

概率潜在语义分析

1. 基本想法

一个文本 (document) 集合包含多个话题, 一个话题由若干单词组成, 话题无法直接观测需要数据估计.

目标: 发现隐变量表示的话题, 即潜在语义

	文本 1	2	3	4
单词 1	2	2	4	3
2	—	—	—	—
3	—	—	—	—
4				
5				

语义相近的单词, 语义相近的文本会被聚向相同的“软的类别”中

↓
话题

2. 生成模型

数据

(1) 单词集合 $W = \{w_1, w_2, \dots, w_M\}$

(2) 文本集合 $D = \{d_1, d_2, \dots, d_N\}$

(3) 话题集合 $Z = \{z_1, z_2, \dots, z_K\}$

观测数据 (word-doc) 矩阵

	$\Lambda(w, d)$	

$(d) \rightarrow (z) \rightarrow (w)$

① $\xrightarrow{P(d)}$ 文本 d

② $\xrightarrow{P(w|d)}$ 话题 z

③ $\xrightarrow{P(w|z)}$ 单词 w

概率表示

单词文本矩阵 T 的生成概率

$$P(T) = \prod_{(w,d)} P(w, d)^{\Lambda(w,d)}$$

$$p(w, d) = p(d) p(w|d) = p(d) \sum_z P(w, z|d)$$

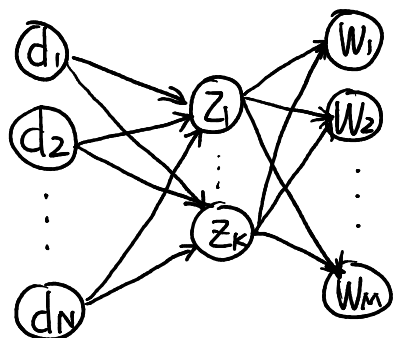
$$= p(d) \sum_z P(z|d) P(w|z)$$

3. 模型参数

(1) 如果认为未知参数是 $P(w, d)$ $M \times N$

(2) 参数个数 \rightarrow 引入话题 $O(M \times K + K \times N)$

其中 $K \ll M, N$, 大大减少了待估参数



4. 概率潜在语义模型求解

对数似然函数

$$\begin{aligned} L(\theta) &= \sum_{i=1}^M \sum_{j=1}^N \Delta(w_i, d_j) \log P(w_i, d_j) \\ &= \sum_{i=1}^M \sum_{j=1}^N \Delta(w_i, d_j) \log \left\{ P(d_j) \sum_k P(w_i | z_k) P(z_k | d_j) \right\} \end{aligned}$$

$r_{ijk} = 1$, 代表 d_j 产生了话题 z_k , z_k 产生了 w_i

$$P(w_i | z_k) P(z_k | d_j)$$

E步: (计算Q函数)

$$Q(\theta, \theta^{(i)}) = \sum_z \log P(Y, z | \theta) \cdot P(z | Y, \theta^{(i)})$$

完全样本似然函数

$$P(w, d, r | \theta) = \prod_{i=1}^N \prod_{j=1}^M P(w_i, d_j, r_{ij1}, r_{ij2}, \dots, r_{ijk})$$

$$= \prod_{i=1}^N \prod_{j=1}^M \left[\prod_{k=1}^K \{ P(w_i | z_k) P(z_k | d_j) \}^{r_{ijk}} \cdot P(d_j) \right]^{\Delta(w_i, d_j)}$$

$$\begin{aligned} \log P(w, d, z | \theta) &= \sum_{i=1}^N \sum_{j=1}^M \left[\Delta(w_i, d_j) \left\{ \log P(d_j) \right. \right. \\ &\quad \left. \left. + \sum_{k=1}^K r_{ijk} \log P(w_i | z_k) \cdot P(z_k | d_j) \right\} \right] \end{aligned}$$

$$r_{ijk}^{(t)} = E[r_{ijk} | w, d, \theta^{(t)}]$$

$$= P(z_k | w, d, \theta) = \frac{P(z_k, w_i | d_j)}{P(w_i | d_j)}$$

$$= \frac{P(w_i | z_k) P(z_k | d_j)}{\sum_k P(w_i | z_k) P(z_k | d_j)} \cdot P(w_i | z_k) P(z_k | d_j) \text{ 用上一步迭代代入}$$

M 步:

$$r_{ijk}^{(t)} = E\{r_{ijk} | w, d, \theta^{(t)}\}$$

$$Q(\theta, \theta^{(i)}) = \sum_{k=1}^K \sum_{i=1}^M \sum_{j=1}^N r_{ijk}^{(t)} \Delta(w_i, d_j) \log\{P(w_i | z_k) P(z_k | d_j)\}$$

$$\text{约束: } \sum_{i=1}^M P(w_i | z_k) = 1, \text{ for } k=1, 2, \dots, K$$

$$\sum_{k=1}^K P(z_k | d_j) = 1, \text{ for } j=1, 2, \dots, M$$

定义 Lagrange 乘数:

$$\begin{aligned} \mathcal{L}(\theta) = Q(\theta, \theta^{(i)}) &+ \sum_{k=1}^K z_k \left\{ 1 - \sum_{i=1}^M P(w_i | z_k) \right\} \\ &+ \sum_{j=1}^N d_j \left\{ 1 - \sum_{k=1}^K P(z_k | d_j) \right\} \end{aligned}$$

$$P(w_i | z_k) = \frac{\sum_{j=1}^N \Delta(w_i, d_j) r_{ijk}}{\sum_{i=1}^M \sum_{j=1}^N \Delta(w_i, d_j) r_{ijk}} \quad P(z_k | d_j) = \frac{\sum_{i=1}^M \Delta(w_i, d_j) r_{ijk}}{\sum_{k=1}^K \sum_{i=1}^M \Delta(w_i, d_j) r_{ijk}}$$