

EM algorithm (Expectation Maximization)

EM 算法. 1977.

Maximum likelihood from incomplete data via the EM model

极大似然估计 ML

已知.

model: 样本服从的分布模型

samples: 随机抽取的样本

ML

未知.

模型参数

样本集 x_1, \dots, x_n

$p(x_i|\theta)$

$$L(\theta) = L(x_1, \dots, x_n|\theta) = \prod_{i=1}^n p(x_i|\theta)$$

$$\hat{\theta} = \operatorname{argmax} L(\theta)$$

求法:

$$\eta(\theta) = \log(L(\theta)) = \sum_{i=1}^n \log p(x_i|\theta)$$

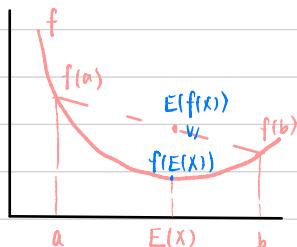
let

$$\eta(\theta) = 0 \quad \text{求解 } \theta.$$

Jesen 不等式.

$f \mapsto \mathbb{R}$, $f'' \geq 0$ (f 是凸函数)

$$E(f(x)) \geq f(E(x)) \quad x \text{ 是常量时, 取等号.}$$



EM algorithm

已知

- ① 分布模型 (several)
- ② 随机抽取的样本

EM

未知

- ① 每个样本属于哪一个分布
- ② 模型参数

样本集 x_1, \dots, x_m

x_i 对应类别 z_i 未知

EM 又被称为数据添加技术, 所添加的数据被称为潜在数据

观测数据 (x_1, \dots, x_m)

隐含数据 (z_1, \dots, z_m)

完整数据 $Y = (X, Z) = \{(x_1, z_1), \dots, (x_m, z_m)\}$

推导 (EM):

$$\begin{aligned} \max H(\theta) &= \mathbb{E} \log P(X^{(i)}, \theta) = \mathbb{E} \log \sum_i P(X^{(i)}, Z^{(i)}; \theta) \\ &= \mathbb{E} \log \sum_i Q_i(Z^{(i)}) \frac{P(X^{(i)}, Z^{(i)}; \theta)}{Q_i(Z^{(i)})} \\ &\geq \mathbb{E} \sum_i \log Q_i(Z^{(i)}) \frac{P(X^{(i)}, Z^{(i)}; \theta)}{Q_i(Z^{(i)})} \quad \downarrow \text{Jensen} \end{aligned}$$

思路:

求 $H(\theta)$ 的下界, 通过 θ, Q 的更新, 使下界逐步提升, 逼近 $\max H(\theta)$, 用下界逼近 $\max H(\theta)$.

E 步骤 (estimation)

固定 θ , 调整 Q 使下界上升.

下界最大取等, 即 $\frac{P(X^{(i)}, Z^{(i)}; \theta)}{Q_i(Z^{(i)})} = C$ (常数)

$$\Rightarrow Q_i(Z^{(i)}) = P(X^{(i)}, Z^{(i)}; \theta) \quad \text{因为要满足 } \sum_i Q_i(Z^{(i)}) = 1$$

M 步骤 (maximization)

固定 Q , 优化 θ . (与 ML 一致).