

## 第四章 Gibbs 抽样和马尔可夫链

# Gibbs 抽样

- Gibbs 抽样是一种迭代抽样算法，随着样本数增加，样本的分布会收敛到目标分布，首先以一个 Bayesian 模型的估计为例介绍如何使用 Gibbs 抽样
- 使用 Bayesian 方法分析数据三要素：
  - ① 模型设定：为数据的抽样分布设定具体形式  $p(\mathbf{y} | \theta)$ ，通常需要引入一些参数  $\theta$
  - ② 设定参数的先验分布  $p(\theta)$ ：一般是主观设定，可以加入参数的先验信息，其样本空间应覆盖参数所有可能的取值
  - ③ 计算参数的后验分布  $p(\theta | \mathbf{y})$  并做统计推断：估计参数的后验期望  $E(\theta | \mathbf{y})$ 、后验方差  $\text{Var}(\theta | \mathbf{y})$ 、置信区间等。参数的后验分布可如下计算：

$$p(\theta | \mathbf{y}) = \frac{p(\theta)p(\mathbf{y} | \theta)}{p(\mathbf{y})} \propto p(\theta)p(\mathbf{y} | \theta) \quad (1)$$

但  $p(\theta | \mathbf{y})$  对应的分布一般很难识别或很难直接对其抽样

## A Bayesian Normal Model

假设数据独立地服从正态分布:

$$Y_i \stackrel{iid}{\sim} N(\mu, \phi^{-1}), \quad i = 1, \dots, n \quad (2)$$

参数的 (共轭) 先验分布:

$$\begin{aligned} \phi &\sim \text{Gam}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right) \\ \mu \mid \phi &\sim N\left(\mu_0, \frac{1}{\kappa_0 \phi}\right) \end{aligned}$$

其中  $\nu_0, \sigma_0^2, \mu_0, \kappa_0^2$  都是确定的常数

# A Bayesian Normal Model

参数  $(\mu, \phi)$  的联合后验分布仍是一个 normal-gamma 分布:

$$\phi \mid y_1, \dots, y_n \sim \text{Gam}\left(\frac{\nu_n}{2}, \frac{S_n}{2}\right)$$
$$\mu \mid \phi, y_1, \dots, y_n \sim \mathcal{N}\left(\mu_n, \frac{1}{\kappa_n \phi}\right)$$

其中

$$\nu_n = \nu_0 + n$$

$$S_n = \nu_0 \sigma_0^2 + \frac{n\kappa_0}{\kappa_0 + n} (\mu_0 - \bar{y})^2 + \sum_{i=1}^n (y_i - \bar{y})^2$$

$$\mu_n = \frac{\kappa_0 \mu_0 + n\bar{y}}{\kappa_0 + n}$$

$$\kappa_n = \kappa_0 + n$$

- 计算  $E(\mu \mid y_1, \dots, y_n)$ , 可以采用如下的 Monte Carlo 方法:

(1) 独立抽样  $\phi^{(s)} \sim \text{Gam}(\nu_n/2, S_n/2)$ ,  $s = 1, \dots, T$

(2) 对每个  $\phi^{(s)}$ , 抽  $\mu^{(s)} \mid \phi^{(s)} \sim \mathcal{N}(\mu_n, (\kappa_n \phi^{(s)})^{-1})$ ,  $s = 1, \dots, T$

则  $E(\mu \mid y_1, \dots, y_n) \approx \sum_{s=1}^T \mu^{(s)} / T$

## A Bayesian Normal Model

如果为模型(2)选取如下的先验分布:

$$\begin{aligned}\mu &\sim N(\mu_0, \tau_0^2) \\ \phi &\sim \text{Gam}(\nu_0/2, \nu_0\sigma_0^2/2)\end{aligned}\tag{3}$$

- 此时  $\phi$  的边际后验分布既不是 Gamma 分布, 也不是任何常见的分布
- 方法一: 使用数值方法得到  $(\mu, \phi)$  的近似联合后验分布

- ▶ 首先为各参数选取足够大的取值范围  $\mu \in [\mu_L, \mu_H]$ ,  $\phi \in [\phi_L, \phi_H] \subseteq (0, \infty)$
- ▶ 然后对区域  $[\mu_L, \mu_H] \times [\phi_L, \phi_H]$  做网格离散
- ▶ 根据(1), 点  $(\mu_i, \phi_j)$  处的后验概率密度为

$$p(\mu_i, \phi_j \mid y_1, \dots, y_n) \propto p(\mu_i, \phi_j)p(y_1, \dots, y_n \mid \mu_i, \phi_j)$$

- ▶ 网格中大部分点的后验概率都很接近 0, 造成计算的浪费, 且该方法只适用于参数较少的情况

# A Bayesian Normal Model

- 方法二：Gibbs 抽样

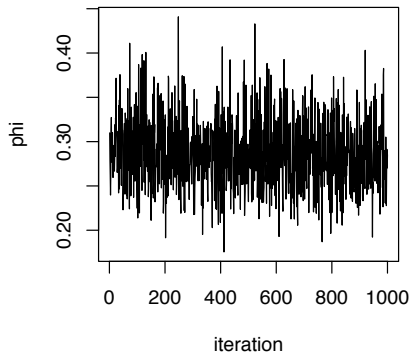
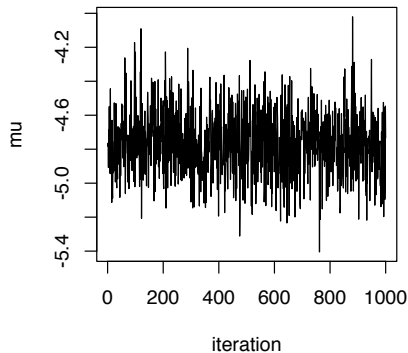
- ▶ 模型(2) - (3) 中  $\mu$  的完全条件分布为： $\mu \mid \phi, y_1, \dots, y_n \sim N(\mu_n, \tau_n^2)$ ，其中  $\mu_n = (\tau_0^{-2}\mu_0 + \phi \sum_{i=1}^n y_i) / (\tau_0^{-2} + n\phi)$ ， $\tau_n^2 = (\tau_0^{-2} + n\phi)^{-1}$
- ▶ 模型(2) - (3) 中  $\phi$  的完全条件分布为： $\phi \mid \mu, y_1, \dots, y_n \sim \text{Gam}(\nu_n/2, S_n/2)$ ，其中  $\nu_n = \nu_0 + n$ ， $S_n = \nu_0\sigma_0^2 + \sum_{i=1}^n (y_i - \mu)^2$
- ▶ 如何利用  $\mu$  和  $\phi$  的完全条件分布得到  $(\mu, \phi)$  的联合后验分布的样本？
  - ★ 假设  $\phi^{(1)}$  是边际后验分布  $p(\phi \mid y_1, \dots, y_n)$  的一个样本，给定  $\phi^{(1)}$ ，从  $\mu$  的完全条件分布抽样：

$$\mu^{(1)} \sim p(\mu \mid \phi^{(1)}, y_1, \dots, y_n)$$

则  $(\mu^{(1)}, \phi^{(1)})$  可以看作联合后验分布  $p(\mu, \phi \mid y_1, \dots, y_n)$  的一个样本

- ▶ 只要给定  $\mu$  或  $\phi$  的一个初始值，然后轮流从  $\mu$  和  $\phi$  的完全条件分布抽样，就可以得到一系列来自联合后验分布  $p(\mu, \phi \mid y_1, \dots, y_n)$  的样本  $\{(\mu^{(s)}, \phi^{(s)}) : s = 1, \dots, T\}$

# A Bayesian Normal Model



样本在参数空间中“移动”的很快，表明样本之间的相关性较小，此时样本均值可以很好地近似后验分布的期望

# A Bayesian Normal Model

- 样本之间的相关性如何影响样本均值对目标期望的近似？
  - ▶ 假设  $\{\theta^{(s)} : s = 1, \dots, T\}$  是由 Gibbs 抽样得到的一系列服从目标分布  $p(\theta)$  的样本，样本均值  $\bar{\theta}$  的方差为

$$\begin{aligned} \text{Var}_G(\bar{\theta}) &= E[(\bar{\theta} - E(\theta))^2] \\ &= \text{Var}_{MC}(\bar{\theta}) + \frac{1}{T^2} \sum_{s=1}^T \sum_{t \neq s} E\left[\left(\theta^{(s)} - E(\theta)\right)\left(\theta^{(t)} - E(\theta)\right)\right] \end{aligned}$$

- ▶  $\text{Var}_G(\bar{\theta}) > \text{Var}_{MC}(\bar{\theta})$ ，且 Gibbs 样本之间的相关性越高， $\text{Var}_G(\bar{\theta})$  就越大，均值的近似效果越差
- 一个衡量 Gibbs 样本相关性的指标是有效样本数 (ESS):

$$ESS = \frac{\text{Var}(\theta)}{\text{Var}_G(\bar{\theta})}$$

为达到与 Gibbs 样本估计量相同精度所需的独立样本个数



## Gibbs 抽样

- 假设模型的参数向量为  $\theta = (\theta_1, \dots, \theta_p)$ , 抽样的目标分布为  $p(\theta)$ , 如果知道每个分量  $\theta_j$  的完全条件分布  $p(\theta_j | \{\theta_k\}_{k \neq j})$ ,  $j = 1, \dots, p$ , 给定初始值  $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_p^{(0)})$ , Gibbs 抽样可以如下从当前样本  $\theta^{(s-1)}$  产生新的样本  $\theta^{(s)}$ :

1. 抽样  $\theta_1^{(s)} \sim p(\theta_1 | \theta_2^{(s-1)}, \theta_3^{(s-1)}, \dots, \theta_p^{(s-1)})$

2. 抽样  $\theta_2^{(s)} \sim p(\theta_2 | \theta_1^{(s)}, \theta_3^{(s-1)}, \dots, \theta_p^{(s-1)})$

3. 抽样  $\theta_3^{(s)} \sim p(\theta_3 | \theta_1^{(s)}, \theta_2^{(s)}, \theta_4^{(s-1)}, \dots, \theta_p^{(s-1)})$

⋮

p. 抽样  $\theta_p^{(s)} \sim p(\theta_p | \theta_1^{(s)}, \theta_2^{(s)}, \dots, \theta_{p-1}^{(s)})$

## Gibbs 抽样

- 不断重复上述过程，Gibbs 抽样可以产生一系列不独立的样本  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(T)}$ ，且每个样本  $\theta^{(s)}$  与之前的样本  $\theta^{(0)}, \dots, \theta^{(s-1)}$  的相关性只取决于  $\theta^{(s-1)}$ ，因此这一列样本是一条**马尔可夫链**
- 在满足一些条件后，从任何初始值  $\theta^{(0)}$  出发，Gibbs 抽样产生的样本  $\theta^{(s)}$  的分布在  $s \rightarrow \infty$  时收敛到目标分布  $p(\theta)$ ，即

$$P(\theta^{(s)} \in A) \rightarrow \int_A p(\theta) d\theta$$

- 此时对任意可积函数  $g$  有

$$\frac{1}{T} \sum_{s=1}^T g(\theta^{(s)}) \rightarrow E[g(\theta)] = \int g(\theta) p(\theta) d\theta, \quad T \rightarrow \infty$$

上述过程被称为**Markov chain Monte Carlo (MCMC)** 方法

# 马尔可夫链

- 给定初始值  $\mathbf{x}^{(0)} \in \mathcal{X}$ , 按照某个条件分布  $p(\mathbf{x} | \mathbf{x}')$  依次抽样

$$\mathbf{x}^{(t)} \sim p(\mathbf{x} | \mathbf{x}^{(t-1)}) \quad \text{且} \quad \mathbf{x}^{(t)} \perp\!\!\!\perp \mathbf{x}^{(t-k)} | \mathbf{x}^{(t-1)}, \quad k > 1$$

称这样得到的序列  $\{\mathbf{x}^{(t)}\}$  为状态空间  $\mathcal{X}$  上的一条**马尔可夫链**

- 在马尔可夫链中, 对任意正整数  $k$ , 定义如下的 **$k$  步转移密度函数**:

- ▶  $p^1(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) = p(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)})$

- ▶  $p^2(\mathbf{x}^{(t+2)} | \mathbf{x}^{(t)}) = \int_{\mathcal{X}} p(\mathbf{x}^{(t+2)} | \mathbf{x}^{(t+1)}) p(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) d\mathbf{x}^{(t+1)}$

- ▶  $p^3(\mathbf{x}^{(t+3)} | \mathbf{x}^{(t)}) = \int_{\mathcal{X}} p(\mathbf{x}^{(t+3)} | \mathbf{x}^{(t+2)}) p^2(\mathbf{x}^{(t+2)} | \mathbf{x}^{(t)}) d\mathbf{x}^{(t+2)}$

⋮

- ▶  $p^k(\mathbf{x}^{(t+k)} | \mathbf{x}^{(t)}) = \int_{\mathcal{X}} p(\mathbf{x}^{(t+k)} | \mathbf{x}^{(t+k-1)}) p^{k-1}(\mathbf{x}^{(t+k-1)} | \mathbf{x}^{(t)}) d\mathbf{x}^{(t+k-1)}$

如果  $p(\mathbf{x} | \mathbf{x}') > 0, \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$ , 那么对任意正整数  $k$ ,  $k$  步转移密度函数  $p^k(\mathbf{x} | \mathbf{x}') > 0, \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$

# 马尔可夫链

- 如果存在一个分布  $\pi(\mathbf{x})$  使马尔可夫链满足

$$\pi(\mathbf{x}) = \int_{\mathcal{X}} p(\mathbf{x} | \mathbf{x}') \pi(\mathbf{x}') d\mathbf{x}'$$

称分布  $\pi(\mathbf{x})$  为该马尔可夫链的一个**平稳分布**

- ▶ 如果马尔可夫链的初始值服从平稳分布, 整个过程的边缘分布会永远“保持”平稳分布

$$\pi(\mathbf{x}) = \int_{\mathcal{X}} p^k(\mathbf{x} | \mathbf{x}') \pi(\mathbf{x}') d\mathbf{x}', \quad \forall k \in \mathbb{N}^+$$

- ▶ Gibbs 抽样产生的马尔可夫链的转移密度函数是

$$p(\boldsymbol{\theta}^{(s)} | \boldsymbol{\theta}^{(s-1)}) = p(\theta_1^{(s)} | \theta_2^{(s-1)}, \dots, \theta_p^{(s-1)}) \cdots p(\theta_p^{(s)} | \theta_1^{(s)}, \dots, \theta_{p-1}^{(s)})$$

如果  $\forall \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^{(s)}, p(\boldsymbol{\theta}^{(s)} | \boldsymbol{\theta}^{(s-1)}) > 0$ , 该过程会收敛到唯一的平稳分布 (参数的后验分布)

# 马尔可夫链

- 如果从任意初始状态  $x' \in \mathcal{X}$  出发，马尔可夫链可以在有限步到达任意其它状态  $x \in \mathcal{X}$ ，称马尔可夫链是**不可约的** (irreducible)

## 定义 (不可约性)

如果  $\forall x, x' \in \mathcal{X}, \exists k < \infty$  使得  $k$  步转移密度  $p^k(x | x') > 0$ ，称该马尔可夫链是不可约的。

- ▶ 如果  $p(x | x') > 0, \forall x, x' \in \mathcal{X}$ ，那么该马尔可夫链是不可约的
- ▶ 一个可约的马尔可夫过程. 令状态空间  $\mathcal{X} = (-1, 1)$  上的一个马尔可夫链有如下的转移分布：

$$p(x | x') = \begin{cases} x \sim U(0, 1), & \text{if } x' \geq 0 \\ x \sim U(-1, 0), & \text{if } x' < 0 \end{cases}$$

# 马尔可夫链

- 如果  $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$ ,  $\gcd\{k : p^k(\mathbf{x} | \mathbf{x}') > 0\} = 1$ , 称该过程具有**非周期性**
  - ▶ 当  $p(\mathbf{x} | \mathbf{x}') > 0$ ,  $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$ , 马尔可夫链是非周期的
  - ▶ 一个**周期性的马尔可夫链**. 令状态空间  $\mathcal{X} = (-1, 1)$  上的一个马尔可夫链有如下的转移分布:

$$p(\mathbf{x} | \mathbf{x}') = \begin{cases} x \sim U(0, 1), & \text{if } x' < 0 \\ x \sim U(-1, 0), & \text{if } x' \geq 0 \end{cases}$$

- 称一个不可约且非周期的马尔可夫链具有**遍历性**(Ergodicity), 具有遍历性的马尔可夫链有以下重要性质:
  - ① 存在唯一的平稳分布  $\pi(\mathbf{x})$
  - ② 从任意初始值出发, 该过程都会收敛到平稳分布  $\pi(\mathbf{x})$ , 即对  $\forall \mathbf{x}' \in \mathcal{X}$ , 当  $k \rightarrow \infty$ ,

$$p^k(\mathbf{x} | \mathbf{x}') \rightarrow \pi(\mathbf{x})$$

# 马尔可夫链

- 对遍历的马尔可夫链, 当  $k \rightarrow \infty$ , 该过程产生的  $\mathbf{x}^{(k)} \sim \pi(\mathbf{x})$
- 此时样本均值会 (almost surely) 收敛到平稳分布的期望, 且对任何可积函数  $g(\cdot)$ , 当  $T \rightarrow \infty$ ,

$$\frac{1}{T} \sum_{t=1}^T g(\mathbf{x}^{(t)}) \rightarrow \int g(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}$$

- 在 Bayesian 分析中通过 Gibbs 抽样产生的马尔可夫链一般是遍历的, 通过产生一条很长的马尔可夫链, 该过程的样本可以近似描述后验分布 (平稳分布)
- 如果初始值选取得不好, 马尔可夫链可能会移动很长时间才收敛到平稳分布概率密度较高的区域, 称这段时间为 burn-in, 用样本均值估计目标期望时通常舍弃 burn-in 阶段的样本