

To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy.

LIGHT ON MATH MACHINE LEARNING

Intuitive Guide to Understanding Word2vec



Thushan Ganegedara  · Follow

Published in Towards Data Science

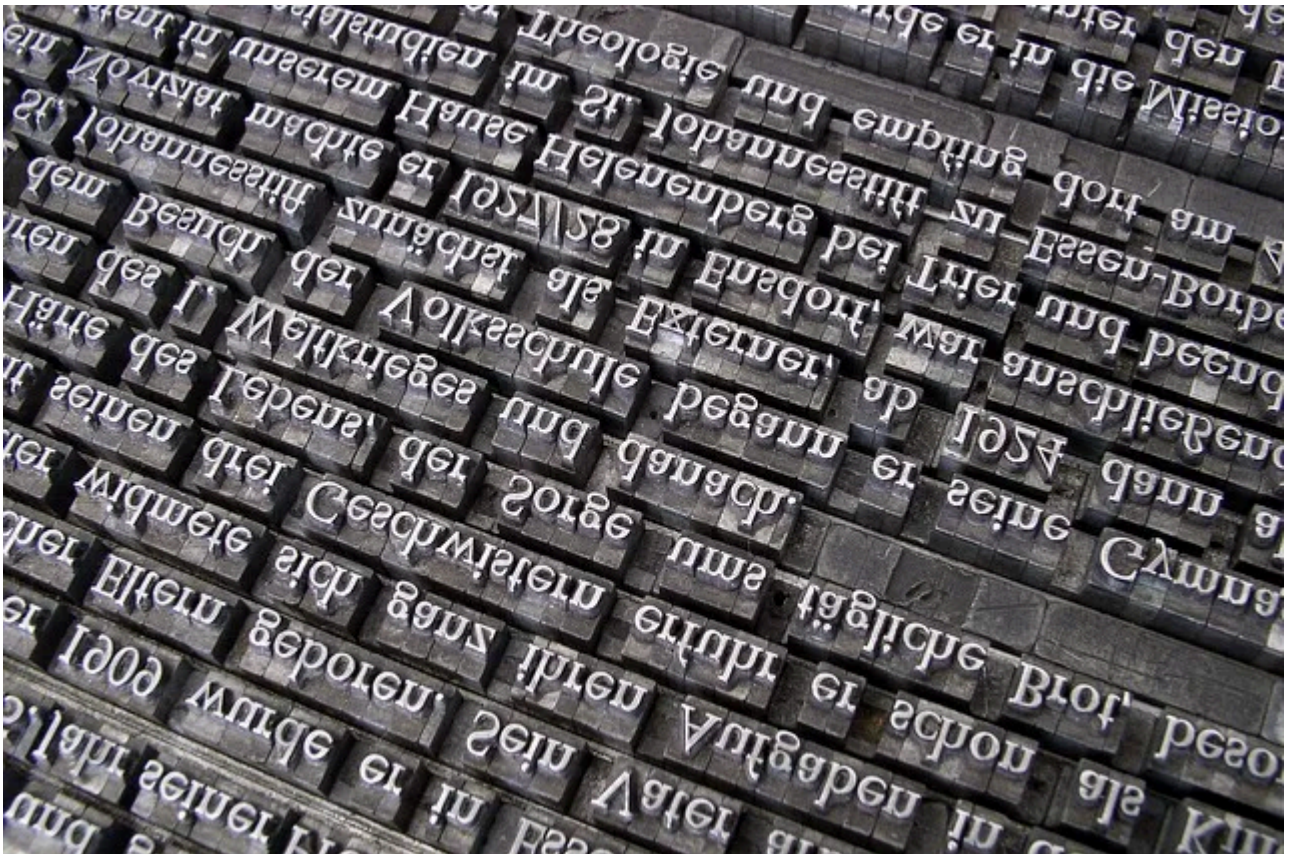
12 min read · Jun 5, 2018



Listen



Share



Here comes the third blog post in the series of *light on math machine learning* A-Z. This article is going to be about Word2vec algorithms. Word2vec algorithms output word vectors. Word vectors, underpin many of the natural language processing (NLP) systems, that have taken the world by a storm (Amazon Alexa, Google translate, etc.). We will talk about the details in the upcoming sections. But first let me spell-out the alphabet with the links to my other blog posts.

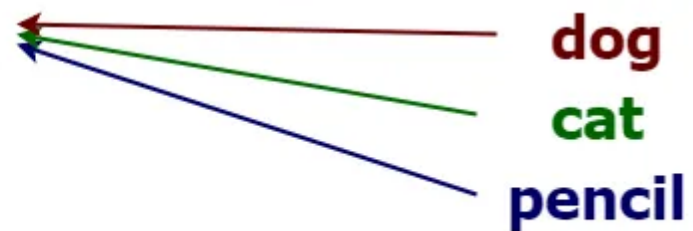
Word vectors, what's the big idea? It's all about that context

Without further due, let us stick our feet in. Word vectors are numerical representations of words, that preserve the semantic relationship between words. For example the vector of the word *cat*, will be very similar to the vector of the word *dog*. However, the vector for *pencil* will be quite different from the word vector of *cat*. And this similarity is defined by the frequency two words in question (i.e. *[cat,dog]* or *[cat,pencil]*), are used in the same context. For example consider the following sentences,

I like to pet my _____

My _____ does not like the postman

dog
cat
pencil



I don't think I need to spell out the odd one in the above sentences, and obviously the ones with *pencil*, as the missing word. Why do you feel it as the odd sentence? The spelling is fine, the grammar is right, then why? It is because the *context*, the word *pencil* used is not correct. This should convince you of the power the context of a word has, on the word itself. Word vector algorithms use the context of the words to learn numerical representations for words, so that words used in the same context have similar looking word vectors.

Impact and implications of Word2vec

To get an idea about the implication of the Word2vec techniques, try the following. Go ahead and bring up [google scholar](#). Type in some NLP related task (e.g. question answering, chatbots, machine translation, etc). Filter the papers published after 2013 (that's when Word2vec methods came out). Get the proportion of the papers using word vectors to total number of papers. I bet this proportion is going to be quite high. To make the statement concrete, word vectors are used for,

- Language modelling
- Chatbots

- Machine translation

- Question answering To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy.

- ... and many more

You can see all the exciting NLP frontiers actually depend on word vectors heavily. Now let us discuss what sort of implications word vectors have to make the model better. When using word vectors, semantically close words will appear to be similar computations within the model, where other words will appear to do different-looking computations. This is a desirable property to have, as encoding such information in the input itself, lead to the model performing well, with lesser data.

From raw text to word vectors: high level approach

Now, with a solid intuition in the bag, we're going first to discuss the high-level mechanics of the Word2vec algorithms. We'll polish the details up in later sections, until we can be sure that we know how to implement a Word2vec algorithm. In order to learn word vectors in an unsupervised manner (i.e. without a human labelling data), we have to define and complete certain tasks. Here's a high level list of these tasks.

1. Create data tuples of the format [input word, output word], where each word is represented as one-hot vectors, from raw text
2. Define a model that can take in the one-hot vectors as inputs and outputs, to be trained
3. Define a loss function that predict the correct word, which is actually in the context of the input word, to optimize the model
4. Evaluate the model by making sure similar words have similar word vectors

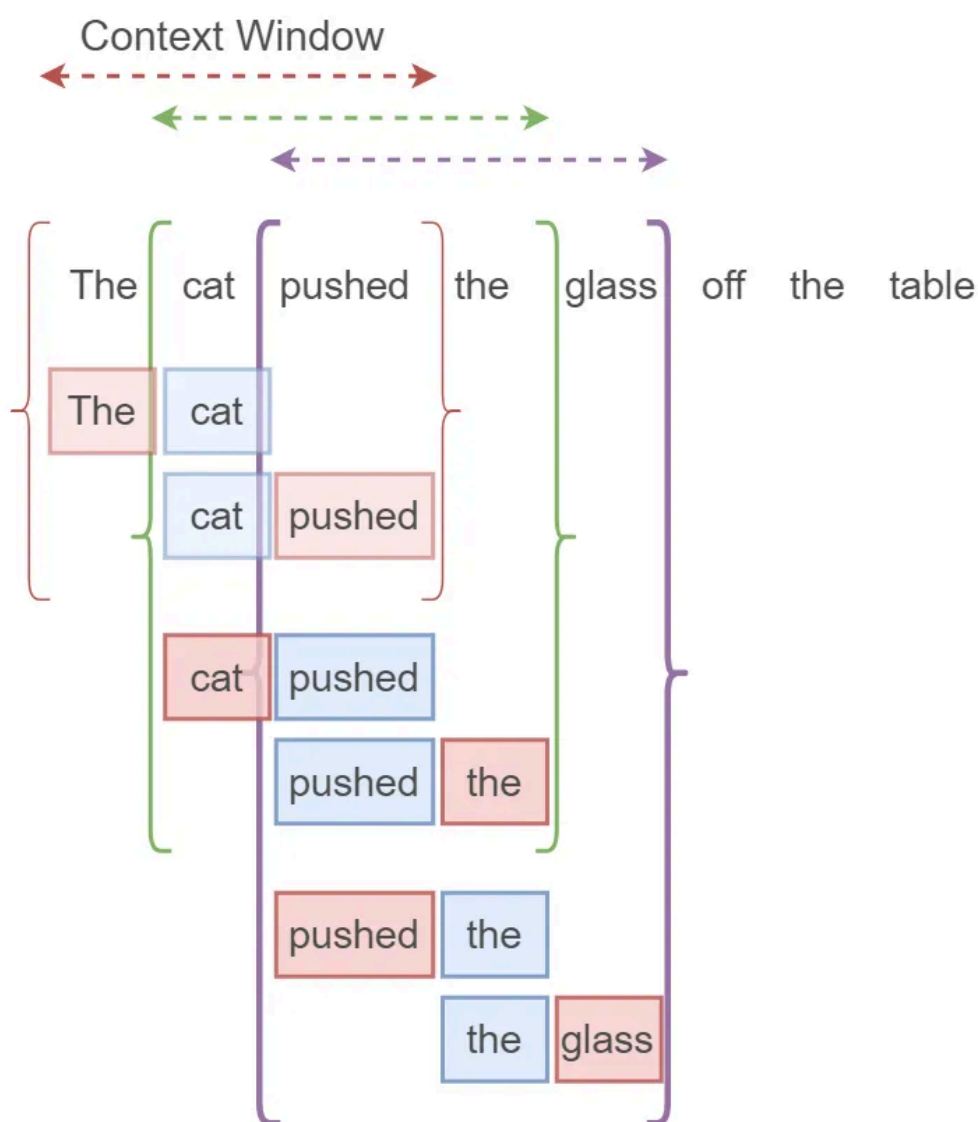
This seemingly simple procedure will lead to learning very powerful word vectors. Let us move onto the nitty-gritties of each step of the above pipeline.

Creating structured data from raw text

This is not a very difficult task. This is simple manipulation of raw text to put that to certain structure. Think of the the following sentence.

The cat pushed the glass off the table

The data created from this sentence would look like as follows. Each row after the sentence represents a context window. To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy. Not input word (*the* red). The one-hot output word (*any word in the context window except the middle word, called context words*). Two data points are created from a single context window. The size of the context window is something defined by the user. Higher the context window size, the better the performance of the model. But when having a large context window size, you'll pay in computational time, as the amount of data increases. Don't confuse the target word with the target (correct output) for the neural network, these are two completely different things.



Creating data for Word2vec

Defining the embedding layer and the neural network

The neural network used to learn from the structured data defined above. However, it comes with a twist! To make clear, you have the following components.

- A batch of inputs represented as one-hot vectors

- A batch To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy. (use)

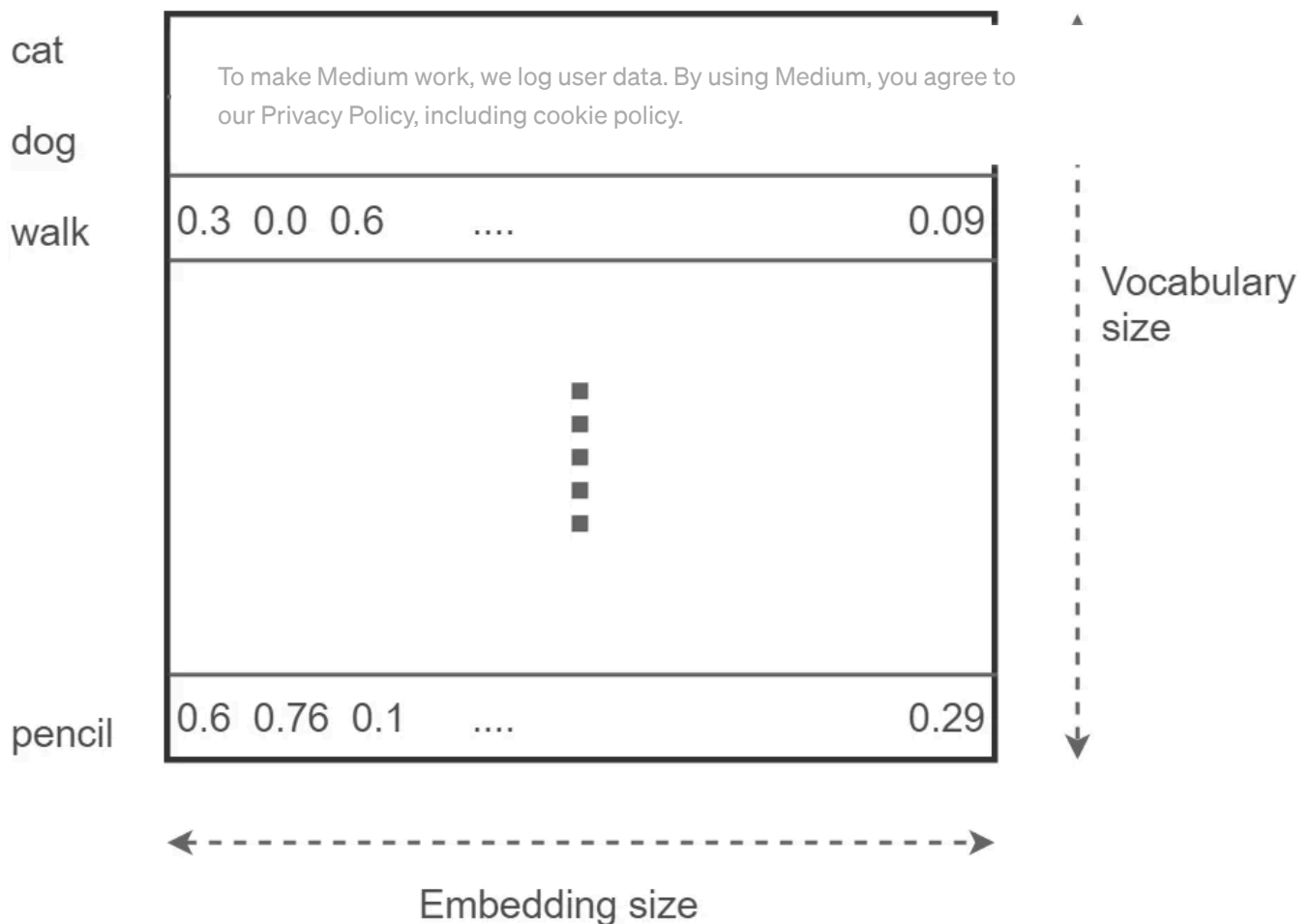
- An embedding layer

- A neural network

No need to be scared if you did not get a feeling of what and how the last two components perform. We will probe each of these components to understand what they do.

Embedding layer: stores all the word vectors

First in our agenda is the *embedding layer*. The embedding layer stores the word vectors of all the words found in the vocabulary. As you can imagine this is an enormous matrix (of size *[vocabulary size x embedding size]*). This embedding size is a user-tunable parameter. The higher it is, the better performing your model will be. But you'll not get much of a jaw-dropping performance/size gain, beyond a certain point (say, an embedding size of 500). This gigantic matrix is initialized randomly (just like a neural network) and is tweaked bit by bit, during the optimization process, to reveal the powerful word vectors. This is what it looks like.



What the embedding layer looks like

Neural network: maps word vectors to outputs

Next in line is the last LEGO block of our model; the neural network. During the training, the neural network takes an input word and attempt to predict the output word. Then using a loss function, we penalize the model for incorrect classifications and reward the model for correct classifications. We will be limiting our conversation to processing a single input and a single output at a time. However in reality, you process data in batches (say, 64 data points). Let's delineate the exact process used during training:

1. For a given input word (the target word), find the corresponding word vector from the embedding layer
2. Feed the word vector to the neural network, then try to predict the correct output word (a context word)
3. By comparing the prediction and true context word, compute the loss
4. Use the loss along with a stochastic optimizer to optimize the neural network and the embedding layer

One thing to note is that when computing the prediction, we use a softmax activation.

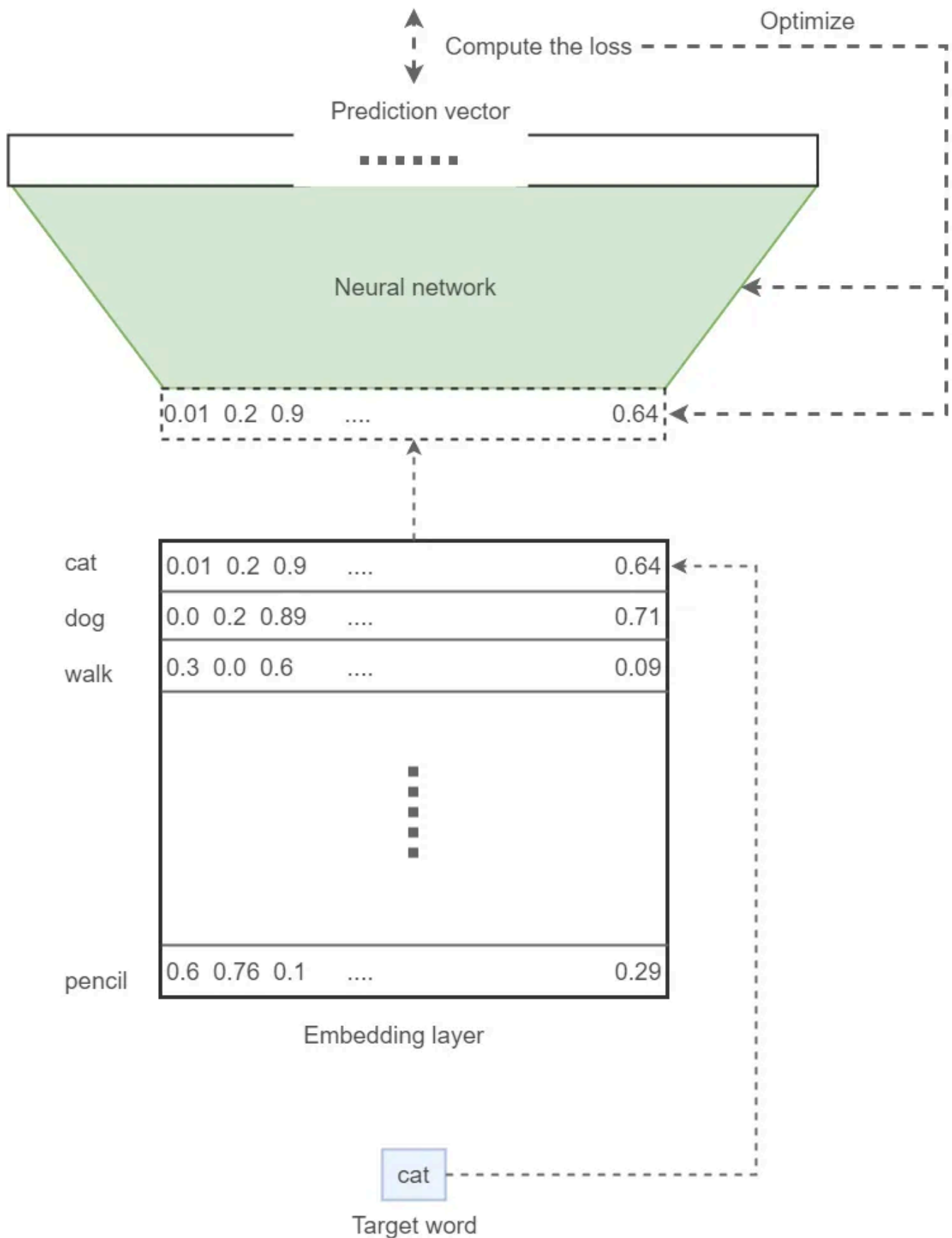
To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy.

Fitting all together: inputs to model to outputs

Knowing all the nuts and bolts of the Word2vec algorithm, we can put all the parts together. So that once this model is trained, all we have to do is *save the embedding layer to the disk*. Then we can enjoy semantic preserved word vectors at any time of the day. Below we see what the full picture looks like.

Context word

To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy.



This particular arrangement of data and the model layout is known as the *skip-gram* algorithm. To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy. us. The other algorithms

Defining a loss function: to optimize the model

One crucial bit of information we haven't discussed so far, but is essential is the loss function. Normally, standard softmax cross entropy loss is a good loss function for a classification task. Using this loss is not very practical for a Word2vec model, as for a simpler task like sentiment analysis (where you have 2 possible outputs: positive or negative). Here things can get funky. In a real word task, that consumes billions of words, the vocabulary size can grow up to 100,000 or beyond easily. This makes the computation of the softmax normalization heavy. This is because the full computation of softmax requires to calculate cross entropy loss with respect to all the output nodes.

So we are going with a smarter alternative, called *sampled softmax loss*. In sampled softmax loss, you do the following. Note that there is quite a lot of changes from the standard softmax cross entropy loss. First, you compute the cross entropy loss between the true context word ID for a given target word and the prediction value corresponding to the true context word ID. Then to that, we add the cross entropy loss of κ negative samples we sampled according to some noise distribution. On a high level, we define the loss as follows:

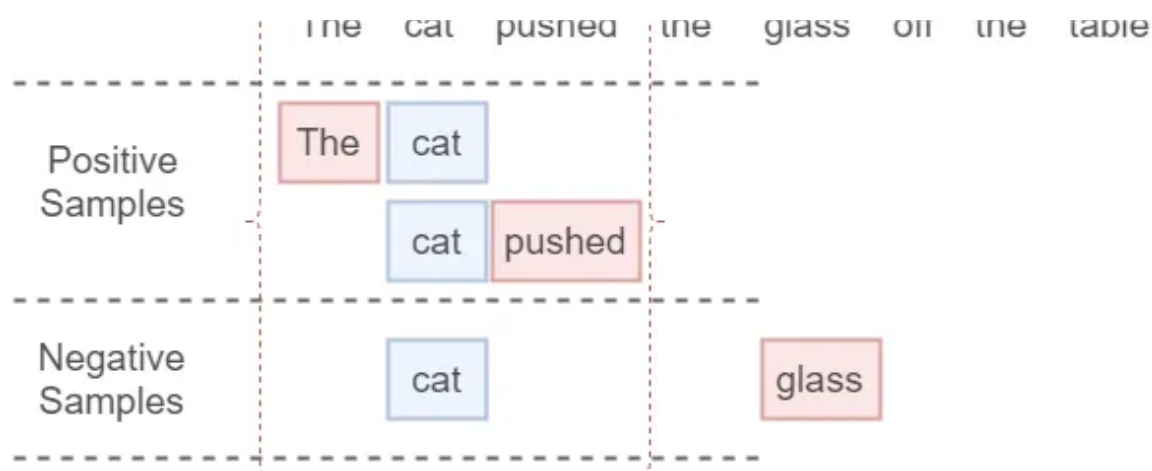
$$\text{Loss} = \text{SigmoidCrossEntropy}(\text{Prediction}, \text{Correct Word}) + \sum_1^K E_{\text{Noise ID}} \text{SigmoidCrossEntropy}(\text{Prediction}, \text{Noise ID})$$

The `SigmoidCrossEntropy` is a loss we can define on a single output node, independent of the rest of the nodes. This makes it ideal for our problem, as our vocabulary can grow quite large. I'm not going to dive into the very details of this loss. You don't need to understand how exactly this is implemented, as these are available as built-in function in TensorFlow. But understanding parameters involved in the loss (e.g. κ) is important. The takeaway message is that, the sampled softmax loss computes the loss by considering two types of entities:

- The index given by the true context word ID in the prediction vector (words within the context window)
- κ indices that indicate word IDs, and are considered to be noise (words outside the context window)

I further visualize this by illustrating an example

To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy, including cookie policy.



Getting positive and negative samples for the sampled softmax layer

TensorFlow implementation: Skip-gram algorithm

Here we will regurgitate what we just discussed into an implementation. This is available as an exercise [here](#). In this section, we're going to implement the following.

- A data generator
- The skip-gram model (with TensorFlow)
- Running the skip-gram algorithm

Data generator

First let us understand how to generate data. We are not going to go into details of this code, as we have already discussed the internal mechanics of the data generation. This is just converting that logic to an implementation.

```
def generate_batch(batch_size, window_size):
    global data_index

    # two numpy arrays to hold target words (batch)
    # and context words (labels)
    batch = np.ndarray(shape=(batch_size), dtype=np.int32)
    labels = np.ndarray(shape=(batch_size, 1), dtype=np.int32)

    # span defines the total window size
    span = 2 * window_size + 1

    # The buffer holds the data contained within the span
    queue = collections.deque(maxlen=span)
```