

1 - hot representation

t-SNE ($300 \text{D} \rightarrow 2\text{D}$) for visualizing high-D embedding
non-linear dimensionality reduction

Why word embedding?

for example, a model have not seen durian, but have seen other fruit such as orange, embedding can draw similarity. This enable process large unlabelled datasets.

Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (1-100B words)

(Or download pre-trained embedding online.)

2. Transfer embedding to new task with smaller training set.
(say, 100k words) $\rightarrow 10,000 \rightarrow 300$

3. Optional: Continue to finetune the word embeddings with new data.

Some similarity with Face Recognition (encoding)

Properties: for example. $\mathbf{Lman} - \mathbf{Lwoman} \approx \mathbf{LKing} - \mathbf{LQ}$

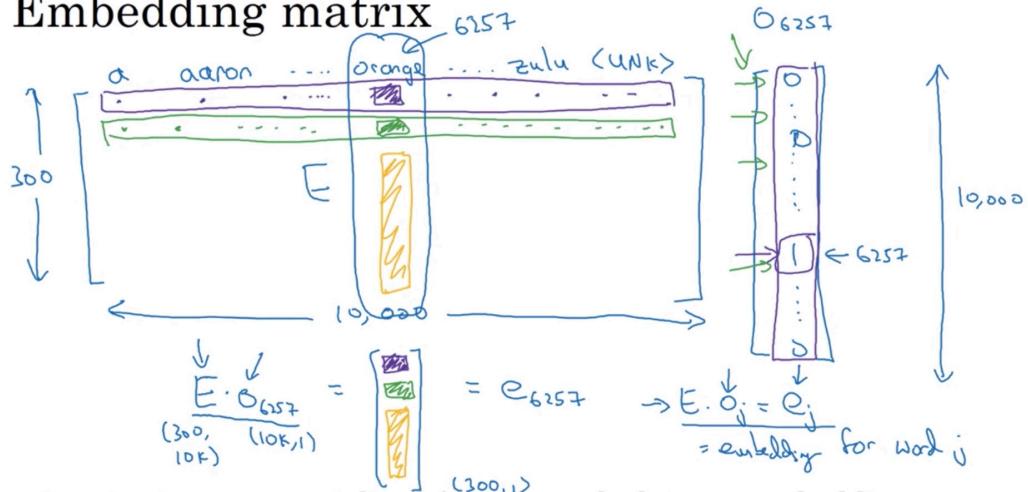
The vector difference is similar for the 2 pair of words

\Rightarrow Find word W : $\arg \max_W \underbrace{\text{similarity}(\mathbf{Lw}, \mathbf{LKing} - \mathbf{Lman} + \mathbf{Lwoman})}_{\text{Cosine Similarity}}$

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

Embedding Matrix

Embedding matrix



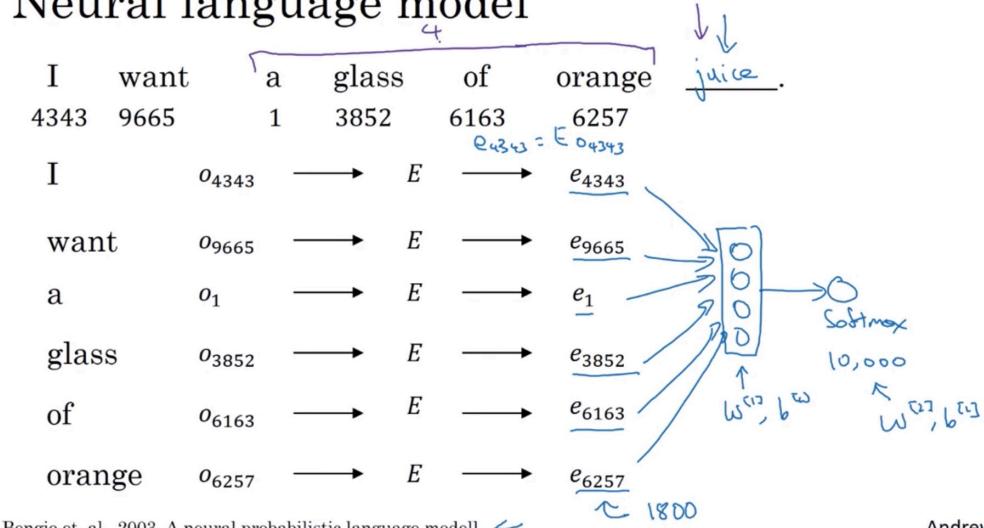
In practice, use specialized function to look up an embedding.

\rightarrow Embedding

Andrew Ng

Because matrix multiplication with sparse matrix inefficient

Neural language model



Bengio et. al., 2003, A neural probabilistic language model]

Andrew

For BRNN, using nearby one word instead of 4 or n-words
is surprising effective as well

Word2Vec

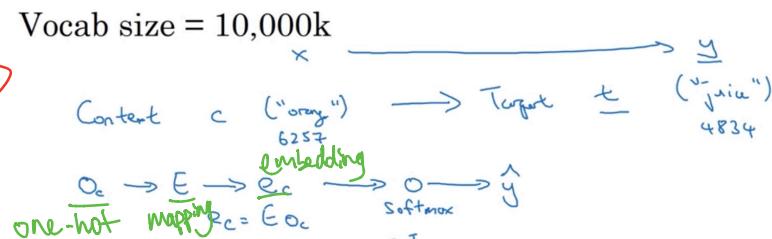
skip gram : Given context word , predict random forward / backward target word

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

Context	Target
orange	juice
orange	glass
orange	my
	↑

Vocab size = 10,000k

Model \Rightarrow



$$\text{Softmax: } p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

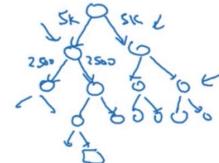
θ_t = parameter associated with output t

extensive, slow

$$L(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4834$$

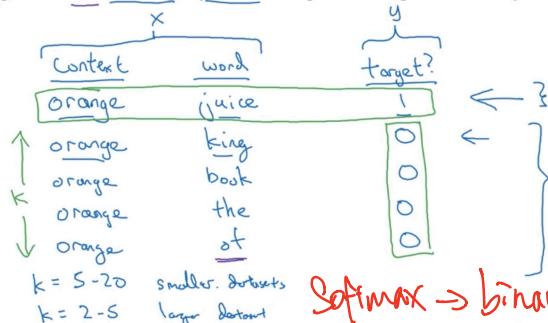
Solution: hierarchical softmax



Negative Sampling

Defining a new learning problem

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



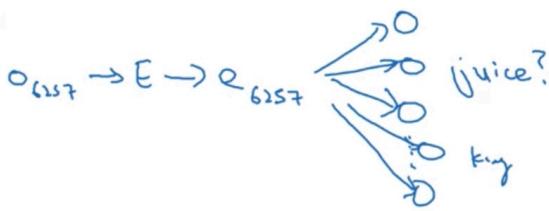
$$\text{Softmax} \quad p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

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

Softmax \rightarrow binary classification

Negative Sampling - continued

Orange
6257



Instead of training 10,000 examples, train on 1 positive example and 4 randomly chosen negative examples.

(0, 000 softmax \rightarrow (0, 000 binary classification \rightarrow k+1 binary)
Therefore this is a lot cheaper to train.

How to Sample?

- ④ $P(w_i)$ \rightarrow getting a lot of "the, of . and , --"
- \hookrightarrow Sample according to Empirical frequency. (How often)
- ⑤ $\frac{1}{|V|}$
 \hookrightarrow uniform random distribution

- ③ The best . combine ④ and ⑤

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

GloVe \rightarrow global vectors for word representation

$X_{ij} =$ # of times i appears in context of j
 i target j context

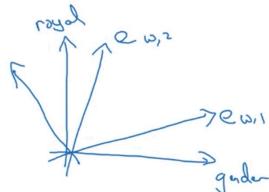
$$\min \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} \underbrace{f(x_{ij})(\theta_i^T \theta_j)}_{\text{weighting term}} - \log X_{ij})^2$$

θ_i, θ_j are symmetric $\frac{\theta_i + \theta_j}{2}$

$f(x_{ij}) = 0$ if $x_{ij} = 0$ \Rightarrow stop words (the, and, is)

A note on the featurization view of word embeddings

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)
Gender	-1	1	-0.95	0.97
Royal	0.01	0.02	0.93	0.95
Age	0.03	0.02	0.70	0.69
Food	0.09	0.01	0.02	0.01



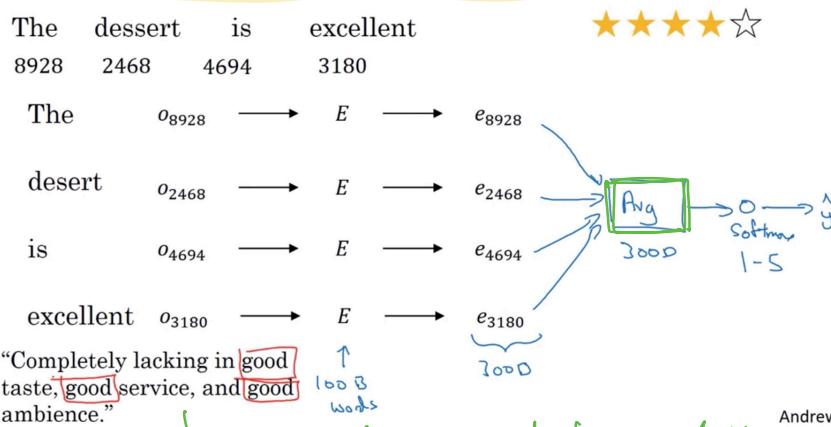
$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i - b_j - \log X_{ij})^2$$

$$[(A\theta)^T (A^T e_j) = \theta^T A^T e_j]$$

why sometimes the component is not human interpretable

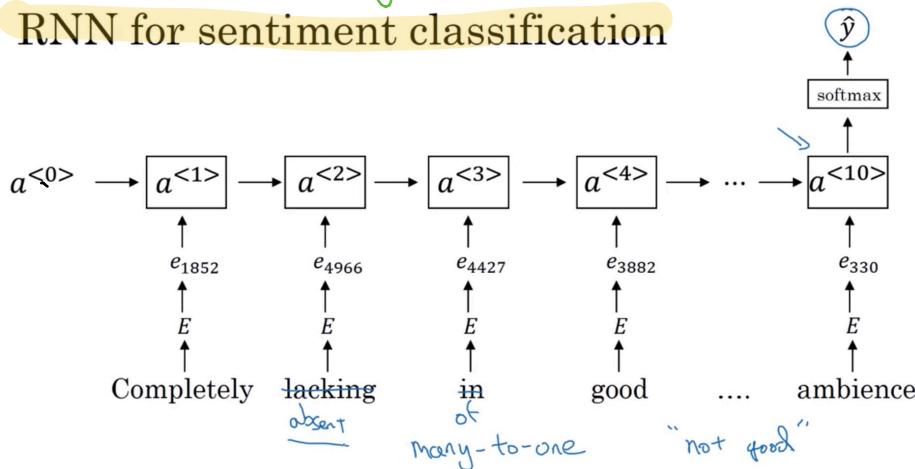
Sentiment classification (lack of labelled training data)

Simple sentiment classification model



↳ average skewed by # of times "good" are mentioned

RNN for sentiment classification



Debiasing Word Embeddings

Eg.

Man:Woman as King:Queen

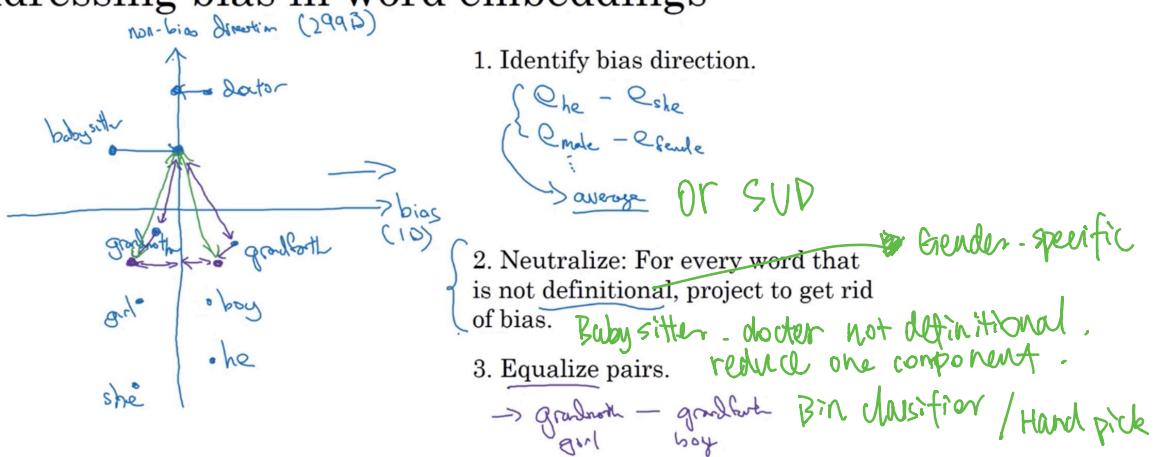
Man:Computer_Programmer as Woman:Homemaker X

Father:Doctor as Mother:Nurse X

Gender Stereotype

Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.

Addressing bias in word embeddings



Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings]

Andrew Ng

1. Suppose you learn a word embedding for a vocabulary of 10000 words. Then the embedding vectors should be 10000 dimensional, so as to capture the full range of variation and meaning in those words.

True
 False

! Incorrect

No, the dimension of word vectors is usually smaller than the size of the vocabulary. Most common sizes for word vectors ranges between 50 and 400.

9. Suppose you have a 10000 word vocabulary, and are learning 500-dimensional word embeddings. The GloVe model minimizes this objective:

$$\min \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij})(\theta_i^T e_j + b_i + b_j' - \log X_{ij})^2$$

Which of these statements are correct? Check all that apply.

θ_i and e_j should be initialized to 0 at the beginning of training.

! This should not be selected

The variables should not be initialized to 0 at the beginning of training.

θ_i and e_j should be initialized randomly at the beginning of training.

X_{ij} is the number of times word i appears in the context of word j.

✓ Correct

The weighting function $f(\cdot)$ must satisfy $f(0) = 0$.

✓ Correct

The weighting function helps prevent learning only from extremely common word pairs. It is not necessary that it satisfies this function.