# Reinforcement Learning for Portfolio Management
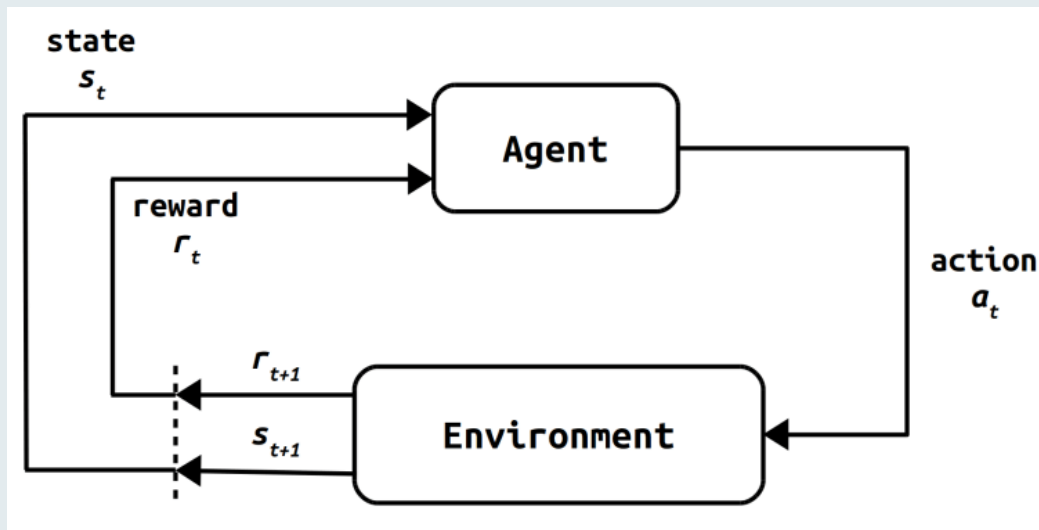
- 环境
- 智能体



R.S. Sutton, A.G. Barto,
Reinforcement Learning: An
Introduction,
MIT Press, Cambridge, MA, 1998

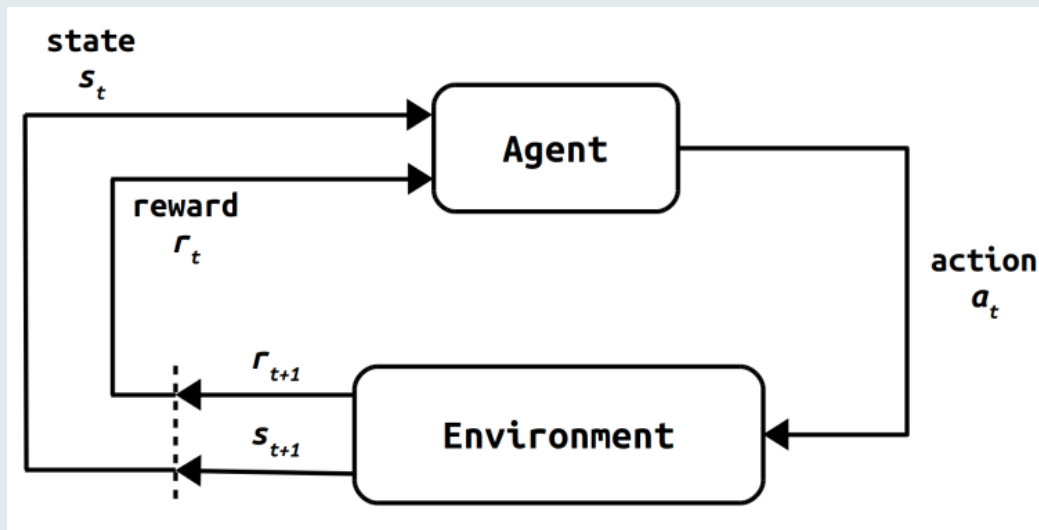- 50*11
- 50: 可投资的股票数目
- 11: 股票的维度：['zopen', 'zhigh', 'zlow', 'zadjcp', 'zclose', 'zd_5', 'zd_10', 'zd_15', 'zd_20', 'zd_25', 'zd_30']
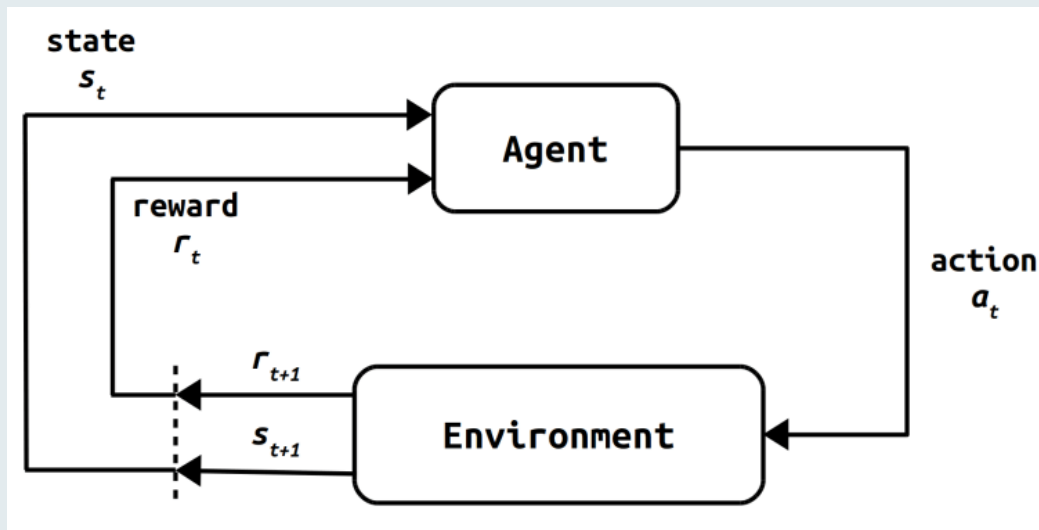
- reward = portfolio_value(t+1) - portfolio_value(t)



R.S. Sutton, A.G. Barto,
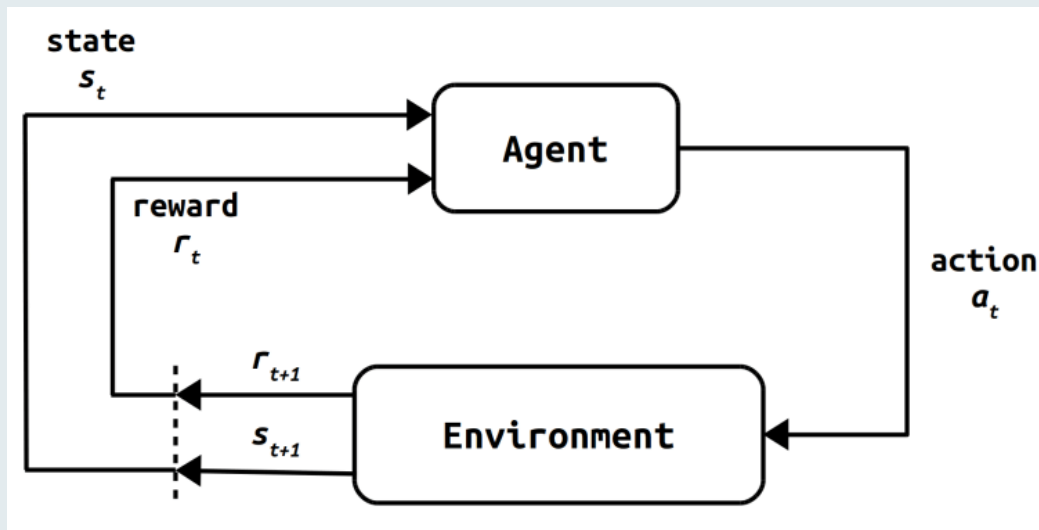Reinforcement Learning: An Introduction,
MIT Press, Cambridge, MA, 1998

- (51,)

- 50: 可投资的股票数目+1:不投资



R.S. Sutton, A.G. Barto,
Reinforcement Learning: An
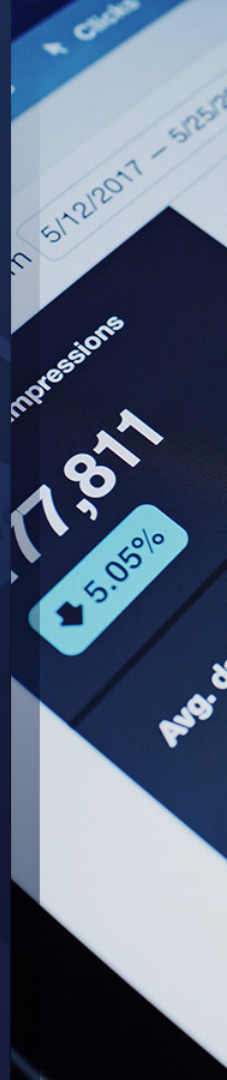Introduction,
MIT Press, Cambridge, MA, 1998

- 状态空间(50,11) ->行动空间 (51,)
- 用反馈的奖励来训练智能体



R.S. Sutton, A.G. Barto,
Reinforcement Learning: An
Introduction,
MIT Press, Cambridge, MA, 1998

# 智能体设计

- Policy Optimization: 策略优化
  - Policy Gradient: 策略梯度
- State Optimization: 状态优化
  - Q-learning: Q学习
- Actor-Critic Methods: 演员-评论家方法

# 随机策略揭示问题

```
def random_agent(env):
    return env.action_space.sample()
```

```
def one_agent(env):
    return np.ones_like(env.action_space_shape)
```

- Fully random vs Evenly distributed



portfolio_value -93.676294%



portfolio_value +323.096709%

- 策略网络结构
  - input状态空间
  - out行动空间
- 训练策略网络
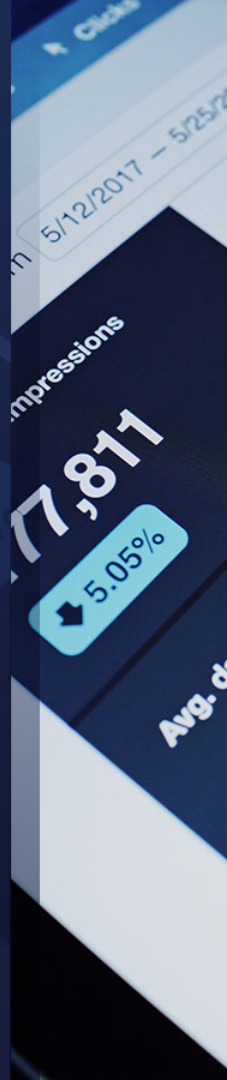  - Step1. MC采样
  - Step2. batch训练

**Algorithm 5:** Model-Carlo Policy Gradient (REINFORCE).

**inputs :** trading universe of $M$-assets
  initial portfolio vector $w_1 = a_0$
  initial asset prices $p_0 = o_0$
  objective function $\mathcal{J}$
  initial agent weights $\theta_0$

**output:** optimal agent policy parameters $\theta_*$

1  initialize buffers: $G, \Delta\theta_c \leftarrow 0$ **repeat**
2      **for** $t = 1, 2, \ldots T$ **do**
3          observe tuple $\langle o_t, r_t \rangle$
4          sample and take action: $a_t \sim \pi_\theta(\cdot | s_t; \theta)$     // portfolio rebalance
5          cache rewards: $G \leftarrow G + r_t$     // (6.19)
6          cache log gradients: $\Delta\theta_c \leftarrow \Delta\theta_c + \nabla_\theta log[\pi_\theta(s, a)]G$ // (6.20)
7      **end**
8      update policy parameters $\theta$ using buffered
9          Monte-Carlo estimates via adaptive optimization   // (6.18), ADAM
10     empty buffers: $G, \Delta\theta_c \leftarrow 0$
11 **until** convergence
12 set $\theta_* \leftarrow \theta$

# 网络设计

- 简单全连接
- Lstm

- 环境:
  - 数据集:
    - 考虑停牌股票
    - 考虑股票的基本面信息

- 智能体:其他的智能体结构
  - 网络结构：可扩展性（现在只能支持50支股票）

- R.S. Sutton, A.G. Barto, Reinforcement Learning: An Introduction, MIT Press, Cambridge, MA, 1998
- Filos, A. (2019). Reinforcement Learning for Portfolio Management. ArXiv [q-Fin.PM]. Retrieved from http://arxiv.org/abs/1909.09571
- Sun, S., Qin, M., Wang, X., & An, B. (2023). PRUDEX-Compass: Towards Systematic Evaluation of Reinforcement Learning in Financial Markets. Transactions on Machine Learning Research. Retrieved from https://openreview.net/forum?id=JjbsIYOuNi

谢谢！