

Word2vec

词向量基础

One-hot-representation

distributed representation

2. CBOW (continuous Bag of words) skip-gram

用于神经网络语言模型

采用三层的神经网络结构(可以多层)

输入
隐藏层
softmax 输出

CBOW: 输入: 词上下文相关词对应词向量

输出: 特定词的词向量

词袋模型, 输入词平等的, 未考虑位置关系

输入 87 词向量, 输出所有词, softmax 概率

期望目标词概率最大

skipgram, 输入: 特定词词向量

输出: softmax 概率排前 8 的 87 词

词表过大时有问题

3. word2vec Huffman 树

用霍夫曼树来代替隐藏层和输出层神经元

霍夫曼树保证高频词节点路径较短, 编码值较短, 反之低频, 常用词编码更短

word2vec 中约定左子树编码为 1, 右子树为 0, 同时约定左子树权重不小于右子树

4 基于 Hierarchical softmax 模型描述

① 输入层到隐藏层, 对所有输入词向量和取平均, 得到一词向量

② $H \rightarrow Output$, 采用霍夫曼树

沿着 Huffman 树走, 未路径到叶子节点, 叶子节点就是词

③ 采用二元逻辑回归:

沿左: 负类 (编码 1)

沿右: 正类 (编码 0)

判断正负:

$$p_H = \sigma(xw) = \frac{1}{1 + e^{-xw}}$$

xw 是当前内部节点词向量 w 的权重

优点: 1) 计算量 $V \rightarrow \log_2 V$

2) 高频词靠近根, 更快被找到, 贪心优化

基于 Hierarchical softmax 模型梯度计算

① 优化描述:

输入词 w , 相关词向量 xw , 相关叶子路径分布

总数 l_w , w 经过的第 j 节点表示为 p_j^w , 对 $j=1, 2, 3, \dots, l_w$

Huffman 编码 $d_j^w \in \{0, 1\}$ $j \in 2, 3, \dots, l_w$

节点序号 $\theta_j^w = 1, 2, \dots, l_w$

w 经过 Huffman 树某节点 j 的概率 $\sigma(xw \theta_j^w)$

$$P(d_j^w | xw, \theta_j^w) = \begin{cases} \sigma(xw \theta_j^w) & d_j^w = 0 \\ 1 - \sigma(xw \theta_j^w) & d_j^w = 1 \end{cases}$$

对某目标输出词 w , 最大似然

$$\prod_{j=2}^{l_w} P(d_j^w | xw, \theta_j^w) = \prod_{j=2}^{l_w} [\sigma(xw \theta_j^w)]^{1-d_j^w} [1 - \sigma(xw \theta_j^w)]^{d_j^w}$$

每次只使用一个样本更新梯度 (随机梯度上升)

$$L = \log \prod_{j=2}^{l_w} P(d_j^w | xw, \theta_j^w) = \sum_{j=2}^{l_w} [(1-d_j^w) \log \sigma(xw \theta_j^w) + d_j^w \log (1 - \sigma(xw \theta_j^w))]$$

$$(\sigma x) = \sigma(x) (1 - \sigma(x))$$

$$\frac{\partial L}{\partial \theta_j^w} = (1-d_j^w) \frac{\sigma(xw \theta_j^w) (1 - \sigma(xw \theta_j^w))}{\sigma(xw \theta_j^w)^2} xw$$

$$- d_j^w \frac{\sigma(xw \theta_j^w) (1 - \sigma(xw \theta_j^w))}{1 - \sigma(xw \theta_j^w)^2} xw$$

$$= [1-d_j^w - \sigma(xw \theta_j^w)] xw$$

同理

$$\frac{\partial L}{\partial xw} = \sum_{j=2}^{l_w} (1-d_j^w - \sigma(xw \theta_j^w)) \theta_j^w$$

② 基于该方法的 CBOW

a) 建立 Huffman 树

2) 随机初始 θ_j, z_{θ_j}

3) a) $z_{\theta_j} = \frac{1}{2c} \sum_{i=1}^c x_i$

b) for $j=2$ to l_w 计算

$$f = \sigma(xw \theta_j^w)$$

$$g = (1-d_j^w - f) \eta$$

$$e = e + g \theta_j^w$$

$$\theta_j^w = \theta_j^w + g xw$$

$$x_i = x_i + e$$

② 基于该方法的 skip-gram 模型

我们期望 $P(x_i | x_w)$ $i=1,2,\dots,2c$ 最大。上下文相关的，在期望 $P(x_i | x_w)$ 最大同时， $P(x_w | x_i)$ 也最大。
选择后者，在 γ 迭代窗口内，不只更新 γ_i ，而 $2c$ 个词，整体迭代更新。
所以 skip-gram 不对输入更新，而是 $2c$ 输出。

a) 建立 Huffman 树

b) 随机初始化

c) 1) for $i=1$ to $2c$

2) 对 $j=1$ to $2c$

计算：
 $f = \sigma(x_w^T \theta_j^w)$
 $g = (1 - d_j^w - f) \eta$
 $e = e + g \theta_j^w$
 $\theta_j^{w+1} = \theta_j^w + g x_i$

iii) $x_i = x_i + e$

缺点：虽然提高了效率，但对于训练词频，需要很多迭代。

6. 基于负采样的模型

① 中心词是 w ， $\text{context}(w)$ 有 $2c$ 个， w 是 γ 正例。负采样得到 $\text{neg } \gamma$ 和 w 不同 i 中心词，利用 γ 正 $\text{neg } \gamma$ 负二元逻辑回归。

② 如何通过 γ 正例 γ 负例二元逻辑回归

正： $P(\text{context}(w), w_i) = \sigma(x_w^T \theta^{w_i})$ $y_i=1, i=0$

负： $P(\text{context}(w_0), w_i) = 1 - \sigma(x_{w_0}^T \theta^{w_i})$ $y_i=0, i=1, \dots, \text{neg}$

$\prod_{i=0}^{\text{neg}} P(\text{context}(w_0), w_i) = \sigma(x_{w_0}^T \theta^{w_0}) \prod_{i=1}^{\text{neg}} (1 - \sigma(x_{w_0}^T \theta^{w_i}))$

似然函数：

$\prod_{i=0}^{\text{neg}} \sigma(x_{w_0}^T \theta^{w_i})^{y_i} (1 - \sigma(x_{w_0}^T \theta^{w_i}))^{1-y_i}$

对数似然：

$L = \sum_{i=0}^{\text{neg}} y_i \log(\sigma(x_{w_0}^T \theta^{w_i})) + (1-y_i) \log(1 - \sigma(x_{w_0}^T \theta^{w_i}))$

$\frac{\partial L}{\partial \theta^{w_i}} = (y_i - \sigma(x_{w_0}^T \theta^{w_i})) x_{w_0}$ $\frac{\partial L}{\partial x^{w_i}} = \sum_{i=0}^{\text{neg}} (y_i - \sigma(x_{w_0}^T \theta^{w_i})) \theta^{w_i}$

③ 如何进行负采样？

$$\text{len}(w) = \frac{\text{count}(w)^{\frac{3}{4}}}{\sum \text{count}(w)^{\frac{3}{4}}}$$

将长度为 $\text{len}(w)$ 的词段作为 $\text{neg } \gamma$ 位置。
从 $\text{neg } \gamma$ 位置中采样 $\text{neg } \gamma$ 位置。
每个位置所在词频就是采样的词 (负例)。
 $\text{neg}(\text{default}) = 10$