

(1)

Naive-Softmax

$$= -\log \hat{y}_0 \quad (\text{简写符号})$$

$$= -\log P(O=0 | C=c)$$

月 日 星期

$J_{\text{naive-softmax}} = -\log \frac{\exp(u_0^T v_c)}{\sum \exp(u_w^T v_c)}$ 对 $v_c, u_{w \neq 0}, u_0$ 求偏导.

$$\begin{aligned} \textcircled{1} \frac{\partial J}{\partial v_c} &= \frac{\partial}{\partial v_c} (\log \sum \exp(u_w^T v_c) - \log \exp(u_0^T v_c)) \\ &= \frac{\partial}{\partial v_c} (\log \sum \exp(u_w^T v_c)) - \frac{\partial}{\partial v_c} (u_0^T v_c) \\ &= \frac{1}{\sum \exp(u_w^T v_c)} \cdot \frac{\partial}{\partial v_c} \sum [e^{u_w^T v_c}] - u_0 \\ \text{Const 可移入} \sum &= \frac{1}{\sum \exp(u_w^T v_c)} \cdot \sum [u_w \exp(u_w^T v_c)] - u_0 \\ &= \sum_w \left[\frac{u_w \exp(u_w^T v_c)}{\sum \exp(u_w^T v_c)} \right] - u_0 \\ &= \sum_w [u_w P(O=w | C=c)] - u_0 \\ &= \sum_w [u_w \hat{y}_w] - u_0 \quad \text{用简写符号代入, 公式简洁} \end{aligned}$$

$$\begin{aligned} \textcircled{2} \frac{\partial J}{\partial u_{w \neq 0}} &= \frac{\partial}{\partial u_{w \neq 0}} (\log \sum \exp(u_w^T v_c) - \log \exp(u_0^T v_c)) \\ &= \frac{1}{\sum \exp(u_w^T v_c)} \cdot \sum \left[\frac{\partial \exp(u_w^T v_c)}{\partial u_{w \neq 0}} \right] - \frac{\partial}{\partial u_{w \neq 0}} u_0^T v_c \\ &= \sum_{w \neq 0} \left[\frac{\exp(u_w^T v_c)}{\sum \exp(u_w^T v_c)} \cdot v_c \right] - 0 \\ &= \sum_{w \neq 0} (\hat{y}_w \cdot v_c) = \hat{y}_w \cdot v_c \quad (\because w \text{ is specific}) \end{aligned}$$

$$\begin{aligned} \textcircled{3} \frac{\partial J}{\partial u_{w=0}} &= \frac{\partial J}{\partial u_0} = \frac{\partial}{\partial u_0} (\log \sum \exp(u_w^T v_c) - u_0^T v_c) \\ &= \frac{1}{\sum \exp(u_w^T v_c)} \cdot \sum \left[\frac{\partial \exp(u_w^T v_c)}{\partial u_0} \right] - \frac{\partial}{\partial u_0} u_0^T v_c \\ &= \frac{\exp(u_0^T v_c)}{\sum \exp(u_w^T v_c)} \cdot v_c - v_c \\ &= (\hat{y}_0 - 1) \cdot v_c \end{aligned}$$

Sigmoid 的快速求导 $\delta(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$

$$\begin{aligned} \frac{\partial}{\partial x} \delta(x) &= \frac{-\frac{\partial}{\partial x} (1+e^{-x})}{(1+e^{-x})^2} = -\frac{1}{(1+e^{-x})^2} \cdot (-1) \cdot e^{-x} \\ &= \frac{e^{-x}}{(1+e^{-x}) \cdot (1+e^{-x})} = \frac{(e^x)^2 \cdot e^{-x}}{(e^x)^2 \cdot (1+e^{-x})(1+e^{-x})} \\ &= \frac{e^x}{(e^x+1) \cdot (e^x+1)} = \delta(x) \cdot \frac{1}{e^x+1} \end{aligned}$$

$$\begin{aligned} &= \delta(x) \cdot \frac{e^x+1-e^x}{e^x+1} = \delta(x) \left(\frac{e^x+1}{e^x+1} - \frac{e^x}{e^x+1} \right) \\ &= \delta(x) (1 - \delta(x)) \end{aligned}$$

(2)

Neg-sampling.

月 日 星期

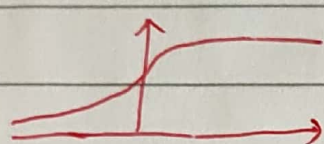
negative-sampling 还是 skip-gram

考虑中心词 c 与上下文 o : 令 $P(D=1 | o, c)$ 为 (o, c) 来自语料库, $P(D=0 | o, c)$ 表示 (o, c) 不来自语料库.

$$P(D=1 | w, c) = \sigma(v_c^T u_w)$$

$$P(D=0 | w, c) = 1 - P(D=1 | w, c)$$

$v_c^T u_w$ 表示了其相似性, 越相似, 数越大, 不相似就很小, 可能会负.

用 sigmoid 约束到 $(0, 1)$.

我们要最大化正确上下文的 $P(D=1)$, 最小化错误上下文的 $P(D=1)$ (即最大化错误上下文的 $P(D=0)$)

$$\max \prod_{(w,c) \in D} P(D=1 | w, c) \prod_{(w,c) \notin D} P(D=0 | w, c)$$

在语料库中 不在语料库中

$$\begin{aligned} \text{取 log: } \max \sum_{(w,c) \in D} \log P(D=1) + \sum_{(w,c) \notin D} \log (1 - P(D=1)) \\ = \max_{\min} \sum_{(w,c) \in D} \log \sigma(v_c^T u_w) + \sum_{(w,c) \notin D} \log (1 - \sigma(v_c^T u_w)) \end{aligned}$$

固定一个来自语料库的 w , $(w, c) \in D$, $w = c - m + j$, 窗口大小 $2m$

$$J = -\log \sigma(u_{c-m+j}^T v_c) - \sum_{\text{sampled}} \log \sigma(-\tilde{u}_k^T v_c)$$

上式中 $\{\tilde{u}_k | k=1 \dots k\}$ 是从 D 中抽样, 抽样概率为实际概率 $3/4$ 次方: $0.9^{3/4} \approx 0.92$, $0.09^{3/4} \approx 0.16$, 可以提高罕见词被抽中的概率.

对 neg-sampling loss 求导:

$$J_{\text{neg-sampling}}(v_c, u, U) = -\log(\delta(u_0^T v_c)) - \sum_{k=1}^K \log(\delta(-u_k^T v_c))$$

$$\begin{aligned} (1) \frac{\partial J}{\partial v_c} &= -\frac{1}{\delta(u_0^T v_c)} \cdot \frac{\partial}{\partial v_c} \delta(u_0^T v_c) - \sum_{k=1}^K \frac{1}{\delta(-u_k^T v_c)} \cdot \frac{\partial}{\partial v_c} \delta(-u_k^T v_c) \\ &= \frac{-1}{\delta(u_0^T v_c)} \cdot \delta(u_0^T v_c) \cdot (1 - \delta(u_0^T v_c)) \cdot u_0 - \sum_{k=1}^K \frac{1}{\delta(-u_k^T v_c)} \cdot \delta(-u_k^T v_c) \cdot (1 - \delta(-u_k^T v_c)) \cdot (-u_k) \end{aligned}$$

✱ You should know $1 - \delta(x) = \delta(-x)$

$$\begin{aligned} &= \delta(u_0^T v_c) \cdot \delta(-u_0^T v_c) \cdot \left(-\frac{u_0}{\delta(u_0^T v_c)} + \sum_{k=1}^K \frac{v_c}{\delta(-u_k^T v_c)} \right) \\ &= -u_0 (1 - \delta(u_0^T v_c)) - \sum_{k=1}^K -u_k (1 - \delta(-u_k^T v_c)) \end{aligned}$$

$$= u_0 (\delta(u_0^T v_c) - 1) - u_0 \delta(-u_0^T v_c) + \sum_{k=1}^K \delta(u_k^T v_c) \cdot u_k$$

$$\begin{aligned} (2) \frac{\partial J}{\partial u_0} &= -v_c \delta(-u_0^T v_c) + v_c \sum_{k=1}^K \delta(u_k^T v_c) \Big|_{u_k=u_0, \text{ but } u_k \neq u_0} \\ &= -v_c \delta(-u_0^T v_c) + v_c \delta(u_0^T v_c) = 0 \\ &= -v_c \delta(-u_0^T v_c) \end{aligned}$$

$$(3) \frac{\partial J}{\partial u_k} = v_c \delta(u_k^T v_c)$$

Naive-softmax

对 neg-sampling 求 skip-gram 的导数

设中心词 $c = w_t$, 窗口大小 m : $[w_{t-m}, \dots, w_{t+m}]$

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

中心词对所有窗口词检索一遍

$$(1) \frac{\partial J_{\text{skip-gram}}}{\partial v_c} = \sum \frac{\partial}{\partial v_c} J_{\text{neg-sampling/naive-softmax}}$$

$$(2) \frac{\partial J_{\text{skip-gram}}}{\partial v_w \neq c} = \frac{\partial}{\partial v_w} J = 0, w \neq c$$

$$(3) \frac{\partial J_{\text{skip-gram}}}{\partial U} = \frac{\partial}{\partial u_0} J + \sum_{k \neq 0} \frac{\partial J}{\partial u_k} = \sum \frac{\partial J}{\partial u}$$