

# Stanford CS224n Assignment 1.

## 1. Softmax.

a) Q. proof that  $\forall x \in \mathbb{R}^n, \forall c \in \mathbb{R} : \text{softmax}(x+c) = \text{softmax}(x)$

A.  $\text{softmax}(x+c) = \text{softmax}(\langle x_1+c; x_2+c; \dots; x_i+c; \dots x_n+c \rangle)$

let  $\text{softmax}(x+c)_i$  be the  $i$ th component of  $\text{softmax}(x+c)$

$$\text{softmax}(x+c)_i = \frac{e^{x_i+c}}{\sum_{j=1}^n e^{x_j+c}}$$

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

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

(can simplify because  $e^c$  is never  $= 0 \forall c$ )

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

$$= \text{softmax}(x)_i$$

$$\Rightarrow \boxed{\text{softmax}(x+c) = \text{softmax}(x)}$$

## 2) Neural Network Basics

a)  $\sigma(x) = \frac{1}{1 + \exp(-x)}$

$$\Rightarrow \frac{d\sigma}{dx}(x) = \frac{0(1 + e^{-x}) - 1 \cdot (-e^{-x})}{(1 + e^{-x})^2}$$

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

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

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

$$= \sigma(x) - \sigma^2(x)$$

$$\frac{d\sigma}{dx}(x) = \sigma(x)(1 - \sigma(x))$$

b)  $\nabla_{\theta} CE(y, \hat{y}) = \nabla_{\theta} \left( -\sum_i y_i \log(\hat{y}_i) \right)$   
 $= \nabla_{\theta} (-\log(\hat{y}_k))$  | Where  $k$  is the correct class

$$= -\frac{1}{\hat{y}_k} \cdot \nabla_{\theta} \hat{y}_k \quad (1)$$

$$\nabla_{\theta} \hat{y}_k = \nabla_{\theta} \frac{e^{\theta_k}}{\sum_i e^{\theta_i}} = \frac{e^{\theta_k} \cdot \sum_i e^{\theta_i} - e^{\theta_k} \cdot e^{\theta_k}}{\sum_i e^{\theta_i}^2}$$

$$= \frac{\left( \frac{1}{\sum_i e^{\theta_i}} \cdot \hat{y}_k \right) - \left( \hat{y}_k \right)^2}{\left( \frac{1}{\sum_i e^{\theta_i}} \right)^2}$$

$\nabla_{\theta} \hat{y}_k = ?$

let's distinguish two case

(a)  $\frac{\partial \hat{y}_k}{\partial \theta_k} = \frac{\partial}{\partial \theta_k} \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} = \frac{e^{\theta_k} \cdot \sum_j e^{\theta_j} - e^{\theta_k} \cdot e^{\theta_k}}{\sum_j e^{\theta_j}^2} = \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} \cdot \left( 1 - \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} \right)$

$= \hat{y}_k (1 - \hat{y}_k)$  (2)

2)

$$b) \text{CE}(y, \hat{y}) = -\sum_j y_j \log(\hat{y}_j) \quad \hat{y}_j = \text{softmax}(z)_j$$

$$\begin{aligned} \nabla_z \text{CE}(y, \hat{y}) &= \nabla_z \left( -\sum_j y_j \log(\hat{y}_j) \right) \\ &= -\nabla_z \log(\hat{y}_k) \quad \text{where } k \text{ is the correct label} \\ &= -\nabla_z \log\left(\frac{e^{z_k}}{\sum_j e^{z_j}}\right) \\ &= -\nabla_z (\log(e^{z_k})) + \nabla_z \log\left(\sum_j e^{z_j}\right) \\ &= -\nabla_z e^{z_k} \\ &= -\nabla_z (z_k) + \nabla_s \log(s) \otimes \nabla_z (s) \quad \left| s = \sum_j e^{z_j} \right. \end{aligned}$$

- let's simplify the first part (1)

$$\nabla_z (z_k) = \begin{bmatrix} \frac{\partial z_k}{\partial z_1} \\ \vdots \\ \frac{\partial z_k}{\partial z_k} \\ \vdots \\ \frac{\partial z_k}{\partial z_n} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 1, \text{ row } = k \\ \vdots \\ 0 \end{bmatrix} = y \text{ (one-hot)} \quad (3)$$

- now let's handle the second part (2)

$$* \nabla_s \log(s) = \frac{1}{s} = \frac{1}{\sum_j e^{z_j}} \quad (4)$$

$$* \nabla_z (s) = \nabla_z \left( \sum_j e^{z_j} \right) =$$

$$\begin{bmatrix} \frac{\partial}{\partial z_1} (e^{z_1} + e^{z_2} + \dots + e^{z_n}) \\ \frac{\partial}{\partial z_2} (e^{z_1} + e^{z_2} + e^{z_3} + \dots + e^{z_n}) \\ \vdots \\ \frac{\partial}{\partial z_n} (e^{z_1} + e^{z_2} + \dots + e^{z_n}) \end{bmatrix}$$

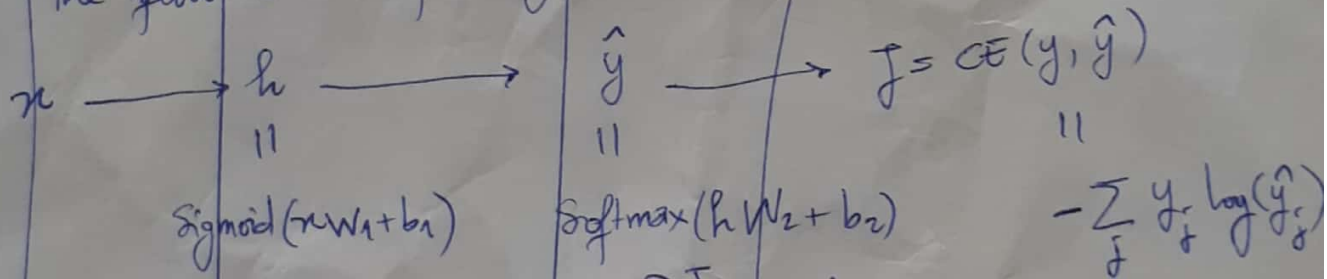
$$= \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_n} \end{bmatrix} = \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_n} \end{bmatrix}$$

$$\begin{aligned}
 (3), (4), (5) &\Rightarrow \nabla_z CE(y, \hat{y}) \\
 &= -y + \left( \begin{bmatrix} z_i \\ \vdots \end{bmatrix} \cdot \frac{1}{\sum_j z_j} \right) \\
 &= -y + \hat{y}
 \end{aligned}$$

$$\Rightarrow \boxed{\nabla_z CE(y, \hat{y}) = \hat{y} - y}$$

c)  $\nabla_n \mathcal{J} = \nabla_n CE(y, \hat{y}) = ?$

The flow is the following



with the previous question we found  $\frac{\partial \mathcal{J}}{\partial h} = \hat{y} - y$

$$\Rightarrow \frac{\partial \mathcal{J}}{\partial n} = \frac{\partial \mathcal{J}}{\partial h} \cdot \frac{\partial h}{\partial n} = (\hat{y} - y) \cdot \frac{\partial h}{\partial n}$$

so we need to find  $\frac{\partial h}{\partial n}$

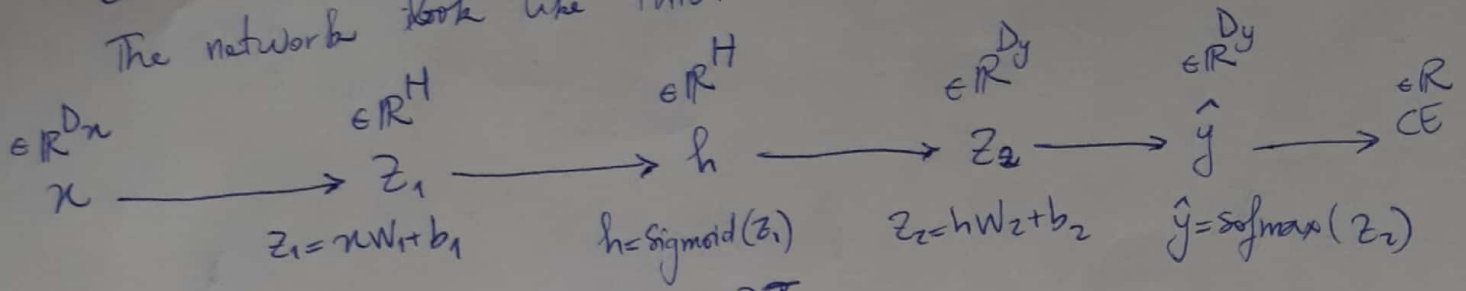
$$\begin{aligned}
 \frac{\partial h}{\partial n} &= \frac{\partial \text{Sigmoid}(nW_1 + b_1)}{\partial n} = \sigma'(nW_1 + b_1) \cdot \frac{\partial (nW_1 + b_1)}{\partial n} \\
 &= \sigma'(nW_1 + b_1)
 \end{aligned}$$

$$(n, D_y) \cdot (D_y, H) \cdot (H, D_h) \cdot (D_h, D_n) \quad (4)$$



2)  $\frac{\partial \mathcal{F}}{\partial x} = \frac{\partial \text{CE}(y, \hat{y})}{\partial x} = ?$

The network look like this:



in the previous question we found  $\frac{\partial \mathcal{F}}{\partial z_2} = \hat{y} - y$

$$- \frac{\partial \mathcal{F}}{\partial x} = \frac{\partial \mathcal{F}}{\partial z_2} \cdot \frac{\partial z_2}{\partial x} = (\hat{y} - y) \cdot \frac{\partial z_2}{\partial x} \quad (1)$$

$$- \frac{\partial z_2}{\partial x} = \frac{\partial z_2}{\partial h} \cdot \frac{\partial h}{\partial x} \quad \left| \quad \frac{\partial z_2}{\partial h} = \frac{\partial (hW_2 + b_2)}{\partial h} = W_2^T \right.$$

$$= W_2^T \cdot \frac{\partial h}{\partial x} \quad (2)$$

$$- \frac{\partial h}{\partial x} = \frac{\partial h}{\partial z_1} \cdot \frac{\partial z_1}{\partial x} = \sigma'(z_1) \cdot \frac{\partial z_1}{\partial x} \quad (3)$$

$$- \frac{\partial z_1}{\partial x} = \frac{\partial (xW_1 + b_1)}{\partial x} = W_1^T \quad (4)$$

$$(1), (2), (3), (4) \Rightarrow \frac{\partial \mathcal{F}}{\partial x} = (\hat{y} - y) \cdot W_2^T \cdot \sigma'(z_1) \cdot W_1^T$$

$\in \mathbb{R}^{D_n} \quad (1 \times D_y) \cdot (D_y \times H) \cdot (H \times H) \cdot (H \times D_n) = (1 \times D_n) \checkmark$

d)  $W_1: D_n \times H \quad b_1: 1 \times H$

$W_2: \cancel{D_n \times H} \quad b_2: 1 \times D_y$   
 $H \times D_y$

total number of parameters =  $(D_n \times H) + (H \times D_y) + H + D_y$

### 3) Word2Vec

a) skip-gram model  $\hat{y} = P(o | c) = \frac{\exp(u_o^T v_c)}{\sum_w \exp(u_w^T v_c)}$

$J = J_{\text{softmax-CE}}(o, v_c, U) = \text{CE}(y, \hat{y})$  | Where  $o$  is the idx of the predicted ~~context~~ word

$\frac{\partial J}{\partial v_c} = ?$

$\text{CE}(y, \hat{y}) = - \sum_w y_w \log(\hat{y}_w) = - \log(\hat{y}_o) = -u_o^T v_c + \log \sum_w e^{u_w^T v_c}$

$\Rightarrow \frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c} (-u_o^T v_c + \log \sum_w \exp(u_w^T v_c))$

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

(element-wise multiplication)  
 $\frac{\partial}{\partial v_c} \left( \sum_w \exp(u_w^T v_c) \right) = \sum_w \frac{\partial}{\partial v_c} (e^{u_w^T v_c}) = \sum_w u_w e^{u_w^T v_c} \quad (2)$

$= \frac{\partial}{\partial v_c} \left( \sum_w u_w e^{u_w^T v_c} \right)$

(1) and (2)  $\Rightarrow$

$\frac{\partial J}{\partial v_c} = -u_o + \frac{\sum_w u_w e^{u_w^T v_c}}{\sum_{w'} e^{u_{w'}^T v_c}}$

$= -u_o + \sum_w u_w \cdot \frac{e^{u_w^T v_c}}{\sum_{w'} e^{u_{w'}^T v_c}}$

$\frac{\partial J}{\partial v_c} = -u_o + \sum_w u_w \cdot \hat{y}_w \rightarrow (F1)$

3) b)  $\frac{\partial F}{\partial u_w} = \frac{\partial}{\partial u_w} (-u_o \cdot v_c + \log \sum_w \exp(u_w^T \cdot v_c))$

$$= \frac{\partial}{\partial u_w} (-u_o \cdot v_c) + \frac{\partial}{\partial u_w} \left[ \log \sum_w \exp(u_w^T \cdot v_c) \right]$$

$$= \begin{cases} 0 & \text{if } o \neq w \\ -v_c & \text{if } o = w \end{cases} + \frac{1}{\sum_{w'} \exp(u_{w'}^T \cdot v_c)} \otimes \frac{\partial}{\partial u_w} \sum_{w'} \exp(u_{w'}^T \cdot v_c)$$

$$= \begin{cases} 0 & \text{if } o \neq w \\ -v_c & \text{if } o = w \end{cases} + \frac{1}{\sum_{w'} \exp(u_{w'}^T \cdot v_c)} \otimes \sum_{w'} v_c \exp(u_{w'}^T \cdot v_c)$$

$$= \begin{cases} 0 & \text{if } o \neq w \\ -v_c & \text{if } o = w \end{cases} + \frac{1}{\sum_{w'} \exp(u_{w'}^T \cdot v_c)} \otimes v_c \cdot \left( \frac{\sum_{w'} \exp(u_{w'}^T \cdot v_c)}{\sum_{w'} \exp(u_{w'}^T \cdot v_c)} = 1 \right)$$

$$\frac{\partial F}{\partial u_w} = \begin{cases} v_c & \text{if } o \neq w \\ 0 & \text{if } o = w \end{cases}$$

$$= \begin{cases} 0 & \text{if } o \neq w \\ -v_c & \text{if } o = w \end{cases}$$

$$+ \frac{1}{\sum_{w'} \exp(u_{w'}^T \cdot v_c)} \otimes v_c \cdot \left[ \frac{\partial}{\partial u_w} \exp(u_w^T \cdot v_c) \right]$$

$$= \begin{cases} 0 & \text{if } o \neq w \\ -v_c & \text{if } o = w \end{cases}$$

$$+ v_c \cdot \frac{\exp(u_w^T \cdot v_c)}{\sum_{w'} \exp(u_{w'}^T \cdot v_c)} \quad \leftarrow \text{softmax}_w$$

$$\frac{\partial F}{\partial u_w} = \begin{cases} v_c \cdot \hat{y}_w & \text{if } o \neq w \\ v_c (1 - \hat{y}_w) & \text{if } o = w \end{cases} \rightarrow (F2)$$

7



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

$$\begin{aligned} \text{c-1)} \quad \frac{\partial J}{\partial v_c} &= -\frac{1}{\sigma(u_o^T v_c)} \otimes \frac{\partial}{\partial v_c} \sigma(u_o^T v_c) - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \otimes \frac{\partial}{\partial v_c} \sigma(-u_k^T v_c) \\ &= \cancel{\sigma'} - \frac{1}{\sigma(u_o^T v_c)} \otimes \left( \sigma'(u_o^T v_c) \cdot \frac{\partial}{\partial v_c} u_o^T v_c \right) - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \otimes \left( \sigma'(-u_k^T v_c) \frac{\partial}{\partial v_c} (-u_k^T v_c) \right) \end{aligned}$$

$$= -\frac{1}{\sigma(u_o^T v_c)} \otimes \left( \sigma'(u_o^T v_c) \cdot U_o \right) - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \otimes \left( \sigma'(-u_k^T v_c) \cdot (-1) U_k \right)$$

$$= -\frac{\sigma(u_o^T v_c)(1-\sigma(u_o^T v_c))}{\sigma(u_o^T v_c)} \otimes U_o + \sum_{k=1}^K \frac{\sigma(-u_k^T v_c)(1-\sigma(-u_k^T v_c))}{\sigma(-u_k^T v_c)} \otimes U_k$$

$$\frac{\partial J}{\partial v_c} = - (1-\sigma(u_o^T v_c)) U_o + \sum_{k=1}^K (1-\sigma(-u_k^T v_c)) U_k \rightarrow (F3)$$

$$\text{c-2)} \quad \frac{\partial J}{\partial u_w} = \begin{cases} 0 & \text{if } w \neq o \\ - (1-\sigma(u_o^T v_c)) v_c + \sum_{k=1}^K (1-\sigma(-u_k^T v_c)) v_c & \text{if } w = k \\ 0 & \text{if } w \neq k \end{cases}$$

$$\Rightarrow \frac{\partial J}{\partial u_w} = \begin{cases} 0; & \text{if } w \neq o \text{ and } w \neq \text{all negative samples } k \\ - (1-\sigma(u_o^T v_c)) v_c; & \text{if } w = o, \text{ the output word} \\ (1-\sigma(-u_k^T v_c)) v_c; & \text{if } w = k, \text{ one of the negative samples} \end{cases}$$

(FA) ←



3)

c) The cost function with the negative sample loss is much more efficient to compute because we don't have to look at all the words in the vocabulary to compute the gradient at each step.

d) d-1) SkipGram  $v \rightarrow \begin{cases} U \end{cases}$

$$\frac{\partial J_{SG}}{\partial v_k} = \sum_{-m, m}^j \frac{\partial}{\partial v_k} F(U_{c+j}, v_c) = \begin{cases} 0 & \text{if } k \neq c \\ \sum_{-m, m}^j \frac{\partial}{\partial v_c} F(U_{c+j}, v_c) \end{cases}$$

$$\frac{\partial J_{SG}}{\partial v_k} = \begin{cases} 0 & \text{if } k \neq c \\ \sum_{-m, m}^j \left[ -U_{c+j} + \sum_w^W U_w \hat{y}_w \right] & \text{if } k=c \text{ and softmax-loss} \\ \sum_{-m, m}^j \left[ (\sigma(U_{c+j}^T v_c) - 1) U_{c+j}^T + \sum_{n=1}^N (1 - \sigma(-U_n^T v_c)) U_n^T \right] & \text{if } k=c \text{ and neg-} \end{cases}$$

Here we change  $k$  previously denoting loss negative samples to  $n \rightarrow (1, -N)$  to break confusion without  $k$

~~$$\frac{\partial J_{SG}}{\partial U_k} = \sum_{-m, m}^j \frac{\partial}{\partial U_k} F(U_{c+j}, v_c) = \begin{cases} 0 & \text{if } k \notin [c-m, c, c+m] \\ \sum_{-m, m}^j \frac{\partial}{\partial U_k} F(U_{c+j}, v_c) & \text{if } k \in E = [c-m, c, c+m] \end{cases}$$~~

~~$$\frac{\partial J_{SG}}{\partial U_k} = \begin{cases} \oplus 0 & \text{if } k \notin E = [c-m, c+m] \text{ and softmax} \\ \oplus \sum_{-m, m}^j v_c y_{c+j} & \text{if } k \in E \text{ and } k \neq c \text{ and softmax} \\ \oplus \sum_{-m, m}^j v_c (1 - y_{c+j}) & \text{if } k \in E \text{ and } k = c \text{ and softmax} \\ \oplus \sum_{-m, m}^j (\sigma(-U_n^T v_c)) v_c & \text{if } k \notin E \text{ and } k = n \text{ and neg-sample} \\ \oplus \sum_{-m, m}^j (\sigma(- \end{cases}$$~~

(this  $k$  is different from the  $k$  used previously to denote neg-samples)

$n = \text{one of the neg-samples}$

$$\frac{\partial J_{SG}}{\partial u_k} = \sum_{-m, m} F(u_{c+j}, v_c)$$

- the output word is neither the predicted or center  
- the output

~~note~~ we

- For the softmax-loss we have 2 cases

- ~~$u_k$  is neither the center nor context word (C1)~~
- $u_k$  is a Context word (C2)
- $u_k$  is the Center (C3)

- For the negsampling-loss we have 3 cases

- $u_k$  is neither neg-sample, nor context, nor center word (C1)
- $u_k$  is neg-sample word (C2)
- $u_k$  is context word (C3)

~~$u_k$  is center word (C4)~~

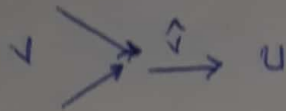
So for softmax-loss:

$$\frac{\partial J_{SG}}{\partial u_k} = \begin{cases} 0 & \text{if } k \notin [c-m, c+m] = E \quad (C1) \\ \sum_{-m, m} v_c \cdot \hat{y}_{c+j} & \text{if } k \in E \text{ and } k \neq c \quad (C2) \\ \sum_{-m, m} v_c (1 - \hat{y}_{c+j}) & \text{if } k \in E \text{ and } k = c \quad (C3) \end{cases}$$

for negative-sampling loss:

$$\frac{\partial J_{SG}}{\partial u_k} = \begin{cases} 0 & \text{if } k \notin [c-m, c+m] = E \text{ and } u_k \text{ is not neg-sample} \quad (C1) \\ \sum_{-m, m} v_c (1 - \sigma(-u_k^T v_c)) & \text{if } k \notin E \text{ and } u_k \text{ is neg-sample} \quad (C2) \\ \sum_{-m, m} v_c (\sigma(u_{c+j}^T v_c) - 1) & \text{if } k \in E \quad (C3) \end{cases}$$

3)

d) d-2) CBDW

$$\frac{\partial \mathcal{F}}{\partial v_k} = \frac{\partial}{\partial v_k} F(u_c, \hat{v})$$

$$\frac{\partial \mathcal{F}}{\partial v_k} = \begin{cases} 0; & \text{if } k \notin E = [c-m, c+m] \\ -u_c + \sum_w u_w \hat{y}_w; & \text{if } k \in E \text{ and softmax} \\ (\sigma(u_c^T \hat{v}) - 1)u_c + \sum_n (1 - \sigma(-u_n^T \hat{v}))u_n & \text{if } k \in E \text{ and neg-loss} \end{cases}$$

$$\frac{\partial \mathcal{F}}{\partial u_k} = \frac{\partial}{\partial u_k} F(u_c, \hat{v})$$

$$\frac{\partial \mathcal{F}}{\partial u_k} = \begin{cases} \hat{v} y_k & \text{if } k \neq c \\ \hat{v} (1 - y_k) & \text{if } k = c \end{cases}$$

Softmax

$$\frac{\partial \mathcal{F}}{\partial u_k} = \begin{cases} 0; & \text{if } k \neq c \text{ and } u_k \neq \text{all negative samples } n \\ (1 - \sigma(-u_n^T \hat{v}))\hat{v}; & \text{if } k \neq c, \text{ and } k = n \text{ a negative sample} \\ (\sigma(u_k^T \hat{v}) - 1)\hat{v}; & \text{if } k = c \end{cases}$$

Neg. sampling



- 3) g) What I can see on the 93-word-vectors.png:
- 1) articles (a, the) are more close to each other than they are to other words
  - 2) semantically <sup>antonym</sup> similar words (boring, waste, worth) (wonderful, great, amazing) are cluster next to each other  $\Rightarrow$  they appear in similar context

4) b) We want to add regularization to our training procedure to keep the weights of our model low and hence fight overfitting; which ultimately will lead to more generalization power when confronted to unseen data

- d) the pretrained works better because:
- trained on more data  $\Rightarrow$  more generalization power
  - hyperparameters were chosen more carefully for pretrained vectors
  - Maybe GloVe ~~is~~ <sup>is</sup> better for this task.
- e) - After a threshold ( $10^{-2}$  here) decreasing the regularization factor doesn't increase the accuracy that much
- train accuracy is lower than dev accuracy once the model has "converged" which is strange
  - the previous point may suggest we may be underfitting.
- f) - globally the model tend to avoid very extreme reviews (+, -)
- hence those classes are poorly classified
  - the model favored (-, +) classes a bit too much

g)

observations for -your vectors since glove url is broken