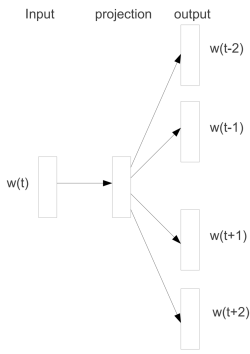


SKIP-GRAM REFRESHER

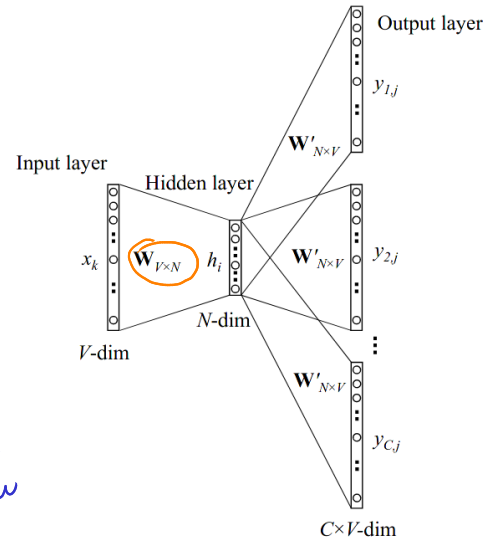


This is the original skip-gram model architecture. It uses a NN with one hidden layer. It receives a one-hot vector as input (center word vector) and outputs a window-size output representing the probabilities of the surrounding context words of the input center word.

The output probabilities are subproducts of the trained model. What we really are interested on are the two matrices learned during training. These matrices are W_{input} and W_{output}

THE SKIP-GRAM NN ARCHITECTURE

- We decide to represent each word with a vector of N dimensions
- We have a corpus of text containing T words and a vocabulary of length V
- Our corpus of text is represented in two matrices: W_{input} , W_{output} . Each word appears in both matrices. After training we use the W_{input} as our word embeddings
- W_{input} represent each word when it acts as a center word
- W_{output} represents each word when it acts as a context word
- Each vector of W_{input} is called v_c
- Each vector of W_{output} is called u_o
- Each input we feed into the NN is called x_k where k is the position that word occupies in the vocabulary. This input is a vector
- During training we select a center word $v_c \rightarrow x_k$ and a context window of size C
- The model's output is a set of probabilities $y_{c,j}$; one per each context word. This probability indicates how likely that word acts as a context word for the input center word.



THE MODEL IN DETAIL

"The man who passes the sentence should swing the sword."
- Ned Stark
 $N=3$ $T=10$ $V=8$ $C=1$

$\Theta = [W_{input} \ W_{output}]$
During backprop. Θ is updated for each v_c

This is a one-hot vector since for the given $v_c = \text{'passes'}$, we know that 'the' acts as a context word

How likely the words in W_{output} act as context words given v_c

$y_{true} (V\text{-dim})$ for c_1

0
0
0
0
0
1
0

$y_{pred} - y_{true} = e_c$

Prediction error of two "c" context words of the input sentence X

$y_{pred} (V\text{-dim})$ for x

0.12
0.21
0.10
0.06
0.13
0.09

$y_{pred} - y_{true} = e_c$

Prediction error of two "c" context words of the input sentence X

$y_{true} (V\text{-dim})$ for c_2

0
0
0
0
0
0
1

$y_{pred} (V\text{-dim})$ for x

0.12
0.21
0.10
0.06
0.13
-0.91

Prediction Error - a.k.a function loss

A given center word \rightarrow man passes sentence should

A context window of 2 words \rightarrow the who

Input Layer one-hot encoded vector

Word-Embedding matrix - a.k.a "Lookup table"

Hidden (Projection) Layer for center word (passes)

Word-Embedding matrix for context words (the, who)

Softmax Output Layer of range [0, 1] Sum = 1

$\hat{y}_{who} = P(\text{who} | \text{passes}) = \text{softmax}(u_{who} \cdot v_{passes}) = 0.09$

This is the way. $F(x) \rightarrow J(v_o, v_c, W_{output})$

When we feed the NN with a v_c we are also given it the two context words for that word "c". This context words have a y one-hot true vector for each one. The model outputs \hat{y} per each context word. The NN is learning to predict the context words, as a byproduct generates the embeddings

This is used during backpropagation!

CS224N - ASSIGNMENT #2

a) Demonstrate: For a training pass, given a v_c we have one \hat{y} containing all the probabilities of the vocab words acting as context words!

$$-y_o \cdot \log \hat{y}_o = -\sum_{w \in V \text{ of } Y} y_w \log(\hat{y}_w) = -\log(\hat{y}_o) = -\log[P(o|c)] = J(u_o, v_c, u)$$

Word p = "passes"
Context word o = "who"

just an index of vector y

not a dot product!

Notice we also have one y for each word in the context window

this happens at index $w=0$ of y

$$-\cancel{0.12} \log \cancel{0.12} + \cancel{0.21} \log \cancel{0.21} + \dots + 1 \cdot \log 0.09 = -\log(0.09)$$

b) Compute the partial derivative of $J_{\text{naive-sfsoftmax}}(v_c, o, u)$ w.r.t. v_c

objective function
outside matrix (size $N \times \text{Woutput}$)
current outside word (in the context window)
current center word being trained

For a single pair of words the objective function J is:

$$J_{ns}(v_c, o, u) = -\log P(o=o | C=c) = -\log P(o|c) = -\log \text{softmax}(u_o, v_c) = -\log \frac{e^{u_o^T v_c}}{\sum_{x \in V} e^{u_x^T v_c}} = -\log \hat{y}_o$$

$$J_{ns}(v_c, o, u) = -\sum_{w \in V} y_w \log(\hat{y}_w) \text{ when } w=o, \text{ meaning } o \text{ is acting as context word of } c$$

$$J_{ns}(v_c, o, u) = -\sum_{w \in V} y_w \log \left[\frac{e^{u_o^T v_c}}{\sum_{x \in V} e^{u_x^T v_c}} \right] = -\sum_{w \in V} y_w [\log e^{u_o^T v_c} - \log \sum_{x \in V} e^{u_x^T v_c}]$$

$P(o|c)$

We know that when $w=0$ $y_w=1$ and 0 for other indices given the center word c and context word o we are currently analyzing

$$J_{ns}(v_c, o, u) = -\sum_{w \in V} y_w [u_o^T v_c - \log \sum_{x \in V} e^{u_x^T v_c}]$$

$$J_{ns}(v_c, o, u) = -y_o [u_o^T v_c - \log \sum_{x \in V} e^{u_x^T v_c}]$$

-1

$$\frac{\partial J_{ns}}{\partial v_c} = -\frac{\partial}{\partial v_c} u_o^T v_c + \frac{\partial}{\partial v_c} \log \sum_{x \in V} e^{u_x^T v_c}$$

$$\downarrow \text{Shape is } 1 \times N$$

$$= -u_o + \log e \times \frac{1}{\sum_{x \in V} e^{u_x^T v_c}} \times \frac{\partial}{\partial v_c} \sum_{x \in V} e^{u_x^T v_c}$$

$$= -u_o + \frac{\sum_{x \in V} \frac{\partial}{\partial v_c} e^{u_x^T v_c}}{\sum_{x \in V} e^{u_x^T v_c}} = -u_o + \sum_{x \in V} \frac{e^{u_x^T v_c} u_x}{\sum_{j \in V} e^{u_j^T v_c}} = -u_o + \sum_{x \in V} P(x|c) u_x = -u_o + \sum_{x \in V} \hat{y}_x u_x$$

current context word o has an y_o

$$= \sum_{x \in V} \hat{y}_x u_x - u_o = \bar{U} \hat{y} - \bar{U} y_o = \bar{U}^T (\hat{y} - y_o)$$

Dot product!

$$\hat{y} \cdot u_1 + \hat{y}_2 \cdot u_2 + \dots + \hat{y}_v \cdot u_v = \bar{U}^T \cdot \hat{y}$$

vectors when processing word c

$$\bar{U}^T \cdot \hat{y} = u_1^T \hat{y}_1 + u_2^T \hat{y}_2 + \dots = \bar{U}^T (\hat{y} - y_o)$$

$$\bar{U}^T = [u_1^T \ u_2^T \ \dots \ u_v^T] \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_v \end{bmatrix}$$

Not sure about moving this sum out.

"inquiries" start here

Shape is vector of $1 \times N \Rightarrow$ follows the shape convention

ALTERNATIVE CALCULATION 3

$$\frac{\partial}{\partial v_c} - \log \frac{e^{u_0^T v_c}}{\sum_w e^{u_w^T v_c}} = - \frac{\partial}{\partial v_c} \left[\log e^{u_0^T v_c} - \log \sum_{w \in V} e^{u_w^T v_c} \right]$$

composite function

$$= \underbrace{- \frac{\partial}{\partial v_c} \log e^{u_0^T v_c}}_{\text{?}} + \frac{\partial}{\partial v_c} \log \sum_{w \in V} e^{u_w^T v_c} \quad \textcircled{2} \text{ See next page}$$

?
 $\log_e e = 1$
 $\log e = 0.43$
 probably \log in
 this context is \ln

Recall: (A)

$$\text{if } F(x) = \frac{f(x)}{g(x)} \Rightarrow -\log F(x) = -\log f(x) + \log g(x)$$

Recall: if s is a function of x ; and
 (B) $y = \log_b s$; then $\frac{dy}{dx} = (\log_b e) \frac{ds}{s}$

Recall: If u is a function of x ; $\frac{d}{dx} e^u = e^u \frac{du}{dx}$

→ This is chain rule!

$$u_0 = [u_{01} \ u_{02} \ \dots \ u_{0n}]$$

$$u_0^T v_c = u_{01}v_{c1} + u_{02}v_{c2} + \dots + u_{0n}v_{cn}$$

$$\frac{\partial}{\partial v_{c1}} u_0^T v_c = u_{01} + 0 + 0 + \dots + 0$$

$$\frac{\partial}{\partial v_{c2}} u_0^T v_c = 0 + u_{02} + 0 + \dots + 0$$

$$\frac{\partial}{\partial v_{c3}} u_0^T v_c = 0 + 0 + u_{03} + 0 + \dots + 0$$

$$\Rightarrow \frac{\partial}{\partial v_c} u_0^T v_c = u_0$$

$$= -\log e \times \frac{1}{e^{u_0^T v_c}} \times \frac{\partial}{\partial v_c} e^{u_0^T v_c}$$

$$= -\frac{1}{e^{u_0^T v_c}} \times e^{u_0^T v_c} \times \frac{\partial u_0^T v_c}{\partial v_c}$$

$$= -u_0$$

Partial
 derivatives
 of v_c

$$\textcircled{2} \frac{\partial}{\partial v_c} \log \sum_{w \in V} e^{u_w^T v_c} = \log e \times \frac{1}{\sum_{w \in V} e^{u_w^T v_c}} \times \frac{\partial}{\partial v_c} \sum_{w \in V} e^{u_w^T v_c}$$

it is needed to change variables!

composite A function of v_c

Apply (B)

$$\log e \times \frac{1}{\sum_{w \in V} e^{u_w^T v_c}} \times \sum_{x \in V} \frac{\partial}{\partial v_c} e^{u_x^T v_c}$$

$$\log e \times \frac{\sum_{x \in V} e^{u_x^T v_c} u_x}{\sum_{w \in V} e^{u_w^T v_c}}$$

why?
 Not sure

I guess just for the softmax

Some people uses \log for \ln ; it seems this is the case here

$$\sum_{x \in V} \frac{e^{u_x^T v_c} u_x}{\sum_{w \in V} e^{u_w^T v_c}}$$

observed representation of the context word

$$\frac{\partial J(v_c, o, u)}{\partial v_c} = -u_0 + \sum_{x \in V} P(x|c) u_x = -u_0 + \sum_{x \in V} \hat{y}_x u_x = -u_0 + U^T \hat{y}$$

one-hot vector!

The expected context word \Rightarrow

what our model thinks what the context should look like

$$\frac{\partial J(v_c, o, u)}{\partial v_c} = -y + U^T \hat{y} = -U^T y + U^T \hat{y} = U^T (\hat{y} - y)$$

A slope in the multidimensional space from the

slope meaning

This slope of this is \mathbb{R}^n , the same of v_c

$$U = \begin{bmatrix} u_0 & u_1 & \dots & u_n \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \end{bmatrix} \quad d\text{-dim.}$$

$$y = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

for a given context word c , y is a hot vector with 1 in the x position where that word is a true outside word

one of the context words of the given center word c

a row in matrix U (U_{output})

This means when $u_w = u_o \Rightarrow$

(c) (5 points) Compute the partial derivatives of $J_{\text{naive-softmax}}(v_c, o, U)$ with respect to each of the 'outside' word vectors, u_w 's. There will be two cases: when $w = o$, the true 'outside' word vector, and $w \neq o$, for all other words. Please write your answer in terms of y , \hat{y} , and v_c . In this subpart, you may use specific elements within these terms as well, such as (y_1, y_2, \dots) .

$$\frac{\partial J_{\text{naive-softmax}}(v_c, o, U)}{\partial u_w \rightarrow \text{shape is } 1 \times N} = \frac{\partial}{\partial u_w} -\log \frac{e^{u_o^T v_c}}{\sum_{w \in V} e^{u_w^T v_c}} = \frac{\partial}{\partial u_w} -\left[\log e^{u_o^T v_c} - \log \sum_{w \in V} e^{u_w^T v_c} \right]$$

When $u_w \neq 0$:

For the given center word "c", any word in $U(\text{output})$ that is not acting as context word

$$= -\frac{\partial}{\partial u_w} \log e^{u_o^T v_c} + \frac{\partial}{\partial u_w} \log \sum_{w \in V} e^{u_w^T v_c}$$

$$= 0 + \log e \times \frac{1}{\sum e^{u_w^T v_c}} \times \frac{\partial \sum_{w \in V} e^{u_w^T v_c}}{\partial u_w} \quad \text{Chain rule!}$$

$$= \frac{\sum_{w \in V} e^{u_w^T v_c} v_c}{\sum_{w \in V} e^{u_w^T v_c}} \quad \text{As } u_w \text{ shape} = 1 \times N \Rightarrow \text{because of shape convention } \frac{\partial J}{\partial u_w} \text{ should also has a shape equals to } 1 \times N.$$

Why does the sum disappears?

\hookrightarrow virginios "

Perhaps rank we have a fixed word w can c

$$= \sum_{w \in V} P(w|c) v_c$$

$$= \hat{y}_{w \neq o} v_c = \hat{y}_{w \neq o} v_c - 0 \approx \hat{y}_{w \neq o} v_c - y_w v_c = v_c (\hat{y}_{w \neq o} - y_w) \quad \text{as } w \text{ is not a context word, the true vector is zero}$$

When $u_w = 0$: for the outside matrix vectors corresponding to the context words o of v_c in the context window

$$\frac{\partial J_{\text{naive-softmax}}(v_c, o, U)}{\partial u_o} = -\frac{\partial}{\partial u_o} \log e^{u_o^T v_c} + \frac{\partial}{\partial u_o} \log \sum_{o \in V} e^{u_o^T v_c}$$

$$= -\log e \times \frac{1}{e^{u_o^T v_c}} \times \frac{\partial}{\partial u_o} e^{u_o^T v_c} + \log e \times \frac{1}{\sum_{o \in V} e^{u_o^T v_c}} \times \frac{\partial \sum_{o \in V} e^{u_o^T v_c}}{\partial u_o}$$

$$= -\frac{e^{u_o^T v_c} v_c}{e^{u_o^T v_c}} + \frac{\sum_{o \in V} e^{u_o^T v_c} v_c}{\sum_{w \in V} e^{u_w^T v_c}}$$

$$= -v_c + \sum_{o \in V} P(o|c) v_c \quad \text{The resulting value is a vector } 1 \times N, \text{ same shape of } u_o$$

$$= -v_c + \sum_{o \in V} \hat{y}_o v_c \quad \text{This is just an index! } \equiv \hat{y}_{u_o} [0]$$

$$= v_c (\hat{y}_o - 1) = v_c (\hat{y}_o - y_o)$$

$\frac{\partial J}{\partial U}$ is a matrix!

- (d) (1 point) Compute the partial derivative of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to U . Please write your answer in terms of $\frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial u_1}, \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial u_2}, \dots, \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial u_{|V_{\text{ocab}}|}}$. The solution should be one or two lines long.

$$\frac{\partial}{\partial U} J(\mathbf{v}_c, o, U) = \left[\frac{\partial J}{\partial u_1} \quad \frac{\partial J}{\partial u_2} \quad \dots \quad \frac{\partial J}{\partial u_o} \quad \dots \quad \frac{\partial J}{\partial u_v} \right] = \begin{bmatrix} \hat{y}_u v_c & \hat{y}_u v_c & \dots & (\hat{y}_o - 1) v_c & \dots & \hat{y}_u v_c \end{bmatrix}$$

→ The result should be a matrix!

- (e) (3 Points) The sigmoid function is given by Equation 4:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (4)$$

Please compute the derivative of $\sigma(x)$ with respect to x , where x is a scalar. Hint: you may want to write your answer in terms of $\sigma(x)$.

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \rightarrow \text{I'm going to choose this}$$

$$\frac{\partial}{\partial x} \sigma(x) = \frac{(e^x + 1)e^x - e^x(\partial/\partial x e^x + 1)}{(e^x + 1)^2} \quad \left\{ \text{Quotient rule} \right.$$

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

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

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

$$= \sigma(x) \frac{e^x + 1}{e^x + 1} - \sigma(x)^2 = \sigma(x) \left[\frac{e^x + 1}{e^x + 1} - \sigma(x) \right] = \sigma(x) [1 - \sigma(x)] //$$

- (f) (4 points) Now we shall consider the Negative Sampling loss, which is an alternative to the Naive Softmax loss. Assume that K negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as w_1, w_2, \dots, w_K and their outside vectors as $\mathbf{u}_1, \dots, \mathbf{u}_K$. For this question, assume that the K negative samples are distinct. In other words, $i \neq j$ implies $w_i \neq w_j$ for $i, j \in \{1, \dots, K\}$. Note that $o \notin \{w_1, \dots, w_K\}$. For a center word c and an outside word o , the negative sampling loss function is given by:

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \quad (5)$$

for a sample w_1, \dots, w_K , where $\sigma(\cdot)$ is the sigmoid function.⁴

Please repeat parts (b) and (c), computing the partial derivatives of $J_{\text{neg-sample}}$ with respect to \mathbf{v}_c , with respect to \mathbf{u}_o , and with respect to a negative sample \mathbf{u}_k . Please write your answers in terms of the vectors \mathbf{u}_o , \mathbf{v}_c , and \mathbf{u}_k , where $k \in [1, K]$. After you've done this, describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss. Note, you should be able to use your solution to part (e) to help compute the necessary gradients here.

Current
→ outside word in context of v_c
given center word
→ outside Matrix

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, U) = -\log[\sigma(\mathbf{u}_o^\top \mathbf{v}_c)] - \sum_{k=1}^K \log[\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)]$$

→ K negative samples (words). All are distinct
→ outside vectors; $o \notin \{w_1, \dots, w_K\}$
 $u_k \neq u_o$

Note: in numpy $\text{np.dot}(-u_k, v_c) \equiv \text{np.dot}(-u_k^T, v_c)$

$$\begin{aligned}
 (1) \quad \frac{\partial J}{\partial v_c} &= - \left\{ \log e \cdot \frac{1}{\sigma(u_0^T v_c)} \times \frac{\partial \sigma(u_0^T v_c)}{\partial v_c} \right\} - \sum_{k \in K} \left\{ \log e \cdot \frac{1}{\sigma(-u_k^T v_c)} \times \frac{\partial \sigma(-u_k^T v_c)}{\partial v_c} \right\} \\
 &= - \left\{ \log e \cdot \frac{1}{\sigma(u_0^T v_c)} \times \frac{\partial \sigma(u_0^T v_c)}{\partial v_c} \times \frac{\partial u_0^T v_c}{\partial v_c} \right\} - \sum_{k \in K} \left\{ \log e \cdot \frac{1}{\sigma(-u_k^T v_c)} \times \frac{\partial \sigma(-u_k^T v_c)}{\partial v_c} \times \frac{\partial (-u_k^T v_c)}{\partial v_c} \right\} \\
 &= - \left\{ \log e \cdot \frac{1}{\sigma(u_0^T v_c)} \times \cancel{\sigma(u_0^T v_c)} [1 - \sigma(u_0^T v_c)] \times u_0 \right\} - \sum_{k \in K} \left\{ \log e \cdot \frac{1}{\sigma(-u_k^T v_c)} \times \cancel{\sigma(-u_k^T v_c)} [1 - \sigma(-u_k^T v_c)] \times -u_k \right\} \\
 \frac{\partial J}{\partial v_c} &= -u_0 [1 - \sigma(u_0^T v_c)] + \sum_{k \in K} u_k [1 - \sigma(-u_k^T v_c)] = \boxed{u_0 [\sigma(u_0^T v_c) - 1] - \sum_{k \in K} u_k [\sigma(-u_k^T v_c) - 1]}
 \end{aligned}$$

$$\begin{aligned}
 (2) \quad \frac{\partial J}{\partial u_0} &= - \left\{ \log e \cdot \frac{1}{\sigma(u_0^T v_c)} \times \frac{\partial \sigma(u_0^T v_c)}{\partial u_0} \right\} - \sum_{k \in K} \frac{\partial}{\partial u_0} \log [\sigma(-u_k^T v_c)] \\
 &= - \left\{ \frac{1}{\sigma(u_0^T v_c)} \times \frac{\partial \sigma(u_0^T v_c)}{\partial u_0} \times \frac{\partial u_0^T v_c}{\partial u_0} \right\} - 0 \\
 &= - \left\{ \frac{1}{\sigma(u_0^T v_c)} \times \cancel{\sigma(u_0^T v_c)} [1 - \sigma(u_0^T v_c)] \times v_c \right\} \\
 &= -v_c [1 - \sigma(u_0^T v_c)] = \boxed{v_c [\sigma(u_0^T v_c) - 1]}
 \end{aligned}$$

$$\begin{aligned}
 (3) \quad \frac{\partial J}{\partial u_k} &= - \left\{ 0 \right\} - \sum \frac{\partial}{\partial u_k} \log [\sigma(-u_k^T v_c)] \\
 &= - \sum \log e \cdot \frac{1}{\sigma(-u_k^T v_c)} \times \frac{\partial \sigma(-u_k^T v_c)}{\partial u_k} \\
 &= - \sum \frac{1}{\sigma(-u_k^T v_c)} \times \cancel{\sigma(-u_k^T v_c)} [1 - \sigma(-u_k^T v_c)] \times -v_c \\
 &= - \sum_{k \in K} -v_c [1 - \sigma(-u_k^T v_c)] = - \sum_{k \in K} v_c [\sigma(-u_k^T v_c) - 1] \\
 &\quad \text{OR } \rightarrow -v_c [\sigma(-u_k^T v_c) - 1], \forall k \in [1, K] \\
 &\quad \rightarrow \text{This notation is useful when inside a for loop of the } k\text{-samples.}
 \end{aligned}$$

Why this loss function is much more efficient to compute than the naive soft-max loss?

\Rightarrow With the naive soft-max we need to normalize the probability by multiplying all the words vectors with the center word. Here we do it just K times.

h) Recall the skip-gram version of word2vec, the total loss for the context window is:

$$J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, u) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(v_c, w_{t+j}, u) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} -\log P(w_{t+j} | w_t)$$

outside word

we can attach here naive soft-max or neg-sampling

Arbitrary loss term for the center word $c = w_t$

$$1) \frac{\partial}{\partial v_c} J_{\text{SG}}(v_c, w_{t+j}, u) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial v_c} J(v_c, w_{t+j}, u)$$

$$2) \frac{\partial}{\partial v_w} J_{\text{SG}}(v_c, w_{t+j}, u) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial v_{w+c}} J(v_c, w_{t+j}, u) = 0; J \text{ is not a function of } v_w$$

Probably they refer to:

$$\frac{\partial}{\partial u_{w+c}} J_{\text{SG}}(v_c, w_{t+j}, u) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial u_w} J(v_c, w_{t+j}, u)$$

$$3) \frac{\partial}{\partial u} J_{\text{SG}}(v_c, w_{t+j}, u) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial u} J(v_c, w_{t+j}, u)$$