

EM算法:

②

混合

① 若数据分布可以建模为方差相同的高斯分布，我可以使用 KMeans

若方差不相等，还想建模成高斯混合模型  $\rightarrow$  就需要用 EM 算法来求参

② 直观理解高斯混合模型 (GMM)

随机变量由  $K$  个高斯分布混合而成，取各个高斯分布的概率为  $\pi_1, \pi_2, \dots, \pi_K$

第  $i$  个高斯分布均值为  $\mu_i$ ，方差为  $\Sigma_i$ ，根据样本  $x_1, x_2, \dots, x_n$  估计参数  $\pi, \mu, \Sigma$

eg 样本

$$\begin{Bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{Bmatrix} = \begin{Bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \end{Bmatrix}$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$   
 $\mu = \{\mu_1, \mu_2, \mu_3, \mu_4\}$

差  $\rightarrow \Sigma_4 = \begin{Bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{Bmatrix}$

元素是任意 2 个特征的协方差

③ 高斯混合模型求解使用 EM 的原因

对数似然函数:  $\ln_{\pi, \mu, \Sigma}(x) = \sum_{i=1}^N \log \left[ \prod_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k) \right]$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

由于隐变量的存在，无法用求导的方法求最大似然，需要用 EM 方法求参

④ EM 算法步骤概要:

问题描述: 随机变量  $x$  由  $K$  个高斯分布混合而成，取各个高斯分布的概率为

$\phi_1, \phi_2, \dots, \phi_K$ ，第  $i$  个高斯分布的均值为  $\mu_i$ ，方差为  $\Sigma_i$

根据样本  $x_1, x_2, \dots, x_n$  估计参数  $\mu, \phi, \Sigma$

E-step: 选取一组参数  $\phi, \mu, \Sigma$ ，计算该参数下，隐变量的条件概率

$$w_j^{(i)} = Q_i(z^{(i)} = j) = P(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma) \quad // \text{for each } i, j$$

$\leftarrow$  M-step 会更新

样本  $i$  属于  
主题  $z_i$   
的条件  
概率

$$= \frac{P(x^{(i)} | z^{(i)} = j; \mu, \Sigma) P(z^{(i)} = j; \phi)}{\sum_{l=1}^K P(x^{(i)} | z^{(i)} = l; \mu, \Sigma) P(z^{(i)} = l; \phi)}$$

M-Step: 结合E步求出的隐变量条件概率  $w_j^{(i)} = Q_i(z^{(i)} = j)$  (对任意的  $i, j$ )

求出似然函数的下界函数(本质上是某个期望函数)的最大值。

<a> 对数似然函数:  $l(\theta) = \sum_{i=1}^m \log P(x; \theta) = \sum_{i=1}^m \log \sum_z P(x, z; \theta)$

<b> 对数似然函数的下界函数:  $\sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$

推导: 令  $Q_i$  是隐变量  $z$  的某一个分布,  $Q_i \geq 0$ , 有:

$$l(\theta) = \sum_{i=1}^m \log \sum_z P(x, z; \theta) = \sum_{i=1}^m \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)$$

$$= \sum_{i=1}^m \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad \leftarrow \text{对 } \square \text{ 求期望再取对数}$$

$$\geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad \leftarrow \text{对 } \square \text{ 取对数再求期望}$$

Jensen 不等式  $\rightarrow$  跟推 ~~不等式~~ 不等式  $m$  个特征,  $m$  个  $\theta$

<c> 求下界函数的最大值

$$\theta := \arg \max_{\theta} \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

求解过程

$$\sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad \text{条件概率展开}$$

$$= \sum_{i=1}^m \sum_{j=1}^K Q_i(z^{(i)} = j) \log \frac{P(x^{(i)} | z^{(i)} = j, \mu, \Sigma) P(z^{(i)} = j; \phi)}{Q_i(z^{(i)} = j)}$$

$$\textcircled{*} = \sum_{i=1}^m \sum_{j=1}^K w_j^{(i)} \log \frac{1}{w_j^{(i)}} \cdot \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_j|^{-\frac{1}{2}}} e^{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)} \cdot \phi_j$$

s.t.  $\phi_1 + \phi_2 + \dots + \phi_K = 1$   
约束条件  $K$  个主题概率之和为 1

对于  $\mu, \Sigma_j$  求解: 对  $\textcircled{*}$  关于  $\mu$  或  $\Sigma_j$  求偏导, 另令偏导为 0 即可

对于  $\phi = \{\phi_1, \phi_2, \dots, \phi_K\}$  求解:  $\textcircled{*}$  式与约束条件构成一个约束条件求极值问题

用拉格朗日乘子法求解。

<d> 求得的结果

$$\mu_j = \frac{\sum_{i=1}^n \omega_j^{(i)} x^{(i)}}{\sum_{i=1}^n \omega_j^{(i)}}$$

$i$  是样本下标

$j$  是特征序号

$$\Sigma_j = \frac{\sum_{i=1}^n \omega_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n \omega_j^{(i)}}$$

$$\Phi_j = \frac{1}{n} \sum_{i=1}^n \omega_j^{(i)}$$

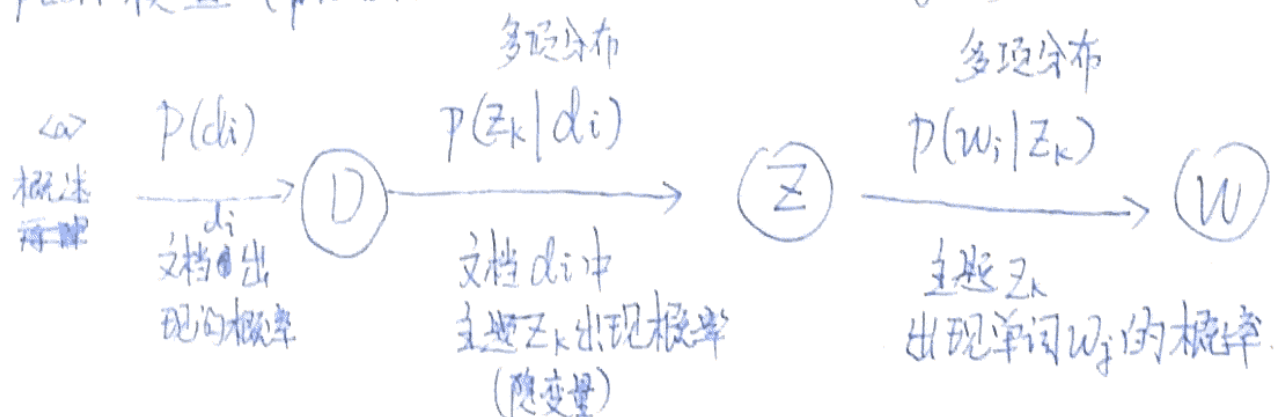
将 M-Step 更新出的  $\mu_j$ ,  $\Sigma_j$ ,  $\Phi_j$  用于 E-Step 重新计算  $\omega_j^{(i)}$

$$\begin{aligned} \text{Estep: } \omega_j^{(i)} &= Q_i(z^{(i)} = j) = P(z^{(i)} = j | x^{(i)}; \Phi, \mu, \Sigma) \\ &= \frac{P(x^{(i)} | z^{(i)} = j; \mu, \Sigma) P(z^{(i)} = j; \Phi)}{\sum_{l=1}^k P(x^{(i)} | z^{(i)} = l; \mu, \Sigma) P(z^{(i)} = l; \Phi)} \end{aligned}$$

直到收敛.



# PLSA 模型 (Probabilistic Latent Semantic Analysis)



与LDA类似, 用文档生成的模拟过程来逆向理解这个算法, 即

以  $P(d_i)$  的概率选中文档  $d_i$ , 以  $P(z_k|d_i)$  的概率选中主题  $z_k$ , 以  $P(w_j|z_k)$  的概率产生一个词  $w_j \rightarrow$  这些词还原出文档的词袋。

公式:

$$P(d_i, w_j) = P(w_j | d_i) P(d_i)$$

文档 词 联合概率

$$P(w_j | d_i) = \sum_{k=1}^K P(w_j | z_k) P(z_k | d_i)$$

目标: 求这两个多项分布。

<c> 似然函数

$$L = \prod_{i=1}^N \prod_{j=1}^M P(d_i, w_j)$$

$n(d_i, w_j)$   $w_j$  在  $d_i$  文档中出现的次数

对数似然

$$l = \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \log P(d_i, w_j)$$

$$= \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \log \sum_{k=1}^K P(w_j | z_k) P(z_k | d_i) P(d_i)$$

未知, 且有隐变量介入  $\rightarrow$  无法求偏导。

<d> 采用类似EM算法的方式求最大似然

E-Step: 假定  $P(z_k | d_i)$ ,  $P(w_j | z_k)$  已知, 求隐变量  $z_k$  的经验概率

$$P(z_k | d_i, w_j) = \frac{P(w_j | z_k) P(z_k | d_i)}{\sum_{k=1}^K P(w_j | z_k) P(z_k | d_i)}$$

M-Step: 在  $(d_i, w_j, z_k)$  已知的前提下, 求关于参数  $P(z_k | d_i)$ ,  $P(w_j | z_k)$  的似然期望的最大值, 从而得到当前阶段  $P(z_k | d_i)$ ,  $P(w_j | z_k)$  的最优解

(方法见下一页)

(d) 变换对数似然函数. 方便求解

$$l = \sum_i \sum_j n(d_i, w_j) \log P(d_i, w_j)$$

$$= \sum_i \sum_j n(d_i, w_j) \log (P(w_j | d_i) \cdot P(d_i))$$

$$= \sum_i \sum_j n(d_i, w_j) (\log P(w_j | d_i) + \log P(d_i))$$

$$= \sum_i \sum_j n(d_i, w_j) \log P(w_j | d_i) + \boxed{\sum_i \sum_j n(d_i, w_j) \log P(d_i)}$$

$$l_{\text{new}} = \sum_i \sum_j n(d_i, w_j) \log P(w_j | d_i)$$

从样本可直接计算  
当作常数去掉

~~$$l_{\text{new}} = \sum_i \sum_j n(d_i, w_j) \log P(w_j | d_i)$$~~

在  $z_k$  上的数学期望

$$\begin{aligned} E(l_{\text{new}}) &= \sum_i \sum_j n(d_i, w_j) \sum_{k=1}^K P(z_k | d_i, w_j) \log P(w_j, z_k | d_i) \\ \text{求最大值} \uparrow &= \sum_i \sum_j n(d_i, w_j) \sum_{k=1}^K P(z_k | d_i, w_j) \log P(w_j | z_k) P(z_k | d_i) \end{aligned}$$

同时, 有约束条件

$$\begin{cases} \sum_{j=1}^M P(w_j | z_k) = 1 \\ \sum_{k=1}^K P(z_k | d_i) = 1 \end{cases}$$

(d2) 这是一个约束条件求极值问题, 用拉格朗日乘子法

~~(d2)~~ 解得 
$$P(w_j | z_k) = \frac{\sum_i n(d_i, w_j) P(z_k | d_i, w_j)}{\sum_{m=1}^M \sum_i n(d_i, w_m) P(z_k | d_i, w_m)}$$

$$P(z_k | d_i) = \frac{\sum_j n(d_i, w_j) P(z_k | d_i, w_j)}{\sum_{k'=1}^K \sum_j n(d_i, w_j) P(z_{k'} | d_i, w_j)}$$

E-Step: M-Step 的结果代回到 E-Step

$$P(z_k | d_i, w_j) = \frac{P(w_j | z_k) P(z_k | d_i)}{\sum_{i=1}^K P(w_j | z_i) P(z_i | d_i)}$$

反复迭代, 直到收敛

## <e> pLSA 与 LDA

pLSA 不需要先验信息就可以完成自学习, 这是它的优势.

如果有先验知识的影响呢? 这时使用 LDA (需要超参数)

LDA 的超参数:

①  $K$ : 主题个数

②  $\vec{\alpha}$ : 主题<sub>先验</sub>集中度

③  $\vec{\beta}$ : 主题突出度 (是否受无义词影响)