

# 连续词袋模型 (CBOW)

与skip-gram 模型相反，CBOW是由背景词推测中心词的模型。

如果输入矩阵是 $I$ 而输出矩阵是 $O$ 。矩阵的形状都是 $K \times H$ ， $K$ 是代表词典，也词典大小， $H$ 是我们设定的中间向量长度。

假设背景词 $t$ 的独热向量 $\hat{t}$ ，则：

$$h = I^T \cdot \hat{t} \quad (H \times 1) \quad (1)$$

$$\hat{o} = O \cdot h \quad (K \times 1) \quad (2)$$

将(1)带入到(2)得到 $O \cdot I^T \cdot \hat{t}$ ，是 $K$ 维向量，记作 $\hat{o}$ ，我们将 $\hat{o}$ 做 softmax 处理，得到背景词对于词典的似然估计 $P(t|K)$ ：

$$P(t|K) = \left[ \frac{e^{\hat{o}_1}}{\sum_{d=1}^K e^{\hat{o}_d}} \quad \cdots \quad \frac{e^{\hat{o}_K}}{\sum_{d=1}^K e^{\hat{o}_d}} \right]^T$$

其中对目标中心词 $r$ 的似然估计，我们通过 $P(t|r) = \hat{r}^T \cdot P(t|K)$ ， $\hat{r}$ 是中心词的独热向量

如果背景词 $t$ 在词典 $K$ 中的索引是 $j$ ，那么 $I^T \cdot \hat{t}$ 就是输入矩阵 $I$ 的第 $j$ 行，也就是说：

$$h = I^T \cdot \hat{t} = I_j^T$$

，而

$$\hat{o} = O \cdot h = [O_1 \cdot I_j^T \quad \cdots \quad O_K \cdot I_j^T]^T$$

，其中的 $O_1 \cdots O_K$ 是矩阵 $O$ 的各行，那么：

$$P(t|K) = \left[ \frac{e^{O_1 \cdot I_j^T}}{\sum_{d=1}^K e^{O_d \cdot I_j^T}} \cdots \frac{e^{O_i \cdot I_j^T}}{\sum_{d=1}^K e^{O_d \cdot I_j^T}} \cdots \frac{e^{O_K \cdot I_j^T}}{\sum_{d=1}^K e^{O_d \cdot I_j^T}} \right]$$

假设中心词 $r$ 在词典中的索引是 $i$ ，则：

$$p(t|r) = \frac{e^{O_i \cdot I_j^T}}{\sum_{d=1}^K e^{O_d \cdot I_j^T}} \quad (3)$$

需要注意这里的 $O_i$ 和 $I_j$ 都是行向量

因为背景词往往是多个(假设中心词是 $r$ ，在文章中的位置是 $p$ ，窗口尺寸是 $w$ ，则背景词分别是 $t_{p-w}, t_{p+1-w} + \cdots, t_{p-1}, t_{p+1}, t_{p+2} + \cdots + t_{p+w}$ 。  $2w$ 个背景词对应的词典索引是 $j_1$ 到 $j_{2w}$ ，对应独热向量分别是： $\hat{t}_{j_1} \cdots \hat{t}_{j_{2w}}$ )，所以我们这里的背景词应该采用背景词的向量的平均值，

也就是说

$$\begin{aligned} I_j^T &= \frac{1}{2w} I^T \cdot (\hat{t}_{j_1} + \dots + \hat{t}_{j_{2w}}) \\ &= \frac{1}{2w} (I_{j_1}^T + \dots + I_{j_{2w}}^T) \end{aligned} \quad (4)$$

我们已知中心词对背景词的似然估计，通过交叉熵损失函数得到损失函数：

$$\begin{aligned} l(\Theta) &= -\log P(t|r) \\ &= -\log \frac{e^{O_i \cdot I_j^T}}{\sum_{d=1}^K e^{O_d \cdot I_j^T}} \end{aligned} \quad (5)$$

将公式（4）带入到公式（5）：

$$\begin{aligned} l(\Theta) &= -\log \frac{e^{\frac{1}{2w} \cdot O_i \cdot (I_{j_1}^T + \dots + I_{j_{2w}}^T)}}{\sum_{d=1}^K e^{\frac{1}{2w} \cdot O_d \cdot (I_{j_1}^T + \dots + I_{j_{2w}}^T)}} \\ &= -\left( \log e^{\frac{1}{2w} \cdot O_i \cdot (I_{j_1}^T + \dots + I_{j_{2w}}^T)} - \log \sum_{d=1}^K e^{\frac{1}{2w} \cdot O_d \cdot (I_{j_1}^T + \dots + I_{j_{2w}}^T)} \right) \end{aligned} \quad (6)$$

我们把公式（6）中  $\frac{1}{2w} \cdot (I_{j_1}^T + \dots + I_{j_{2w}}^T)$  用  $\overline{I_j^T}$  替代，公式（6）简化成：

$$l(\Theta) = -(O_i \cdot \overline{I_j^T} - \log \sum_{d=1}^K e^{O_d \cdot \overline{I_j^T}}) \quad (7)$$

公式（7）是损失函数，我们对其求  $\overline{I_j^T}$  的偏微分，得到  $\overline{I_j^T}$  的梯度：

$$\begin{aligned} \frac{\partial l(\Theta)}{\partial \overline{I_j^T}} &= -\left( O_i - \frac{e^{O_1 \cdot \overline{I_j^T}} \cdot O_1 + \dots + e^{O_K \cdot \overline{I_j^T}} \cdot O_K}{\sum_{d=1}^K e^{O_d \cdot \overline{I_j^T}}} \right) \\ &= -\left( O_i - \sum_{d=1}^K \frac{e^{O_d \cdot \overline{I_j^T}} \cdot O_d}{\sum_{d=1}^K e^{O_d \cdot \overline{I_j^T}}} \right) \\ &= \sum_{d=1}^K P(t|r_d) \cdot O_d - O_i \end{aligned} \quad (8)$$