Computational
Semantics
for
Natural
Language
Processing

Alexis Tabin

16 - 821 - 803

# Assignment 1

## 1. Skip-Gram

$w \in V_w$ , a word , $V_w$ word vocabulary , $w \rightarrow \vec{w} \in \mathbb{R}^d$

$c \in V_c$ , a context , $V_c$ context vocabulary , $c \rightarrow \vec{c} \in \mathbb{R}^d$

$D$ , collection of observed word and context pairs

$\#(w,c)$ , numbers of time $(w,c)$ appears in $D$

$$\#(w) = \sum_{c' \in V_c} \#(w,c') \qquad \# c = \sum_{w' \in V_w} \#(w',c)$$

1. $P(D=1 \mid w,c) = \sigma(\vec{w} \cdot \vec{c}) = \dfrac{1}{1 + \exp(-\vec{w} \cdot \vec{c})}$

$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w,c) \log\left(\sigma(\vec{w} \cdot \vec{c})\right)$$

$$= \sum_{w \in V_w} \sum_{c \in V_c} \frac{\#(w,c)}{1 + \exp(-\vec{w} \cdot \vec{c})}$$

### Answers 1.1

An objective functions tries to optimize something. Here, we want to maximize $P(D=1 \mid w,c)$ for ~~an~~ observed $(w,c)$ pairs, ~~while maximizing~~ ~~$P(D=0)$~~ . We can see that we could use this log likelihood function as the objective function, it would ~~work~~, but it is really heavy on computation It is heavy because for each time we need to compute $\ell$, we need to go through every word, and for each word, we go though every context.

1

.2.

$c_{N_i}$ randomly drawn negative samples

for $(w,c)$ we want to

    1. maximize $P(D=1|w,c)$

    2. while maximizing $P(D=0|w,c)$ for negative samples $\Big|$ i

$k = $ # of samples

$c_N \sim \dfrac{\#(c)}{|D|} = P_D(c)$

## Answers 1.2

Recall that $\quad P(D=1|w,c) = \sigma(\vec{w}\cdot\vec{c}) = \dfrac{1}{1+\exp(-\vec{w}\cdot\vec{c})}$

$\Rightarrow \quad P(D=0|w,c) = 1-P(D=1|w,c) = 1 - \sigma(\vec{w}\cdot\vec{c})$

$$\mathcal{L}(w,c) = P(D=1|W,C) \overset{\#k \leftarrow \text{# of samples}}{\prod_{i=1}} P(D=0|w,c) \quad \text{by } \Big| \text{i}$$

we calculate the log

$$\log \mathcal{L}(w,c) = \log\left(P(D=1|w,c)\prod_{i=1}^{k} P(D=0|w,c)\right) = \log(\sigma(\vec{w}\cdot\vec{c})) + \log\left(\prod_{i=1}^{k}(1-\sigma(\vec{w}\cdot\vec{c}))\right)$$

$$= \log(\sigma(\vec{w}\cdot\vec{c})) + \sum_{i=1}^{k} \log(1-\sigma(\vec{w}\cdot\vec{c}))$$

💡 multiply 2nd term and divise by $k$.

$$= \log(\sigma(\vec{w}\cdot\vec{c})) + \frac{k}{k}\sum_{i=1}^{k} \log(1-\sigma(\vec{w}\cdot\vec{c}))$$

$k \cdot \dfrac{1}{k}$

$$= \mathbb{E}_{c_N \sim P_D}\left(\log(1-\sigma(\vec{w}\cdot\vec{c}))\right)$$

$\boxed{\text{ii}}$

$1 - \sigma(\vec{w}\cdot\vec{c}) = 1 - \dfrac{1}{1+e^{-\vec{w}\cdot\vec{c}}}$

$\qquad = \dfrac{1+e^{-\vec{w}\cdot\vec{c}} - 1}{1+e^{-\vec{w}\cdot\vec{c}}}$

$\qquad \overset{\triangle}{=}$

$= \log$

$1-\sigma(x) = 1 - \dfrac{1}{1+e^x} = \dfrac{1+e^x-1}{1+e^x} = \dfrac{e^{-x}}{1+e^{-x}} = \dfrac{1}{e^x(1+e^{-x})} = \dfrac{1}{e^x+1}$

2

Start Following of 1.2

with ii , we have

$$\log \, h(w,c) = \lg(\sigma(\vec{w}.\vec{c}')) + \frac{k}{\alpha} \, h \cdot \mathbb{E}_{c_N \sim P_D} \left( \log \left( \sigma(-\vec{w}.\vec{c}') \right) \right)$$

So we have for each pair
the objective function we want to minimize :

$$h(\vec{w},\vec{c}) = \log \left( \sigma(\vec{w}.\vec{c}') \right) + h \, \mathbb{E}_{c_N \sim P_D} \left( \log \left( \sigma(\vec{w}.\vec{c}') \right) \right)$$

For all pairs, it gives

$$l = \sum_{w \in V_w} \sum_{c \in V_c} \#(w,c) \left( \log \sigma(\vec{w}.\vec{c}') + h \cdot \mathbb{E}_{c_N \sim P_D} \left[ \log(-\vec{w}.\vec{c}') \right] \right)$$

___

**3.** Assuming
$\vec{w}.\vec{c}'$ independent , $l(w,c) = ?$ , let's rewrite $l$

$$l = \sum_w \sum_c \#(w,c) \left( \log \sigma(\vec{w}.\vec{c}') \right) + \sum_w \sum_c \#(w,c) \left( h \cdot \mathbb{E}_{c_N \sim P_D} \left[ \lg \sigma(-\vec{w}.\vec{c}_N') \right] \right)$$

$$= \sum_w \sum_c \#(w,c) \left( \log \sigma(\vec{w}.\vec{c}') \right) + \sum_{w \in V_w} \#(w) \left( h \cdot \mathbb{E}_{c_N \sim P_D} \left[ \log \sigma(-\vec{w}.c_N') \right] \right)$$

recall the expectation term : $\mathbb{E}_{c_N \sim P_D} \left[ \log(\sigma(\vec{w}.c_N')) \right] = \sum_{c_N \in C} \frac{\#(c_N)}{|D|} \log \sigma(-\vec{w}.c_N)$

$$= \frac{\#(c)}{|D|} \log \sigma(-\vec{w}.\vec{c}') + \sum_{c_N \in V_c \backslash \{c\}} \frac{\#(c_N)}{|D|} \log(-\vec{w}.c_N)$$

$\rightarrow$ we then plug that back in the loss function :

$$l = \sum_w \sum_c \#(w,c) \log \left( \sigma(\vec{w}.\vec{c}') \right) + \sum_w \#w \cdot h \cdot ( --- ) \qquad \underrightarrow{\text{answer}} \qquad \boxed{3}$$

We can see how to express $\ell(w,c)$:

$$\ell(w,c) = \#(w,c)\log(\vec{w}\cdot\vec{c}) + h\cdot\#w\cdot\frac{\#c}{|D|}\log(-\vec{w}\cdot\vec{c})$$

1.4  $x = \vec{w}\cdot\vec{c}$ , we want $x^* = \max_x \ell$

$\rightsquigarrow \ell(x) = \#(w,c)\cdot\sigma(-x)$

$\rightsquigarrow \ell(x) = \#(w,c)\cdot\log(\sigma(x)) + \#(w)\cdot h\cdot\frac{\#c}{|D|}\log(\sigma(-x))$

$x^* = \max_x(\ell) \rightsquigarrow \dfrac{\partial \ell(x)}{\partial x} = 0$

$\rightsquigarrow \#(w,c)\,\sigma'(x)\cdot\dfrac{1}{\sigma(x)} + \underbrace{\dfrac{\#(w)\cdot h\cdot\#(c)}{|D|}}_{\text{as we saw in 1.1}}\cdot\sigma'(-x)\dfrac{1}{\sigma(-x)} = 0$

$$\sigma(x) = \frac{1}{1+e^{-x}} = (1+e^{-x})^{-1}$$

$$\sigma'(x) = -(1+e^{-x})^{-2}(1+e^{-x})'$$

$$= \#(1+e^{-x})^{-2}e^{-x}$$

$$= \frac{1}{e^x + e^{-2x}}$$

$$= \cdots = \sigma(x)(1-\sigma(x))$$

$\rightsquigarrow \#(w,c)\,\underbrace{\dfrac{\sigma(x)(1-\sigma(x))\dfrac{1}{\sigma(x)}}{\sigma(-x)}}$

$+\dfrac{h\#(w)\#(c)}{|D|}\sigma(x) = 0$

$\vdots \rightarrow$ with trivial adjustments

result $\rightarrow$

$$e^{2x} - \left(\dfrac{\#(w,c)}{h\,\#(w)\dfrac{\#(c)}{|D|}} - 1\right)e^x$$

$$-\dfrac{\#(w,c)}{h\,\#(w)\dfrac{\#c}{|D|}} = 0$$

rewritten

4

following af 1.4

$$\to e^{2x} - \left( \frac{\#(w,c)}{h \, \#(w) \, \frac{\#(c)}{|D|}} - 1 \right) e^x - \frac{\#(w,c)}{h \, \#(w) \, \frac{\#c}{|D|}} = 0$$

let's define $y = e^x$

$$\Rightarrow y^2 - by - c = 0$$

$\hookrightarrow$ this equation has 2 solutions :

$$y_1 = -1 \quad \Rightarrow \quad e^x = -1 \quad \Rightarrow \quad \text{not possible}$$

$$y_2 = \frac{\#(w,c)}{h \cdot \#(w) \cdot \frac{\#c}{|D|}}$$

Remember, we have $y = e^x = e^{\vec{w} \cdot \vec{c}}$

$$\Rightarrow \vec{w} \cdot \vec{c} = \log(y)$$

$$x^* =$$

$$\Rightarrow \vec{w} \cdot \vec{c} = \log \left( \frac{\#(w,c) \cdot |D|}{\#(w) \cdot \#(c)} \cdot \frac{1}{a} \right) = \text{argmax } \ell$$

## 5. $k=1$

$$x^* = \log\left(\frac{\#(w,c)\cdot|D|}{\#(w)\cdot\#(c)}\cdot\frac{1}{k}\right)$$

$$= \log\left(\frac{\#(w,c)}{|D|}\cdot\frac{|D|}{\#(w)}\cdot\frac{|D|}{\#(c)}\right)$$

$$= \log\left(\frac{P_D(w,c)}{P_D(c)\,P(w)}\right) = PMI(w,c)$$

by definition

## 6. $M = W^*\cdot C^{*T}$ ?

$$M_{ij}^{SGNS} = W_i\cdot C_j = \vec{w}_i\cdot\vec{c}_j = PMI(w_i,c_j) - \log k$$

For negative-sampling values $k > 1$, SGNS is factorizing a shifted PMI matrix $M^{PMI_k} = M^{PMI} - \log k$