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 , $w \to \vec{w} \in \mathbb{R}^d$

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

$D$ , 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 ~~woh~~ 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.

$c_N$, 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  |i

$k = \#$ of samples

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

## Answers 1.2

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

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

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

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 $k$.

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

|i|

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

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

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

$=$

$= \log$

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

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

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

---

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

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

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

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

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

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

$$\ell = \sum_w \sum_c \#(w,c) \log \left( \sigma(\vec{w}.\vec{c}') \right) + \sum_w \#w \cdot h \cdot ( \cdots )$$

answer $\longrightarrow$

$\overline{\ulcorner 3}$

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

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

$\boxed{1\text{-}4}$ $x = \vec{w} \cdot \vec{c}$ , we want $x^* = \max_x \ell$

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

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

$x^* = \max_x (\ell) \quad \sim \frac{\partial \ell(x)}{\partial x} = 0$

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

$\sim \#(w,c) \, \underbrace{\sigma(x)(1 - \sigma(x)) \frac{1}{\sigma(x)}}_{\sigma(-x)}$

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

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

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

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

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

$\quad = \ldots = \sigma(x)(1 - \sigma(x))$

$\begin{matrix} \circ \\ \circ \\ \circ \end{matrix} \to$ with trivial adjustments

result $\to$

$\Rightarrow \boxed{\begin{array}{l} e^{2x} - \left( \dfrac{\#(w,c)}{h \, \#(w) \frac{\#(c)}{|D|}} - 1 \right) e^x \\[2em] \quad - \dfrac{\#(w,c)}{h \, \#(w) \frac{\#c}{|D|}} = 0 \end{array}}$ rewritten

$\boxed{4}$

following of 1.4

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

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.** $h = 1$

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

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

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

by definition

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

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

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

$$\sqrt{6}$$

## Assignment 1

## 2. Multi-prototype Word Embeddings

- $(w, c)$, pair of a word and a context
- $N_w$, word senses (prototypes)
- $W \in \mathbb{R}^{|V_w| \times N_w \times d}$
- $h_w \in \{1, \dots, N_w\}$, $h_w = h$ ~~means~~ if $w$ means if $h^{th}$ prototype
- $\pi_{wh} = P(h_w = h | w)$

We model
$$P(D=1 | w, c) = \sum_{h=1}^{N_w} P(h_w = h | w) P(D=1 | w, h_w = h, c)$$
$$= \sum_{h=1}^{N_w} \pi_{wh} P(D=1 | w, h_w = h, c)$$

### 2.1. $P(D)$, $L(D) = ?$

likelihood    log-likelihood

$$P(D) = \prod_{(w,c) \in D} P_r\left[ D = 1 | w, c \right]$$

$$= \prod_{(w,c) \in D} \sum_{h=1}^{N_w} P(h_w = h | w) P(D=1 | w, h_w = h, c)$$

$$L(D) = \log(P(D)) = \log \left( \prod_{(w,c)} \sum_{h=1}^{N_w} \pi_{wh} P(D=1 | w, h_w = h, c) \right)$$

$$= \sum_{(w,c) \in D} \log \left( \sum_{h=1}^{N_w} \pi_{wh} P(D=1 | w, h_w = h, c) \right)$$

**2.2** $\quad M_{whc} \in \{0,1\} := \mathbb{1}\{h_w = h\}$

$P(D, \Pi), \quad \log(P(D,\Pi)) = \mathcal{L}(D, \Pi)$

$\triangleright \ P(D, \Pi) = \underbrace{P(\Pi \mid D) P(D)}$

$\left| \ P(\Pi \mid D) = P(\Pi = 0 \mid D) P(\Pi = 0) + P(\Pi = 1 \mid D) P(\Pi = 1) \right.$

$\triangleright \quad \prod_{(w,c) \in D} \sum_{h=1}^{N_w} \Pi_{wh} P_i(D = 1 \mid w, h_w = h, \text{hing } c) \cdot P(\Pi \mid D)$

$$= \prod_{(w,c) \in D} \prod_{h=1}^{N_w} \Pi_{wh} \, P(D = 1 \mid w, h_w = h, c)^{M_{wac}}$$

$\mathcal{L}(D, \Pi) =$

$\triangleright \ \log(P(D, \Pi)) = \log \prod_{(w,c)} \prod_{h=1}^{N_w} \Pi_{wh} \, P(D = 1 \mid w, h_w = h, c)^{M_{wac}}$

$$= \sum_{(w,c)} \sum_{h=1}^{N_w} M_{wac} \log(P(D = 1 \mid w, h_w = h, c)) + \log(\Pi_{wa})$$

☑

**2.3)** E-step.  $Q(\Theta) := \mathbb{E}_{M|D} \, \mathcal{L}(D,M)$

$$\mathbb{E}_{M|D} \, M_{whc} = \mu_{whc}$$

▷ $Q(\Theta) = \mathbb{E}_{M|D} \, \mathcal{L}(D,M) = \mathbb{E}\left[ \sum_{(w,c)\in D} \sum_{h=1}^{N_w} M_{whc} \log\left(P(D=1 | \dots)\right) + \log(\pi_{w...} \right.$

$$= \sum_{(w,c)\in D} \sum_{h=1}^{N_w} \mathbb{E}_{M|D} \, M_{whc} \log\left(P(D=1 | w, h_w=h, c) + \log(\pi_{whc})\right)$$

$$= \sum_{(w,c)\in D} \sum_{h=1}^{N_w} \mu_{whc} \log\left(P(D=1 | w, h_w=h, c) + \log\left(\pi_{wh}\right)\right)$$

parameters $\Theta$ :  $\pi_{wh}$  for $h=1$ to $N_w$

~~and $\pi_{whc}$~~

~~Correction.~~

     parameters $\Theta$ are $\left\{ \pi_{wh'} \right\}_{h'=1}^{N_w}$

## 2.4

write $\mu_{wkc}$ as $\sim \{\pi_{wkc}\}_{k'=1}^{N_w}$

and $P(\underset{=1}{D} \mid \omega, h_w = k', c)$

① $\mu_{wkc} = \mathbb{E}_{M\mid D} M_{wkc} = 1 \cdot P(M = \underset{1}{\cancel{0}}(D)$

$\qquad \qquad \qquad \qquad \underbrace{+ 0 \cdot P(A = 0 \mid D)}$

$M=1$ if $h_w = k$ (2.2)

$\qquad \qquad = P(h_w = k \mid D, \omega) = \dfrac{P(D \mid \omega, h_w = k) \underbrace{P(h_w = k \mid \omega)}_{\pi_{wk}}}{P(D \mid \omega)}$

$\qquad = \dfrac{\pi_{wk} \, P(D \mid \omega, h_w = k)}{P(D \mid \omega)}$

for $D = 1$, we have : $\mu_{wkc} = \dfrac{\pi_{wk} \, P(D \mid \omega, h_w = k)}{\sum\limits_{k=1}^{N_w} P(D = 1 \mid \omega, h_w = k, c)}$

4

## 2-5 ) $\Pi$-step

D  $\Pi^{*}_{wc} := \arg\max_{\Pi_{wc}} Q(\Theta)$      parameter

$$Q(\Theta) \overset{2.3}{=} \sum_{(w,c)\in D} \sum_{a=1}^{N_w} \left( \underbrace{\mu_{wac}\left(\log\left(P(D=1\,|\,w, h_w = h, c)\right)\right) + \log(\Pi_{wac}}_{\text{no } \Pi_{wac}} \right)$$

$$\Rightarrow \arg\max_{\Pi_{wac}} Q(\Theta) = \arg\max_{\Pi_{wac}} \left( \sum_{(w,c)} \sum^{N_w} \mu_{wac} \log(\Pi_{wac}) \right)$$

lagrangian constraint :  $\sum^{N_w} \Pi_{wac} = 1$

$$\Rightarrow \mathcal{L}(\Pi_{wac}; \lambda) = \sum_{(w,c)} \sum^{N_w} \mu_{wac} \log(\Pi_{wac}) + \lambda\left(1 - \sum_{h=1}^{N_w} \Pi_{wac}\right)$$

now we derive the lagrangian

$$\frac{\partial \mathcal{L}}{\partial \Pi_{wac}} = 0 \iff \text{and with a few calculations,}$$
we have

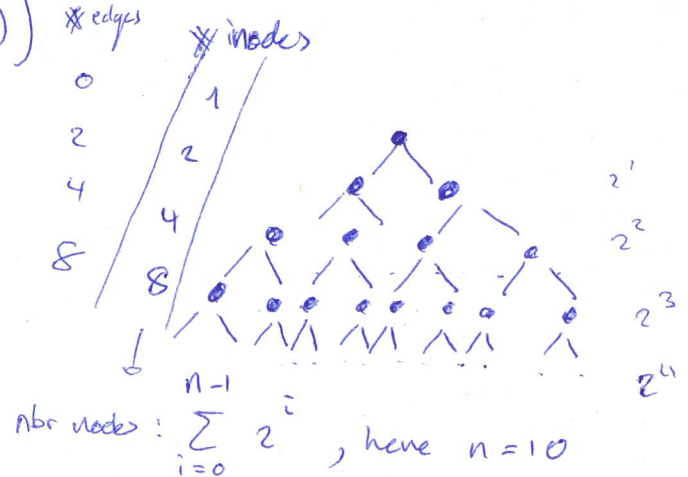$$\Pi^{*}_{wac} = C \cdot \sum_{(w,c)} \mu_{wac} \quad \text{where } c \text{ is just a scaling cst}$$

Computational Semantics for Natural Language Processing | Alexis Tabin
16-821-803

# Assignment 1

## Hierarchical Softmax & Huffmann Coding

$D$, document

$|D|$, # of words

$V$, vocabulary

$|V|$, voc size

$p(w) = \dfrac{\#w}{|D|}$, word frequency

(# of times $w$ appears in $D$)

$$P_w(c) = \frac{\exp\{(s(w,c))\}}{\sum\limits_{c \in V} \exp\{(s(w,c))\}} \quad , \quad \mathcal{O}(|V|) \text{ for each pair } (w,c)$$

approximated by (depth of node $c$ in tree $\frac{\circ}{\bullet}$)

$$P_w(c) = \prod_{j=1}^{L(c)-1} P_w((c,j) \longrightarrow (c,j+1))$$

(node of depth $j$ on road to leaf $c$)

( prob of transition

## 3.1 binary tree, $|V| = 2^{10} = 1024$

To compute $P_w(c)$, we need to go all the way down to the leaf, the way is composed of 9 edges here.

General case :

$$O(\log|V|\frac{\circ}{\bullet} - 1)$$
$$= O(\log|V|)$$

| # edges | # nodes |
|---|---|
| 0 | 1 |
| 2 | 2 |
| 4 | 4 |
| 8 | 8 |
| ↓ | $n-1$ |



$2^1$
$2^2$
$2^3$
$2^4$

nbr nodes : $\sum\limits_{i=0}^{n-1} 2^i$ , here $n = 10$

$\Rightarrow$ | # inner nodes = 1023 |

Total # nodes = inner nodes + leaf nodes
$= 1023 + 1024 = 2047$

# of edge $= 2046$

$|V| = 4$

Intuitively, we have:



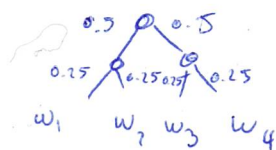$$\mathbb{E}\left[L(c)\right] = 3 \cdot 0.1 + 3 \cdot 0.1 + 2 \cdot 0.3 + 1 \cdot 0.5 = 1,7$$

$\cancel{\mathbb{E}\left[L(c)\right] \geq \mathbb{E}\left[-\log(p(w_i))\right]}$

$$L(w_i) \geq -\log(p(w_i))$$

no clue

- worst case: let take the last example.

  the worst case is this one:

  

  $$\mathbb{E}\left[L(c)\right] = \sum 4 \cdot (0.25 \cdot 2)$$

  $$= 2 = \log_2 4$$

  general case: $\mathbb{E}\left[L(c)\right] = O(\log|V|)$

Assignment 1

4. FastText Embeddings

4.1

a) 40'000

b) 2e6

4.2

a) i) $\dfrac{\partial l}{\partial u_c} = \dfrac{\partial}{\partial u_c}\left( -\log(\sigma(s(w_t, w_c))) - \cdots \right)$

doesn't depend on $u_c$

$= \dfrac{\partial}{\partial u_c} -\log\left( \sigma\left( \sum_g z_g^T u_c \right) \right)$

$= -\dfrac{1}{\sigma\left(\sum_g z_g^T u_c\right)} \dfrac{\partial}{\partial u_c} \sigma\left(\sum_g z_g^T u_c\right)$

$= -\dfrac{1}{\sigma\left(\sum_g z_g^T u_c\right)} \left(1 - \sigma\left(\sum_g z_g^T u_c\right)\right) \sigma\left(\sum_g z_g^T u_c\right) \sum_g z_g$

$= \left(\sigma\left(\sum_g z_g^T u_c\right) - 1\right) \sum_g z_g$

14

a) ii) $\dfrac{\partial \mathcal{L}}{\partial u_n} = \underbrace{\dfrac{\partial}{\partial u_n} - \log(\ldots)}_{=0} - \sum\limits_{n \in N_h} \log\left(\sigma\left(-s(w_t, w_n)\right)\right)$

$= \dfrac{\partial L}{\partial u_n} - \sum\limits_{n \in N_a} \log\left(\sigma\left(-s(w_t, w_n)\right)\right)$

Same as i) but with a negative sign

$\Rightarrow \dfrac{\partial L}{\partial u_n} = \left(1 - \sigma\left(-s(w_t, w_n)\right)\right) \sum\limits_{g} z_g$

b) $\dfrac{\partial \mathcal{L}}{\partial z_g} = \dfrac{\partial L}{\partial u_c} \dfrac{\partial L}{\partial u_n} = \left(1 - \sigma\left(-s(w_t, w_n)\right)\right) z_g$

$+ \sum\limits_{n \in N_a} \left(1 - \sigma\right)$     *sorry for the mess*

$= \sum\limits_{n \in N_a} \left(1 - \sigma\left(-s(w_t, w_n)\right)\right) u_n$

$\quad + \left(\sigma\left(s(w_t, w_c)\right) - 1\right) u_c$