

Word2Vec

背景 *onehot* 用向量表达词意义

两个模型

Skip-gram

$$\max \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w^{(t+j)} | w^{(t)}),$$
 总词数 T 背景词 中心词
 Skip-gram 最大似然估计
 每个词都是独立

等价
log

$$\min - \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w^{(t+j)} | w^{(t)}).$$
 损失函数

套入向量

引入所有的词, 导致
梯度训练很慢

$$P(w_o | w_c) = \frac{\exp(u_o^T v_c)}{\sum_{i \in V} \exp(u_i^T v_c)}$$

$$\log P(w_o | w_c) = u_o^T v_c - \log \left(\sum_{i \in V} \exp(u_i^T v_c) \right)$$
 u_o 是 w_o 的向量表达
 v_c 是 w_c 的向量表达

求微分进行梯度下降

$$\frac{\partial \log P(w_o | w_c)}{\partial v_c} = u_o - \frac{\sum_{j \in V} \exp(u_j^T v_c) u_j}{\sum_{i \in V} \exp(u_i^T v_c)}$$

$$= u_o - \sum_{j \in V} \left(\frac{\exp(u_j^T v_c)}{\sum_{i \in V} \exp(u_i^T v_c)} \right) u_j$$

$$\frac{\partial \log P(w_c | W_o)}{\partial v_{o1}} = \frac{1}{2m} \left(u_c - \sum_{j \in V} \frac{\exp(u_j^T \bar{v}_o) u_j}{\sum_{i \in V} \exp(u_i^T \bar{v}_o)} \right) = \frac{1}{2m} \left(u_c - \sum_{j \in V} P(w_j | W_o) u_j \right).$$

CBOW

$$\max \prod_{t=1}^T P(w^{(t)} | w^{(t-m)}, \dots, w^{(t-1)}, w^{(t+1)}, \dots, w^{(t+m)}).$$

$$\min - \sum_{t=1}^T \log P(w^{(t)} | w^{(t-m)}, \dots, w^{(t-1)}, w^{(t+1)}, \dots, w^{(t+m)}).$$

$$P(u_c | w_{o1}, \dots, w_{o2m}) = \frac{\exp(\frac{1}{2m} u_c^T (v_{o1} + \dots + v_{o2m}))}{\sum_{i \in V} \exp(\frac{1}{2m} u_i^T (v_{o1} + \dots + v_{o2m}))}$$

$$\frac{\partial \log P(w_c | W_o)}{\partial v_{o1}} = \frac{1}{2m} \left(u_c - \sum_{j \in V} \frac{\exp(u_j^T \bar{v}_o) u_j}{\sum_{i \in V} \exp(u_i^T \bar{v}_o)} \right) = \frac{1}{2m} \left(u_c - \sum_{j \in V} P(w_j | W_o) u_j \right).$$

随机梯度

注: 如果用梯度下降的话, 开销很大, 实际很难学习, 所以引出下面训练方法

两个近似的训练方法

负采样

负采样

我们以跳字模型为例讨论负采样。

词典 \mathcal{V} 大小之所以会在目标函数中出现，是因为中心词 w_c 生成背景词 w_o 的概率 $\mathbb{P}(w_o | w_c)$ 使用了softmax，而softmax正是考虑了背景词可能是词典中的任一词，并体现在softmax的分母上。

我们不妨换个角度，假设中心词 w_c 生成背景词 w_o 由以下两个相互独立事件联合组成来近似

- 中心词 w_c 和背景词 w_o 同时出现在该训练数据窗口
- 中心词 w_c 和第1个噪声词 w_1 不同时出现在该训练数据窗口（噪声词 w_1 按噪声词分布 $\mathbb{P}(w)$ 随机生成）
- ...
- 中心词 w_c 和第 K 个噪声词 w_K 不同时出现在该训练数据窗口（噪声词 w_K 按噪声词分布 $\mathbb{P}(w)$ 随机生成）

我们可以使用 $\sigma(x) = 1/(1 + \exp(-x))$ 函数来表达中心词 w_c 和背景词 w_o 同时出现在该训练数据窗口的概率：

公式

$$P(D = 1 | w_c, w_o) = \sigma(\mathbf{u}_o^\top \mathbf{v}_c),$$

$$\max \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(D = 1 | w^{(t)}, w^{(t+j)}).$$

然而，以上模型中包含的事件只考虑了正负样本。这导致当所有词向量相等且值为无穷大时，以上的联合概率才被最大化为1。很明显，这样的词向量毫无意义。也导致难以直接对词向量施加梯度下降目标函数来求最优。设背景词 w_o 出现在中心词 w_c 的一个数据窗口为事件 D ，我们记 $D=1$ 和 $P(w)$ 采样 K 个未出现在该数据窗口中的词，即噪声词。 w_k ($k = 1, \dots, K$) 不出现在中心词 w_c 的数据窗口为事件 N_k 。假设同时含有正负样本和负类样本的事件 P, N_1, \dots, N_K 相互独立，负采样将以上需要最大化的仅考虑正负样本的联合概率改写为

$$\max \log \mathbb{P}(w_o | w_c) = \log[\mathbb{P}(D = 1 | w_o, w_c) \prod_{k=1, w_k \sim \mathbb{P}(w)}^K \mathbb{P}(D = 0 | w_k, w_c)]$$

$$\max \log \mathbb{P}(w_o | w_c) = \log \frac{1}{1 + \exp(-\mathbf{u}_o^\top \mathbf{v}_c)} + \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log[1 - \frac{1}{1 + \exp(-\mathbf{u}_k^\top \mathbf{v}_c)}]$$

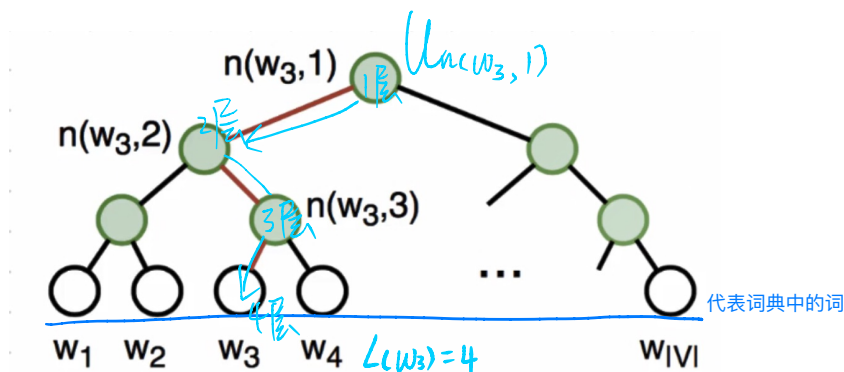
$$\min -\log \mathbb{P}(w_o | w_c) = -\log \frac{1}{1 + \exp(-\mathbf{u}_o^\top \mathbf{v}_c)} - \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log \frac{1}{1 + \exp(\mathbf{u}_k^\top \mathbf{v}_c)}$$

注：计算大小从词数量 V 降到 K

Trick:噪声词设为单词概率的3/4次方（0.99 0.01 各自的3/4次方）

层次Soft max

用了哈夫曼树，根据词频（字频）构造



$$P(w_o | w_c) = \prod_{j=1}^{L(w_o)-1} \sigma \left(\left[n(w_o, j+1) = \text{leftChild}(n(w_o, j)) \right] \cdot \mathbf{u}_{n(w_o, j)}^T \mathbf{v}_c \right),$$

真! 假-1
我的下一步是不是左子树
wc对应向量

例子

step 1, 正 step 2, 负 step, 正

$$P(w_3 | w_i) = \sigma(\mathbf{u}_{n(w_3,1)}^T \mathbf{v}_i) \cdot \sigma(-\mathbf{u}_{n(w_3,2)}^T \mathbf{v}_i) \cdot \sigma(\mathbf{u}_{n(w_3,3)}^T \mathbf{v}_i)$$

$$\sum_{w=1}^V P(w | w_i) = 1$$

Question:

- 1、如何处理无意义的高频词
- 2、如何学习词组
- 3、噪声选3/4平方原因

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

按词频公式去掉

用有限的非叶子结点向量表达了所要生成表达词语的向量

unigram and bigram count
如果某些词语经常在一起就判定他们是短语
输入两个词若大于某值就认定它们在一起

$$\text{score}(w_i, w_j) = \frac{\text{count}(w_i w_j) - \delta}{\text{count}(w_i) \times \text{count}(w_j)}.$$

通常拿2-4个词

在取样的过程中，我们把每个词频的值取四分之三次方的时候效果最好。这样做是因为对频率较大的数进行取3/4次幂，原数字损失的多，反之较小的，去掉的就少，这样做是一种平滑方式。