

# Homework for “Algorithms for Big-Data Analysis”

Beijing International Center for Mathematical Research  
Peking University

February 17, 2024

**Note: Please write up your solutions independently. If you get significant help from others, write down the source of references. A formal mathematical proof for all your claims is required.**

1. 考虑有限情形的MDP  $(S, A, P, R, \gamma)$ ，其中  $S$  是有限个离散状态的集合， $A$  是有限个离散动作的集合， $R$  是奖励函数， $\gamma \in (0, 1)$  是折扣因子。给定时刻  $t$  的状态  $s$  和动作  $a$ ，下一时刻转移到状态  $s'$  的概率是  $P(s' | s, a) = P(s_{t+1} = s' | s_t = s, a_t = a)$ 。令  $V(s)$  为价值函数，定义Bellman 算子  $B$ ：

$$BV(s) = \max_a \left\{ R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) V(s') \right\}, \quad \forall s \in S.$$

- (a) 证明算子  $B$  是压缩映射，即：

$$\|BV - BV'\|_\infty \leq \gamma \|V - V'\|_\infty,$$

其中  $\|V - V'\|_\infty = \max_s |V(s) - V'(s)|$ 。

- (b) 从  $V_0$  开始执行迭代算法： $V_{k+1} = BV_k$ 。对于任意  $k > 0$ ，证明：

$$\|V_{n+k} - V_n\|_\infty \leq \frac{\gamma^n}{1-\gamma} \|V_1 - V_0\|_\infty.$$

2. 令考虑有限步MDP  $(S, A, s_1, P, R, H)$ ，其中  $S$  为状态集合， $A$  为动作集合， $s_1$  为初始状态， $P$  为转移概率矩阵， $R$  为奖励矩阵， $H$  为终止时间， $\gamma$  为折扣因子。定义  $\pi(s, a)$  是在状态  $s$  根据随机策略  $\pi$  执行动作  $a$  的概率。定义  $\tau = (s_1, a_1, s_2, \dots, a_{H-1}, s_H)$  是从状态  $s_1$  出发，执行策略  $\pi$  产生的轨道(状态-动作对)，即  $a_t \sim \pi(s_t, \cdot)$ 。

- (a) 写出轨道  $\tau$  的概率表达式  $D^\pi(\tau)$ 。

- (b) 定义  $\rho(\pi)$  是有限步总奖励的均值，即：

$$\rho(\pi) = \mathbb{E} \left[ \sum_{t=1}^H \gamma^{t-1} r_t | \pi, s_1 \right].$$

令  $R(\tau)$  是在轨道  $\tau$  获得的总奖励。写出  $\rho(\pi)$  关于  $R(\tau)$  的表达式。

(c) 假设 $\pi_\theta(s, a)$ 是参数化之后的策略, 其中 $\theta$ 是参数。证明

$$\nabla_{\theta} \rho(\pi_{\theta}) = \mathbb{E}_{\tau} [R(\tau) \nabla_{\theta} \log(D^{\pi_{\theta}}(\tau))] = \mathbb{E}_{\tau} \left[ R(\tau) \sum_{t=1}^{H-1} \nabla_{\theta} \log(\pi_{\theta}(s_t, a_t)) \right].$$

(d) 给定状态 $s_t$ 和参数 $\theta$ , 假设 $b_t(s_t)$ 条件独立于 $\pi_\theta$ 产生的抽样。证明:

$$\mathbb{E}_{\tau} \left[ \sum_{t=1}^{H-1} b_t(s_t) \frac{\partial}{\partial \theta_j} \log(\pi_{\theta}(s_t, a_t)) \mid \theta, s_1 \right] = 0.$$

(e) 给出满足(d)里条件的一种 $b_t(s_t)$ 。