

Softmax and cross entropy loss

$$\text{softmax}(x+c)_i = \frac{e^{(x_i+c)}}{\sum_j e^{(x_j+c)}} = \quad (1)$$

$$\frac{\sum_j e^c}{\sum_j e^{x_j+c}} = \frac{e^c \cdot e^{x_i}}{e^c \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(x); \quad (a)$$

As a result $\text{softmax}(x) = \text{softmax}(x+c)$
(it's the same math just on multiple coordinates).

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

$$G'(x) = \frac{0 - e^{-x}}{(1+e^{-x})^2} = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{e^{-x}}{(1+e^{-x})^2} =$$

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

$$G(x) - G(x)^2 = G(x)(1-G(x)).$$

$$\text{So: } G'(x) = G(x)(1-G(x)).$$

$$CE(\hat{y}, \hat{y}) = -\sum_i y_i \log(\hat{y}_i)$$

(a) (2)

$$\frac{\partial CE(\hat{y}, \hat{y})}{\partial V_C} = \frac{\partial}{\partial V_C} \left(-\log \left(\frac{\exp(u_0^T v_C)}{\sum_w \exp(u_w^T v_C)} \right) \right) =$$

$$\frac{\partial}{\partial V_C} \left(-\log(\exp(u_0^T v_C)) - \log \left(\sum_w \exp(u_w^T v_C) \right) \right)$$

$$\frac{\partial}{\partial V_C} \left(-u_0^T v_C + \log \left(\sum_w \exp(u_w^T v_C) \right) \right) =$$

$$-u_0 + \frac{1}{\sum_w \exp(u_w^T v_C)} \cdot \sum_w \exp(u_w^T v_C) \cdot u_w^T x =$$

$$\boxed{-u_0 + \sum_x p(x|C) \cdot u_w^T x}$$

Answers (b)

Case 1: Vector y has a 1 at location u_0

$$\frac{\partial}{\partial u_0} \mathbb{E}(y, y) = \frac{\partial}{\partial u_0} \left(-\log \left(\frac{\exp(u_0^\top v_c)}{\sum_w \exp(u_w^\top v_c)} \right) \right)^2$$

$$\frac{\partial}{\partial u_0} \left(-\log(\exp(u_0^\top v_c)) - \log(\sum_w \exp(u_w^\top v_c)) \right)^2$$

$$\begin{aligned} \frac{\partial}{\partial u_0} & \left(-u_0^\top v_c + \log(\sum_w \exp(u_w^\top v_c)) \right)^2 \\ &= -v_c + \frac{1}{\sum_w \exp(u_w^\top v_c)} \cdot v_c \exp(u_0^\top v_c) \end{aligned}$$

$$-v_c + v_c \beta(\alpha_c) = v_c (\beta(\alpha_c) - 1).$$

(Case 2): Vector y has a 1 at location k show that $\nabla \ell(u)$.

$$\frac{\partial}{\partial u_k} \ell(y, \hat{y}) = \frac{\partial}{\partial u_k} \left(-\log \frac{\exp(u_0^T v_c)}{\sum_w \exp(u_w^T v_c)} \right) =$$

$$\frac{\partial}{\partial u_k} \left(-\log (\exp(u_0^T v_c)) - \log \left(\sum_w \exp(u_w^T v_c) \right) \right) =$$

$$\frac{\partial}{\partial u_k} \left(-u_0^T v_c + \log \left(\sum_w \exp(u_w^T v_c) \right) \right) =$$

$$0 + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot v_c \exp(u_k^T v_c) =$$

$$V_c \cdot \frac{\exp(u_k^T v_c)}{\sum_w \exp(u_w^T v_c)} = \boxed{V_c p(k|c)}$$

~~Notwendige Bedingungen~~

(C)

$$\frac{\partial J_{\text{Heg-Schle}}}{\partial v_c} = \frac{\partial}{\partial v_c} \left(-\log(\phi(u_c^\top v_c)) - \sum_{k=1}^K \phi(-u_k^\top v_c) \right) =$$

$$-\frac{1}{\phi(u_c^\top v_c)} \cdot \phi(u_c^\top v_c)(1 - \phi(u_c^\top v_c)) \cdot u_0$$

$$-\sum_{k=1}^K \frac{\phi(-u_k^\top v_c)(1 - \phi(-u_k^\top v_c)) \cdot -u_k}{\phi(-u_k^\top v_c)} =$$

$$u_0(\phi(u_c^\top v_c) - 1) - \sum_{k=1}^K (\phi(-u_k^\top v_c) - 1) u_k$$

$$\frac{\partial}{\partial u_0} \left(J_{\text{log-simle}} \right) =$$

$$\frac{\partial}{\partial u_0} \left(-\log \left(\delta(u_0^\top v_c) \right) \right) - \sum_{c=1}^k \log \left(\delta(-u_0^\top v_c) \right) =$$

∇ has no no

$$\frac{\partial}{\partial u_0} \left(-\log \left(\delta(u_0^\top v_c) \right) \right) =$$

$$-\frac{1}{\delta(u_0^\top v_c)} \cdot \delta(u_0^\top v_c) (1 - \delta(u_0^\top v_c)) \cdot v_c =$$

$$v_c (\delta(u_0^\top v_c) - 1)$$

$$\frac{\partial}{\partial h_u} J_{\text{log-surface}} =$$

$$\frac{\partial}{\partial h_u} \left(-\log(\delta(h_u^\top v_c)) - \sum_{u=1}^k \log(\delta(-h_u^\top v_d)) \right) =$$


 has no h_u

$$\frac{\partial}{\partial h_u} \left(- \sum_{u=1}^k \log(\delta(-h_u^\top v_c)) \right) =$$

$$- \frac{1}{\delta(-h_u^\top v_c)} \delta(-h_u^\top v_c) (1 - \delta(-h_u^\top v_c)) \cdot v_c =$$

$$-v_c (\delta(-h_u^\top v_c) - 1).$$

$$J = \sum_{\substack{-M \leq j \leq M \\ j \neq 0}} F(w_{c+j}, v_c)$$

(d)

$$\frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c} \left(\sum_{\substack{-M \leq j \leq M \\ j \neq 0}} F(w_{c+j}, v_c) \right) =$$

$$\sum_{\substack{-M \leq j \leq M \\ j \neq 0}} \frac{\partial F(w_{c+j}, v_c)}{\partial v_c}$$

$$\frac{\partial J}{\partial w_{c+i}} = \frac{\partial}{\partial w_{c+i}} \left(\sum_{\substack{-M \leq j \leq M \\ j \neq 0}} F(w_{c+j}, v_c) \right) =$$

↑
∂F(w_{c+i}, v_c)

$$\frac{\partial F(w_{c+i}, v_c)}{\partial w_{c+i}}$$

-M ≤ i ≤ M
i ≠ 0

In our picture we can

see that good and bad
wonderful and great
boring and waste
cool and sweet

has received similar values in the 2d
PCA based on the kth CS feature
it makes sense as wonderful and great
are similar many words, good and bad
are opposites, and boring and waste
and cool and sweet could be
used in similar situations.