

更新

蒙特卡洛方法

Policy Gradient

TRPO/PPO

时间差分方法

Actor-Critic

Q-Learning

DDPG

DQN

