

Lesson 2: word embeddings

I

(or how to make sense of

King - Man + Woman \approx Queen
when you are a machine)

In this lesson, we make a huge leap forward in time from 1913 (Markov) to 2013 (Mikolov et al.)

- Reminder et al. \equiv "et alia" meaning "and others".

- Flash article on screen.

This article is the basis of a miraculous algorithm called word2vec (code.google.com/archive/p/word2vec/)

My choice of wording is careful:

"Miraculous" \equiv Incredible results but we do not understand why.

Indeed, as we shall see, the ~~definition~~ justification of the underlying statistical model is not at the mathematical level of rigor yet the technique works incredibly well.

I think it is important to ~~have~~ ^{know} that, while keeping an open mind

General idea: Starting from a (large)

corpus of text, word2vec constructs an

embedding $\hat{\varphi}$: Words $\longrightarrow \mathbb{R}^d$ $d \gg 1$

$$w \longmapsto \varphi(w) = v_w$$

Here \mathbb{R}^d is "the vector space of features"

and $\hat{\varphi}$ is constructed by training a model

which tries to predict a word w according

to its context i.e neighboring words.

The space parameter Θ is the set of maps φ

$$\hat{\varphi} = \underset{\varphi \in \Theta}{\text{Argmax}} \log L(\varphi, \text{Text Corpus})$$

$\longmapsto \log \text{Likelihood}$

More precisely: Although there are variants

and multiple additional tricks, let us

detail the standard CBOW Model (Continuous Bag-of-words)

Given a corpus made of an ordered set

of words (w_1, w_2, \dots, w_T) ; there are multiple assumptions

1/ we assume that the log-likelihood is of the "Context of t-th word"

form $L = \sum_{t=1}^T \log(w_t | \underline{C_t})$

$$C_c = \{ w_{t+k} \mid 1 \leq |k| \leq c \}$$

$$|C_c| = 2c$$

Remark: As we read the text, this ~~is~~ can be seen as if the appearance of words is Markovian with the context as the single explanatory variable

2/ Assume there is a scoring function $S_\varphi(w, C) = s(w, C)$ which indicates the score of

a word w appearing in context C . \triangle This is where φ enters φ
 We take as definition:

$$s(w, C) = \frac{1}{|C|} \sum_{w' \in C} \langle \varphi(w), \varphi(w') \rangle = S_\varphi(w, C)$$

" " "

σ_w $\sigma_{w'}$

3/ We take for $\log p(w|C)$

the following $\left(\begin{array}{l} \rightarrow \text{with score of positive outcomes} \\ \searrow \text{with score of negative outcomes} \end{array} \right)$

$$\log p(w|C) = \log(1 + e^{-s(w, C)}) + \sum_{n \in N_c} \log(1 + e^{-s(w, C)})$$

$N_c =$ set of sampled negative ~~examples~~ candidates

i.e. words which never appear with C

Indeed, as much it is important to teach the machine what are the positive outcomes, it is also important to teach the negative ones!

In the end,

$$\log L(\varphi) = \sum_{t=1}^T \left[\log(1 + e^{s_{\varphi}(w_t, c_t)}) \right.$$

$$\left. + \sum_{m \in N_{c_t}} \log(1 + e^{-s_{\varphi}(w_t, c_t)}) \right]$$

$$\Rightarrow \hat{\varphi} = \underset{\varphi}{\text{Argmax}} \log L(\varphi)$$

L → Learned through gradient descent.