

word2vec

单词 w ;

词典 $D = \{w_1, w_2, \dots, w_N\}$, 由单词组成的集合;

语料库 C , 由单词组成的文本序列;

单词 w 的上下文 $Context(w)$, 由语料库中单词 w 的前 c 个单词和后 c 个单词组成的文本序列。

CBOW模型网络结构

输入层: $\mathbf{v}(Context(w)_1), \mathbf{v}(Context(w)_2), \dots, \mathbf{v}(Context(w)_{2c}) \in \mathbb{R}^m$;

投影层: $\mathbf{x}_w = \sum_{i=1}^{2c} \mathbf{v}(Context(w)_i) \in \mathbb{R}^m$;

输出层: $T_{Huff}(\mathbf{x}_w) = s_{q(\mathbf{x}_w)}, s \in \mathbb{R}^N, q: \mathbb{R}^m \rightarrow \{1, 2, \dots, N\}$ 。

基于Hierarchical softmax的CBOW模型

记

$$p^w = (p_1^w, p_2^w, \dots, p_{l^w}^w)$$

为从根节点出发到达 w 对应的叶子结点的路径。其中, l^w 为路径长度, 即路径中结点数目; p_i^w 为路径中的结点, p_1^w 为根结点, $p_{l^w}^w$ 为 w 对应的叶子结点。

记

$$d^w = (d_2^w, d_3^w, \dots, d_{l^w}^w)$$

为 w 的Huffman编码。其中, $d_i^w \in \{0, 1\}$ 为路径 p^w 中第 i 个结点对应的编码 (根结点对应编码)。

记

$$\theta^w = (\theta_1^w, \theta_2^w, \dots, \theta_{l^w-1}^w)$$

为路径 p^w 中非叶子结点对应的参数向量。其中, $\theta_i^w \in \mathbb{R}^m$ 为路径 p^w 中第 i 个非叶子结点对应的参数向量。

条件概率

$$p(w|Context(w)) = \prod_{j=2}^{l^w} p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w)$$

其中

$$p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w) = \begin{cases} \sigma(\mathbf{x}_w^\top \theta_{j-1}^w), & d_j^w = 0; \\ 1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w), & d_j^w = 1, \end{cases}$$

或者

$$p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w) = [\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{1-d_j^w} \cdot [1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{d_j^w}$$

对数似然函数

$$\begin{aligned}\mathcal{L} &= \sum_{w \in \mathcal{C}} \log \prod_{j=2}^{l^w} p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w) \\ &= \sum_{w \in \mathcal{C}} \sum_{j=2}^{l^w} \{ (1 - d_j^w) \cdot \log [\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] + d_j^w \cdot \log [1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \}\end{aligned}$$

对数似然函数 \mathcal{L} 关于 θ_{j-1}^w 的梯度

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta_{j-1}^w} &= \frac{\partial}{\partial \theta_{j-1}^w} \left\{ \sum_{w \in \mathcal{C}} \sum_{j=2}^{l^w} \{ (1 - d_j^w) \cdot \log [\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] + d_j^w \cdot \log [1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \} \right\} \\ &= (1 - d_j^w) [1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \mathbf{x}_w - d_j^w \sigma(\mathbf{x}_w^\top \theta_{j-1}^w) \mathbf{x}_w \\ &= [1 - d_j^w - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \mathbf{x}_w\end{aligned}$$

θ_{j-1}^w 的更新

$$\theta_{j-1}^w = \theta_{j-1}^w + \eta [1 - d_j^w - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \mathbf{x}_w$$

其中, η 为学习率。

对数似然函数 \mathcal{L} 关于 \mathbf{x}_w 的梯度

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_w} = \sum_{j=2}^{l^w} [1 - d_j^w - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \theta_{j-1}^w$$

$\mathbf{v}(\tilde{w})$ 的更新

$$\mathbf{v}(\tilde{w}) = \mathbf{v}(\tilde{w}) + \eta \frac{\partial \mathcal{L}}{\partial \mathbf{x}_w}$$

其中, $\tilde{w} \in \text{Context}(w)$ 。

Skip-gram模型网络结构

输入层: $\mathbf{v}(w) \in \mathbb{R}^m$

输出层: $T_{Huff}(\mathbf{v}_w) = s_{q(\mathbf{v}_w)}, s \in \mathbb{R}^N, q: \mathbb{R}^m \rightarrow \{1, 2, \dots, N\}$

基于Hierarchical softmax的Skip-gram模型

条件概率

$$p(\text{Context}(w) | w) = \prod_{u \in \text{Context}(w)} p(u | w)$$

其中

$$p(u | w) = \prod_{j=2}^{l^u} p(d_j^u | \mathbf{v}(w), \theta_{j-1}^u)$$

且

$$p(d_j^u | \mathbf{v}(w), \theta_{j-1}^u) = [\sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)]^{1-d_j^u} \cdot [1 - \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)]^{d_j^u}$$

对数似然函数

$$\begin{aligned}\mathcal{L} &= \sum_{w \in C} \log \prod_{u \in \text{Context}(w)} \prod_{j=2}^{l^u} \left\{ \left[\sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right]^{1-d_j^u} \cdot \left[1 - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right]^{d_j^u} \right\} \\ &= \sum_{w \in C} \sum_{u \in \text{Context}(w)} \sum_{j=2}^{l^u} \left\{ (1 - d_j^u) \cdot \log \left[\sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] + d_j^u \cdot \log \left[1 - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] \right\}\end{aligned}$$

对数似然函数 \mathcal{L} 关于 θ_{j-1}^u 的梯度

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta_{j-1}^u} &= \frac{\partial}{\partial \theta_{j-1}^u} \left\{ \sum_{w \in C} \sum_{u \in \text{Context}(w)} \sum_{j=2}^{l^u} \left\{ (1 - d_j^u) \cdot \log \left[\sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] + d_j^u \cdot \log \left[1 - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] \right\} \right\} \\ &= \sum_{w \in C} \left\{ (1 - d_j^u) \left[1 - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] \mathbf{v}(w) - d_j^u \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \mathbf{v}(w) \right\} \\ &= \sum_{w \in C} \left[1 - d_j^u - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] \mathbf{v}(w)\end{aligned}$$

θ_{j-1}^u 的更新

$$\theta_{j-1}^u = \theta_{j-1}^u + \eta \sum_{w \in C} \left[1 - d_j^u - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] \mathbf{v}(w)$$

其中, η 为学习率。

对数似然函数 \mathcal{L} 关于 $\mathbf{v}(w)$ 的梯度

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}(w)} = \sum_{u \in \text{Context}(w)} \sum_{j=2}^{l^u} \left[1 - d_j^u - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] \theta_{j-1}^u$$

$\mathbf{v}(w)$ 的跟新

$$\mathbf{v}(w) = \mathbf{v}(w) + \eta \sum_{u \in \text{Context}(w)} \sum_{j=2}^{l^u} \left[1 - d_j^u - \sigma \left(\mathbf{v}(w)^\top \theta_{j-1}^u \right) \right] \theta_{j-1}^u$$

基于Negative Sampling的CBOW模型

设 $\text{Context}(w)$ 的负样本子集为

$$\text{NEG}(w) \neq \emptyset$$

对于 $\forall \tilde{w} \in D$, 定义

$$L^w(\tilde{w}) = \begin{cases} 1, & \tilde{w} = w \\ 0, & \tilde{w} \neq w \end{cases}$$

表示词 \tilde{w} 的标签, 正样本标签为1, 负样本标签为0。

关于字典 D 的子集 $\{w\} \cup NEG(w)$ 的似然函数

$$g(w) = \prod_{u \in \{w\} \cup NEG(w)} p(u|Context(w)) = \sigma(\mathbf{x}_w^\top \theta^w) \prod_{u \in NEG(w)} [1 - \sigma(\mathbf{x}_w^\top \theta^w)]$$

其中

$$p(u|Context(w)) = \begin{cases} \sigma(\mathbf{x}_w^\top \theta^u), L^w(u) = 1 \\ 1 - \sigma(\mathbf{x}_w^\top \theta^u), L^w(u) = 0 \end{cases}$$

或者

$$p(u|Context(w)) = [\sigma(\mathbf{x}_w^\top \theta^u)]^{L^w(u)} \cdot [1 - \sigma(\mathbf{x}_w^\top \theta^u)]^{1-L^w(u)}$$

\mathbf{x}_w 为 $Context(w)$ 词向量之和, $\theta^u \in \mathbb{R}^m$ 为模型参数。

关于语料库 C 的对数似然函数

$$\begin{aligned} \mathcal{L} &= \log \prod_{w \in C} g(w) = \sum_{w \in C} \log g(w) \\ &= \sum_{w \in C} \log \prod_{u \in \{w\} \cup NEG(w)} \left\{ [\sigma(\mathbf{x}_w^\top \theta^u)]^{L^w(u)} \cdot [1 - \sigma(\mathbf{x}_w^\top \theta^u)]^{1-L^w(u)} \right\} \\ &= \sum_{w \in C} \sum_{u \in \{w\} \cup NEG(w)} \left\{ L^w(u) \cdot \log[\sigma(\mathbf{x}_w^\top \theta^u)] + [1 - L^w(u)] \cdot \log[1 - \sigma(\mathbf{x}_w^\top \theta^u)] \right\} \end{aligned}$$

对数似然函数 \mathcal{L} 关于 θ^u 的梯度

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta^u} &= \frac{\partial}{\partial \theta^u} \left\{ \sum_{w \in C} \sum_{u \in \{w\} \cup NEG(w)} \left\{ L^w(u) \cdot \log[\sigma(\mathbf{x}_w^\top \theta^u)] + [1 - L^w(u)] \cdot \log[1 - \sigma(\mathbf{x}_w^\top \theta^u)] \right\} \right\} \\ &= L^w(u) [1 - \sigma(\mathbf{x}_w^\top \theta^u)] \mathbf{x}_w - [1 - L^w(u)] \sigma(\mathbf{x}_w^\top \theta^u) \mathbf{x}_w \\ &= [L^w(u) - \sigma(\mathbf{x}_w^\top \theta^u)] \mathbf{x}_w \end{aligned}$$

θ^u 的更新

$$\theta^u = \theta^u + \eta [L^w(u) - \sigma(\mathbf{x}_w^\top \theta^u)] \mathbf{x}_w$$

对数似然函数 \mathcal{L} 关于 \mathbf{x}_w 的梯度

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_w} = \sum_{u \in \{w\} \cup NEG(w)} [L^w(u) - \sigma(\mathbf{x}_w^\top \theta^u)] \theta^u$$

$\mathbf{v}(\tilde{w})$ 的更新

$$\mathbf{v}(\tilde{w}) = \mathbf{v}(\tilde{w}) + \eta \frac{\partial \mathcal{L}}{\partial \mathbf{x}_w}$$

其中, $\tilde{w} \in Context(w)$ 。

基于Negative Sampling的Skip-gram模型

关于字典 D 的子集 $\{w\} \cup NEG^{\tilde{w}}(w)$ 的似然函数

$$g(w) = \prod_{\tilde{w} \in Context(w)} \prod_{u \in \{w\} \cup NEG^{\tilde{w}}(w)} p(u|\tilde{w})$$

其中

$$p(u|\tilde{w}) = \begin{cases} \sigma(\mathbf{v}(\tilde{w})^\top \theta^u), L^w(u) = 1 \\ 1 - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u), L^w(u) = 0 \end{cases}$$

或者

$$p(u|\tilde{w}) = [\sigma(\mathbf{v}(\tilde{w})^\top \theta^u)]^{L^w(u)} \cdot [1 - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)]^{1-L^w(u)}$$

$NEG^{\tilde{w}}(w)$ 为处理词 \tilde{w} 时生成的负样本子集。

关于语料库 C 的对数似然函数

$$\begin{aligned} \mathcal{L} &= \log \prod_{w \in C} g(w) = \sum_{w \in C} \log g(w) \\ &= \sum_{w \in C} \log \prod_{\tilde{w} \in Context(w)} \prod_{u \in \{w\} \cup NEG^{\tilde{w}}(w)} \left\{ [\sigma(\mathbf{v}(\tilde{w})^\top \theta^u)]^{L^w(u)} \cdot [1 - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)]^{1-L^w(u)} \right\} \\ &= \sum_{w \in C} \sum_{\tilde{w} \in Context(w)} \sum_{u \in \{w\} \cup NEG^{\tilde{w}}(w)} \{ L^w(u) \cdot \log[\sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] + [1 - L^w(u)] \cdot \log[1 - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] \} \end{aligned}$$

对数似然函数 \mathcal{L} 关于 θ^u 的梯度

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta^u} &= \frac{\partial}{\partial \theta^u} \left\{ \sum_{w \in C} \sum_{\tilde{w} \in Context(w)} \sum_{u \in \{w\} \cup NEG^{\tilde{w}}(w)} \{ L^w(u) \cdot \log[\sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] + [1 - L^w(u)] \cdot \log[1 - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] \} \right\} \\ &= L^w(u) [1 - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] \mathbf{v}(\tilde{w}) - [1 - L^w(u)] \sigma(\mathbf{v}(\tilde{w})^\top \theta^u) \mathbf{v}(\tilde{w}) \\ &= [L^w(u) - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] \mathbf{v}(\tilde{w}) \end{aligned}$$

θ^u 的更新

$$\theta^u = \theta^u + \eta [L^w(u) - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] \mathbf{v}(\tilde{w})$$

对数似然函数 \mathcal{L} 关于 $\mathbf{v}(\tilde{w})$ 的梯度

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}(\tilde{w})} = \sum_{u \in \{w\} \cup NEG^{\tilde{w}}(w)} [L^w(u) - \sigma(\mathbf{v}(\tilde{w})^\top \theta^u)] \theta^u$$

$\mathbf{v}(\tilde{w})$ 的更新

$$\mathbf{v}(\tilde{w}) = \mathbf{v}(\tilde{w}) + \eta \frac{\partial \mathcal{L}}{\partial \mathbf{v}(\tilde{w})}$$

负采样算法

设词典 D 中词 w_i 对应线段 $l(w_i)$ ，长度为

$$len(w_i) = \frac{counter(w_i)}{\sum_{u \in D} counter(u)}$$

其中， $counter(\cdot)$ 为词在语料 C 中的出现次数。可将线段 $l(w_1) \cdots l(w_N)$ 拼接为长度为1的单位线段。

记

$$l_0 = 0$$

$$l_k = \sum_{j=1}^k \text{len}(w_j), k = 1, 2, \dots, N$$

则以 $\{l_j\}_{j=0}^N$ 为剖分点可得到区间 $[0, 1]$ 上的一个非等距剖分

$$I_i = (l_{i-1}, l_i], i = 1, 2, \dots, N$$

在区间 $[0, 1]$ 上以剖分点 $\{m_j\}_{j=0}^M$ 做等距剖分，其中 $M \gg N$ 。

将等距剖分的内部点 $\{m_j\}_{j=1}^{M-1}$ 投影到非等距剖分。则可建立 $\{m_j\}_{j=1}^{M-1}$ 与区间 $\{I_j\}_{j=1}^N$ 的映射，进一步建立与词 $\{w_j\}_{j=1}^M$ 之间的映射

$$w_k = \text{Table}(i), m_i \in I_k, i = 1, 2, \dots, M-1$$