

## 2 Understanding word2vec

2.a

$$\sigma(\mathbf{x} + c)_j = \frac{e^{z_j + c}}{\sum_{k=1}^K e^{z_k + c}} = \frac{e^c e^{z_j}}{e^c \sum_{k=1}^K e^{z_k}} = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} = \sigma(\mathbf{x})_j$$

*Quod.Erat.Demonstrandum*

2.b

$$\begin{aligned} & - \sum_{w \in Vocab} y_w \log(\hat{y}_w) = \\ & -y_o \log(\hat{y}_o) - \sum_{w \in Vocab, w \neq o} y_w \log(\hat{y}_w) = \\ & -\log(\hat{y}_o) \end{aligned}$$

2.c

$$\begin{aligned} & \frac{\partial J(v_c, o, U)}{\partial v_c} = \\ & -\frac{\partial(u_o^T v_c)}{\partial v_c} + \frac{\partial(\log(\sum_w \exp(u_w^T v_c)))}{\partial v_c} = \\ & -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \frac{\partial(\sum_w \exp(u_w^T v_c))}{\partial v_c} = \\ & -u_o + \sum_w \frac{\exp(u_w^T v_c) u_w}{\sum_w \exp(u_w^T v_c)} = \\ & -u_o + \sum_w p(O = w | C = c) u_w = \\ & -u_o + \sum_w \hat{y}_w u_w = \\ & U(\hat{y} - y) \end{aligned}$$

2.d

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

*Either :*

$w = 0 :$

$$\begin{aligned} \frac{\partial J(v_c, o, U)}{\partial u_w} &= \\ -v_c + p(O = o | C = c)v_c &= \\ \hat{y}_w v_c - v_c &= \\ (\hat{y}_w - 1)v_c \end{aligned}$$

$w \neq 0 :$

$$\begin{aligned} \frac{\partial J(v_c, o, U)}{\partial u_w} &= \\ 0 + p(O = w | C = c)v_c &= \\ \hat{y}_w v_c \end{aligned}$$

*Eventually :*

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

2.e

$$\begin{aligned} \frac{\partial \sigma(x_i)}{\partial x_i} &= \frac{1}{(1 + \exp(-x_i))^2} \exp(-x_i) = \sigma(x_i)(1 - \sigma(x_i)) \\ \frac{\partial \sigma(x)}{\partial x} &= \left[ \frac{\partial \sigma(x_j)}{\partial x_i} \right]_{d \times d} = \begin{bmatrix} \sigma'(x_1) & 0 & \dots & 0 \\ 0 & \sigma'(x_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma'(x_d) \end{bmatrix} = \text{diag}(\sigma'(x)) \end{aligned}$$

2.f

For  $v_c$  :

$$\begin{aligned} \frac{\partial J_{\text{neg-sample}}}{\partial v_c} &= \\ (\sigma(u_o^T v_c) - 1)u_o + \sum_{k=1}^K (1 - \sigma(-u_k^T v_c)) u_k &= \\ (\sigma(u_o^T v_c) - 1)u_o + \sum_{k=1}^K \sigma(u_k^T v_c) u_k \end{aligned}$$

For  $u_o : o \notin [w_i]_{i=1}^K$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_o} = (\sigma(u_o^T v_c) - 1)v_c$$

For  $u_k : k = [1, K] \in \mathbb{N}$

$$\frac{\partial J}{\partial \mathbf{u}_k} = -(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1) \mathbf{v}_c = \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{v}_c,$$

For naive softmax loss function :

$$\begin{aligned} \frac{\partial J(v_c, o, U)}{\partial U} &= v_c(\hat{y} - y)^T \\ \frac{\partial J(v_c, o, U)}{\partial v_c} &= U(\hat{y} - y) \end{aligned}$$

For negative sampling loss function :

$$\frac{\partial J}{\partial \mathbf{v}_c} =$$

$$(\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{u}_o + \sum_{k=1}^K \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{u}_k =$$

$$\sigma(-\mathbf{u}_o^\top \mathbf{v}_c) \mathbf{u}_o + \sum_{k=1}^K \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{u}_k$$

$$\frac{\partial J}{\partial \mathbf{u}_k} = \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{v}_c : k = [1, K] \in \mathbb{N}$$

$$\frac{\partial J}{\partial \mathbf{u}_o} = (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{v}_c = \sigma(-\mathbf{u}_o^\top \mathbf{v}_c) \mathbf{v}_c$$

*With this loss function the computation is not related with going over all of the words in vocabulary,  $V$ , rather than in the number of samples,  $K$ .*

2.g.i

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} =$$

$$\sum_{-m \leq j \leq m, j \neq 0} \frac{\frac{\partial J(v_c, w_{t+j}, U)}{\partial U}}{\partial U}$$

2.g.ii

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_c} =$$

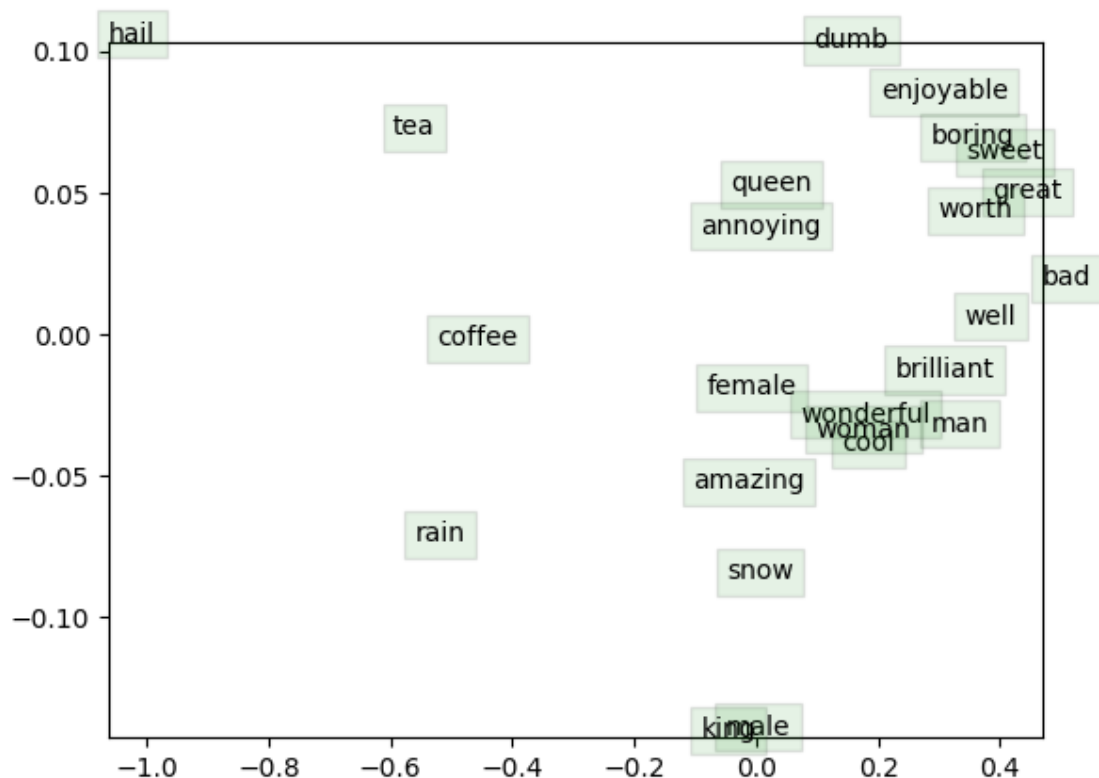
$$\sum_{-m \leq j \leq m, j \neq 0} \frac{\frac{\partial J(v_c, w_{t+j}, U)}{\partial v_c}}{\partial v_c}$$

2.g.iii

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_w} = 0$$

### 3 Implementing word2vec

3.e



*We observe the words  $\{\text{sweet}, \text{great}\}$  clustered together, and it's reasonable that they can appear in proximity.*

*The words  $\{\text{coffee}, \text{tea}\}$  are not clustered together, however this is surprising because they in many cases used interchangeably, or along with each other.*

*However the clustering didn't catch that the words  $\{\text{man}, \text{male}\}$  can have similar meanings, as they are afar from each other.*