

时序差分

Sarsa

无模型强化学习，不知道env的转移矩阵P

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

α 是学习率，此处本来应该是 $1/N(s)$

这是对 $Q(s,a)$ 做一阶矩估计的变种

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma r_{t+1} + \dots + \gamma^n Q(s_{t+n}, a_{t+n}) - Q(s_t, a_t)]$$

多步Sarsa

在线策略

区别

离线策略

QLearning

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[R_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

不一定下一个动作就是a