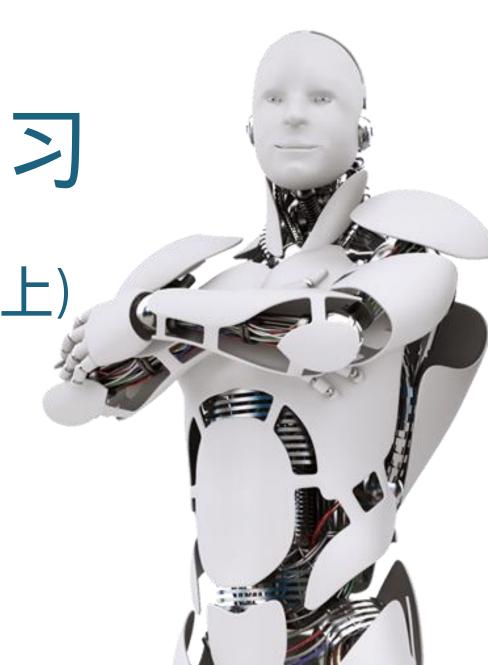
彪哥带你学强化学习

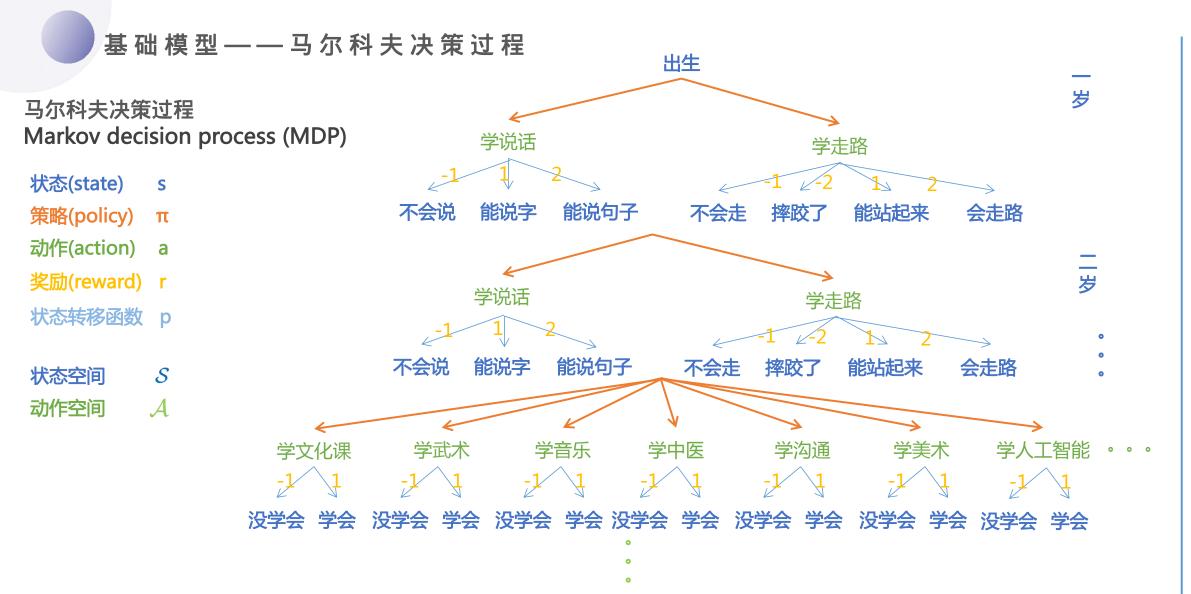
二、强化学习基础概念(上)

DEEPLY UNDERSTAND REINFORCEMENT LEARNING

讲师: 韩路彪







医生 教师 工程师 科学家 律师 教授 公务员 企业家 中层领导 高层领导 销售 顾问 。。。

基础模型——马尔科夫决策过程

马尔科夫决策过程
马尔科夫奖励过程 + 策略
马尔科夫过程 + 奖励



具体状态

基础模型——马尔科夫过程Markov Process(MP)



动作变量 At 象棋第一步棋 具体动作 at 单头炮

18



$$egin{aligned} p(s_{t+1}|s_1,a_1,s_2,a_2,\ldots,s_t,a_t) &= p(s_{t+1}|s_t,a_t) \ p(s_{t+1}|s_1,s_2,s_3,\ldots,s_t) &= p(s_{t+1}|s_t) \end{aligned}$$

概率转移矩阵

$$\mathcal{P} = egin{pmatrix} P(s_1|s_1) & \cdots & P(s_n|s_1) \ dots & \ddots & dots \ P(s_1|s_n) & \cdots & P(s_n|s_n) \end{pmatrix}$$

	饿了	饱了
饿了	0.1	0.9
饱了	0.9	0.1



基础模型——马尔科夫奖励过程Markov reward process(MRP)



回报

$$G_t = R_t + R_{t+1} + R_{t+2} + \dots$$

 $g_t = r_t + r_{t+1} + r_{t+2} + \dots$

折扣回报

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots \ g_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$



基础模型——马尔科夫决策过程Markov decision process(MDP)

