Some more complicated algorithms

I want a glass of orange <span style="color:red">juice</span>
4343  9665  1  3852  6163  6257

I        $O_{4343} \rightarrow E \rightarrow e_{4343}$
want   $O_{9665} \rightarrow E \rightarrow e_{9665}$    Concatenate 连接.

$\begin{bmatrix} a \\ glass \\ of \\ orange \end{bmatrix}$    $\begin{matrix} e_1 \\ e_{3852} \\ e_{6163} \\ e_{6257} \end{matrix}$   $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  $\rightarrow$  $\bigcirc$  softmax.
$\underline{\frac{10,000}{\uparrow}}$
vocab size.

$300 \times 6$
$= 1800.$

$\rightarrow$ 4 words  $300 \times 4 = 1200.$

or  only look at the
previous 4 words

use a fixed history ~~size~~

~~en~~ enables you to work with
arbitrarily long sentences.

This will do a descent job.

other context / target pairs.

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

4 words on left and right.   a glass of orange ___ to go with .

✱ | Last 1 word |    nearby one word.
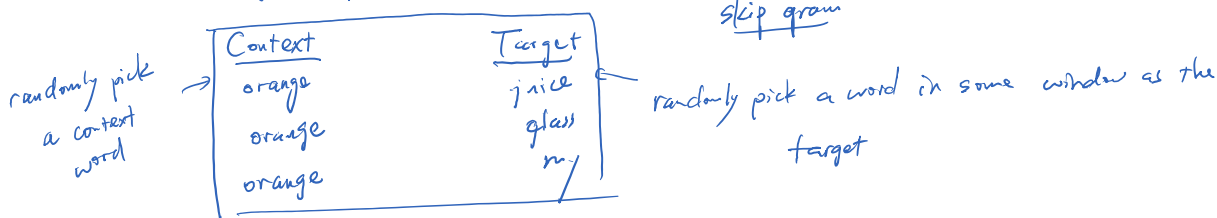      ⌐ skip gram model

simpler ~~a~~ algorithm — just one word.
if your main goal is just word embedding,  one word is enough!

| Word 2 Vec Algorithm |   Thomas Mikolov. ~~et al~~ <span style="color:red">Kai Chen, Greg Corrado, Jeff Dean</span>

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

skip gram

randomly pick → | Context | Target |
a context     | orange  | juice  |
word          | orange  | glass  |     randomly pick a word in some window as the
              | orange  | my     |                    target

The goal is good word embedding.    $O = $ one-hot encoding

Model.   vocab size = 10,000 k.
        Context  $c(\text{"orange"}) \longrightarrow$ Target  $t(\text{"juice"}).$
        $O_c \rightarrow E \rightarrow e_c \rightarrow O \rightarrow \hat{y}$
                                    softmax
Softmax    $p(t|c) = \dfrac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$    $\theta_t$ parameter associated with output $t$.
        <span style="color:red">开始? identifiable?</span>
        <span style="color:red">如何求?</span>   $y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$  $\longleftarrow$ 4834

$$\sum_{j=1}^{10000} e^{\theta_j^T e_c}$$

$\text{何以 ? identifiable?}$ (learn 到の $W_t$.) $\quad y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \longleftarrow 4834$

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

Problem : Computationally taxing . — Sum over 10,000 in softmax

Solution : "hierarchical softmax"    ✓ not real really symmetric

balanced

≥/₀ .  tree   in applications.

P(c).

read the original paper        alternative version

2 versions    ① surrounding words → middle word

method

Negative Sampling  → computationally cheap method version for skip gram

Mikolov, Sutskever, Chen, Corrado. Dean.

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

$\overset{\times}{\phantom{x}}$   $\overset{y}{\phantom{y}}$

| Context | Target | Target ? |
|---|---|---|
| orange | juice | 1 |
| | king | 0 |
| | book | 0 |
| | the | 0 |
| | of | 0 |

$k$

Which k to choose?

k = 5 – 20 for small data set

k = 2 – 5 for larger datasets
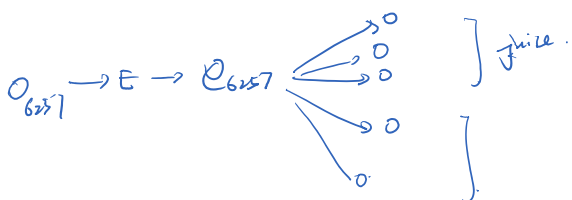
① pick a context word, pick a target word. — positive example

② for k time pick random words in the dictionary. label as 0.

It is ok if one of the randomly picked word is actually in the window.

Softmax : $p(t|c) = \dfrac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$ ← 10000 way softmax

$p(y=1 | c,t) = \sigma(\theta_t^T e_c)$.    Sigmoid.   — with negative sampling

$O_{6257} \to E \to e_{6257}$  ⟨ juice.    10,000 binary classification problem

use computational power !

Selecting negative example?

Selecting negative example?

$\left[\begin{array}{l} \cdot \text{ proportional to word freq} \longrightarrow \text{end up with lots of "a" "the"} \\ \cdot \text{ equal prob? Not good either} \end{array}\right.$

use <u>something in between</u>.

折中 (↑)

$$p(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$$

$\boxed{\text{GloVe word vector}}$ — not used as much as

· word2vec

skip gram

Global Vector for word representation.

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

$c, t$.

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

$\underset{c}{\underset{\swarrow\searrow}{t}}$      $\underset{t}{\hat{\uparrow}}$      $\underset{c}{\uparrow}$

for GloVe define $X_{ij} = X_{j,i}$

$\underline{X_{ij}}$ a count of how much $i$ and $j$ appear with each other.

$\boxed{\text{Minimize}}$ $\sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) \left( \underset{\underset{\Theta_t^T e_t}{\uparrow}}{\Theta_i^T e_j} + b_i + \underset{\underset{c}{\downarrow}}{b_j'} - \log X_{ij} \right)$   $\underset{\downarrow}{t}$   $\underset{i}{c}$

Weighting term $f(x_{ij}) = 0$ if $x_{ij} = 0$   "$0 \log 0 = 0$" 防止计算 $\log 0$.

$\curvearrowright$ <u>frequent</u> word. this, is, a.

<u>infrequent</u> word. <u>durian</u>

weight — not give frequent words too much weight

not give infrequent words too little weight.

$\underline{\underline{\Theta_i}} \ \underline{e_j}$ are <u>symmetric</u>

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

About the featurization : You <u>cannot</u> guarantee that each dimension of featurization is interpretable.    royal$\nearrow$   e.v.2.

Abou... ... ... ... ...

is interpretable.

You can't guarantee the features are human-interpretable.

royal

$e_{v,2}$

$e_{v,1}$

gender