

Machine Learning and Financial Applications

Lecture 10

Reinforcement Learning for Portfolio Optimization

Liu Peng

liupeng@smu.edu.sg

Video tutorial

- https://youtu.be/Yd3H_h-7C58

Introducing portfolio optimization

- The foundations of portfolio optimization can be traced back to the seminal work of Harry Markowitz, who introduced the concept of Modern Portfolio Theory (MPT) in his 1952 seminal paper Markowitz (1952)
- The essence of MPT lies in the quantification of the trade-off between risk and return, both at portfolio level.
- This leads to the formation of the efficient frontier, which represents a set of optimal portfolios offering the highest expected return for a given level of risk, or equivalently, the lowest risk for a target expected return.

Expected
return and
risk of the
portfolio

$$E[R_p] = \mathbf{w}^T \boldsymbol{\mu}$$

$$\sigma_p^2 = \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w},$$

$$\mathbf{w} = [w_1, w_2, \dots, w_n]^T$$

$$r_{t+1} = \frac{P_{t+1} - P_t}{P_t} = \frac{P_{t+1}}{P_t} - 1.$$

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_{nn} \end{bmatrix},$$

$$\mu_i = \frac{1}{T} \sum_{t=1}^T r_{i,t}, \quad \text{Cov}(r_i, r_j) = \sigma_{ij} = \frac{1}{T-1} \sum_{t=1}^T (r_{i,t} - \mu_i)(r_{j,t} - \mu_j),$$

Understanding
asset return
and covariance
matrix

Mean-variance optimization

$$\begin{aligned} & \underset{\boldsymbol{w}}{\text{maximize}} && \boldsymbol{w}^T \boldsymbol{\mu} \\ & \text{subject to} && \boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w} \leq \sigma_p^2, \\ & && \boldsymbol{w}^T \mathbf{1} = 1, \end{aligned}$$

$$\begin{aligned} & \underset{\boldsymbol{w}}{\text{minimize}} && \boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w} \\ & \text{subject to} && \boldsymbol{w}^T \boldsymbol{\mu} = \mu_p, \\ & && \boldsymbol{w}^T \mathbf{1} = 1, \end{aligned}$$

DRL for portfolio optimization

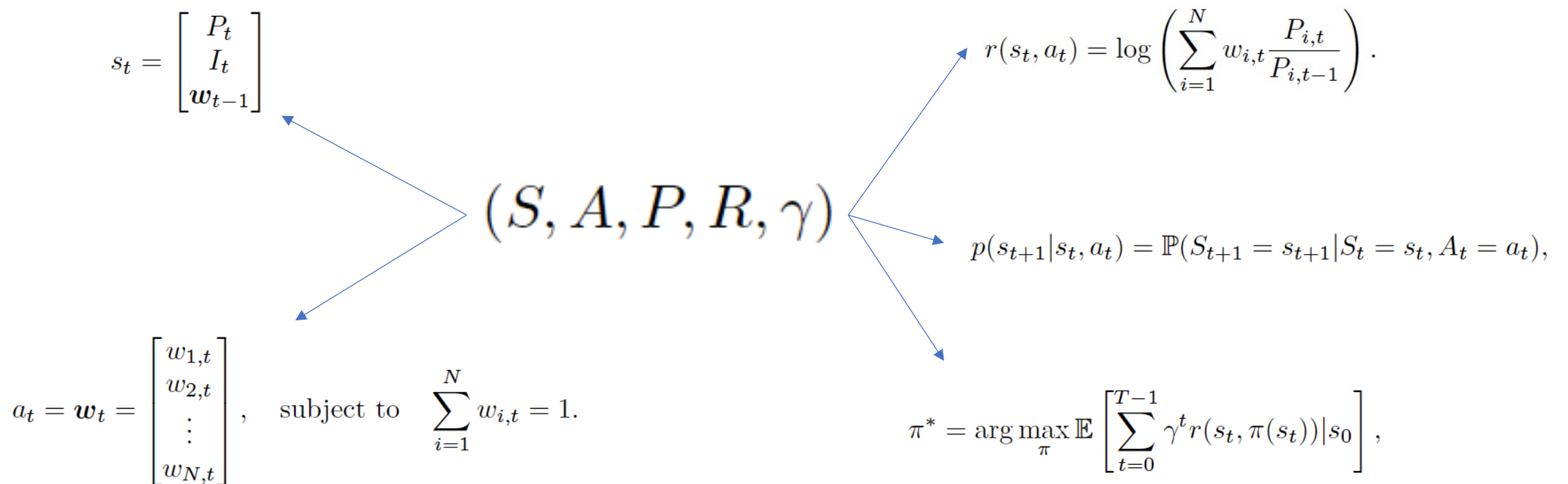
- Model-free learning
- Adaptive strategies
- Complex decision-making
- Scalability



Exploration versus exploitation

Markov decision process

- An MDP is a mathematical framework that provides a formal description of an environment for RL.



The learning environment cont'd

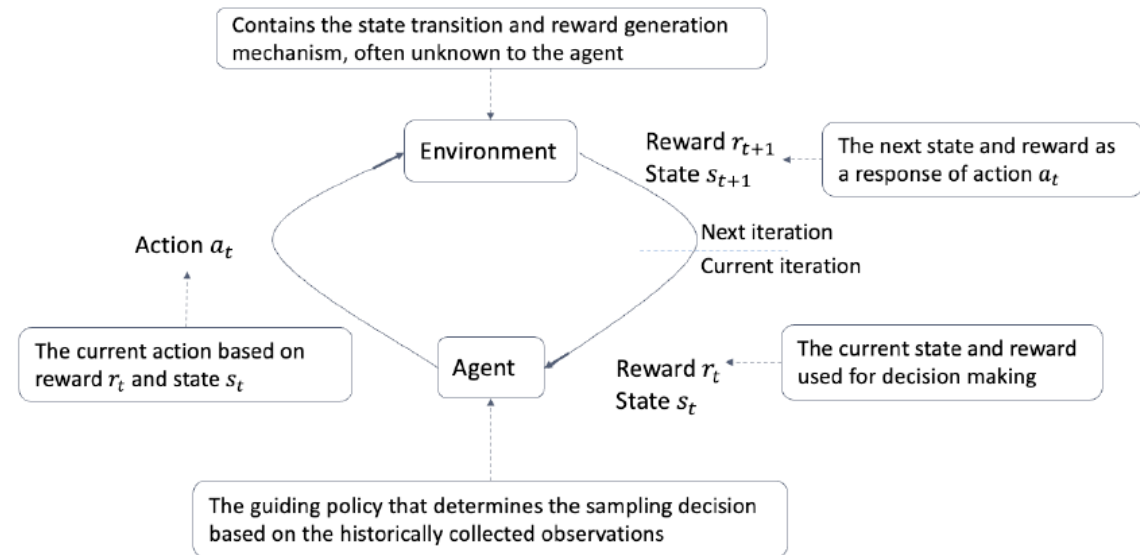


Figure 1.1: Iterative interaction between the agent and the environment.



Value
functions

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$V^\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$$

$$V^\pi(s) = \mathbb{E}_\pi[R_{t+1} | S_t = s] + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s]$$

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V^\pi(s')]$$

$$Q^\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

$$Q^\pi(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a' | s') Q^\pi(s', a')]$$