

# Doc2vec 模型结构和训练

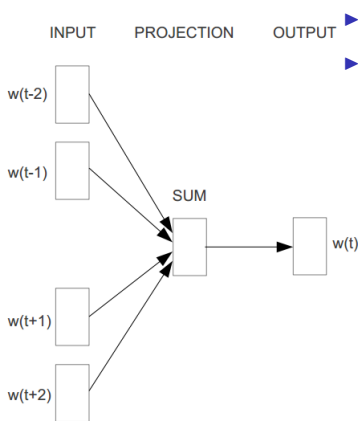
方敏

驰声研发

2016 年 10 月

- ▶ 背景知识: word2vec 模型
- ▶ Doc2vec 模型结构介绍
- ▶ 模型训练与 doc 向量获取
- ▶ Optimizing Computational Efficiency
- ▶ 其它模型

# 背景知识: word2vec 模型



► 预测  $p(w_t | w_{t-k}, \dots, w_{t+k})$

► 目标函数如下

$$\frac{1}{T} \sum_{t=k}^{T-k} \log p(w_t | w_{t-k}, \dots, w_{t+k})$$

$$p(w_t | w_{t-k}, \dots, w_{t+k}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}}$$

$$y_{w_t} = b + Uh(w_{t-k}, \dots, w_{t+k}) \quad (1)$$

$h$  表示 INPUT 到 PROJECTION 的转化关系, 如求和

图: CBOW(Continuous Bag-of-Words Model)

# Doc2vec 模型结构介绍

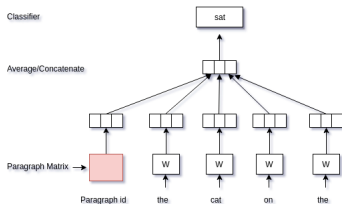


图: Doc2vec PV-DM

- ▶ 预测  $p(w_{t,j} | w_{t-k,j}, \dots, w_{t+k,j}, d_j^j)$   
 $w_{t,j}$  表示第  $j$  篇 doc 的第  $t$  个词,  $d_j^j$  表示第  $j$  篇 doc,  $k$  为  $1/2$  窗长
- ▶ 分类器

$$p(w_{t,j} | w_{t-k,j}, \dots, w_{t+k,j}) = \frac{e^{y_{w_{t,j}}}}{\sum_i e^{y_i}}$$
$$y_{w_{t,j}} = b + Uh(w_{t-k,j}, \dots, w_{t+k,j}, d_j) \quad (2)$$

目标函数：

$$\frac{1}{J} \sum_{j=1}^J \frac{1}{N_j - 2k} \sum_{t=k}^{N_j-k} \log p(w_{t,j} | w_{t-k,j}, \dots, w_{t+k,j}, d_j) \quad (3)$$

其中， $N_j$  表示第  $j$  篇 doc 的字符长度， $J$  表示语料中包含的 doc 数目

若不考虑使用输出反馈：

$$h_t = f(W_{xh}x_t + W_{hh}h_{t-1}) \quad (4)$$

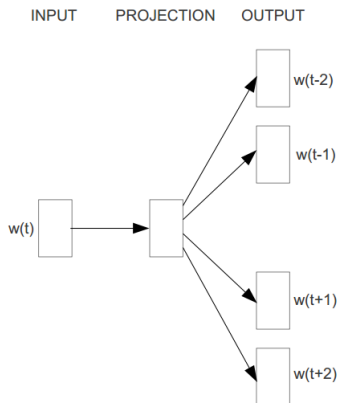
$$y_t = g(W_{hy}h_t) \quad (5)$$

# Optimizing Computational Efficiency

- ▶ Negative Sampling
- ▶ Hierarchical Softmax

$$\begin{aligned} \frac{1}{T} \sum_{t=k}^{T-k} \log p(w_t | w_{t-k}, \dots, w_{t+k}) \\ p(w_t | w_{t-k}, \dots, w_{t+k}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}} \\ y_{w_t} = b + Uh(w_{t-k}, \dots, w_{t+k}) \end{aligned} \tag{6}$$

# 其它模型



► 预测  $p(w_t | w_{t-k}, \dots, w_{t+k})$

► 目标函数如下

$$\frac{1}{T} \sum_{t=k}^{T-k} \sum_{-k \leq l \leq k, l \neq 0} \log p(w_{t+l} | w_t)$$

$$p(w | w_t) = \frac{e^{y_w}}{\sum_i e^{y_i}}$$

$$y_w = b + Uh(w_t)$$

(7)

图: CBOW(Continuous Bag-of-Words Model)



# 其它模型

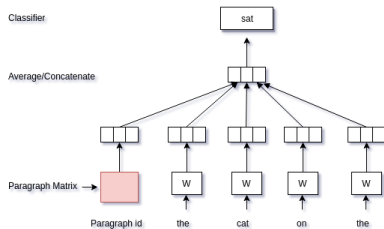


图: Doc2vec PV-DM

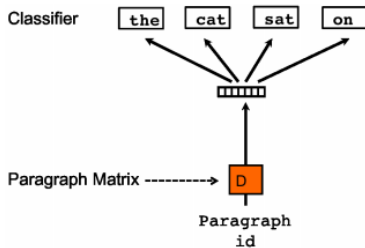


图: Doc2vec PV-DBOW