

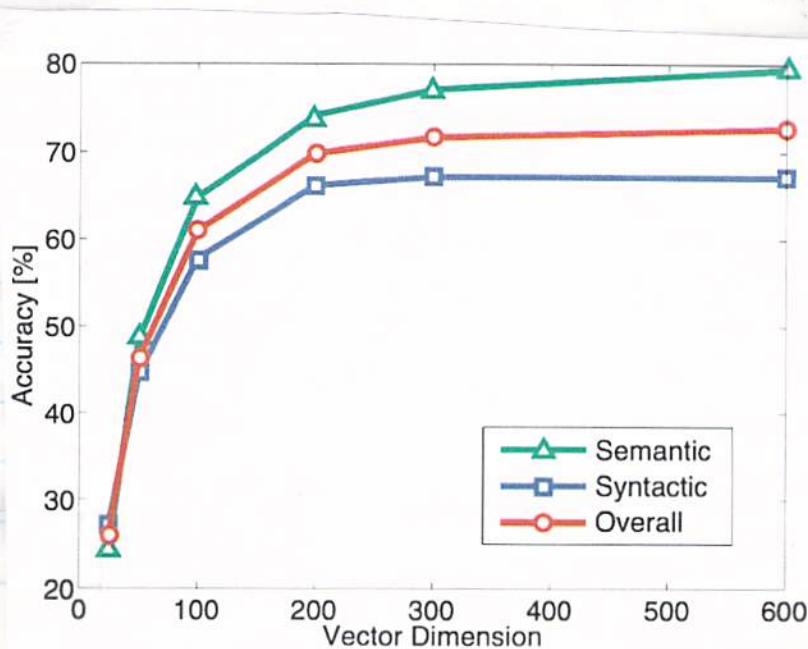
ECE 1786 Lecture #2

Last Day: Word Embedding Properties & Meaning Extraction.

Work-in-Flight: - Assignment 1 due next week.
- Fill out survey if ~~needed~~ have not - new registrants discuss adding.

Today: How Word Embeddings are Created.

- Recall:
- Word embeddings represent meaning of words in numbers.
 - helps neural networks deal with ambiguity in language
 - ~~some~~ embeddings that are close; mean that the associated words are close in meaning
 - also other associations Queen - King = Woman - Man
- You should have done Assign 1 parts 1 & 2 now.
- last week we discussed (a little) how big ^{size} dimension should be; diminishing returns beyond 300 according to GloVe paper. by Pennington et. al.



How to create (train) vectors?

- a clever + complex idea that begins with Bergio → Mikolov.

Big picture of method: We train a neural network to make a prediction based on the meaning of a word. Inside the neural network that meaning ^{of each word} will be encoded.

= the word vector / embedding

The training of the network will cause the encoding to be learned.

That encoding ~~is~~ is the embedding / vector.

Capologies for now using 3 terms that are the same thing: encoding / embedding / vector.
→ relates to auto-encoders

LINK TO AUTO ENCODERS

So, what is that prediction, ^{task} and where does the data come from? Where do the labels come from?

* The Prediction Task: ^{Given} ~~Does~~ word A "belong" ^{with} ~~word B?~~ ^{it?}
i.e. ~~are they related somehow?~~ ^{what works?}

Let's contemplate this by ^{→ give them often used together?} ^{→ predict a word that is often used together with word A}

Considering these three sentences:

1. The mathematician ran to the store.
2. The engineer ran to the store.
3. The mathematician solved the problem. ^{which is what?}

Sentences 1 & 2 imply similarity ^{between} of engineer & mathematician.

⇒ Not surprising that this is a sensible sentence

4. "The engineer solved the problem"

- this is an example of the Distributional Hypothesis ^{of other words} \circ
 "Words appearing in similar contexts are related" ^{are closer in meaning}

- since we want to predict which words are ~~related~~ ^{appear near}, we have a ready made dataset of examples and labels: every sentence ever written!
 - an enormous dataset! yippeee!!

- let's consider some simple sentences taken from the SimpleCorpus.txt in Assignment 1. ¹³ in this course

- looks like this: I hold a dog. She holds a dog. He holds a dog. I rub the dog. She rubs the cat...

- these ^{even} simple sentences give some clues to the meaning of the words used, and which are related \rightarrow what are they?

- we can generate training examples of words that are ^{appear together} ~~related~~ by selecting words that are 'near' each other in ^{NN} correct (valid) written sentences.
 - near \triangleq within a few ⁽¹⁻³⁾ words (on either side) of target word

- for example, consider "I hold a dog"; if we take 'near' to be words \pm ~~2~~ ^{from each side} 2 words away, then the training examples of ^{pairs of} related words are

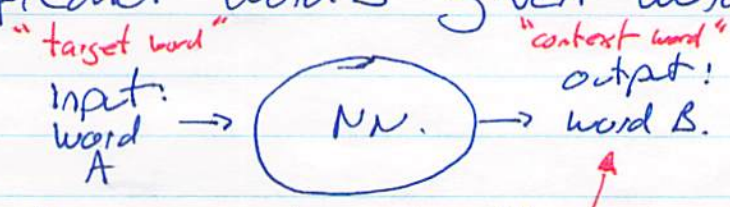
	(I, hold)	(hold, I)	(a, I)	(dog, hold)
	(I, a)	(hold, a)	(a, hold)	(dog, a)
Target		(hold, dog)	(a, dog)	
is word	I	hold	a	dog

- clearly can make a lot of these.
 + context word.

1-hot 4 word vocabulary: hello goodbye begin end. 2-4 ~~2-4~~

hello = $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ goodbye = $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ begin = $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ end = $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

- want to build a NN predictor that given a word pair (wordA, wordB) can predict wordB given wordA:



- ~~for now~~, assume that words can only come from a specific, limited vocabulary → define output

- a key part of this is how the input is represented; what is the simplest way to represent all the words in a vocabulary of size V ? - 1 bit
- discuss 1-hot - word B is 1-hot encoded. eg $|V|$
- we know that we want to represent the words by a vector of size \ll vocabulary size (that's what 1st ~~lect~~ lecture was about) $|V|$

- let's consider a simple example similar to Assignment 1, part 3

- let's say the vocabulary size $|V| = 10$ (10 word: typical vocab = 2K-6K)
- let's assume the embedding size has $\dim = 4$ (typically 5 or 1)

- so that implies we've ^{want} got 10 embeddings of size 4 each;
- these are typically stored in a matrix like so

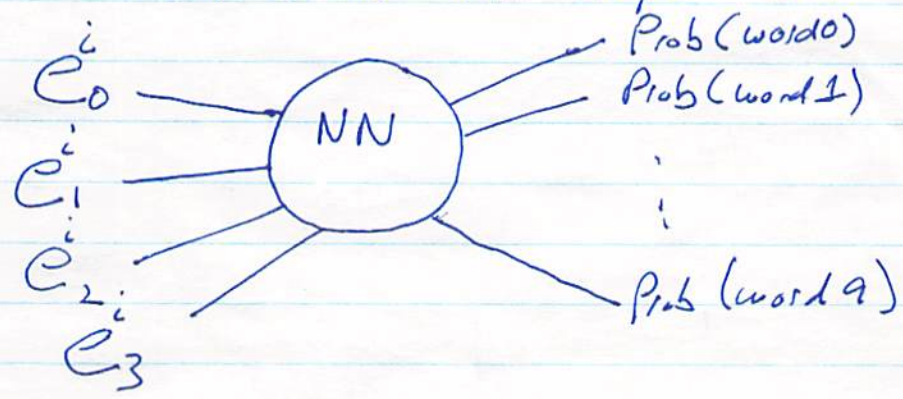
$$\begin{bmatrix} e_0^0 & e_0^1 & e_0^2 & \dots & e_0^9 \\ e_1^0 & e_1^1 & e_1^2 & \dots & e_1^9 \\ e_2^0 & e_2^1 & e_2^2 & \dots & e_2^9 \\ e_3^0 & e_3^1 & e_3^2 & \dots & e_3^9 \end{bmatrix} \rightarrow \dim \times |V|$$

"embedding matrix"

"she" "he" "hold" "dog"

- these will represent word A.
- Key: these values will be randomly initialized and learned through training like weights & biases (can pose this so that these are exactly weights, but didn't)

- so the input to the NN will be these 4 values ($e_0 \rightarrow e_3$) associated with word A in the embedding matrix
- but what should the output be?

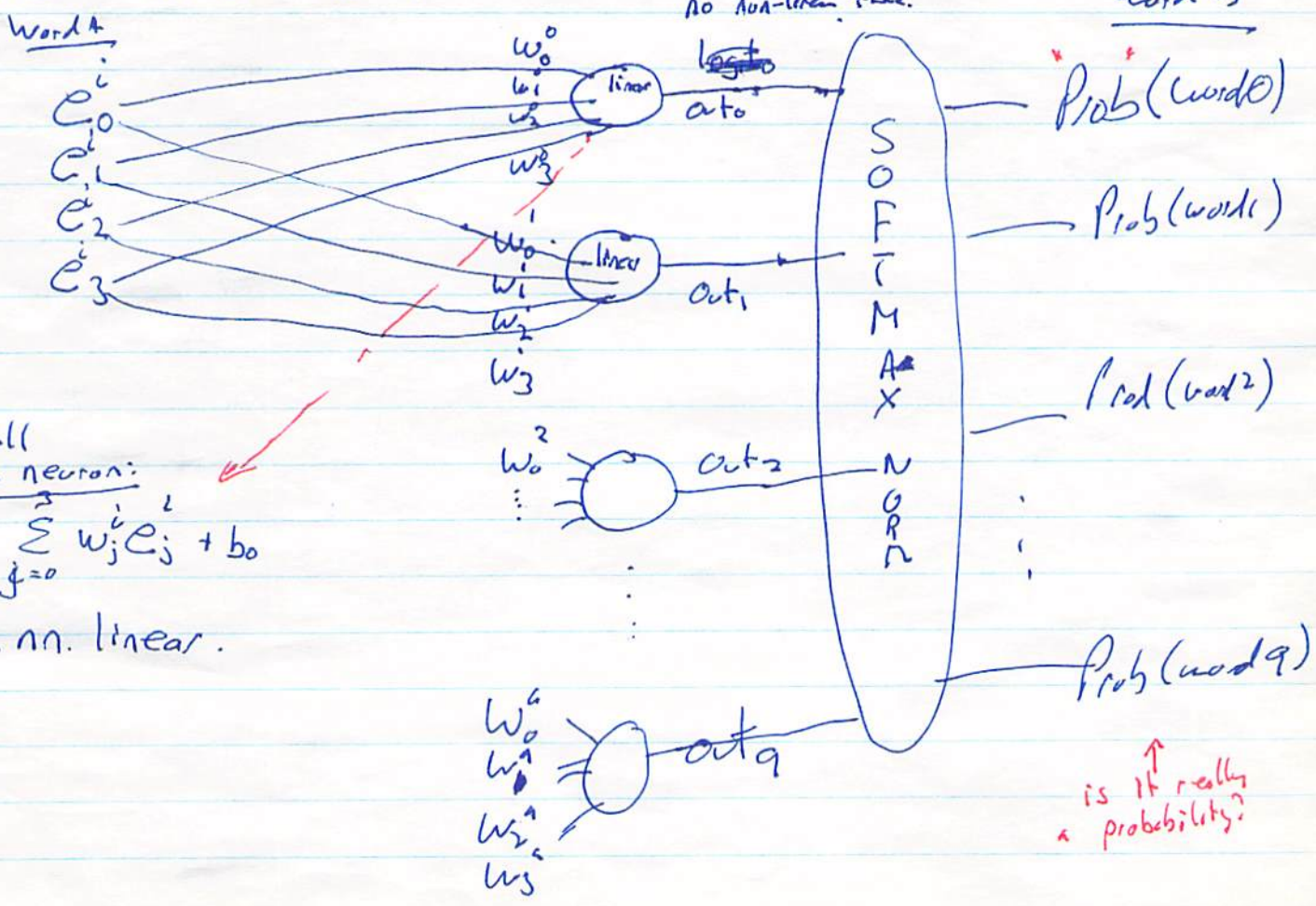


Predict which specific word is word B

- can't force it to be
- 1-hot
- so use probabilities

- where i is the word number of word A (aka 'target' word)

here is the ¹⁰ ~~lets draw what~~ the 10-neuron neural network ~~looks like:~~ that will do the prediction task



recall linear neuron:

$$out_0 = \sum_{j=0}^3 w_j^i e_j^i + b_0$$

→ torch.nn.Linear.

is it really a probability?

- So to train this network, you would present many examples of (wordA, wordB) pairs.

→ wordA is the e_j^i target

→ wordB, the label, is 1-hot encoded in a vector of size $|V| = 10$.

→ Use "cross entropy" loss to train. (the $-\log$ of correct answer)

ask → discuss.

note PyTorch combines softmax + cross-entropy into 1 function

- Crucial, to repeat: the values of e_j^i are also learned through gradient descent.

for numerical stability

- Nuanced improvement: observe the w_j^i - these are different word embeddings for each word.

⇒ we don't really need 2 embeddings, so can save computation \odot ($\frac{1}{2}$ # parameters) by just using 1

⇒ i.e. ^{con} replace w_j^i with e_j^i (the same parameter appears in 2 places in the network)

~~→ do this in assignment 1, part 2~~

→ a different way to think about this:

→ to perform the prediction task, "do these two words relate?", just compute the dot product of their word vectors

→ the higher it is, the more similar they are

→ just like a convolution kernel.

Assignment 1 Section 3 asks you to train this prediction task / embedding creation

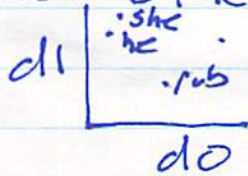
① Using a vocabulary of size 11:

I, she, he, rub, hold, dog, cat, and, the, a

② On a small corpus (smallsimplecorpus.txt)

③ with an embedding size 2

⇒ want you to see if similar words end up 'near' each other on a 2-D plot



- An important detail: words in the corpus are first reduced to their "root" in a process called lemmatization:

e.g. holds → hold
rubs → rub
rub → rub.

- so we only make embedding for the root

- Section 3 will help you get back to familiarity with software for training a neural net; you'll be training the neural network "from scratch" (ie not pre-trained)

→ (called Skip Gram)

Next: This method of training embeddings is slow in general because there are many outputs to compute → equal to the size of the vocabulary; faster method called Skip Gram with Negative Sampling

~~At~~ - also is Assignment 1 Section 4:

2-8

Skip-Gram with Negative Sampling

- rather than predict which of the |V| words in the vocabulary are associated ~~related~~ to the target word, (a multi-class classification) ↗ which is log > 5, 10 → 3000!
- instead, we make a binary prediction as to whether a given pair of words (target word, context word) ~~are related~~ or not belong together
- to do so, we need positive examples of related words - exactly the ^{same} as ^{the} SKIP-GRAM, ^{method} above, but also need negative examples of word pairs that are not related. (to properly train a binary classifier)
- so, to train, create the ^{same} positive examples as above. (with a given window size, etc)
- to create the negative samples, for each target word, we randomly sample words from the entire corpus ^{→ omitting the word itself.} - why is this OK? [i.e. won't some positive examples be included?]
- this sampling gives the method its name: skip-gram with negative sampling.
- described in Jurafsky, ^{text} Section 6.8.
- How many negative samples?
 - often 2x as many negative than positive.

So, the binary prediction for this method is the answer, yes or no, to the question "do these two input belong together or not?"

Given two input words WordA and WordB, expressed as word embeddings of size d , predict 1 if yes, they do belong together, and 0 if the answer is no, they do not.



WordA's embedding is $(a_0, a_1, a_2, \dots, a_d)$, Word B is $(b_0, b_1, b_2, \dots, b_d)$

In a similar way to the skip-gram method, which essentially compares a given input wordA to every possible wordB in the vocabulary, through a dot product, we just use one dot product here to compute the similarity between wordA and WordB. It is interesting, again, to note that this operation is the same as the basic linear neuron operation. So the NN becomes:

$$\text{WordA} \cdot \text{WordB} = \sum_{i=0}^{d-1} (a_i \times b_i) = \text{Out}$$

We convert Out to a "probability" (P) using the sigmoid function (really just force it to be a number between 0 and 1):

$$\text{i.e. } P = \sigma(\text{Out})$$

This probability is turned into the neural network training loss function using binary cross entropy, i.e.

if label = 1 (yes, related)

$$\text{loss} = -\log(P)$$

if label = 0 (no, not related)

$$\text{loss} = -\log(1-P)$$

} recall from previous course
- why is it a logarithm?

→ really pushes hard against wrong answer

The full skip gram with negative sampling method has some nuances discussed in the assignment.

When randomly sampling words for the negative examples and the positive examples, avoid the high-frequency words (such as "the", "and" ...) as they appear near many words, and so don't add much information. Reducing their appearance in the labeled dataset makes training more efficient.

Other notes about Assignment 1, Section 4, on Skip Gram with Negative Sampling:

Because this method is more efficient, we can train larger sized vectors than Section 3, use a much larger vocabulary, and train from a much larger corpus (the provided `LargeCorpus.txt`).

In Section 3 we just used 'lemmatization' as mentioned above as the 'tokenization' method. Tokenization is the process of converting input words into known, specific inputs. The tokenization process given to you in Section 4 is a little more complex, mainly just removing punctuation as well as lemmatizing.

Since you're asked to use a dimension size of 8 for the embeddings in Section 4, to visualize these, you need to reduce the dimensionality down to 2. Dimensionality reduction would have been covered in a previous course, but the code for doing so is done with principle component analysis (one of several ways to do this), and the code is given to you to do this. This allows you to visualize the embeddings with the same 2-D plot as used in Section 3 of the assignment.