

CS224n: Deep Learning for NLP¹

Lecture Notes: Part VII ²

Winter 2017

¹ Course Instructor: Richard Socher

² Authors: Francois Chaubard, Richard Socher

Keyphrases: RNN, Recursive Neural Networks, MV-RNN, RNTN

This set of notes discusses and describes the many variants on the RNN (Recursive Neural Networks) and their application and successes in the field of NLP.

1 Recursive Neural Networks

In these notes, we introduce and discuss a new type of model that is indeed a superset of the previously discussed Recurrent Neural Network. Recursive Neural Networks (RNNs) are perfect for settings that have nested hierarchy and an intrinsic recursive structure. Well if we think about a sentence, doesn't this have such a structure? Take the sentence "A small crowd quietly enters the historical church". First, we break apart the sentence into its respective Noun Phrase, Verb Phrase, "A small crowd" and "quietly enters the historical church", respectively. But there is a noun phrase, verb phrase within that verb phrase right? "quietly enters" and "historical church". etc, etc. Seems pretty recursive to me.

The syntactic rules of language are highly recursive. So we take advantage of that recursive structure with a model that respects it! Another added benefit of modeling sentences with RNN's is that we can now input sentences of arbitrary length, which was a huge head scratcher for using Neural Nets in NLP, with very clever tricks to make the input vector of the sentence to be of equal size despite the length of the sentences not being equal. (see Bengio et al., 2003; Henderson, 2003; Collobert & Weston, 2008)

Let's imagine our task is to take a sentence and represent it as a vector in the same semantic space as the words themselves. So that phrases like "I went to the mall yesterday", "We went shopping last week", and "They went to the store", would all be pretty close in distance to each other. Well we have seen ways to train unigram word vectors, should we do the same for bigrams, trigrams, etc. This very well may work but there are two major issues with this thinking. 1) There are literally an infinite amount of possible combinations of words. Storing and training an infinite amount of vectors would just be absurd. 2) Some combinations of words while they might be completely reasonable to hear in language, may never be represented in our training/dev corpus. So we would never learn them.

We need a way to take a sentence and its respective words vectors,

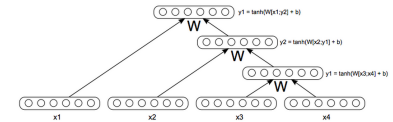


Figure 1: A standard Recursive Neural Network

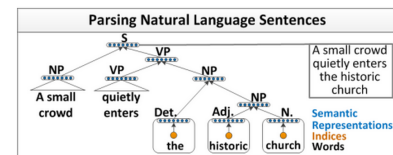


Figure 2: Parse Tree of a sentence.

and derive what the embedded vector should be. Now let's first ask a very debated question. Is it naive to believe that the vector space that we used to represent all words, is sufficiently expressive to also be able to represent all sentences of any length? While this may be unintuitive, the performance of these models suggests that this is actually a reasonable thing to do.

Let's first discuss the difference between semantic and grammatical understanding of a sentence. Semantic analysis is an understanding of the meaning of a sentence, being able to represent the phrase as a vector in a structured semantic space, where similar sentences are very nearby, and unrelated sentences are very far away. The grammatical understanding is one where we have identified the underlying grammatical structure of the sentence, which part of the sentence depends on which other part, what words are modifying what other words, etc. The output of such an understanding is usually represented as a parse tree as displayed in Figure 2.

Now for the million dollar question. If we want to know the semantic representation, is it an advantage, nay, required, to have a grammatical understanding? Well some might disagree but for now we will treat this semantic composition task the following way. First, we need to understand words. Then, we need to know the way words are put together, then, finally, we can get to a meaning of a phrase or sentence by leveraging these two previous concepts.

So let's begin with our first model built on this principle. Let's imagine we were given a sentence, and we knew the parse tree for that sentence, such as the one displayed in Figure 2, could we figure out an encoding for the sentence and also perhaps a sentiment score just from the word vectors that are in the sentence? We observe how a Simple RNN can perform this task.

1.1 A simple single layer RNN

Let's walk through the model displayed in Figure 3 above. We first take a sentence parse tree and the sentence word vectors and begin to walk up the tree. The lowest node in the graph is Node 3, so we concatenate L_{29} and L_{430} to form a vector $\in \mathbb{R}^{2d}$ and feed it into our network to compute:

$$h^{(1)} = \tanh(W^{(1)} \begin{bmatrix} L_{29} \\ L_{430} \end{bmatrix} + b^{(1)}) \quad (1)$$

Since $W^{(1)} \in \mathbb{R}^{d \times 2d}$ and $b^{(1)} \in \mathbb{R}^d$, $h^{(1)} \in \mathbb{R}^d$. We can now think of $h^{(1)}$ as a point in the same word vector space for the bigram "this assignment", in which we did not need to learn this representation separately, but rather derived it from its constituting word vectors.

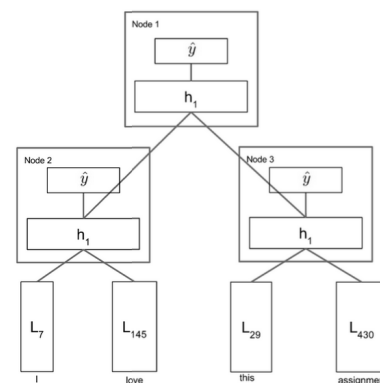


Figure 3: An example standard RNN applied to a parsed sentence "I love this assignment"

We now take $h^{(1)}$ and put it through a softmax layer to get a score over a set of sentiment classes, a discrete set of known classes that represent some meaning. In the case of positive/negative sentiment analysis, we would have 5 classes, class 0 implies strongly negative, class 1 implies negative, class 2 is neutral, class 3 is positive, and finally class 4 is strongly positive.

Now we do the same thing with the "I" and "love" to produce the vector $h^{(1)}$ for the phrase "I love". Again, we compute a score over the semantic classes again for that phrase. Finally, for the most interesting step, we need to merge the two phrases "I love" and "this assignment". Here we are concatenating word phrases, rather than word vectors! We do this in the same manner, concatenating the two $h^{(1)}$ vectors and compute

$$h^{(1)} = \tanh(W^{(1)} \begin{bmatrix} h_{Left}^{(1)} \\ h_{Right}^{(1)} \end{bmatrix} + b^{(1)}) \tag{2}$$

Now we have a vector in the word vector space that represents the full sentence "I love this assignment". Furthermore, we can put this $h^{(1)}$ through the same softmax layer as before, and compute sentiment probabilities for the full sentence. Of course the model will only do this reliably once trained.

Now lets take a step back. First, is it naive to think we can use the same matrix W to concatenate all words together and get a very expressive $h^{(1)}$ and yet again use that same matrix W to concatenate all phrase vectors to get even deeper phrases? These criticisms are valid and we can address them in the following twist on the simple RNN.

1.2 Syntactically Untied SU-RNN

As we discussed in the criticisms of the previous section, using the same W to bring together a Noun Phrase and Verb Phrase and to bring together a Prepositional Phrase and another word vector seems intuitively wrong. And maybe we are bluntly merging all of these functionalities into too weak of a model.

What we can do to remedy this shortcoming is to "syntactically untie" the weights of these different tasks. By this we mean, there is no reason to expect the optimal W for one category of inputs to be at all related to the optimal W for another category of inputs. So we let these W 's be different and relax this constraint. While this for sure increases our weight matrices to learn, the performance boost we gain is non-trivial.

As in Figure 4 above, we notice our model is now conditioned upon what the syntactic categories of the inputs are. Note, we deter-

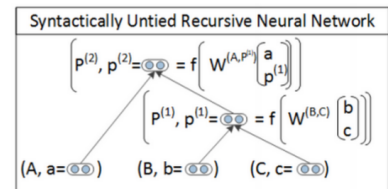


Figure 4: Using different W 's for different categories of inputs is more natural than having just one W for all categories

mine what the categories are via a very simple Probabilistic Context Free Grammar (PCFG) which is more or less learned by computing summary statistics over the Penn Tree Bank to learn rules such as "The" is always a DT, etc, etc. No deeper understanding of this part is really necessary, just know its really simple.

The only major other difference in this model is that we initialize the W 's to the identity. This way the default thing to do is to average the two word vectors coming in. Slowly but surely, the model learns which vector is more important and also any rotation or scaling of the vectors that improve performance. We observe in Figure 5 that the trained weight matrices learn actual meaning! For example, the DT-NP rule or Determiner followed by a Noun Phrase such as "The cat" or "A man", puts more emphasis on the Noun Phrase than on the Determiner. (This is obvious because the right diagonals are red meaning higher weights). This is called the notion of soft head words, which is something that Linguists have long observed to be true for sometime, however the model learned this on its own just by looking at data. Pretty cool!

The SU-RNN does indeed outperform previously discussed models but perhaps it is still not expressive enough. If we think of modifying words, such as adverbs like "very", any interpolation with this word vector and the following one, is definitely not what the understood nature of "very" is.

As an adverb, it's literal definition is "used for emphasis". How can we have a vector that emphasizes any other vector that is to follow when we are solely performing a linear interpolation? How can we construct a vector that will "scale" any other vector this way? Truth is we can not. We need to have some form of multiplication of word on another word. We uncover two such compositions below that enable this. The first utilizes word matrices and the other utilizes a Quadratic equation over the typical Affine.

1.3 MV-RNN's (Matrix-Vector Recursive Neural Networks)

We now augment our word representation, to not only include a word vector, but also a word matrix! So the word "very" will have a word vector $v_{very} \in \mathbb{R}^d$ but also $V_{very} \in \mathbb{R}^{d \times d}$. This gives us the expressive ability to not only embed what a word means, but we also learn the way that words "modify" other words. The word matrix enables the latter. To feed two words, a and b , into a RNN, we take their word matrices A and B , to form our input vector x as the concatenation of vector Ab and Ba . For our example of "very", V_{very} could just be the identity times any scalar above one. Which would scale any neighboring word vector by that number! This is the type of expres-

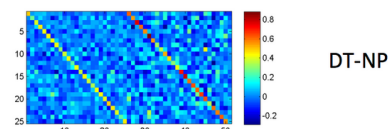


Figure 5: The learnt W weights for DT-NP composition match Linguists theory

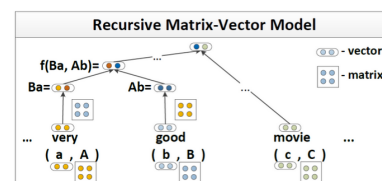


Figure 6: An example MV-RNN

sive ability we desired. While the new word representation explodes our feature space, we can express much better the way words modify each other.

By observing the errors the model makes, we see even the MV-RNN still can not express certain relations. We observe three major classes of mistakes.

First, Negated Positives. When we say something positive but one word turns it negative, the model can not weigh that one word strong enough to flip the sentiment of the entire sentence. Figure 7 shows such an example where the swap of the word "most" to "least" should flip the entire sentiment of the sentence, but the MV-RNN does not capture this successfully.

The second class of mistakes is the Negated Negative case. Where we say something is not bad, or not dull, as in Figure 8. The MV-RNN can not recognize that the word "not" lessens the sentiment from negative to neutral.

The final class of errors we observe is the "X but Y conjunction" displayed in Figure 9. Here the X might be negative BUT if the Y is positive then the model's sentiment output for the sentence should be positive! MV-RNNs struggle with this.

Thus, we must look for an even more expressive composition algorithm that will be able to fully capture these types of high level compositions.

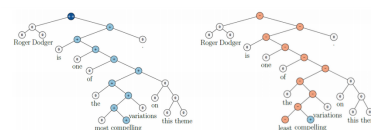


Figure 7: Negated Positives



Figure 8: Negated Negatives

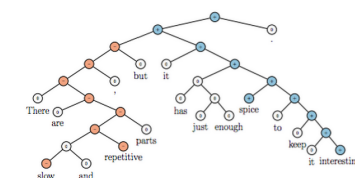


Figure 9: Using a Recursive Neural Net can correctly classify the sentiment of the contrastive conjunction X but Y but the MV-RNN can not

1.4 RNTNs (Recursive Neural Tensor Network)

The final RNN we will cover here is by far the most successful on the three types of errors we left off with. The Recursive Neural Tensor Network does away with the notion of a word matrix, and furthermore, does away with the traditional affine transformation pre-tanh/sigmoid concept. To compose two word vectors or phrase vectors, we again concatenate them to form a vector $\in \mathbb{R}^{2d}$ but instead of putting it through an affine function then a nonlinear, we put it through a quadratic first, then a nonlinear, such as:

$$h^{(1)} = \tanh(x^T V x + W x) \tag{3}$$

Note that V is a 3rd order tensor in $\in \mathbb{R}^{2d \times 2d \times d}$. We compute $x^T V [i] x \forall i \in [1, 2, \dots, d]$ slices of the tensor outputting a vector $\in \mathbb{R}^d$. We then add $W x$ and put it through a nonlinear function. The quadratic shows that we can indeed allow for the multiplicative type of interaction between the word vectors without needing to maintain and learn word matrices!

As we see in Figure 11, the RNTN is the only model that is capable of succeeding on these very hard datasets.

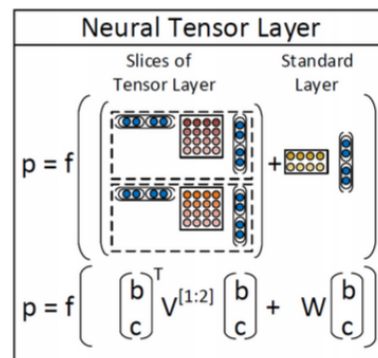


Figure 10: One slice of a RNTN. Note there would be d of these slices.

Model	Accuracy	
	Negated Positive	Negated Negative
biNB	19.0	27.3
RNN	33.3	45.5
MV-RNN	52.4	54.6
RNTN	71.4	81.8

Figure 11: Comparing performance on the Negated Positive and Negated Negative data sets.

We will continue next time with a model that actually outperforms the RNTN in some aspects and it does not require an input parse tree! This model is the Dynamic Convolutional Neural Network, and we will talk about that soon.

2 CNNs (Convolutional Neural Networks)

Back to sentence representation, what if you did NOT know the parse tree? The RecNNs we have observed in this set of lecture notes depend on such an initial parsing. What if we are using Recurrent Networks and the context of the phrase is on its right side? If we only applied Softmax on the last time-step, then the last few words would have a disproportionately large impact on the output, e.g. Sentiment Classification.

2.1 Why CNNs?

Convolutional Neural Networks take in a sentence of word vectors and first create a phrase vector for all subphrases, not just grammatically correct phrases (as with Recursive Neural Network)! And then, CNNs group them together for the task at hand.

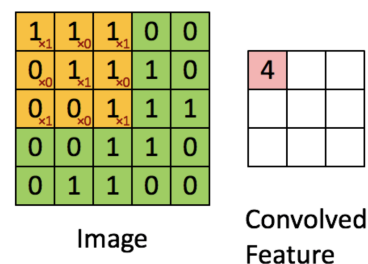
2.2 What is Convolution?

Let's start with the 1D case. Consider two 1D vectors, f and g with f being our primary vector and g corresponding to the **filter**. The convolution between f and g , evaluated at entry n is represented as $(f * g)[n]$ and is equal to $\sum_{m=-M}^M f[n-m]g[m]$.

Figure 12 shows the 2D convolution case. The 9×9 green matrix represents the primary matrix of concern, f . The 3×3 matrix of red numbers represents the filter g and the convolution currently being evaluated is at position $[2, 2]$. Figure 12 shows the value of the convolution at position $[2, 2]$ as 4 in the second table. Can you complete the second table?

2.3 A Single-Layer CNN

Consider word-vectors $x_i \in R^k$ and the concatenated word-vectors of a n -word sentence, $x_{1:n} = x_1 \oplus x_2 \dots \oplus x_n$. Finally, consider a Convolutional filter $w \in R^{hk}$ i.e. over h words. For $k = 2$, $n = 5$ and $h = 3$, Figure 13 shows the Single-Layer Convolutional layer for NLP. We will get a single value for each possible combination of three consecutive words in the sentence, "the country of my birth". Note, the filter w is itself a vector and we will have $c_i = f(w^T x_{i:i+h-1} + b)$



Stanford UFLDL wiki

Figure 12: Convolution in the 2D case

to give $\mathbf{c} = [c_1, c_2 \dots c_{n-h+1}] \in \mathbb{R}^{n-h+1}$. For the last two time-steps, i.e. starting with the words "my" or "birth", we don't have enough word-vectors to multiply with the filter (since $h = 3$). If we necessarily need the convolutions associated with the last two word-vectors, a common trick is to pad the sentence with $h - 1$ zero-vectors at its right-hand-side as in Figure 14.

2.4 Pooling

Assuming that we don't use zero-padding, we will get a final convolutional output, \mathbf{c} which has $n - h + 1$ numbers. Typically, we want to take the outputs of the CNN and feed it as input to further layers like a Feedforward Neural Network or a RecNN. But, all of those need a fixed length input while our CNN output has a length dependent on the length of the sentence, n . One clever way to fix this problem is to use max-pooling. The output of the CNN, $\mathbf{c} \in \mathbb{R}^{n-h+1}$ is the input to the max-pooling layer. The output of the max-pooling layer is $\hat{c} = \max\{\mathbf{c}\}$, thus $\hat{c} \in \mathbb{R}$.

We could also have used min-pooling because typically we use ReLU as our non-linear activation function and ReLU is bounded on the low side to 0. Hence a min-pool layer might get smothered by ReLU, so we nearly always use max-pooling over min-pooling.

2.5 Multiple-Filters

In the example above related to Figure 13, we had $h = 2$, meaning we looked only at bi-gram with a single specific combination method i.e. filter. We can use multiple bi-gram filters because each filter will learn to recognize a different kind of bi-gram. Even more generally, we are not restricted to using just bi-grams, we can also have filters using tri-grams, quad-grams and even higher lengths. Each filter has an associated max-pool layer. Thus, our final output from the CNN layers will be a vector having length equal to the number of filters.

2.6 Multiple-Channels

If we allow gradients to flow into the word-vectors being used here, then the word-vectors might change significantly over training. This is desirable, as it specializes the word-vectors to the specific task at hand (away from say GloVe initialization). But, what about words that appear only in the test set but not in the train set? While other semantically related word vectors which appear in the train set will have moved significantly from their starting point, such words will still be at their initialization point. The neural network will be specialized for inputs which have been updated. Hence, we will get low

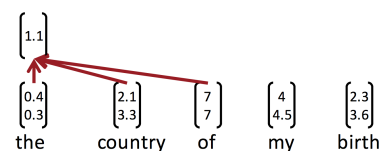


Figure 13: Single-Layer Convolution: one-step

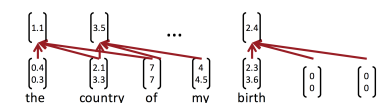


Figure 14: Single-Layer Convolution: all-steps

performance on sentences with such words (words that are in test but not in train).

One work-around is to maintain two sets of word-vectors, one 'static' (no gradient flow into them) and one 'dynamic', which are updated via SGD. Both are initially the same (GloVe or other initialization). Both sets are simultaneously used as input to the neural network. Thus, the initialized word-vectors will always play a role in the training of the neural network. Giving unseen words present in the test set a better chance of being interpreted correctly.

There are several ways of handling these two channels, most common is to simply average them before using in a CNN. The other method is to double the length of the CNN filters.

2.7 CNN Options

1. Narrow vs Wide

Refer to Figure 15. Another way to ask this is should we not (narrow) or should we (wide) zero-pad? If we use Narrow Convolution, we compute a convolution only in those positions where all the components of a filter have a matching input component. This will clearly not be the case at the start and end boundaries of the input, as in the left side network in Figure 15. If we use Wide Convolution, we have an output component corresponding to each alignment of the convolution filter. For this, we will have to pad the input at the start and the end with $h - 1$ zeros.

In the Narrow Convolution case, the output length will be $n - h + 1$ and in the Wide Convolution case, the length will be $n + h - 1$.

2. k -max pooling

This is a generalization of the max pooling layer. Instead of picking out only the biggest (max) value from its input, the k -max pooling layer picks out the k biggest values. Setting $k = 1$ gives the max pooling layer we saw earlier.

3 Constituency Parsing

Natural Language Understanding requires being able to extract meaning from large text units from the understanding of smaller parts. This extraction requires being able to understand how smaller parts are put together. There are two main techniques to analyze the syntactic structure of sentences: **constituency parsing** and **dependency parsing**. Dependency parsing was covered in the previous lectures (see *lecture notes 4*). By building binary asymmetric relations between a word and its dependents, the structure shows which word

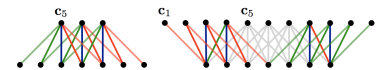


Figure 15: Narrow and Wide Convolution (from Kalchbrenner et al. (2014))

depends on which other words. Now we focus on constituency parsing, that organizes words into nested constituents.

Constituency Parsing is a way to break a piece of text (e.g. one sentence) into sub-phrases. One of the goals of constituency parsing (also known as "phrase structure parsing") is to identify the constituents in the text which would be useful when extracting information from text. By knowing the constituents after parsing the sentence, it is possible to generate similar sentences syntactically correct.

3.1 Constituent

In syntactic analysis, a constituent can be a single word or a phrases as a single unit within a hierarchical structure. A phrase is a sequence of two or more words built around a head lexical item and working as a unit within a sentence. To be a phrase, a group of words should come together to play a specific role in the sentence. In addition, the group of words can be moved together or replaced as a whole, and the sentence should remain fluent and grammatical.

For instance, the following sentence contains the noun phrase: "wonderful CS224N".

- I want to be enrolled in the wonderful CS224N!

We can rewrite the sentence by moving that whole phrase to the front as below.

- The wonderful CS224N I want to be enrolled in!

Or the phrase could be replaced with an constituent of similar function and meaning, like "great CS course in Stanford about NLP and Deep Learning".

- I want to be enrolled in the great CS course in Stanford about NLP and Deep!

For constituency parsing, the basic clause structure is understood as a binary division of the clause into subject (noun phrase NP) and predicate (verb phrase VP), expressed as following rule. The binary division of the clause results in a one-to-one-or-more correspondence. For each element in a sentence, there are one or more nodes in the tree structure.

- S -> NP VP

In fact, the process of parsing illustrates certain similar rules. We deduce the rules by beginning with the sentence symbol S, and applying the phrase structure rules successively, finally applying replacement rules to substitute actual words for the abstract symbols.

We interpret large text units by **semantic composition** of smaller elements. These smaller elements can be changed while keeping a similar meaning, like in the following examples.

Based on the extracted rules, it is possible to generate similar sentences. If the rules are correct, then any sentence produced in this way should be syntactically correct. However, the generated sentences might be syntactically correct but semantically nonsensical, such as the following well-known example:

- Colorless green ideas sleep furiously

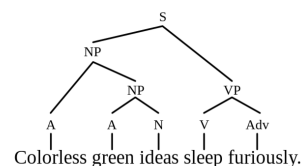


Figure 16: Constituency Parse Tree for 'Colorless green ideas sleep furiously'

3.2 Constituency Parse Tree

Interestingly, in natural language, the constituents are likely to be nested inside one another. Thus a natural representation of these phrases is a tree. Usually we use a consistency parse tree to display the parsing process. The constituency-based parse trees of constituency grammars distinguish between terminal and non-terminal nodes. Non-terminals in the tree are labeled as types of phrases (e.g. Noun Phrase), the terminals are the exact words in the sentence.

Take 'John hits the ball' as an example, the syntactic structure of the English sentence is shown as below.

We have a parse tree starting from root S, which represents the whole sentence, and ending in each leaf node, which represents each word in a sentence. We use the following abbreviation:

- S stands for sentence, the top-level structure.
- NP stands for noun phrase including the subject of the sentence and the object of the sentence.
- VP stands for verb phrase, which serves as the predicate.
- V stands for verb.
- D stands for determiner, such as the definite article "the"
- N stands for noun

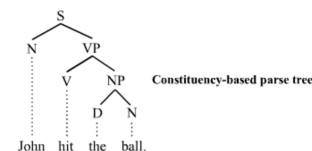


Figure 17: Consistency Parse Tree for 'John hits the ball'