

word2vec 公式推导

(Objective function + 梯度更新公式)

skip-gram model:

$$\max J(\theta) = \prod_{t=1}^T \prod_{\substack{m \leq j \leq m \\ j \neq 0}} \mathcal{P}(w_{t+j} | w_t; \theta)$$

最大似然  $\Rightarrow \theta$ : word vectors

除了 outside word 其他位置为 0  
实际为最小化交叉熵.

$$\Leftrightarrow \text{minimize } J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{m \leq j \leq m \\ j \neq 0}} \log \mathcal{P}(w_{t+j} | w_t; \theta) - \mathcal{P} \log(\mathcal{P})$$

$$\mathcal{P}(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} = \left[ \begin{array}{l} V: \text{vocabulary size} \\ u_o: \text{word vector for "outside" words} \\ v_c: \text{word vector for "center" word} \end{array} \right]$$

梯度下降

$$A: \frac{\partial}{\partial v_c} \log \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} = \frac{\partial}{\partial v_c} \log \exp(u_o^T v_c) - \frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c)$$

$$\left[ \textcircled{1} (a^x)' = a^x \ln a \quad \textcircled{2} (e^x)' = e^x \quad \textcircled{3} (\ln x)' = \frac{1}{x} \right]$$

$$\frac{\partial}{\partial v_c} \log \exp(u_o^T v_c) = \frac{\partial}{\partial v_c} (u_o^T v_c) = u_o$$

$$\begin{aligned} \frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c) &= \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \sum_{w=1}^V \frac{\partial}{\partial v_c} \exp(u_w^T v_c) \\ &= \sum_{w=1}^V \left( \exp(u_w^T v_c) \cdot u_w \right) \end{aligned}$$

Therefore:

$$\begin{aligned} A &= u_o - \frac{\sum_{w=1}^V (\exp(u_w^T v_c) \cdot u_w)}{\sum_{w=1}^V \exp(u_w^T v_c)} = u_o - \sum_{w=1}^V \frac{\exp(u_w^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot u_w \\ &= u_o - \sum_{w=1}^V P(w|c) \cdot u_w \end{aligned}$$

$$\begin{aligned}
 B: \frac{\partial}{\partial u_0} \log \frac{\exp(u_0^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} &= \frac{\partial}{\partial u_0} \log \exp(u_0^T v_c) - \frac{\partial}{\partial u_0} \log \sum_{w=1}^V \exp(u_w^T v_c) \quad (\ln x)' = \frac{1}{x} \\
 &= \frac{\partial}{\partial u_0} u_0^T v_c - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \frac{\partial}{\partial u_0} \left( \sum_{w=1}^V \exp(u_w^T v_c) \right) \\
 &= v_c - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \frac{\partial}{\partial u_0} \exp(u_0^T v_c) = v_c - \frac{\exp(u_0^T v_c) \cdot v_c}{\sum_{w=1}^V \exp(u_w^T v_c)} \\
 &= v_c - p(u_0 | v_c) \cdot v_c
 \end{aligned}$$

△交叉熵：机器/深度学习中所用来描述目标与预测值差距，即定义目标函数

**信息量**：一个事件发生的概率越低，获取到的信息量就越大

A：巴西队进入2018年世界杯决赛 B：中国队进入2018年世界杯决赛

显然<sup>低</sup>发生的概率<sup>低</sup>，因此<sup>B</sup>的信息量更大  $I(X=A) = -\log(p(A))$



**熵**：所有信息量的期望  $H(X) = -\sum_{i=1}^n p(x_i) \log(p(x_i)) \geq 0$

**相对熵**：也称KL散度 (Kullback-Leibler divergence)

用于衡量两个分布的差异

机器学习中， $P \Rightarrow$  真实分布  $[1, 0, 0]$ ； $Q \Rightarrow$  模型预测分布  $[0.7, 0.2, 0.1]$

显然 $Q$ 用来描述样本分类不够完美，信息量不足，需要额外的“信息增量”

$$D_{KL}(p||q) = \sum_{i=1}^n p(x_i) \log \frac{p(x_i)}{q(x_i)} \geq 0$$

用 $p$ 比用 $q$ 多出来的信息增量

$$\begin{aligned}
 \text{交叉熵 } D_{KL}(p||q) &= \sum_{i=1}^n p(x_i) \log p(x_i) - \sum_{i=1}^n p(x_i) \log q(x_i) \\
 &= -H(p(x)) - \sum_{i=1}^n p(x_i) \log q(x_i)
 \end{aligned}$$

机器/深度学习中此项为固定值，只需优化  $-\sum_{i=1}^n p(x_i) \log q(x_i)$  使其最小。