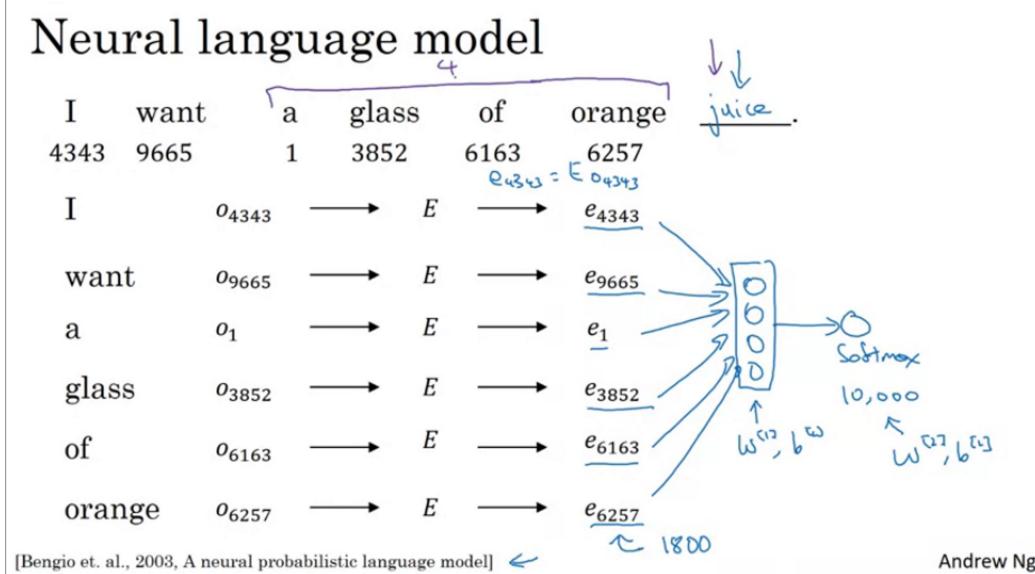


# Learning Word Embeddings: Word2vec & GloVe

Wednesday, September 30, 2020 3:11 PM

## ① learning word embeddings

- The newer models are simpler than initial models
- We will learn the initial models first to gain intuition
- Neural language model



$$\rightarrow e_j = E \cdot o_j$$

$\rightarrow$  The embeddings are passed to NN layer followed by softmax layer

$\rightarrow$  Input size is  $(300 \times 6, 1)$  if all 6 previous words are used (window size = 6)

- Other context / target pairs

Eg: I want a glass of orange juice to go along with my cereal

Context: last 4 words  $\rightarrow$  a glass of orange  
4 words both sides  $\rightarrow$  a glass of orange  
to go along with

Last 1 word  $\rightarrow$  orange  
Nearby 1 word  $\rightarrow$  glass

$\hookrightarrow$  skip gram model  
(works remarkably well)

Nearby 1 word → glass  
 ↳ skip gram model  
 (works remarkably well)  
 (for learning embeddings)

## ② Word2Vec

- Skip grams

→ We choose context and randomly a target within a specific window (eg: +10/-10 words)

context	target
Orange	juice (next word)
orange	glass (2 words back)
orange	mug (6 words after)

→ Now this is used as a supervised learning problem

- Word2vec model

→ Vocabulary size = 10000 words

→ we need to learn context ( $c$ ) to target( $t$ )

→  $e_c = E \cdot O_c$   
 → soft max is used to get  $P(t|c) = \hat{y}$



$$\text{softmax: } P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{i=1}^{10000} e^{\theta_i^T e_c}} \quad \theta_t = \text{parameter associated with } t$$

→ Cross entropy loss is used

$$L(\hat{y}, y) = -\sum_{i=1}^{10000} y_i \log \hat{y}_i$$

$$y = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \text{one hot vector}$$

- Problems with softmax classifier

→ The softmax layer

$$P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{i=1}^{10000} e^{\theta_i^T e_c}}$$

is very computationally intensive as we are summing through our entire vocabulary

→ This can be solved by using a hierarchical softmax classifier which uses a tree like structure with more common words at the top and rare words in the deeper nodes

- To sample context  $c$ 
  - can be chosen randomly but frequent words like a, the, my.. will dominate
  - In reality, it is chosen randomly but not uniformly to balance the common words and non-common words
- 2 versions of Word2Vec are used  
skip grams & CBOW (continuous bag of words)

### ③ Negative Sampling

- Negative Sampling is similar to skip gram model but much more efficient algorithm
- For the example  
I want a glass of orange juice to go along with my cereal  
The sampling would be

context	word	target	
orange	juice	1	→ positive example
orange	king	0	
orange	book	0	
orange	the	0	
orange	of	0	

↳ it is fine if some negative samples appear in same sentence

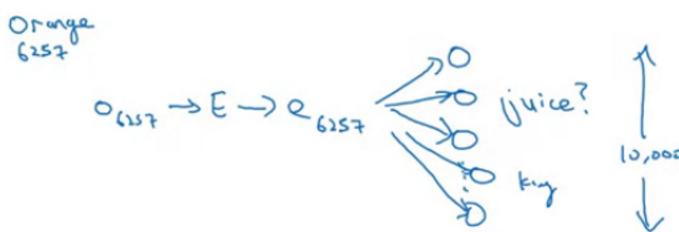
- Positive example is chosen by skip-gram technique
- Negative examples ( $K \rightarrow 5$  to 20 acc. to original paper) are picked randomly from vocabulary

- The model is defined as

→ target is the probability that a pair of words are a context target pair

→ we use logistic regression model

$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$



10,000 binary classification problems  
as we train  
 $k+1$  classifiers  
each iteration

- To select negative examples

→ we can sample according to word frequencies but results in frequent words like the, of, and ...

→ Hence, we use (acc. to original paper)

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{1000} f(w_j)^{3/4}} \quad f \rightarrow \text{freq.}$$

## ④ GloVe word vectors

- Used less frequently, known for simplicity stands for global vectors
- We will choose context and target similar to previous models.
- We define

$$x_{ct} = \underset{C}{\text{No. of time } t \text{ appears in context}}$$

- $x_{ct} = x_{tc}$  is also possible if we choose a window pair

- The model is defined as

$$\text{minimize} \rightarrow \sum_{i=1}^{100000} \sum_{j=1}^{100000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})^2$$

→  $f(x_{ij})$  is weighting term used to avoid calculations if  $x_{ij} = 0$ , hence for solving the log0 problem and to not give too much weight to frequent words and not too little to infrequent words

→  $\theta$  &  $e$  are symmetrical and swappable unlike previous models

Model

## Mimimize

$$\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\Theta_i^T e_j + b_i + b_j) - \log \underline{x_{ij}}^2$$

"weighting term"

$f(x_{ij}) = 0$  at  $x_{ij} = 0$ .      " $0 \log 0 = 0$ "

this, is, at, a, ...

derivation

$\Theta_w^{(\text{final})} = \frac{\Theta_w + \Theta_w}{2}$