

Assignment-2 | NLP

Q:-1 (a) :- given $\rightarrow \text{softmax}(x) = \frac{e^{x_i}}{\sum_j e^{x_j}}$ — (1)

R.H.S. = $\text{softmax}(x+c)$

$$= \frac{e^{(x_i+c)}}{\sum_j e^{(x_j+c)}} \quad (\text{from — (1)})$$

$$= \frac{e^c \cdot e^{x_i}}{\sum_j e^{x_j} \cdot e^c}$$

$$= \frac{e^c \cdot e^{x_i}}{e^c \sum_j e^{x_j}} \quad (\because c = \text{constant})$$

$$= \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$= \text{softmax}(x) = \text{L.H.S.}$$

Hence, it is proved that softmax is invariant to constant offset in the input.

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

$$\Rightarrow \sigma'(x) = \frac{\partial}{\partial x} \left(\frac{1}{1+e^{-x}} \right)$$

$$= \frac{(-1)}{(1+e^{-x})^2} \cdot \frac{\partial}{\partial x} (1+e^{-x})$$

$$= \left(\frac{-1}{(1+e^{-x})^2} \right) \cdot (0 - e^{-x})$$

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

$$= \frac{1}{e^x} \cdot \frac{e^{2x}}{(e^x+1)^2}$$

$$= \frac{e^x \cdot e^x}{e^x \cdot (e^x+1)^2}$$

$$= \left(\frac{e^x}{e^x+1} \right) \left(\frac{1}{e^x+1} \right)$$

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

$$= \left(\frac{1}{1+e^{2c}} \right) \left(1 - \left(\frac{e^{2c}+1}{e^{2c}+1} \right) \right)$$

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

$$= \boxed{\sigma(x)(1 - \sigma(x))} \quad (\because \text{from } \text{---} (1))$$

Hence, \square is proved!
given

Q-2(a) :- $\frac{\partial}{\partial v_c} J_{CE}(0, v_c, U) = ?$

$$J_{CE} = - \sum_i y_i \log(\hat{y}_i) \quad \text{--- (1)}$$

- It is given that y_i is one-hot encoded vector and value for \hat{y} is also given,

Hence

$$J_{CE} = - \log \left\{ \frac{\exp(u_o^T v_c)}{\sum_w \exp(u_w^T v_c)} \right\}$$

$$= - \left[\log(\exp(\psi_0^T \psi_c)) - \log \left\{ \sum_w \exp(\psi_w^T \psi_c) \right\} \right]$$

$$-y_0^T v_c + \log \sum_w \exp(y_w^T v_c) \quad \text{--- (2)}$$

\therefore From (2)

$$\frac{\partial J_{CE}}{\partial v_c} = \left(\frac{1}{\sum_w \exp(u_w^T v_c)} \right) \left(\exp(u_w^T v_c) \cdot u_w^T - u_w^T \right)$$

~~$$= \sum_w u_w^T \left\{ \frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} - 1 \right\}$$~~

$$= \sum_{j=1}^W u_j^T \left\{ \frac{\exp(u_j^T v_c)}{\sum_{i=1}^W \exp(u_i^T v_c)} - 1 \right\}$$

$$\boxed{\frac{\partial J_{CE}}{\partial v_c} = \sum_{j=1}^W u_j^T (P(u_j | v_c) - 1)} \quad (3)$$



Scanned with
CamScanner

(3) is the required solution

Q-2(b):- $\frac{\partial J_{CE}}{\partial u_w} = ?$

- Same as Q-2(a), We can derive,

$$J_{CE} = -u_w^T v_c + \log \sum_w \exp(u_w^T v_c) \quad \text{--- (1)}$$

→ Now,

$$\frac{\partial J_{CE}}{\partial u_w} = -v_c + \left\{ \frac{1}{\sum_w \exp(u_w^T v_c)} \right\} (\exp(u_w^T v_c)) (v_c)$$

$$\left[\frac{\partial J_{CE}}{\partial u_w} = v_c \left\{ \frac{\exp(u_w^T v_c)}{\sum_{j=1}^w \exp(u_j^T v_c)} - 1 \right\} \right] \quad \text{--- (2)}$$

Which is required derivation:

$$\left[\frac{\partial J_{CE}}{\partial u_w} = v_c \left\{ p(u_w | v_c) - 1 \right\} \right]$$



#1

Q:- 2(c)!

$$\frac{\partial}{\partial v_c} J_{\text{neg-sample}} = ?$$

#2

$$\frac{\partial}{\partial w} J_{\text{neg-sample}} = ?$$

⇒ Given,

$$J_{\text{neg-sample}} = -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c))$$

①

$$\begin{aligned} \Rightarrow \therefore \frac{\partial}{\partial v_c} J_{\text{neg-sample}} &= (-1) \left\{ \frac{\sigma(u_0^T v_c) \cdot (1 - \sigma(u_0^T v_c))}{\sigma(u_0^T v_c)} \right\} \cdot u_0^T \\ &\quad - \sum_{k=1}^K \left(\frac{1}{\sigma(-u_k^T v_c)} \right) \cdot (\sigma(-u_k^T v_c)) \cdot (1 - \sigma(-u_k^T v_c)) \cdot (-u_k^T) \end{aligned}$$

$$= \left[u_0^T \cdot (\sigma(u_0^T v_c) - 1) - \sum_{k=1}^K u_k^T (\sigma(-u_k^T v_c) - 1) \right] \quad \text{②}$$

Which is required solution #1

$$\Rightarrow \therefore \frac{\partial}{\partial w} J_{\text{neg-sample}} = \left(\frac{-1}{\sigma(u_0^T v_c)} \right) \left\{ \sigma(u_0^T v_c) \cdot (1 - \sigma(u_0^T v_c)) \right\} \cdot v_c^T$$

~~Which is required solution #2~~



Scanned with
CamScanner

∴ $0 \notin \{1, \dots, K\}$, hence 2nd-part derivative
w.r.t. $w = 0$

$$= \left[\mathbf{v}_c^T \left(\mathbf{y}^T \mathbf{v}_c \right) - 1 \right] \quad \text{--- (3)}$$

Which is required Solution #2

Q1-2(d):- For skip-gram, the cost for a context centered around c is:

$$J_{\text{skip-gram}} = \sum_{-w \leq j \leq w, j \neq 0} F(\mathbf{w}_{c+j}, \mathbf{v}_c) \quad \text{--- (1)}$$

\therefore derivatives would be,

$$\frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{w}} = \left[\sum_{-w \leq j \leq w, j \neq 0} \frac{\partial F}{\partial \mathbf{w}_{c+j}} \right] \quad \text{--- (2)}$$

$$\frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{v}_c} = \left[\sum_{-w \leq j \leq w, j \neq 0} \frac{\partial F}{\partial \mathbf{v}_c} \right] \quad \text{--- (3)}$$



required!

Nearest Neighbors:

Nearest Words to the word **great**:

great
juan
ringside
toast
kenneth

Nearest Words to the word **brilliant**:

stephen
brilliant
liberating
undermines
ugly-duckling

Nearest Words to the word **female**:

female
illuminates
dancing
fear-inducing
skateboarding

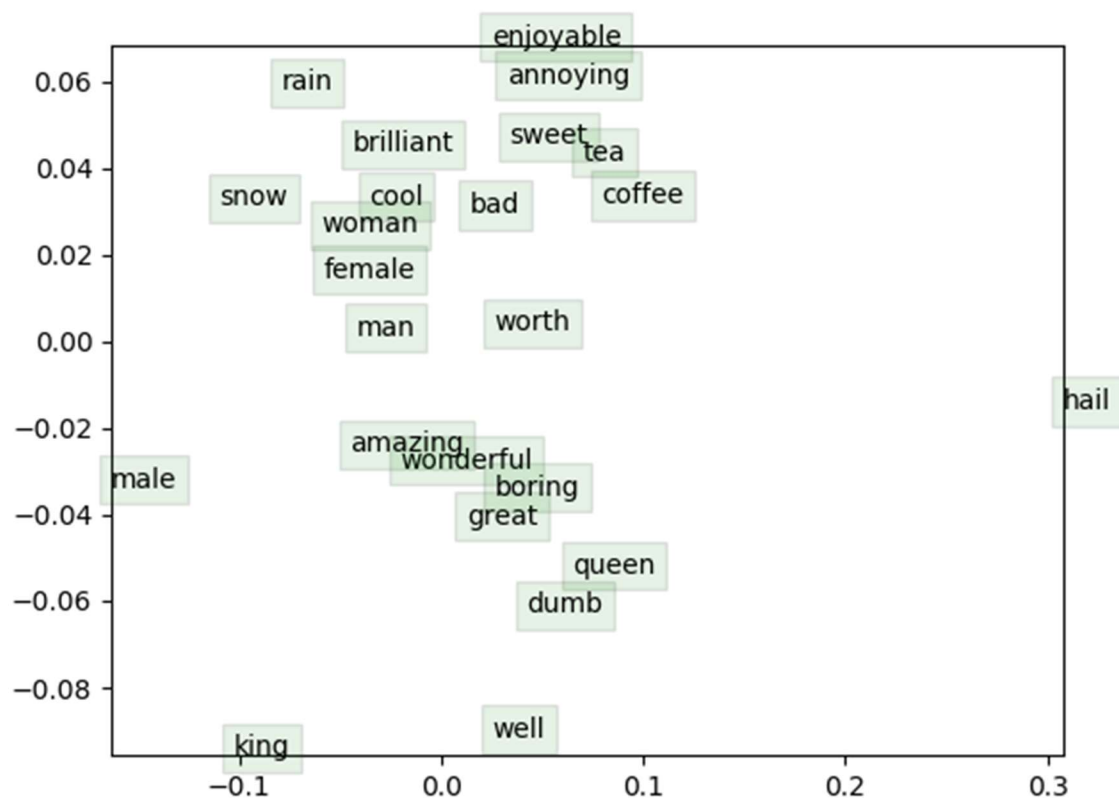
Nearest Words to the word **king**:

instinct
king
xtc
composition
lector

Nearest Words to the word **snow**:

snow
presents
savvy
compensate
squandering

Vectors plot:



Analysis:

- By comparing results from kNN and word vectors plot, we can see the drastic difference in word similarities resulted from both.
- By looking at 5 nearest words for the given word, we can see that kNN is only able to give same word as the closest similar word to given word itself. Other 4 words are not even close to being similar. The reason is the way kNN calculates similarity. In this case kNN only depends on cosine distance to calculate similarity, which doesn't take context of the word into account. It only tries to calculate distance based on simple dot product.
- On other hand, word vectors are designed to capture the context of the word given its neighboring words. As we can see in the plot, similar words are grouped together. For example: (Woman, female), (amazing, wonderful, great), (tea, coffee) etc.
- Only advantage of kNN is its speed. Because, word vectors heavily rely on many calculations to capture proper context of the word, word vectors are slow to calculate.