knowledge (SICK) dataset (Marelli et al., 2014), and the Stanford Natural Language Inference (SNLI) corpus (Bowman et al., 2015), both of which involve a 3-way classification for textual entailment. SICK was created for SemEval 2014 based on image caption data. The premises and hypotheses are automatically generated from the original captions and so contain some unintentional systematic patterns. Most approaches to SICK involve hand-engineered features (Lai and Hockenmaier, 2014)

or large collections of entailment rules (Beltagy et al., 2015). SNLI is the largest textual entailment dataset by several orders of magnitude. It was created to provide enough entailment data to train neural network models for textual entailment. The premises in SNLI are captions from the Flickr30K corpus (Young et al., 2014), and the hypotheses (entailed, contradictory, or neutral in relation to the premise) were solicited from workers on Mechanical Turk. Bowman et al. (2015) illustrate the effectiveness of LSTMs on SNLI. More recent approaches include sentence embedding models (Liu et al., 2016; Munkhdalai and Yu, 2016a), neural attention models (Rocktäschel et al., 2016; Parikh et al., 2016), and neural tree-based models (Munkhdalai and Yu, 2016b; Chen et al., 2016), but we will show that Bowman et al.’s model can be significantly improved by adding a single feature based on our predicted denotational probabilities.

3 Denotational Similarities

In contrast to traditional distributional similarities that are intended to capture which expressions occur in similar contexts, Young et al. (2014) introduced the concept of “denotational similarities” to capture which expressions can be used to describe similar situations, and estimate these similarities from Flickr30K, a corpus of 30,000 images, each paired with five descriptive captions. Borrowing the notion of the denotation of a declarative sentence s , $\llbracket s \rrbracket$, as the set of possible worlds in which the sentence is true from classical truth-conditional semantics, Young et al. first define the *visual denotation* of a sentence (or phrase) s , $\llbracket s \rrbracket$, as the (sub)set of images that s can describe. In order to compute these visual denotations from their corpus, they define a set of normalization and reduction rules (e.g. lemmatization, dropping modifiers, replacing nouns with their hypernyms, dropping PPs, extracting NPs) that augment the original Flickr30K captions with a large number of shorter, more generic phrases that are each associated with a subset of the Flickr30K images. This results in a large subsumption hierarchy over phrases, which they call a denotation graph (see Figure 1). The structure of the denotation graph is similar to the idea of an entailment graph (Berant et al., 2012). Each node in the denotation graph corresponds to a phrase s , associated with its denotation $\llbracket s \rrbracket$, i.e. the set of images for which this

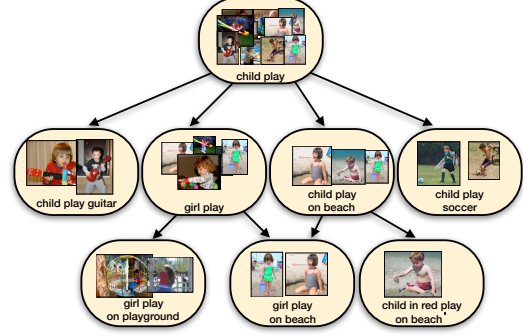


Figure 1: The denotation graph is a subsumption hierarchy over phrases associated with images

phrase could be derived from the original captions. For example, the denotation of a phrase “*woman jog on beach*” is the set of images in the corpus that depict a woman jogging on a beach. Note that the denotation of a node (e.g. “*woman jog on beach*”) is always a subset of the denotations of any of its ancestors (e.g. “*woman jog*”, “*person jog*”, “*jog on beach*” or “*beach*”).

The denotational probability of a phrase s , $P_{\llbracket \rrbracket}(s)$, is a Bernoulli random variable that corresponds to the probability that a randomly drawn image can be described by s . Given a denotation graph over N images, $P_{\llbracket \rrbracket}(s) = \frac{|\llbracket s \rrbracket|}{N}$. The joint denotational probability of two phrases x and y , $P_{\llbracket \rrbracket}(x, y) = \frac{|\llbracket x \rrbracket \cap \llbracket y \rrbracket|}{N}$, indicates how likely it is that a situation can be described by both x and y . Young et al. propose to use pointwise mutual information scores (akin to traditional distributional similarities), and conditional probabilities $P_{\llbracket \rrbracket}(x|y) = \frac{|\llbracket x \rrbracket \cap \llbracket y \rrbracket|}{|\llbracket y \rrbracket|}$ as so-called denotational similarities. The latter are what we will work with here. They are intended to capture entailment-like relations (what is the probability that x is true, given that y can be said about this situation) that can hold due to commonsense knowledge, hyponymy, etc. In an ideal scenario, if the premise p entails the hypothesis h , then the conditional probability $P(h|p)$ is 1 (or close to 1). Similarly, if h contradicts p , then the conditional probability $P(h|p)$ is close to 0. We therefore stipulate that learning to predict the conditional probability of one phrase h given another phrase p would be helpful in predicting textual entailment. We also note that y entails x , and $P_{\llbracket \rrbracket}(x|y) = 1$, if x is an ancestor of y in the graph.

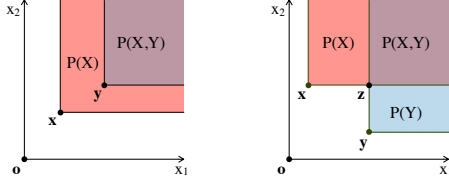


Figure 2: Order embedding interpreted as probabilities.

4 An order embedding for probabilities

Young et al. (2014) and Lai and Hockenmaier (2014) show that denotational probabilities can be at least as useful as traditional distributional similarities for tasks that require semantic inference such as entailment or textual similarity recognition. However, their systems only use those denotational probabilities that can be directly read off of the denotation graph.

Here, we present a model that learns to predict denotational probabilities $P_{\square}(x)$ and $P_{\square}(x|y)$ even for phrases it has not seen during training. Our model is inspired by Vendrov et al. (2016), who observed that a partial ordering \preceq over the vector representations of phrases can be used to express entailment relations. They induce a so-called order embedding for words and phrases such that the vector \mathbf{x} corresponding to phrase x is smaller than the vector \mathbf{y} , i.e. $\mathbf{x} \preceq \mathbf{y}$, for phrases y entailed by x , where \preceq corresponds to the reversed product order on \mathbb{R}_+^N ($\mathbf{x} \preceq \mathbf{y} \Leftrightarrow x_i \geq y_i \forall i$). But although they use their model to predict entailment labels between pairs of sentences, their model is only capable of making binary decisions (entailed or not entailed). We generalize this idea to learn an embedding space that expresses not only the binary relation that phrase x is entailed by phrase y , but also the probability that phrase x is true given phrase y . Specifically, we learn a mapping from phrases x to N -dimensional vectors $\mathbf{x} \in \mathbb{R}_+^N$ such that the vector $\mathbf{x} = (x_1, x_2, \dots, x_N)$ corresponding to phrase x defines the denotational probability of x as $P_{\square}(x) = \exp(-\sum_i x_i)$. The origin (the zero vector) therefore has probability $\exp(0) = 1$. Any vector \mathbf{x} that lies above the origin (i.e. $\forall_i x_i > 0$) has probability less than 1, and a vector \mathbf{x} that is farther from the origin than a vector \mathbf{y} represents a phrase x that has a smaller denotational probability than phrase y . We can visualize this as each phrase vector occupying a region in the embedding space that is proportional to the denotational probability of the phrase. Figure 2 illustrates this

in two dimensions. The zero vector at the origin has a probability proportional to the entire region of the positive orthant, while other points in the space correspond to smaller regions and thus probabilities less than 1. Since the joint probability $P_{\square}(x, y)$ in this embedding space should be proportional to the size of the intersection of the regions of \mathbf{x} and \mathbf{y} , we define the joint probability of two phrases x and y to correspond to the point \mathbf{z} that is the element-wise maximum of \mathbf{x} and \mathbf{y} : $z_i = \max(x_i, y_i)$. This allows us to compute the conditional probability $P_{\square}(x|y)$ as follows:

$$\begin{aligned} P_{\square}(x|y) &= \frac{P_{\square}(x, y)}{P_{\square}(y)} \\ &= \frac{\exp(-\sum_i z_i)}{\exp(-\sum_i y_i)} \\ &= \exp(\sum_i y_i - \sum_i z_i) \end{aligned}$$

Shortcomings We note that this order embedding does not allow us to represent the negation of x as a vector. We also cannot learn an embedding for two phrases with completely disjoint denotations. The best we can do is learn an embedding that assumes that any two phrases are independent. For any pair $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$, $P_{\square}(X, Y) \geq P_{\square}(X)P_{\square}(Y)$, (equality holds if the non-zero coordinates of \mathbf{x} and \mathbf{y} are disjoint):

$$\begin{aligned} P_{\square}(X, Y) &= \exp(-\sum_i \max(x_i, y_i)) \\ &\geq \exp(-\sum_i x_i - \sum_i y_i) \\ &= P_{\square}(X)P_{\square}(Y) \end{aligned}$$

Nevertheless, we will see below that we will come very close to learning that some phrase pairs have very low denotational similarities, while others are highly similar.

5 Our model for $P_{\square}(x)$ and $P_{\square}(x, y)$

We train a neural network model to predict $P_{\square}(x)$, $P_{\square}(y)$, and $P_{\square}(x|y)$ for new phrases x and y . This model consists of an LSTM (Hochreiter and Schmidhuber, 1997) that outputs a 512d vector which is passed through an additional 512d layer. We use 300d GloVe vectors (Pennington et al., 2014) trained on 840B tokens as the word embedding input to the LSTM. We use the same model to represent both phrases x and y regardless of which phrase is the premise or the hypothesis. Thus, we

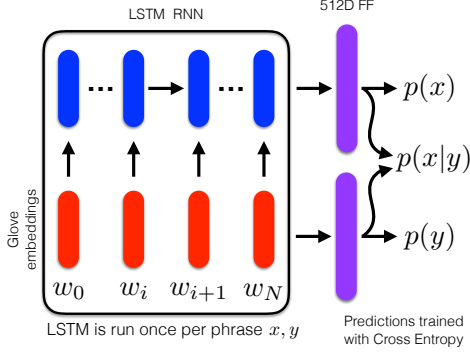


Figure 3: Our probability model architecture.

pass the sequence of word embeddings for phrase x through the model to get \mathbf{x} , and we do the same for phrase y to get \mathbf{y} . As previously described, we sum the elements of \mathbf{x} and \mathbf{y} to get the probabilities. From \mathbf{x} and \mathbf{y} we can compute the joint vector \mathbf{z} , which we use to compute $P_{\square}(x|y)$. Figure 3 illustrates the structure of our model. Our training data consists of ordered phrase pairs $\langle x, y \rangle$. We train this model to predict both the denotational probabilities of each phrase ($P_{\square}(x)$ and $P_{\square}(y)$) as well as the conditional probability $P_{\square}(x|y)$ of x given y (typically the pair $\langle y, x \rangle$ will also appear in the training data). Our per-example loss is the sum of the cross entropy losses for $P_{\square}(x)$, $P_{\square}(y)$, and $P_{\square}(x|y)$:

$$L = - \left(P_{\square}(x) \log Q(x) + (1 - P_{\square}(x)) \log(1 - Q(x)) \right) \\ - \left(P_{\square}(y) \log Q(y) + (1 - P_{\square}(y)) \log(1 - Q(y)) \right) \\ - \left(P_{\square}(x|y) \log Q(x|y) + (1 - P_{\square}(x|y)) \log(1 - Q(x|y)) \right)$$

We use the Adam optimizer with a learning rate of 0.001, and a dropout rate of 0.5. These parameters were tuned on the development data.

Numerical issues In section 4, we described the probability vectors \mathbf{x} as being in the positive orthant. However, in our implementation, we use unnormalized log probabilities. This puts all of our vectors in the negative orthant instead, but it prevents the gradients from becoming too small during training. To ensure that the vectors are in \mathbb{R}_+^N , we clip the values of the elements of \mathbf{x} so that $x_i \leq 0$. To compute $\log(P_{\square}(x))$, we sum the elements of \mathbf{x} and clip the sum to the range $(\log(10^{-10}), -0.0001)$ in order to avoid errors caused by passing $\log(0)$ values to the loss function. The conditional log probability is simply $\log P_{\square}(x|y) = \log P_{\square}(x, y) - \log P_{\square}(y)$, where

$\log P_{\square}(x, y)$ is now the element-wise minimum:

$$\log P_{\square}(x, y) = \sum_i \min(x_i, y_i) \quad (1)$$

5.1 Training regime

To train our model, we use phrase pairs $\langle x, y \rangle$ from the denotation graph generated on the training split of the Flickr30K corpus (Young et al., 2014). We consider all 271,062 phrases that occur with at least 10 images in the training split of the graph, to ensure that they are frequent enough that their probabilities are reliable.

We include all phrase pairs where the two phrases have at least one image in common. These constitute 22 million phrase pairs $\langle x, y \rangle$ with $P_{\square}(x|y) > 0$. To train our model to predict $P_{\square}(x|y) = 0$, we include those phrase pairs $\langle x, y \rangle$ that have no images in common if $N \times P_{\square}(x)P_{\square}(y) \geq N^{-1}$ (N is the total number of images), meaning the phrases occur frequently enough that we would expect them to co-occur at least once. This yields 1.9 million pairs where $P_{\square}(x|y) = 0$. For additional examples of $P_{\square}(x|y) = 1$, we include phrase pairs that have an ancestor-descendant relationship in the denotation graph. We include all ancestor-descendant pairs where each phrase occurs with at least 2 images, for an additional 2.8 million phrase pairs.

For evaluation purposes, we first assign out 5% of the phrases to the development pool and 5% to the test pool. The actual test data then consists of all phrase pairs where at least one of the two phrases comes from the test pool. The resulting test data contains 10.6% unseen phrases by type and 51.2% unseen phrases by token. All phrase pairs in the test data contain at least one phrase that was unseen in the training or development data. The development data was created the same way.

This dataset will be available at <http://XXX>.

We train our model on all 42 million phrase pairs with batch size 512 for 10 epochs and use the mean KL divergence on the conditional probabilities of the phrase pairs in the development data to select the best model. Since $P_{\square}(x|y)$ is a Bernoulli distribution, we compute the KL divergence for each phrase pair $\langle x, y \rangle$ as

$$D_{KL}(P||Q) = P_{\square}(x|y) \log \frac{P_{\square}(x|y)}{Q(x|y)} \\ + (1 - P_{\square}(x|y)) \log \frac{1 - P_{\square}(x|y)}{1 - Q(x|y)} \quad (2)$$

where $Q(x|y)$ is the conditional probability predicted by our model.

6 Predicting denotational probabilities

6.1 Prediction on new phrase pairs

We evaluate our model using both the KL divergences $D_{KL}(P||Q)$ of the gold individual and conditional probabilities $P_{\square}(x)$ and $P_{\square}(x|y)$ against the corresponding predicted probabilities Q , and the Pearson correlation r , which expresses the correlation between two variables (the per-item gold and predicted probabilities) as a value between -1 (total negative correlation) and 1 (total positive correlation). As described above, we compute the KL divergence on a per-item basis, and report the mean over all items in the test or development set.

Table 1 shows the performance of our trained model on unseen test data. The full test data consists of 4.6 million phrase pairs, all of which contain at least one phrase that was not observed in either the training or development data. Our model does reasonably well at predicting these conditional probabilities, reaching a correlation of $r = 0.949$ on the complete test data. On the subset of test phrase pairs where both phrases are previously unseen, the model’s predictions are almost as good at $r = 0.920$.

We also analyze our model’s accuracy on phrase pairs where the gold $P_{\square}(x|y)$ is either 0 or 1. The latter case reflects an important property of the denotation graph, since $P_{\square}(x|y) = 1$ when x is an ancestor of y . More generally, we can interpret $P_{\square}(h|p) = 1$ as a confident prediction of entailment, and $P_{\square}(h|p) = 0$ as a confident prediction of contradiction. Figure 4 shows the distribution of probability predictions when gold $P_{\square}(h|p) = 0$ (top) and $P_{\square}(h|p) = 1$ (bottom). Our model does nearly as well at predicting these probabilities when both phrases in the pair are unseen in the training data (gray bars) as it does on the full test data (black bars).

6.2 Prediction on longer sentences

We evaluate the performance of our probability model on two existing textual entailment datasets, SICK and SNLI.

Our model up to this point has only been trained on short phrases, since conditional probabilities in the denotation graph are only reliable for phrases that occur with multiple images (see Figure 5 for

| | $P(x)$ | | $P(x y)$ | |
|----------------|--------|-------|----------|-------|
| | KL | r | KL | r |
| Training data | 0.0003 | 0.998 | 0.017 | 0.974 |
| Full test data | 0.001 | 0.979 | 0.031 | 0.949 |
| Unseen pairs | 0.002 | 0.837 | 0.048 | 0.920 |

Table 1: Our model predicts the probability of unseen phrase pairs with high correlation to the gold probabilities.

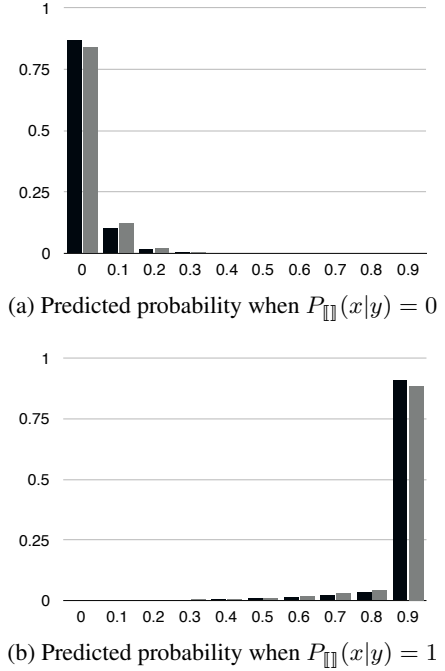


Figure 4: Predicted probabilities on denotational phrase test data when $P_{\square}(x|y)$ is 0 or 1. Black is the full test data and grey is the subset of pairs where both phrases are unseen. Frequency is represented as a percentage of the size of the data.

the distribution of phrase lengths in the training data). To improve our model’s performance on longer sentences, we now train it from scratch on both the 42 million denotation graph training phrase pairs (which are lemmatized) with their original denotational probabilities, and on all 550,000 lemmatized sentence pairs from the SNLI training data (the SNLI training data has a mean sentence length of 11 words). We do not train on SICK because the corpus is much smaller and also has a different distribution of phenomena, including explicit negation. We augment the SNLI data with approximate gold denotational probabilities by assigning a probability $P_{\square}(S) = s/N$ to any sentence S that occurs s times in the N train-

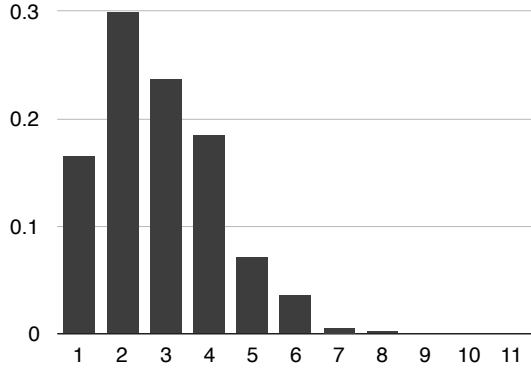


Figure 5: Distribution of phrase lengths as a fraction of the data size on the denotation graph phrase training data.

ing sentences. We approximate the gold conditional probabilities for each sentence pair $\langle p, h \rangle$ according to the entailment label: if p entails h , then $P(h|p) = 0.9$. If p contradicts h , then $P(h|p) = 0.001$. Otherwise, $P(h|p) = 0.5$.

Figure 6 shows predicted probabilities on SNLI sentence pairs. The top row shows predictions on the SNLI test data of our previously described model, which was trained only on short denotational phrases. We observe that the median probabilities increase from contradiction to neutral to entailment, even though this model was only trained on short phrases with a much smaller vocabulary. Given the training data, we can’t expect these probabilities to align cleanly with the entailment labels, but even so, there is already some information here to distinguish between entailment classes.

The bottom row shows that training on both denotational phrases and SNLI sentence pairs with approximate conditional probabilities improves our representation of longer sentences. This model’s predicted conditional probabilities align much more closely with the entailment class labels. Entailing sentence pairs have high probability (median 0.72), neutral sentence pairs have mid-range probabilities (median 0.46), and contradictory sentence pairs have probabilities approaching 0 (median 0.19).

7 Predicting textual entailment

In Section 6.2, we trained our probability model on both short phrase pairs for which we had gold probabilities and longer SNLI sentence pairs for which we estimated probabilities. We now eval-

| | Premise | Hypothesis | G | P |
|----|----------------------------------|------------------------|-----|-----|
| 1 | person walk on trail in woods | in forest | 1.0 | 1.0 |
| 2 | group of person bike | group of person ride | 0.9 | 0.8 |
| 3 | adult sing while play instrument | adult play guitar | 0.8 | 0.8 |
| 4 | person sit on bench outside | on park bench | 0.4 | 0.4 |
| 5 | tennis player hit ball | person swing | 0.2 | 0.2 |
| 6 | girl sleep | on pillow | 0.1 | 0.2 |
| 7 | man practice martial art | person kick person | 0.1 | 0.3 |
| 8 | person skateboard on ramp | man ride skateboard | 0.2 | 0.2 |
| 9 | busy intersection | city street | 0.3 | 0.2 |
| 10 | person dive into swim pool | person fly through air | 0.1 | 0.1 |
| 11 | sit at bench | adult read book | 0.1 | 0.1 |
| 12 | person leap into air | jump over obstacle | 0.0 | 0.0 |
| 13 | person talk on phone | man ride skateboard | 0.0 | 0.0 |

Table 2: Gold and predicted conditional probabilities from the denotational phrase development data.

| Premise | Hypothesis | Gold | Pred |
|----------------------|-------------|------|------|
| skier on snowy hill | athlete | 1.00 | 0.99 |
| pitcher throw ball | mound | 0.53 | 0.84 |
| golf ball | athlete | 0.53 | 0.66 |
| person point | man point | 0.48 | 0.41 |
| in front of computer | person look | 0.36 | 0.21 |

Table 3: Gold and predicted conditional probabilities from unseen pairs in the denotational phrase development data.

uate the effectiveness of this model for textual entailment, and demonstrate that these predicted probabilities are informative features for predicting entailment on both SICK and SNLI.

Model We first train an LSTM similar to the 100d LSTM that achieved the best accuracy of the neural models in Bowman et al. (2015). It produces 100d sentence vectors for the premise and hypothesis based on the words’ GloVe vectors. The concatenated 200d sentence pair representation from the LSTM passes through three 200d tanh layers and a softmax layer for entailment classification. We train the LSTM on SNLI with batch size 512 for 10 epochs. We use the Adam optimizer with a learning rate of 0.001 and a dropout rate of 0.85, and use the development data to select the best model.

Next, we take the output vector produced by the LSTM for each sentence pair and append our predicted $P_{\square}(h|p)$ value (the probability of the hypothesis given the premise). We train another classifier that passes this vector through two tanh layers with a dropout rate of 0.5 and a final softmax

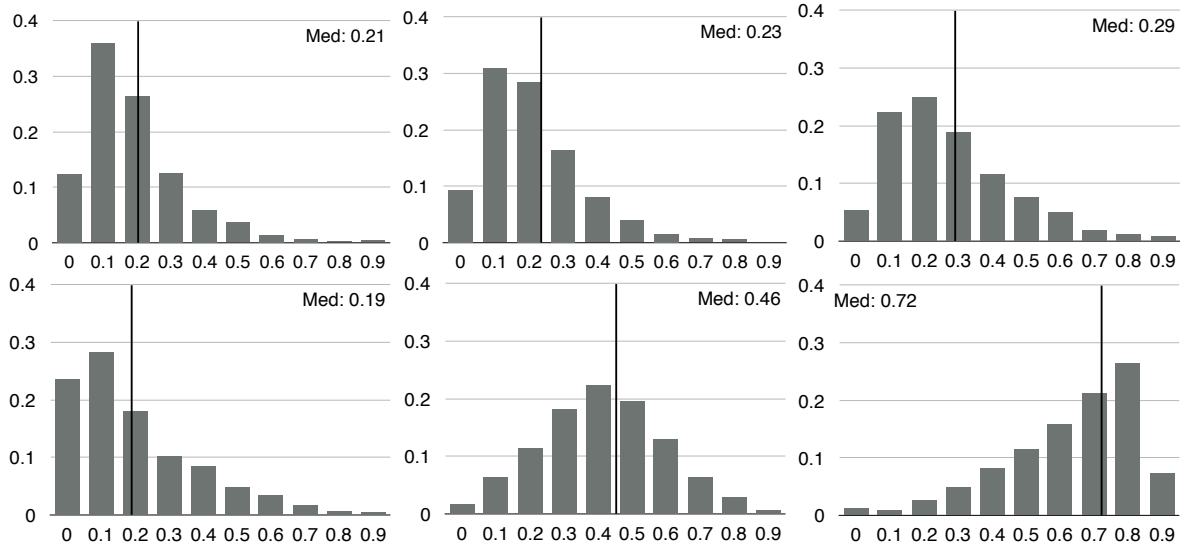


Figure 6: Predicted probabilities for SNLI sentence pairs (test data) by entailment label, as a percentage of pairs with that label. Top row: predictions from the model trained only on short denotational phrases. Bottom row: predictions from the model trained on both short denotational phrases and SNLI (training data).

classification layer. Holding the parameters of the LSTM fixed, we train this model for 10 epochs on SNLI with batch size 512.

Results Table 5 contains our results on SNLI. Our baseline LSTM achieves the same 77.2% accuracy as reported by Bowman et al. (2015), whereas a classifier that combines the output of this LSTM with only a single feature from our probability model improves performance to 78.2% accuracy.

We use the same approach to evaluate the effectiveness of our predictions on SICK (Table 6). SICK does not have enough data to train an LSTM, so we combine the SICK and SNLI training data to train both the LSTM and the final model. When we add the predicted conditional probability as a single feature for each SICK sentence pair, performance increases from 81.5% to 82.7% accuracy. This approach outperforms the transfer learning approach of Bowman et al. (2015), which was also trained on both SICK and SNLI.

8 Discussion

Section 6 has demonstrated that we can successfully learn to predict denotational probabilities both for phrases that we have not encountered during training and for longer sentences. Section 7 has illustrated the utility of these predicted probabilities by showing that a single feature based

on the conditional denotational probabilities predicted by our model improves the accuracy of an LSTM on the SICK and SNLI entailment data sets by 1% or more. Although we were not able to evaluate the impact on more complex, recently proposed models that outperform this LSTM, this improvement is quite encouraging, in particular since we only have accurate denotational probabilities for the short phrases from the denotation graph (mostly 6 words or fewer) that have a much more limited vocabulary than the full SNLI data (there are 5263 word types in the denotation graph training data, while the lemmatized SNLI training data has a vocabulary size of 31,739).

Examining the predicted conditional probability of specific pairs of phrases and sentences demonstrates what our model does well and where it could improve. Table 2 has example predictions from the denotation phrase development data. Table 3 has examples for phrase pairs where both phrases were unseen. Our model correctly learns to predict a high probability for entailed phrases even when there is no direct hypernym involved, as in example 2, and for closely related phrases that are not strictly entailing, as in example 3. Our model also learns to predict reasonable probabilities for phrases that describe events that frequently co-occur but are not required to do so, such as example 7. In examples 10 and 11, our model predicts low probabilities for occasionally co-occurring events, which are still more likely

| | Premise | Hypothesis | CPR |
|---------------|---|---|------|
| Entailment | 1 A person rides his bicycle in the sand beside the ocean. | A person is on a beach. | 0.88 |
| | 2 Two women having drinks and smoking cigarettes at the bar. | Two women are at a bar. | 0.86 |
| | 3 A senior is waiting at the window of a restaurant that serves sandwiches. | A person waits to be served his food. | 0.61 |
| | 4 A man with a shopping cart is studying the shelves in a supermarket aisle. | There is a man inside a supermarket. | 0.47 |
| | 5 The two farmers are working on a piece of John Deere equipment. | John Deere equipment is being worked on by two farmers. | 0.16 |
| Neutral | 6 A group of young people with instruments are on stage. | People are playing music. | 0.86 |
| | 7 Two doctors perform surgery on patient. | Two doctors are performing surgery on a man. | 0.56 |
| | 8 Two young boys of opposing teams play football, while wearing full protection uniforms and helmets. | boys scoring a touchdown | 0.30 |
| | 9 Two men on bicycles competing in a race. | Men are riding bicycles on the street. | 0.24 |
| Contradiction | 10 Two women having drinks and smoking cigarettes at the bar. | Three women are at a bar. | 0.79 |
| | 11 A man in a black shirt is playing a guitar. | The man is wearing a blue shirt. | 0.47 |
| | 12 An asian woman sitting outside an outdoor market stall. | A woman sitting in an indoor market. | 0.22 |
| | 13 A white dog with long hair jumps to catch a red and green toy. | A white dog with long hair is swimming underwater. | 0.09 |
| | 14 Two women are embracing while holding to go packages. | The men are fighting outside a deli. | 0.06 |

Table 4: Examples of predicted conditional probability on sentence pairs from the SNLI development data.

| Model | Test Acc. |
|---------------------------|-------------|
| Our LSTM | 77.2 |
| Our LSTM + CPR | 78.2 |
| Bowman et al. (2015) LSTM | 77.2 |

Table 5: Entailment accuracy on SNLI (test).

| Model | Test Acc. |
|-------------------------------|-------------|
| Our LSTM | 81.5 |
| Our LSTM + CPR | 82.7 |
| Bowman et al. (2015) transfer | 80.8 |

Table 6: Entailment accuracy on SICK (test).

than the improbable co-occurrences of examples 12 and 13.

Table 4 has example predictions for sentence pairs from SNLI development data. Some cases of entailment are straightforward and predicting high probability is easy, such as example 2, which simply involves dropping words from the premise. In other cases, our model correctly predicts high probability for an entailed hypothesis that does not have such obvious word-to-word correspondence with the premise, e.g. example 1. Our model’s predictions are less accurate when the sentence structure differs between premise and hypothesis or when there are many unknown words, as in ex-

ample 5. For neutral pairs, our model usually predicts mid-range probabilities, but there are some exceptions like example 6. In this example, it is not certain that the people are playing music, but it is a reasonable assumption from the premise, more so than other neutral examples. Our model cannot reason about numbers and quantities, as example 10 shows. It also fails to predict in example 11 that a man wearing a black shirt is probably not also wearing a blue shirt. However, our model does predict low probabilities for some contradictory examples that have reasonably high word overlap, as in example 13.

9 Conclusion

We have presented a framework for representing denotational probabilities in a vector space and demonstrate that we can successfully learn to predict these probabilities for new phrases with a neural network model. We have demonstrated that when trained on sentences with approximate probabilities in addition to short phrases, our model can learn reasonable representations for longer sentences. We have shown that our model’s predicted probabilities are useful for textual entailment and provide additional gains in performance when added to existing competitive textual entailment classifiers. Future work will examine whether the embeddings it learns can be used directly by these classifiers, and how to incorporate negation into our model.

References

- Islam Beltagy, Stephen Roller, Pengxiang Cheng, Katrin Erk, and Raymond J. Mooney. 2015. Representing meaning with a combination of logical form and vectors. *CoRR*, abs/1505.06816.
- Jonathan Berant, Ido Dagan, and Jacob Goldberger. 2012. Learning entailment relations by global graph structure optimization. *Computational Linguistics*, 38(1):73–111, March.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642. Association for Computational Linguistics.
- Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, and Hui Jiang. 2016. Enhancing and combining sequential and tree LSTM for natural language inference. *CoRR*, abs/1609.06038.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The PASCAL recognising textual entailment challenge. In *Proceedings of the First International Conference on Machine Learning Challenges: Evaluating Predictive Uncertainty Visual Object Classification, and Recognizing Textual Entailment*, MLCW’05, pages 177–190, Berlin, Heidelberg. Springer-Verlag.
- Ido Dagan, Dan Roth, Mark Sammons, and Fabio Massimo Zanzotto. 2013. Recognizing textual entailment: Models and applications. *Synthesis Lectures on Human Language Technologies*, 6(4):1–220.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November.
- Alice Lai and Julia Hockenmaier. 2014. Illinois-LH: A denotational and distributional approach to semantics. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 329–334. Association for Computational Linguistics.
- Yang Liu, Chengjie Sun, Lei Lin, and Xiaolong Wang. 2016. Learning natural language inference using bidirectional LSTM model and inner-attention. *CoRR*, abs/1605.09090.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. 2014. A SICK cure for the evaluation of compositional distributional semantic models. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*. European Language Resources Association (ELRA).
- Tsendsuren Munkhdalai and Hong Yu. 2016a. Neural semantic encoders. *CoRR*, abs/1607.04315.
- Tsendsuren Munkhdalai and Hong Yu. 2016b. Neural tree indexers for text understanding. *CoRR*, abs/1607.04492.
- Ankur P. Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. 2016. A decomposable attention model for natural language inference. *CoRR*, abs/1606.01933.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomas Kocisky, and Phil Blunsom. 2016. Reasoning about entailment with neural attention. In *ICLR*.
- Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 2016. Order-embeddings of images and language. In *Proceedings of ICLR*.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations. *Transactions of the Association of Computational Linguistics – Volume 2, Issue 1*, pages 67–78.