Learning (Training) Word Embedding

method 1: Bengio et. al., 2003, A neural probabilistic languge model

I    want    a    glass    of    orange    ———.    (total words: 10,000)
4343   9665   1   3852   6163   6257

Given a few words before a blank, predict the blank.

I        $O_{4343}$ → E → $e_{4343}$
(denoted a one-hot)

want    $O_{9665}$ → E → $e_{9665}$ →

a        $O_1$ → E → $e_1$ ↗

glass    $O_{3852}$ → E → $e_{3852}$ ↗

of        $O_{6163}$ → E → $e_{6163}$ ↗

orange  $O_{6257}$ → E → $e_{6257}$ ↗

$\begin{bmatrix} O \\ O \\ O \\ O \\ \vdots \\ O \end{bmatrix}$ → O
softmax
10,000

method 2: Mikolov et. al., 2013 Efficient estimation of word representation in vector space
(Word2Vec)

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

skip grams: randomly choose target words within a given window size.

| context | target |
|---------|--------|
| orange  | juice  |
| orange  | glass  |
| orange  | my     |
| ⋮       | ⋮      |

supervised learning problem
context → target

How to sample context?
   need to find a balance between common words and less common words.

Total vocab size = 10,000 ~~k~~

$$O_{context} \rightarrow E \rightarrow O_{context} \rightarrow O_{softmax} \rightarrow \hat{y}$$

Softmax:

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum\limits_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

⟩ problems with softmax classification

$$\sum\limits_{j=1}^{10,000} e^{\theta_j^T e_c} \longleftarrow O(n) \quad n: \text{number of vocabs}$$

Hierarchial Softmax

$$\log^{|n|}$$

method 3: Mikolov et. al., 2013, Distributed representation of words and
phrases and their compositionality
(Negative Sampling)

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

context word : a word in the sentence

positive word pair : pair the context word with a random chosen
word within the window size

negative word pair : pair the context word with a random
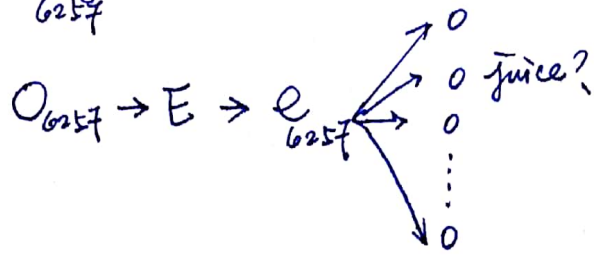chosen word from the dictionary

ex: orange , juice    label
                        1
    orange , king       0
    orange , book       0

pick 1 context word          K: 5-20
generate 1 positive pair     smaller dataset
& generate k negative pairs  K: 2-5
                             larger dataset

## Model ( # vocab = 10,000 )

| context word | | $y$ target |
|---|---|---|
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

orange
6257

$$O_{6257} \rightarrow E \rightarrow e_{6257} \to \begin{matrix} 0 \\ 0 \text{ juice?} \\ 0 \\ \vdots \\ 0 \end{matrix}$$

10,000 binary
logistic regression problems
but each iteration only train (k+1) of them

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

How to sample negative pairs?

$p(w_i)$
sample
according to the empirical frequency
( $a \cdot the \cdot of \ldots$ tend to be picked )

$\longleftrightarrow$

$$\dfrac{f(w_i)^{3/4}}{\sum\limits_{v=1}^{10000} f(w_i)^{3/4}}$$

$\longleftrightarrow$

$$\dfrac{1}{|V|}$$

uniformly
sample
( not representative )

Method 4: GloVe (global vectors for word representation)

(Pennington et. al. 2014 GloVe: Global vectors for word representation)

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

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

target
word

context
word

(How often word $i$ and word $j$ appear together)

Model

$\rightarrow$ $\theta_i$ $e_j$ are symmetric

$\Rightarrow$ $e_w^{(final)} = \dfrac{e_w + \theta_w}{2}$

$$\text{minimize } \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(X_{ij}) \left( \theta_i^T e_j + b_i + b_j' - \log X_{ij} \right)^2$$

a weighting term

- $f(X_{ij}) = 0$ if $X_{ij} = 0$
- used to find a balance between common words and non-common ones