

# 潜在语义分析 (LSA)

① 无监督学习方法：一种文本降维方法

② 通过矩阵分解发现：文本-单词-话题关系

③ 在1990年提出，在推荐系统、文本搜索、图像处理中都有应用

## 1. 单词向量空间与话题向量空间

### 1.1. 单词向量空间

word vector space model: 给定一个文本，用一个向量表示文本“语义”  
向量每一维对应一个单词，取值是频率或权重，向量内积表示文本相似度。

单词-文本矩阵：

(1) 文本集合  $D = \{d_1, d_2, \dots, d_n\}$ ，单词集合  $W = \{w_1, w_2, \dots, w_m\}$

(2) 单词-文本矩阵可表示成  $X = (x_{ij}) \in \mathbb{R}^{m \times n}$ ，其中  $x_{ij}$  表示单词  $w_i$  在文本  $d_j$  中出现的频数或权重

权重通常用 TF-IDF (单词频率-逆文本频率) 表示

$$\text{TF-IDF}_{ij} = \frac{\text{tf}_{ij}}{\text{tf}_j} \log \frac{df}{df_i}$$

其中， $\text{tf}_{ij}$  为单词  $w_i$  在文本  $d_j$  中出现的频数， $\text{tf}_j = \sum_i \text{tf}_{ij}$

$df_i$  为包含单词  $w_i$  的文本数， $df = n$  是总文本数

单词文本矩阵的  $j$  列  $x_j = (x_{1j}, x_{2j}, \dots, x_{mj})^T$  表示  $d_j$  的信息， $d_i, d_j$  之间的相似度为  $x_i \cdot x_j$  或  $\frac{x_i \cdot x_j}{\|x_i\| \|x_j\|}$  (余弦)

(1) 模型优点：模型简单，计算效率高

(2) 模型局限：不能准确刻画语义相似度，自然语言有一词多义和多词一义。

## 1.2 话题向量空间

文本的语义相似度可以用两者之间的“话题”相似度表示

假设文本共有  $K$  个话题 ( $K \ll m$ ) 每个话题由定义在单词集合  $W$  上的  $n$  维向量表示:  $t_l = (t_{l1}, t_{l2}, \dots, t_{ln})^T$ ,  $t_{li}$  表示单词  $w_i$  在话题  $t_l$  上的权重

得到  $K$  个话题  $t_1, t_2, \dots, t_K$  张成一个  $K$  维话题向量空间.

文本  $d_j$  在单词向量空间是  $x_j$ , 将  $x_j$  投影到话题向量空间  $T$  中得到

$y_j = (y_{1j}, \dots, y_{Kj})^T \in \mathbb{R}^K$ ,  $y_{lj}$  代表文本  $d_j$  在话题  $t_l$  上的权重

话题-文本矩阵  $Y = (y_1, y_2, \dots, y_n) \in \mathbb{R}^{K \times n}$

$x_j \approx y_{1j} \cdot t_1 + y_{2j} \cdot t_2 + \dots + y_{Kj} \cdot t_K$  (即  $x_j$  可以通过  $K$  个话题进行线性近似)

用单词-话题矩阵  $T$  以及话题-文本矩阵  $Y$  近似单词-文本矩阵  $X$ :

$$X \approx TY \quad T \in \mathbb{R}^{m \times K}, Y \in \mathbb{R}^{K \times n}$$

文本相似度 ① 单词  $x_i \cdot x_j$

② 话题空间  $y_i \cdot y_j$

## 2. 潜在语义分析算法

### 2.1 矩阵奇异值分解

$U$  左特征向量矩阵

$$X = U \Sigma V^T \quad X \in \mathbb{R}^{m \times n} \quad n \leq m \quad \Sigma = \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_n\}$$

$$= \sum_{i=1}^n \sigma_i u_i v_i^T$$

$V$  右特征向量矩阵

根据话题个数  $K$  对  $X$  进行截断的奇异值分解

$$X \approx U^{(K)} \Sigma^{(K)} V^{(K)T}, \text{ 则 } T = U^{(K)}, Y = \Sigma^{(K)} V^{(K)T}$$

### 2.2 非负矩阵分解算法

给定一个非负  $X \geq 0$  (所有元素非负) 找到两个非负矩阵  $W \in \mathbb{R}^{m \times K} \geq 0$  与  $H \in \mathbb{R}^{K \times n} \geq 0$

$X \approx WH$ , 由于  $K < \min\{m, n\}$ , 非负矩阵分解是对原数据的压缩

## 1. 损失函数

(1) 平方损失  $\|A-B\|_F^2 = \sum_{ij} (a_{ij} - b_{ij})^2$

(2) 散度 (divergence) 散度损失函数

① 不对称, 在  $A=B$  时取下界

$$D(A\|B) = \sum_{ij} (a_{ij} \log \frac{a_{ij}}{b_{ij}} - a_{ij} + b_{ij}) \quad \text{② } 0 \log 0 = 0$$

当  $\sum_j a_{ij} = \sum_j b_{ij} = 1$  时, 即 K-L 散度,  $D(A\|B) = \sum_{ij} a_{ij} \log \frac{a_{ij}}{b_{ij}}$

注: K-L 散度用于度量两个概率分布之间的差异  $D_{KL}(P\|Q) = E_{x \sim P}(\log \frac{P(x)}{Q(x)})$

非负矩阵分解转换为带约束的优化问题  $\min_{W, H} \|X - WH\|_F^2$  or  $\min_{W, H} D(X\|WH)$

## 2. 算法

非负矩阵分解的优化算法, 迭代对  $W$  与  $H$  分别进行优化

定理: 平方损失对以下乘法更新规则  $H_{ij} \leftarrow H_{ij} \frac{(W^T X)_{ij}}{(W^T W H)_{ij}}$   $W_{il} \leftarrow W_{il} \frac{(X H^T)_{il}}{(W H H^T)_{il}}$

是非增的, 当且仅当  $W$  和  $H$  是损失的稳定点时函数损失不变

$$J(W, H) = \frac{1}{2} \|X - WH\|_F^2 = \frac{1}{2} \sum_{ij} (x_{ij} - (WH)_{ij})^2$$

$$\frac{\partial J(W, H)}{\partial w_{il}} = - \sum_j \{x_{ij} - (WH)_{ij}\} H_{ij} = - \{(X H^T)_{il} - (W H H^T)_{il}\}$$

$$\frac{\partial J(W, H)}{\partial h_{ij}} = - \{(W^T X)_{ij} - (W^T W H)_{ij}\}$$

$$W_{il} \leftarrow W_{il} + \lambda_{il} \{(X H^T)_{il} - (W H H^T)_{il}\}, \quad \lambda_{il} = \frac{W_{il}}{(W H H^T)_{il}}$$

$$H_{ij} \leftarrow H_{ij} + u_{ij} \{(W^T X)_{ij} - (W^T W H)_{ij}\}, \quad u_{ij} = \frac{H_{ij}}{(W^T W H)_{ij}}$$

(1) 在算法更新时, 选取初始矩阵非负, 更新中可保证  $W, H$  非负

(2) 每次迭代对  $W$  的列进行归一化为单位向量