

# 采样——MCMC

Mia Feng

2018 年 4 月 20 日

## 1 概述

MCMC: 粗暴的采样模拟方式, 用于模拟直接计算困难的分布。用于采样, 数值积分等等。

求解目标: 用多次采样得到的频率分布近似原概率分布。

$$y = \arg \max_{c_k} P(Y = c_k) \prod_{j=1}^n P(X_j = x_j | Y = c_k) \quad (1)$$

求解思路: 微积分思想 (recall: 学习微积分的时候, 用无数个划分的小矩形的面积来近似面积, 但是此时小矩形是来自均匀分布的)。实际上, 来自均匀分布的可能性很小, 此时需要求解方法: Markov Chain, 蒙特卡洛积分。

### 1.1 推导

推导 取  $I$  为示性函数。 $a_l$  表示  $X$  的第  $l$  个特征。样本有  $n$  个。类标  $m$  个, 特征  $s$  个。

$$P(Y = c_k) = \frac{\sum_{i=1}^n I(y_i = c_k)}{n}, k = 1, 2, \dots, m \quad (2)$$

$$P(X_j = a_{jl} | Y = c_k) = \frac{\sum_{i=1}^n I(x_i^j = a_{jl}, y_i = c_k)}{\sum_{i=1}^n I(y_i = c_k)} \quad (3)$$

其中,  $i = 1, 2, \dots, n, l = 1, 2, \dots, s, k = 1, 2, \dots, m$

改进 为了避免分母为 0，进行了拉普拉斯平滑，即在分母上加了类数目。

$$P(X_j = a_{jl}|Y = c_k) = \frac{\sum_{i=1}^n I(x_i^j = a_{jl}, y_i = c_k) + 1}{\sum_{i=1}^n I(y_i = c_k) + m}$$

(4)

2 算法实现

注意实现时取了拉普拉斯平滑，见公式 (4)，且为了防止下溢取对概率值取了对数。[? ? ]

3 Implementation

MCMC

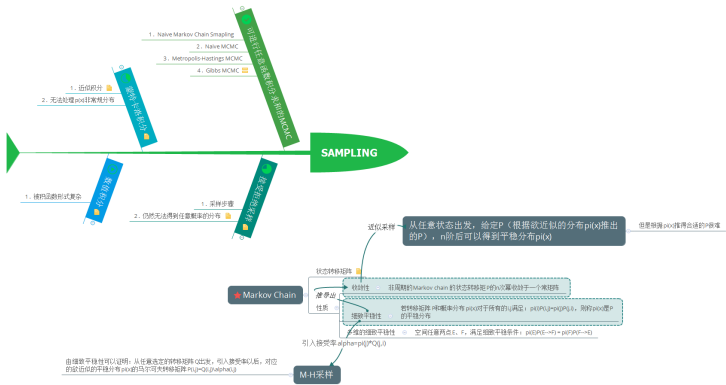


图 1: MCMC 思维导图