



中国石油大学
CHINA UNIVERSITY OF PETROLEUM

随机过程作业

题 目	Bayes 的极限分布与极大似然
学 院	信息科学与工程学院
专 业	控制科学与工程
姓名学号	孟庆鑫 2018 312 032
	刘芳正 2019 310 704
	韩亚楠 2019 310 708
任课教师	严彦文

Bayes 的极限分布与极大似然

孟庆鑫

2018 312 032

刘芳正

2019 310 704

韩亚楠

2019 310 708

摘 要: 类比于 Markov 链的状态转移算子思路, 本文构造了序贯 Bayes 分布列, 以不太严谨的方式证明了 Bayes 的极限分布与极大似然的关系: 对于存在极大似然估计的似然函数、任意有界正先验分布下, Bayes 极限分布弱收敛于取值在极大似然估计值的 Dirac δ 分布. 共轭先验方法可用于构造特殊的显式序贯 Bayes 分布列.

1 基本定义

一个 Markov 链 \mathcal{M} 由其初始分布 $\pi(0)$ 和状态转移矩阵 P 完全确定. 若 \mathcal{M} 存在极限分布 $\lim_{n \rightarrow \infty} \pi(n) = \pi$, 则 π 满足

$$\pi P = \pi, \quad (1)$$

π 是平稳分布, 刚好是 P 的特征值 1 对应的特征向量. 一个 Markov 链, 不论初始分布 $\pi(0)$ 是什么, 如果能达到平稳分布, 那么平稳分布由状态转移矩阵 P 确定.

记 Θ 为分布参数, \mathcal{D} 为样本对应的随机变量, $\mathbb{P}(\Theta)$ 是 Θ 的先验分布, 根据 Bayes 公式

$$\underbrace{\mathbb{P}(\Theta|\mathcal{D})}_{\text{后验分布}} = \frac{\mathbb{P}(\mathcal{D}|\Theta)\mathbb{P}(\Theta)}{\mathbb{P}(\mathcal{D})} \propto \underbrace{\mathbb{P}(\mathcal{D}|\Theta)}_{\text{似然函数}} \underbrace{\mathbb{P}(\Theta)}_{\text{先验分布}}, \quad (2)$$

若 $\mathbb{P}(\Theta|\mathcal{D}), \mathbb{P}(\Theta)$ 有相同的形式, 则称 $\mathbb{P}(\Theta)$ 是 $\mathbb{P}(\mathcal{D}|\Theta)$ 关于 Θ 的共轭先验分布^[1].

2 Bayes 的极限分布

将 Markov 链的状态转移矩阵 P 看作算子, 等式 (1) 的 π 即为算子 P 的不动点. 一个可达平稳分布的 Markov 链可以看作算子 P 连续作用下趋近不动点的过程.

同样的, 等式 (2) 的先验分布 $\mathbb{P}(\Theta)$ 可以看作是来自于其它数据集 \mathcal{D}' 上的经验, 即 $\mathbb{P}(\Theta) = \mathbb{P}(\Theta|\mathcal{D}')$, 为使先验尽可能符合当前样本, 取 $\mathcal{D}' = \mathcal{D}$, 如果把

$$\frac{(\bullet)\mathbb{P}(\mathcal{D}|\Theta)}{\sum_{\Omega_{\Theta}}(\bullet)\mathbb{P}(\mathcal{D}|\Theta)}$$

看作算子 D , 这就与 Markov 链很相似了: 算子 D 连续作用在先验分布上, 极限分布趋近于分布空间的某个不动点. 下面通过例子来说明这一观点.

2.1 二项分布的例子

考虑 n 次独立重复试验并观测到一组样本, 事件出现的次数记为随机变量 X , $X \sim B(n, \theta)$,

$$\mathbb{P}(X = x|\theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x},$$

将 θ 也视为随机变量, $\theta \in (0, 1)$, 但 θ 的分布不清楚, 我们先假设 $\theta \sim \text{Uniform}(0, 1)$,

$$\mathbb{P}(\theta|x) = \frac{\mathbb{P}(x|\theta)\mathbb{P}(\theta)}{\int_0^1 \mathbb{P}(x|\theta)\mathbb{P}(\theta) d\theta} = \frac{\Gamma(n+2)}{\Gamma(x+1)\Gamma(n-x+1)} \theta^{(x+1)-1} (1-\theta)^{(n-x+1)-1},$$

即 $\theta|x \sim \text{Beta}(x+1, n-x+1)$. 注意到 $\text{Uniform}(\theta|0, 1) = \text{Beta}(\theta|1, 1)$, 实际上我们选取均匀分布恰好是共轭先验的一个特例. 为此, 重新假设 $\theta \sim \text{Beta}(\alpha, \beta)$, 同理可得 $\theta|x \sim \text{Beta}(x+\alpha, n-x+\beta)$. 再将此 posterior 分布作为先验分布, 如此, 可构造分布列 $\{\text{Beta}(kx+\alpha, k(n-x)+\beta)\}_{k=0}^{+\infty}$, 并有 $k \rightarrow +\infty$,

$$\text{Beta}(kx+\alpha, k(n-x)+\beta) \rightarrow \delta\left(\frac{x}{n}\right), \quad (3)$$

δ 是 Dirac 函数. (3) 式说明 $\forall \varepsilon > 0$,

$$\mathbb{P}\left(\left|\theta - \frac{x}{n}\right| \leq \varepsilon\right) = 1,$$

因此可确定 $\hat{\theta} = \frac{x}{n}$.

上面的讨论是在共轭先验构成的分布列上做的, 我们再重新假设 $\theta \sim \mathbb{P}(\theta)$, 其中 $\mathbb{P}(\theta)$ 满足

$$\forall \theta \in (0, 1), \mathbb{P}(\theta) > 0, |\mathbb{P}(\theta)| < \infty, \text{ 且 } \int_0^1 \mathbb{P}(\theta) d\theta = 1, \quad (4)$$

仍有 $k \rightarrow +\infty$,

$$\frac{1}{Z_k(x)} \mathbb{P}(\theta) (\theta^x (1-\theta)^{n-x})^k \rightarrow \delta\left(\frac{x}{n}\right),$$

其中 $Z_k(x)$ 是归一化系数. 这个结论说明, 不论先验分布如何选, 只要满足 (4) 式, 如此构造的序贯贝叶斯的极限分布都是 $\delta\left(\frac{x}{n}\right)$, 即参数 θ 在极限分布上只能取 $\frac{x}{n}$, 这个值恰好是 θ 的极大似然估计. 这里需要注意, 有界分布序列收敛到了无界可积分布 Dirac δ .

2.2 Poisson 分布的例子

假设有一组观测样本 $n_i \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda)$, $i = 1, 2, \dots, m$, 记 $\mathbf{n} = (n_1, n_2, \dots, n_m)$, 即

$$\mathbb{P}(\mathbf{n}|\lambda) = \prod_{i=1}^m \frac{e^{-\lambda} \lambda^{n_i}}{\Gamma(n_i + 1)},$$

选取参数 λ 的共轭先验^[2]

$$\text{Gamma}(\lambda|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} e^{-\beta\lambda} \lambda^{\alpha-1},$$

可以得到后验分布 $\mathbb{P}(\lambda|\mathbf{n}, \alpha, \beta) = \text{Gamma}(\lambda|\alpha + \sum \mathbf{n}, \beta + m)$, 再将此分布作为先验分布, 如此, 可构造分布列 $\{\text{Gamma}(\lambda|\alpha + k \sum \mathbf{n}, \beta + km)\}_{k=0}^{+\infty}$, 并有 $k \rightarrow +\infty$,

$$\text{Gamma}(\lambda|\alpha + k \sum \mathbf{n}, \beta + km) \rightarrow \delta(\bar{\mathbf{n}}).$$

选取先验 $\lambda \sim \mathbb{P}(\lambda)$ 满足

$$\forall \lambda \in \mathbb{R}^+, \mathbb{P}(\lambda) > 0, |\mathbb{P}(\lambda)| < \infty, \text{ 且 } \int_{\mathbb{R}^+} \mathbb{P}(\lambda) d\lambda = 1,$$

仍有 $k \rightarrow +\infty$,

$$\frac{1}{Z_k(\mathbf{n})} \mathbb{P}(\lambda) (e^{-\lambda} \lambda^{\bar{\mathbf{n}}})^{mk} \rightarrow \delta(\bar{\mathbf{n}}),$$

其中 $Z_k(\mathbf{n})$ 是归一化系数. 再次注意到, $\bar{\mathbf{n}}$ 恰好是 λ 的极大似然估计.

2.3 正态分布的例子

假设一组观测样本 $x_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, $i = 1, 2, \dots, n$, 其中 μ 是未知的, σ^2 是已知的, 记 $\mathbf{x} = (x_1, x_2, \dots, x_n)$, 则似然为

$$\mathbb{P}(\mathbf{x}|\mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2},$$

假设参数 μ 的共轭先验为 $\mathcal{N}(\mu|\mu_k, \sigma_k^2)$ ^[3], 则可计算后验分布为 $\mathcal{N}(\mu|\mu_{k+1}, \sigma_{k+1}^2)$, 满足递推关系

$$\frac{1}{\sigma_{k+1}^2} = \frac{1}{\sigma_k^2} + \frac{n}{\sigma^2}, \quad \frac{\mu_{k+1}}{\sigma_{k+1}^2} = \frac{\mu_k}{\sigma_k^2} + \frac{n\bar{x}}{\sigma^2},$$

从而有

$$\frac{1}{\sigma_k^2} = \frac{1}{\sigma_0^2} + \frac{kn}{\sigma^2}, \quad \frac{\mu_k}{\sigma_k^2} = \frac{\mu_0}{\sigma_0^2} + \frac{kn\bar{x}}{\sigma^2},$$

因此

$$\mathcal{N}(\mu|\mu_k, \sigma_k^2) = \mathcal{N}\left(\mu \left| \frac{\mu_0\sigma^2 + kn\bar{x}\sigma_0^2}{\sigma^2 + kn\sigma_0^2}, \left(\frac{1}{\sigma_0^2} + \frac{kn}{\sigma^2}\right)^{-1} \right.\right) \xrightarrow{k \rightarrow +\infty} \delta(\bar{x}).$$

对 μ 的任意先验不再特别讨论, 而是统一的放在指数家族分布里说明.

2.4 指数家族的例子

考虑指数家族分布^[1]

$$\mathbb{P}(X|\eta) = h(X) e^{\eta^T T(X) - A(\eta)}$$

及一组观测样本 $\mathbf{x} = (x_1, x_2, \dots, x_n)$, 似然函数为

$$\mathbb{P}(\mathbf{x}|\eta) = e^{\sum_{i=1}^n \eta^T T(x_i) - nA(\eta)} \prod_{i=1}^n h(x_i), \quad (5)$$

选取先验 $\eta \sim \mathbb{P}(\eta)$ 满足

$$\forall \eta \in \Omega_\eta, \mathbb{P}(\eta) > 0, |\mathbb{P}(\eta)| < \infty, \text{ 且 } \int_{\Omega_\eta} \mathbb{P}(\eta) d\eta = 1,$$

注意 Bayes 分布列

$$\left\{ \frac{1}{Z_k(\mathbf{x})} \mathbb{P}(\eta) \left(e^{\eta^T \overline{T(\mathbf{x})} - A(\eta)} \right)^{nk} \right\}_{k=0}^{\infty},$$

记 $f(\eta) = e^{\eta^T \overline{T(\mathbf{x})} - A(\eta)}$ 的最大值在 η^* 处取得, 一方面, η^* 是 (5) 式的极大似然估计, 一方面, 考虑 $f \in C(\Omega_\eta)$, $\exists \zeta > 0$, 使得 $\forall \varepsilon > 0, \varepsilon < \zeta, \forall \xi \in B(\eta^*, \varepsilon), \forall \eta \notin B(\eta^*, \varepsilon), f(\xi) > f(\eta)$, 则

$$\begin{aligned} \lim_{k \rightarrow +\infty} \mathbb{P}\left(\frac{\mathbb{P}(\eta) f^{nk}(\eta)}{Z_k(\mathbf{x})} \in B(\eta^*, \varepsilon)\right) &= \lim_{k \rightarrow +\infty} \int_{B(\eta^*, \varepsilon)} \frac{\mathbb{P}(\eta) f^{nk}(\eta)}{\int_{\Omega_\eta} \mathbb{P}(\eta) f^{nk}(\eta) d\eta} d\eta \\ &\stackrel{\exists \vartheta \in B(\eta^*, \varepsilon)}{=} \lim_{k \rightarrow +\infty} \frac{\int_{B(\eta^*, \varepsilon)} \mathbb{P}(\eta) d\eta}{\int_{\Omega_\eta \setminus B(\eta^*, \varepsilon)} \mathbb{P}(\eta) \left(\frac{f(\eta)}{f(\vartheta)}\right)^{nk} d\eta + \int_{B(\eta^*, \varepsilon)} \mathbb{P}(\eta) d\eta} = 1, \end{aligned}$$

这说明 $k \rightarrow +\infty$,

$$\frac{1}{Z_k(\mathbf{x})} \mathbb{P}(\eta) f^{nk}(\eta) \rightarrow \delta(\eta^*),$$

于是得到结论: 指数家族分布的 Bayes 极限分布弱收敛于 Dirac $\delta(\eta^*)$ 分布, η^* 恰好是极大似然估计值.

3 结论

将 Bayes 后验分布当作先验分布, 构造了类似 Markov 链的分布列, 对于存在极大似然估计的似然函数, 任意有界正先验分布, 其 Bayes 极限分布弱收敛于 Dirac $\delta(\theta)$ 分布, 其中 θ 恰好是唯一的极大似然估计值, 即 $\forall \theta \in \Omega_\Theta, \mathbb{P}(\theta) > 0, |\mathbb{P}(\theta)| < \infty$, 且 $\sum_{\Omega_\Theta} \mathbb{P}(\theta) = 1$, $\exists! \theta = \arg \max_{\Theta} \mathbb{P}(\mathcal{D}|\theta)$, 记 $f(\theta) = \mathbb{P}(\mathcal{D}|\theta)$, $f \in C(\Omega_\Theta)$, $\exists \zeta > 0$, 使得 $\forall \varepsilon > 0, \varepsilon < \zeta, \forall \xi \in B(\theta, \varepsilon), \forall \eta \notin B(\theta, \varepsilon), f(\xi) > f(\eta)$, 则有

$$\left(\frac{(\cdot) \mathbb{P}(\mathcal{D}|\theta)}{\sum_{\Omega_\Theta} (\cdot) \mathbb{P}(\mathcal{D}|\theta)} \right)^k \mathbb{P}(\theta) \xrightarrow{k \rightarrow +\infty} \delta(\theta),$$

且有推论

$$\arg \max_{\Theta} \left(\frac{(\cdot) \mathbb{P}(\mathcal{D}|\theta)}{\sum_{\Omega_\Theta} (\cdot) \mathbb{P}(\mathcal{D}|\theta)} \right)^k \mathbb{P}(\theta) \xrightarrow{k \rightarrow +\infty} \theta.$$

证明过程类似于指数家族分布例子的证明,

$$\lim_{k \rightarrow +\infty} \mathbb{P} \left(\left(\frac{(\cdot) \mathbb{P}(\mathcal{D}|\theta)}{\sum_{\Omega_\Theta} (\cdot) \mathbb{P}(\mathcal{D}|\theta)} \right)^k \mathbb{P}(\theta) \in B(\theta, \varepsilon) \right) = \lim_{k \rightarrow +\infty} \mathbb{P} \left(\frac{\mathbb{P}(\theta) f^k(\theta)}{Z_k(\mathcal{D})} \in B(\theta, \varepsilon) \right) = 1,$$

以及 $\exists K$, 使 $k > K$,

$$\arg \max_{\Theta} \left(\frac{(\cdot) \mathbb{P}(\mathcal{D}|\theta)}{\sum_{\Omega_\Theta} (\cdot) \mathbb{P}(\mathcal{D}|\theta)} \right)^k \mathbb{P}(\theta) \in B(\theta, \varepsilon),$$

再由 ε 的任意性可得结论.

本文基于以下想法: 当 Markov 链到达极限分布的时候, 极限分布就是平稳分布, 并只由状态转移矩阵决定, 亦可看作状态转移矩阵的特征, 与初始分布无关; 从 Bayes 估计的角度看, 先验经过似然算子反复作用到达极限分布, 那么极限分布也由似然算子决定, 亦可看作似然算子的特征, 与初始先验无关, 似然由样本生成, 这意味着经过无穷次数据作用得到的后验, 将“抹去”先验的痕迹.

从结论来看, 任何 Dirac δ 分布都是似然算子的不动点, 这意味着初始先验分布不能选择 Dirac δ 分布, 选取有界正先验的意义在于, 先验是依赖有限次数据作用后得到, 并且对随机变量的任何取值均有概率, 而 Dirac δ 分布意味着经历了无限的数据作用, 变量已不再随机.

本文在举例中使用了共轭先验, 该方法可用于构造特殊的显式序贯 Bayes 分布列, 尤其是正态分布的例子, 可以直接看出极限分布. 对于非共轭先验分布, 先验后验之间一般不具备良好的解析性质.

References

- [1] C. Bishop and S. Yip, *Pattern Recognition and Machine Learning*, vol. 1. 01 2006.
- [2] “Gamma distribution.” https://en.wikipedia.org/wiki/Gamma_distribution.
- [3] K. Murphy, “Conjugate bayesian analysis of the gaussian distribution,” pp. 445–470, 11 2007.