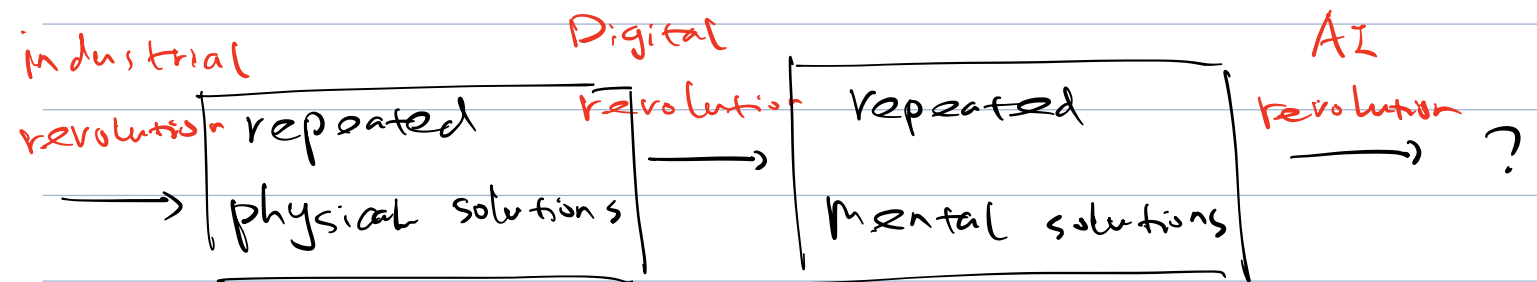


Reinforcement Learning

(RL)

— What is RL ?



RL: interacting the environment

~~X~~ Learning to make decisions from interactions

- ① active
- ② sequential interactions

Goal:

- ① find previously unknown solutions
- ② find solutions online for unforeseen circumstances

RL vs ML

Supervised
signal



reward

{ label
action

how good something
is compared to something
else

1. No supervision, only a reward signal

2. feedback can be delayed.

Not instantaneous

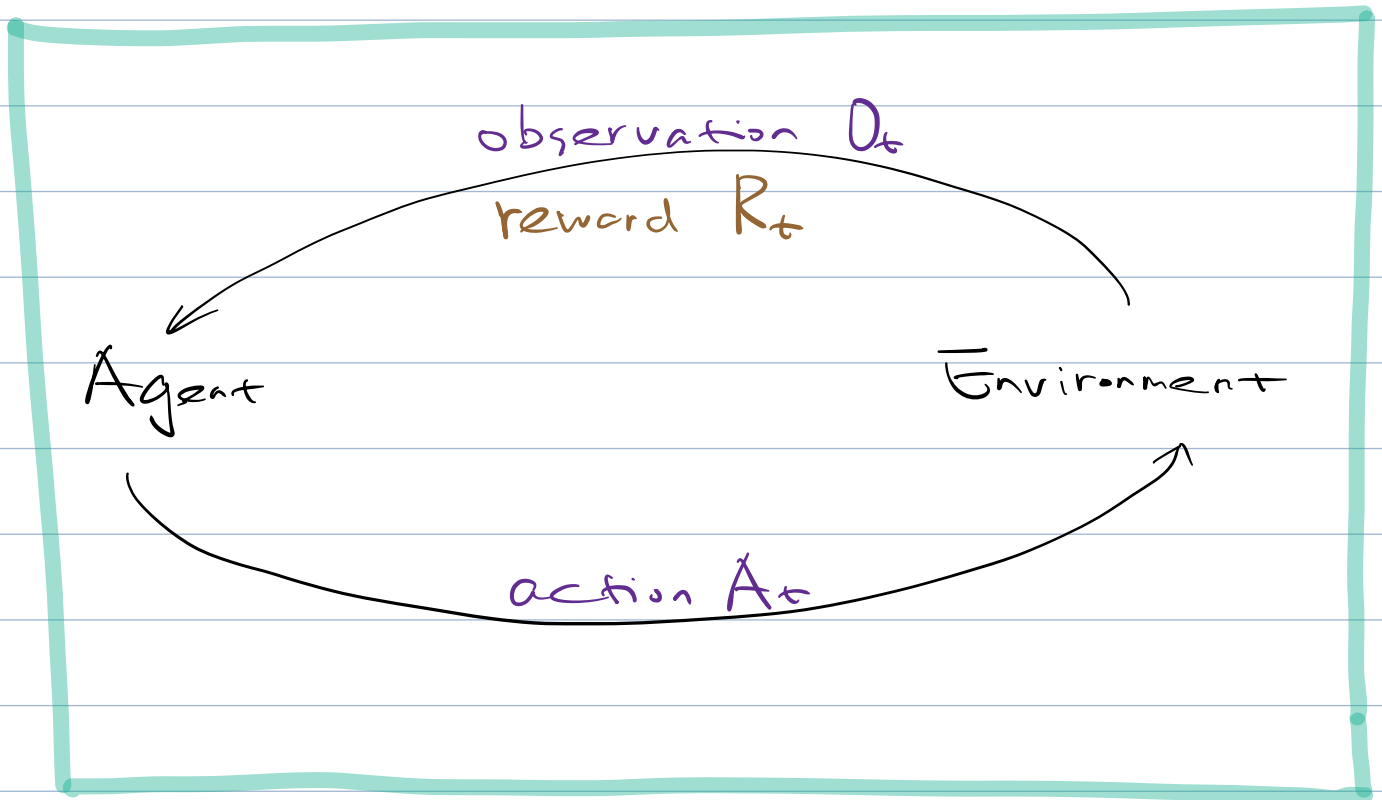
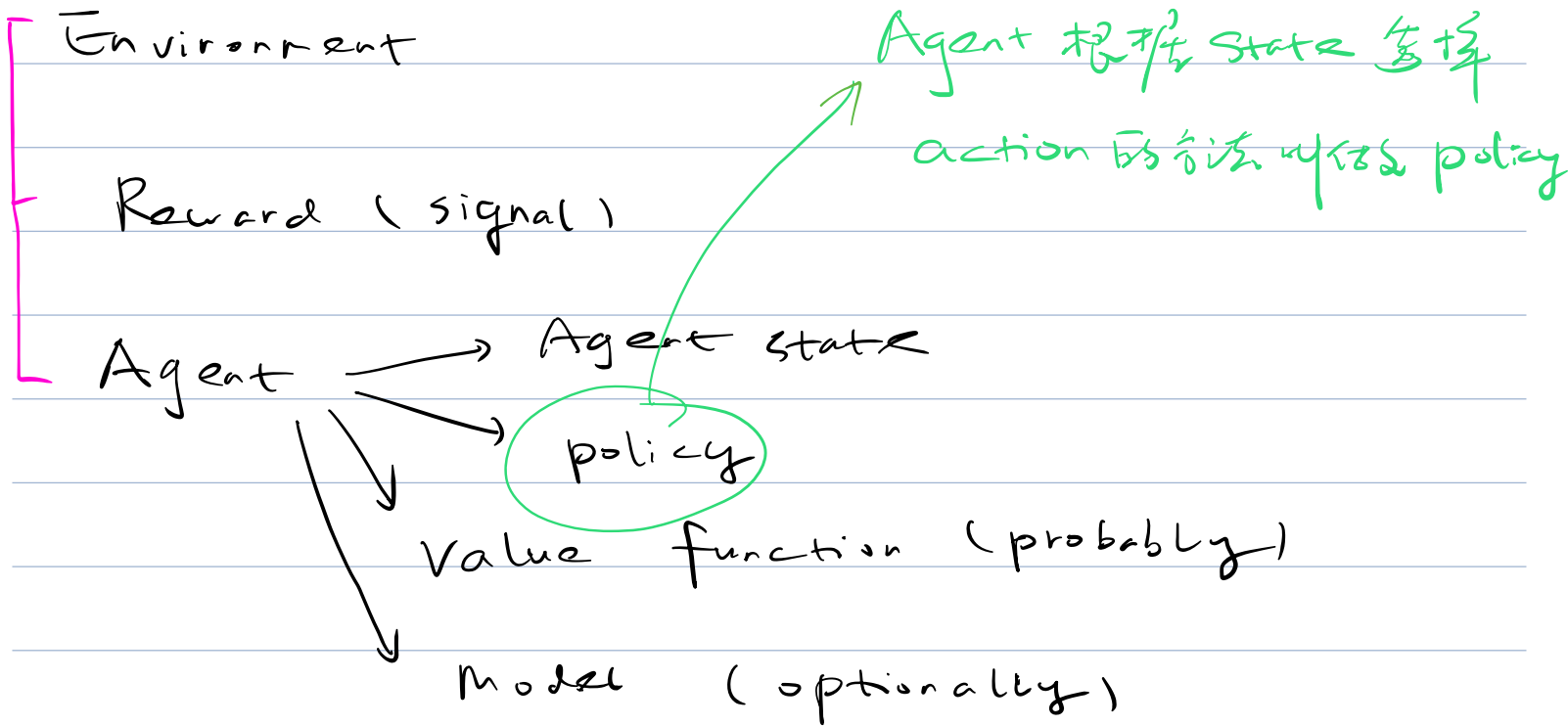
3. time matters

4. earlier decisions affect later interactions

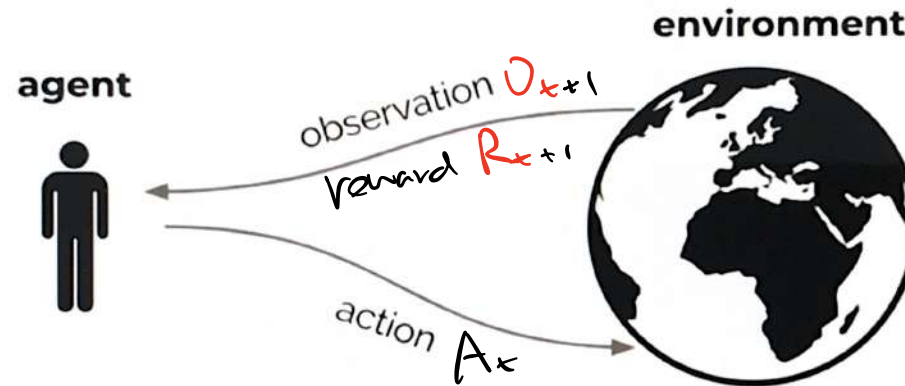
RL is learning what to do, how to map situations
to actions, so as to maximize a numerical
reward signals.

学习一个根据当前 state 选择一个能 最大化全局收益
的 action.

Core concepts of RL



Agent and Environment



- ▶ At each step t the agent:
 - ▶ Receives observation O_t (and reward R_t)
 - ▶ Executes action A_t
- ▶ The environment:
 - ▶ Receives action A_t
 - ▶ Emits observation O_{t+1} (and reward R_{t+1})

Rewards R_t (reward hypothesis)

R_t : scalar feedback signal (reward/penalty)

indicates how well agent is doing at step t
(define its goal)

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

Maximize

Return (cumulative reward)

Reward Hypothesis:

Any goal can be formalized as the outcome
of maximizing a cumulative reward.

Values: expected/mean cumulative reward
from a state s

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

$$= \mathbb{E}[R_{t+1} + R_{t+2} + R_{t+3} + \dots | S_t = s]$$

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

Maximize value
by picking suitable actions

Define Desirability of a state or action

(No supervised feedback)

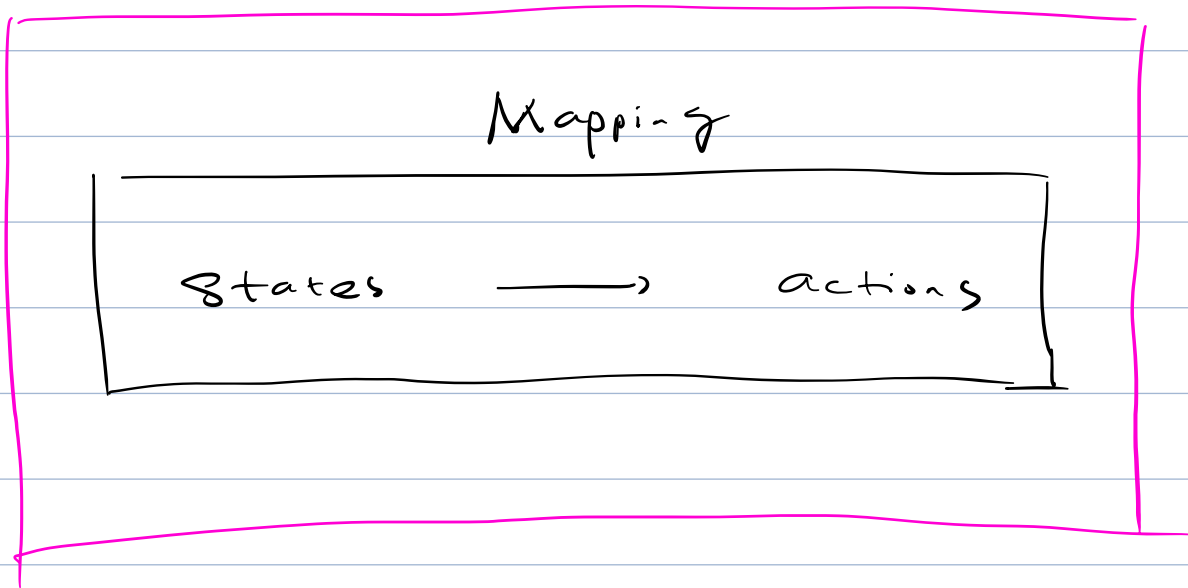
Goal: select actions to maximize value

problem: actions \rightarrow long-term consequences

reward \rightarrow delayed

Solution: sacrifice immediate reward to gain
long-term reward

policy



Action values (condition the value on action):

$$Q(s, a) = \mathbb{E}[G_t \mid S_t = s, A_t = a]$$

$$= \mathbb{E}[R_{t+1} + R_{t+2} + R_{t+3} + \dots \mid S_t = s, A_t = a]$$

Agent components

Agent state: $S_{t+1} = f(S_t, A_t, R_{t+1}, O_{t+1})$

policy: $A = \pi(s)$; $\pi(A|s) = P(A|s)$

value function: $V(s) = \mathbb{E}[G_t | S_t = s]$

$$= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s, \pi]$$

Model

P.S. Environment State: H_t

(Full observable / Partially observable)

Environment State

(observation)

Environment State : environment's internal state

Markov decision process

history : a sequence of observations.

actions

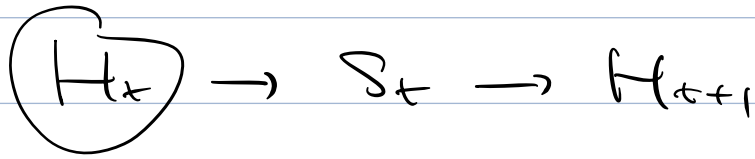
rewards

$$H_t = O_0, A_0, R_0, O_1, \dots, O_{t-1}, A_{t-1}, R_t, O_t$$

(used to construct an agent state S_t)

MDPs: Markov decision Processes

$$p(r, s | s_t, A_t) = p(r, s | H_t, A_t)$$



Markov

full observability

partial observability: agent gets partial information

(Partially observable MDPs)

Agent State

function of history

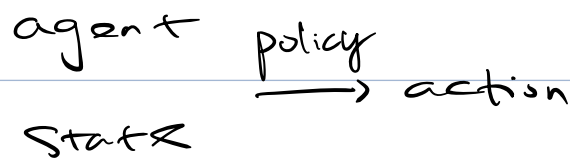
$$S_{t+1} = f(S_t, A_t, R_{t+1}, O_{t+1})$$

f : state update function

P.S. agent state \ll environment state

Policy

define the agent's behaviour:



Deterministic policy: $A = \pi(S)$

Stochastic policy: $\pi(A|S) = p(A|S)$

Value Function

$$V_{\pi}(s) = \mathbb{E}[G(t) \mid S_t = s, \pi]$$

$$= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s, \pi]$$

discount
factor

$$\gamma \in [0, 1]$$

trades off importance of
immediate vs long-term
rewards

Bellman Equation (Bellman 1957)

Recursive form:

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

$$V_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t \sim \pi(s)]$$

$$= \mathbb{E}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s, A_t \sim \pi(s)]$$

$a \sim \pi(s)$: a is chosen by policy π in state s

(even if π is deterministic)

$$V_{*}(s) = \max_a \mathbb{E}[R_{t+1} + \gamma V_{*}(S_{t+1}) \mid S_t = s, A_t = a]$$

Hamilton-jacobi (bellman) equation

Value function approximations

① State space might be too big

② (Sample) for expectations

Model

Predict what the environment will do next.

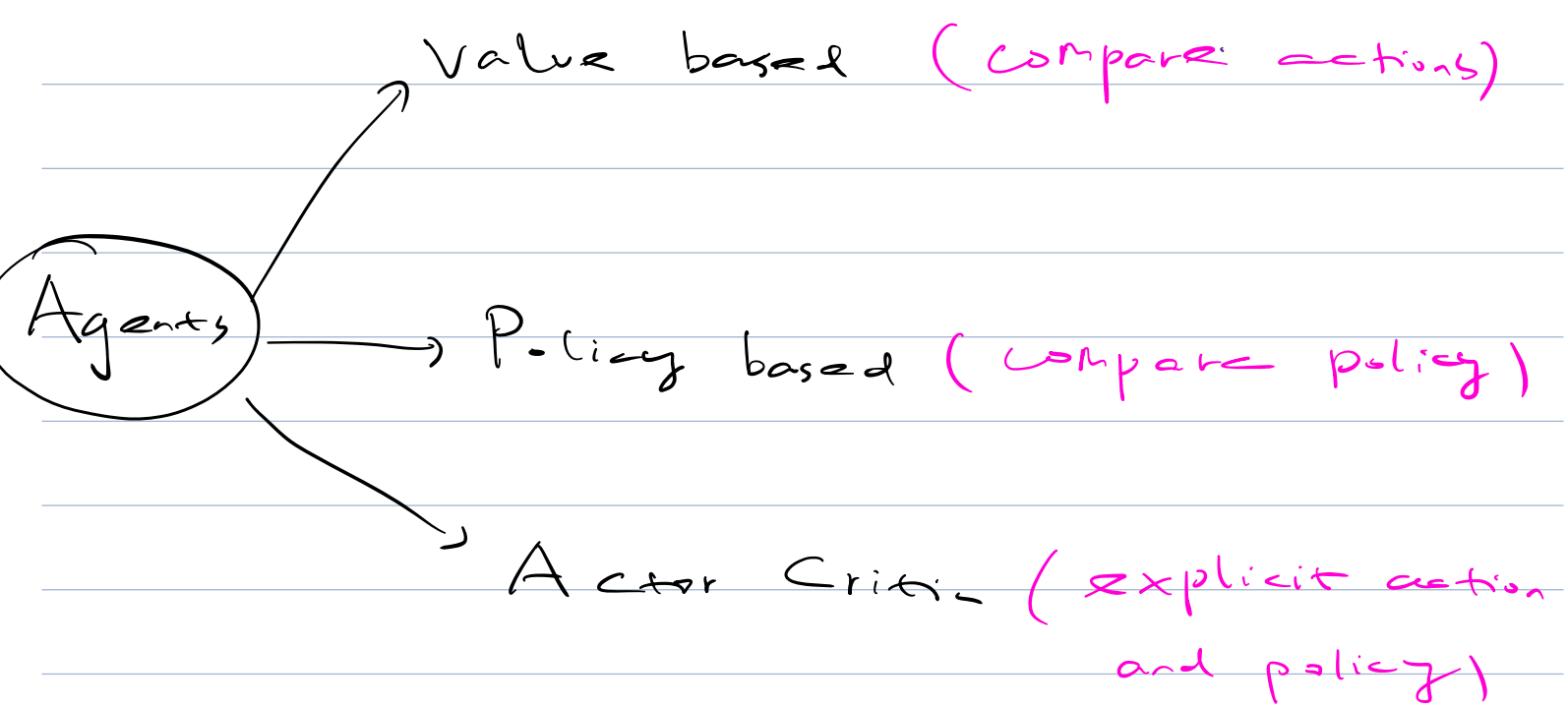
P predicts the next state:

$$P(s, a, s') \approx P(S_{t+1} = s' | S_t = s, A_t = a)$$

R predicts the next (immediate) reward

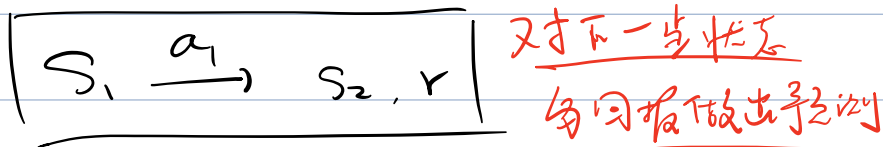
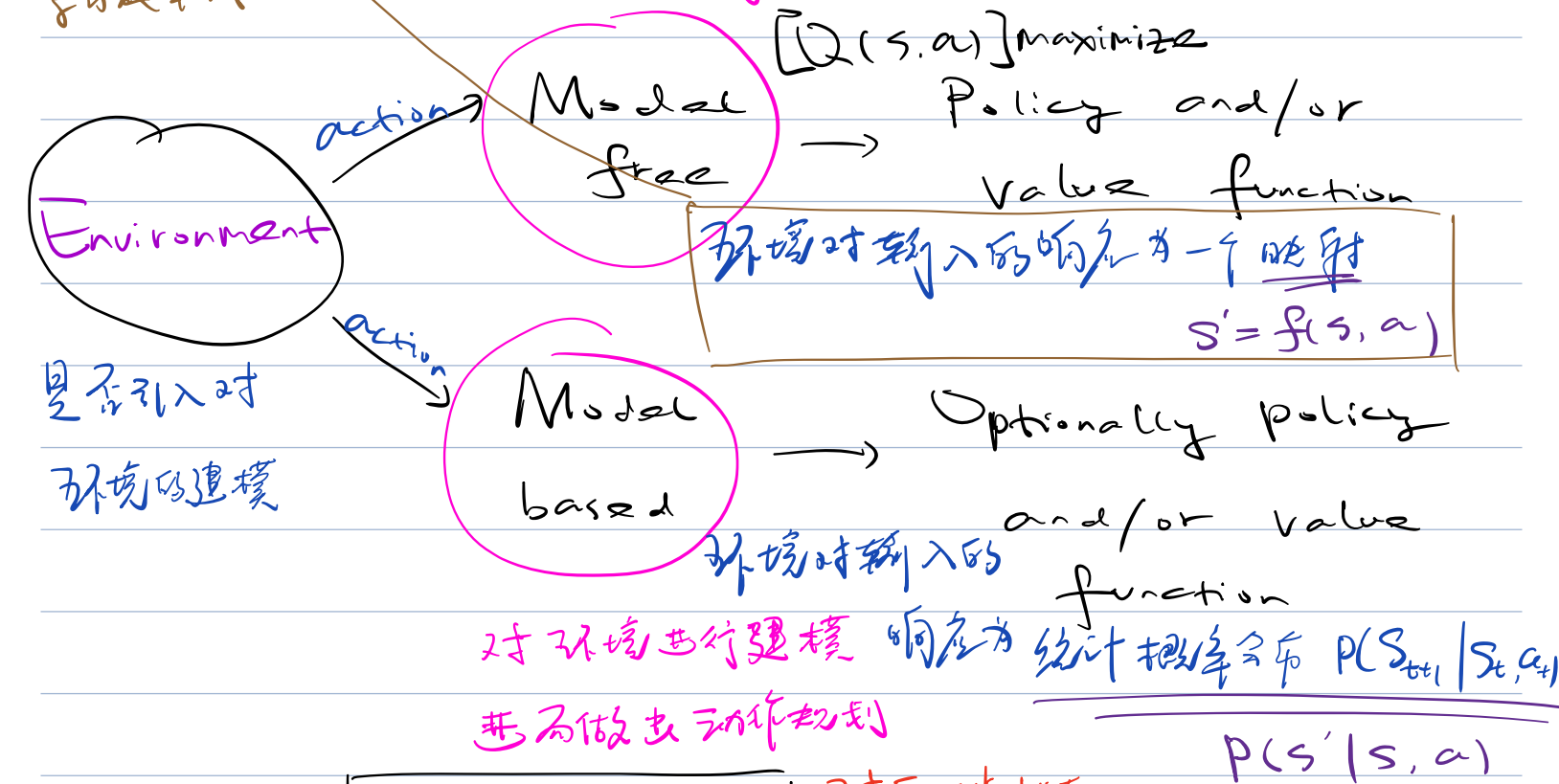
$$R(s, a) \approx \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

Stochastic / generative models



需要大量样本去学习
学习成本低

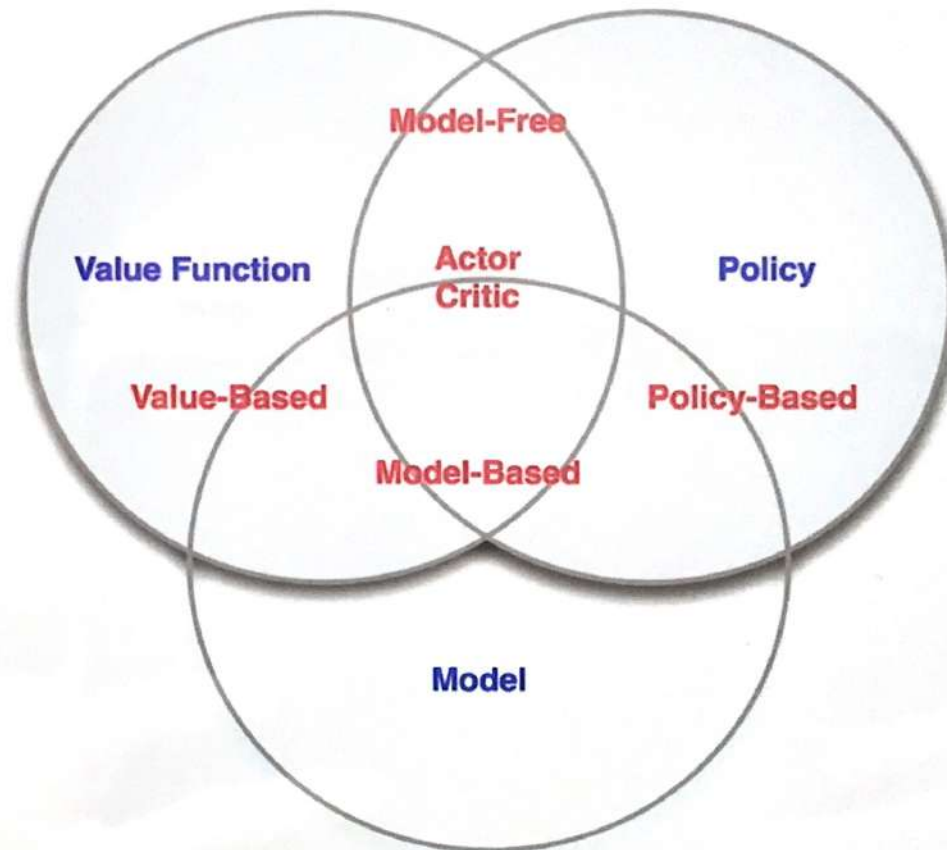
Value-based 单点迭代, 不是GD. (最小化拟合误差)
Policy-based 是GD
不需对环境进行建模 → 最优策略



状态转移概率: $T(s_2 | s_1, a_1)$

奖励函数: $R(s_1, a_1)$

Agent Taxonomy



based extra critic you could have a
model based value based Asians and the

Challenges

1. Learning:

① Environment initially unknown

② agent $\xleftrightarrow{\text{interact}}$ environment

2. planning:

① A model of the environment is given

② The agent plans in the model
(without external interaction)

Predicting

prediction:

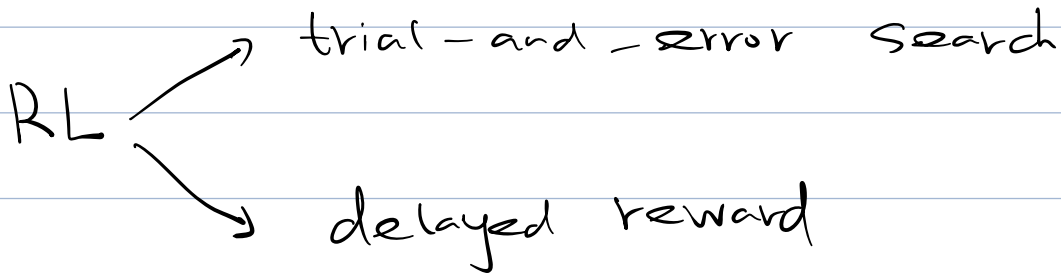
evaluate the future
(for a given policy)

Control:


optimize the future
(for a given action)

optimizing

$$\pi_*(s) = \operatorname{argmax}_{\pi} V_{\pi}(s)$$



Learning the components of an agent

- 
- ▶ All components are functions
 - ▶ Policies map states to actions
 - ▶ Value functions map states to values
 - ▶ Models map states to states and/or rewards
 - ▶ State updates map states and observations to new states
 - ▶ We could represent these functions as neural networks, then use deep learning methods to optimize these
 - ▶ Take care: we often violate assumptions from supervised learning (iid, stationarity)
 - ▶ Deep reinforcement learning is a rich and active research field
 - ▶ (Current) neural networks are not always the best tool (but they often work well)

any subset of those or superjet and
state updates which

IID: identically and independently
distributed