

word2vec 优化方式推导

一. Negative Sampling (负采样)

(k+1)

选择一个正样本和K个负样本，训练一个二分类器。

1. CBOW 的公式推导

假设中心词为w，周围词为context(w)，负样本集 NEG(w)

定义 Indicator Function. $I^w(\tilde{w}) = \begin{cases} 1, & \tilde{w} = w \\ 0, & \tilde{w} \neq w \end{cases}$

用来表示词对 [context(w), \tilde{w}] 的标签， $\text{context}(w) = \frac{j+0}{\sum_j u_{j+w}} \text{ (周围词平均向量)}$
 $= \vec{x}_w$

似然函数 likelihood

$$= \prod_{\tilde{w} \in \{w\} \cup \text{NEG}(w)} p(\tilde{w} | \text{context}(w)) \quad \textcircled{1}, \text{令目标词 } \tilde{w} \text{ 的词向量为 } \theta^{\tilde{w}}$$

$$\text{其中 } p(\tilde{w} | \text{context}(w)) = \begin{cases} \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}}), & L^w(\tilde{w}) = 1 \\ 1 - \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}}), & L^w(\tilde{w}) = 0 \end{cases}$$

$$\therefore \textcircled{1} \text{式} = \sigma(\vec{x}_w^T \cdot \theta^w) \cdot \prod_{\tilde{w} \in \text{NEG}(w)} (1 - \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})) = g(w)$$

给定语料库，我们的目标 $\max_{w \in C} \prod g(w) \Leftrightarrow \max_{w \in C} \log \prod g(w)$

$$\Rightarrow \max_{w \in C} \sum \left[\log \sigma(\vec{x}_w^T \cdot \theta^w) + \sum_{\tilde{w} \in \text{NEG}(w)} \log (1 - \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})) \right] \quad (1 - \sigma(x) = \sigma(-x))$$

$$\text{令 } L(w, \tilde{w}) = I^w(\tilde{w}) \log \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}}) + (1 - I^w(\tilde{w})) \log (1 - \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})) \quad \textcircled{2}$$

代表每个训练样本的 log-likelihood，使用梯度上升法

$$\frac{\partial L}{\partial \theta^{\tilde{w}}} = I^w(\tilde{w}) (1 - \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})) \cdot \vec{x}_w + (1 - I^w(\tilde{w})) (-\sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})) \cdot \vec{x}_w$$

$$= (I^w(\tilde{w}) - I^w(\tilde{w}) \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})) \cdot \vec{x}_w + (I^w(\tilde{w}) \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}}) - \cancel{I^w(\tilde{w}) \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})}) \cdot \vec{x}_w$$

$$= [I^w(\tilde{w}) - \sigma(\vec{x}_w^T \cdot \theta^{\tilde{w}})] \cdot \vec{x}_w$$

$$6'(x) = 6(x)(1 - 6(x))$$

$$(\log 6(x))' = 1 - 6(x)$$

$$[\log(1 - \sigma(x))]' = -\sigma(x)$$

由于 θ^w 和 x_w^T 在②式中对称,

$$\frac{\partial L}{\partial \theta^w} \text{ 换成 } \frac{\partial L}{\partial x_w^T} = [I^w(\tilde{w}) - \sigma(\vec{x}_w^T \cdot \vec{\theta}^{\tilde{w}})] \cdot \vec{\theta}^{\tilde{w}}$$

对于 $K+1$ 个样本, 可使用下述公式更新梯度

$$\vec{\theta}^{\tilde{w}}: \vec{\theta}^w + \eta [I^w(\tilde{w}) - \sigma(\vec{x}_w^T \cdot \vec{\theta}^{\tilde{w}})] \cdot \vec{x}_w^T$$

$$x_w^*: x_w^* + \bar{z} \eta [I^w(\tilde{w}) - \sigma(\vec{x}_w^T \cdot \vec{\theta}^{\tilde{w}})] \cdot \vec{\theta}^{\tilde{w}}$$

$\tilde{w} \in \{w\} \cup \text{NEG}(w)$

2. skipgram 公式推导

$$\text{似然函数} \Rightarrow \prod_{w \in C} \prod_{u \in \text{context}(w)} \prod_{\tilde{w} \in \{u\} \cup \text{NEG}(u)} P(\tilde{w} | w)$$

$$\text{其中 } P(\tilde{w} | w) = \begin{cases} \sigma(\vec{v}_w^T \cdot \vec{u}_{\tilde{w}}), & I^w(\tilde{w}) = 1 \\ 1 - \sigma(\vec{v}_w^T \cdot \vec{u}_{\tilde{w}}), & I^w(\tilde{w}) = 0 \end{cases}$$

$$\text{Log-likelihood} = \sum_{w \in C} \sum_{u \in \text{context}(w)} \sum_{\tilde{w} \dots} \left[I^w(\tilde{w}) \log \sigma(\vec{v}_w^T \cdot \vec{u}_{\tilde{w}}) + (1 - I^w(\tilde{w})) \log (1 - \sigma(\vec{v}_w^T \cdot \vec{u}_{\tilde{w}})) \right]$$

由于公式与(1)一样, 我们可以得到

$$\frac{\partial L}{\partial \vec{u}_{\tilde{w}}} = \cancel{I^w(\tilde{w})} [I^w(\tilde{w}) - \sigma(\vec{v}_w^T \cdot \vec{u}_{\tilde{w}})] \cdot \vec{v}_w^T$$

$$\vec{u}_{\tilde{w}}: \vec{u}_{\tilde{w}} + \eta \frac{\partial L}{\partial \vec{u}_{\tilde{w}}}$$

$$\text{同理: } \frac{\partial L}{\partial \vec{v}_w^T} = [I^w(\tilde{w}) - \sigma(\vec{v}_w^T \cdot \vec{u}_{\tilde{w}})] \cdot \vec{u}_{\tilde{w}}$$

对于 $K+1$ 个样本:

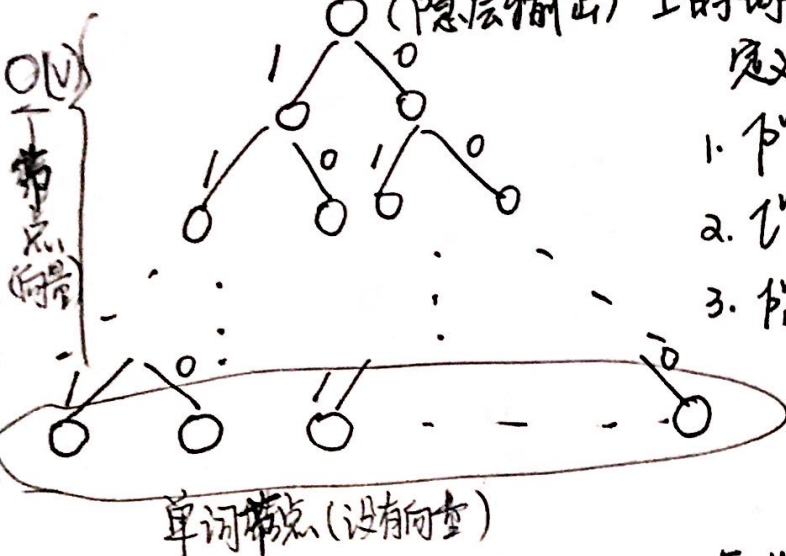
$$\vec{v}_w: \vec{v}_w + \bar{z} \eta \frac{\partial L}{\partial \vec{v}_w^T}$$

$\tilde{w} \in \{u\} \cup \text{NEG}(u)$

二. Hierarchical Softmax (层次 softmax) (哈夫曼树)

叶子结点代表的是词典中的所有词 $|V|$ 。通往目标词的路径上的点

(隐层输出) 上的向量是我们优化的。



定义:

1. p^w : 根结点到单词结点(叶子结点)的路径

2. l^w : 路径 p^w 中结点的个数

3. $p_1^w, p_2^w, \dots, p_{l^w}^w$: 路径上的点, p_1^w 根结点, $p_{l^w}^w$: 单词结点

4. $d_2^w, d_3^w, \dots, d_{l^w}^w \in \{0, 1\}$: 词 w 的哈夫曼编码, 由 $(l^w - 1)$ 位二进制编码构成

5. $\theta_1^w, \theta_2^w, \dots, \theta_{l^w}^w$: 对应 p^w 中非叶子结点的向量

1. CBOW 的公式推导

$$\max \prod_{w \in C} \prod_{j=2}^{l^w} p(d_j^w | x_w, \theta_{j-1}^w)$$

(x_w 为隐层输出)
 $\vec{x}_w = \frac{\sum \text{context}(w)}{|\text{context}(w)|}$, 有 $p = \begin{cases} \sigma(\vec{x}_w^T \cdot \theta_{j-1}^w), & d_j^w = 0 \\ 1 - \sigma(\vec{x}_w^T \cdot \theta_{j-1}^w), & d_j^w = 1 \end{cases}$

对数似然 $\Rightarrow \log \textcircled{1} = \sum_{w \in C} \sum_{j=2}^{l^w} ((1 - d_j^w) \log \sigma(\vec{x}_w^T \cdot \theta_{j-1}^w) + d_j^w \log (1 - \sigma(\vec{x}_w^T \cdot \theta_{j-1}^w)))$

$= L(j, w)$

对于每个样本

$$\frac{\partial L(j, w)}{\partial \theta_{j-1}^w} = (1 - d_j^w) (1 - \sigma(\vec{x}_w^T \cdot \theta_{j-1}^w)) \cdot \vec{x}_w + d_j^w (-\sigma(\vec{x}_w^T \cdot \theta_{j-1}^w)) \cdot \vec{x}_w$$

$$= (1 - d_j^w - \sigma(\vec{x}_w^T \cdot \theta_{j-1}^w)) \cdot \vec{x}_w$$

$$\theta_{j-1}^w = \theta_{j-1}^w + \eta \frac{\partial L}{\partial \theta_{j-1}^w}$$

对称性 $\Rightarrow \frac{\partial L}{\partial x_w} = (1 - \sigma(\vec{x}_w^T \cdot \theta_{j-1}^w) - d_j^w) \cdot \theta_{j-1}^w$

根据 x_w 更新 $v_w \Rightarrow v_m = v_m + \eta \sum_{j=2}^{l^w} \frac{\partial L}{\partial x_w}$

2. skipgram 的公式推导

非常类似 CBOW, 多了一个遍历 context(w)

$$\max_w \prod_{u \in \text{context}(w)} \prod_{j=2}^L p(d_j^w | x_w, \theta_{j-1}^w) \quad (x_w \text{ 隐层输出})$$

对数似然 ~~$L = \sum_w \sum_{u \in \text{context}(w)} \sum_{j=2}^L \log p(d_j^w | x_w, \theta_{j-1}^w)$~~

$$= \sum_w \sum_{u \in \text{context}(w)} \underbrace{\sum_{j=2}^L \left((1-d_j^w) \log \sigma(x_w^T \cdot \theta_{j-1}^w) + d_j^w \log (1 - \sigma(x_w^T \cdot \theta_{j-1}^w)) \right)}_{L(w, j)}$$

参考上节: $\frac{\partial L}{\partial \theta_{j-1}^w} = (1-d_j^w - \sigma(x_w^T \cdot \theta_{j-1}^w)) \cdot x_w$

$$\theta_{j-1}^w : \theta_{j-1}^w + \eta \frac{\partial L}{\partial \theta_{j-1}^w}$$

$$\frac{\partial L}{\partial x_w} = (1-d_j^w - \sigma(x_w^T \cdot \theta_{j-1}^w)) \cdot \theta_{j-1}^w$$

$$x_w^* : x_w + \eta \sum_{u \in \text{context}(w)} \sum_{j=2}^L \frac{\partial L}{\partial x_w}$$