

1. 环境

gym 搭建格子环境

① reset() 结束一个 Episode 重新训练

② step() action 对动作改变环境

互动, 返回 next-state reward done

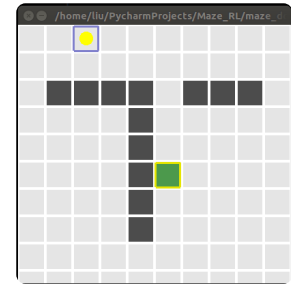
$$\text{reward} \begin{cases} \text{goal} = 10 \\ \text{back} = -1 \\ \text{move} = -0.1 \end{cases}$$

next-state agent 位置

goal 位置

agent 和障碍物之间距离 位置信息

每个场景输入状态不一致
Model 无法统一
状态如何改变呢



2. QAgent

① Neural Network

relu(Dense) + relu(Dense) + Dense

② Exploitation - Exploration

$$\epsilon_{\text{decay}} = \left(\frac{\epsilon_{\text{min}}: 0.01}{\epsilon: 1} \right)^{\frac{1}{\text{MaxEpisode}}}$$

$$\epsilon *= \epsilon_{\text{decay}}$$

③ train

for Episode:

reset()

while not done:

a_0 , 随机策略选择动作.

self.step(a_0)

④ 经验回放

更新

$$Q_{k+1}(S_t, a_t) = Q_k(S_t, a_t) + \alpha (R_{k+1}(S_t, a_t) + \gamma \max_{a_{t+1}} Q_k(S_{t+1}, a_{t+1}) - Q_k(S_t, a_t))$$

⑤ tensorboard 记录

loss reward 变化.