

## 1. Basic Concept

$$H_t = O_1, R_1, A_1, \dots, O_t, R_t, A_t$$

$$S_t = f(H_t)$$

$\&$  : used to determine what happens next

$S_t^e$  : invisible to agent always

$S_t^a$  类比: 信息用于选下一步做什么  
用于

$$S_t^a = f(H_t) \xrightarrow{\text{action.}}$$

information state : Markov State

Markov 状态

$$P(S_{t+1} | S_t) = P(S_{t+1} | S_1, \dots, S_t)$$

i.e.  $S_t$  sufficient statistics

for future

↓

$O_t = S_t^a = S_t^e =$

1 full observable environment

Markov decision process

2. Partial observable environment

$$O_t = S_t^a \neq S_t^e$$

partially observable POMDP  
Markov decision process

$$S_t^a = 1^{\#} H_t \text{ complete history}$$

for  $S_t^e$

$$= 2^{\#} (P[S_t^e = s'], \dots \Rightarrow \text{概率分布})$$

Beliefs of  $S_t^e$

$$P(S_t^e = s^n)$$

$$= 3^{\#} S_t^a = G(S_{t-1}^a W_3 + O_t W_0)$$

observation.

recurrent  
neural  
network

linear 旧信息  $\rightarrow$  新  
combination

2.

RL Agent

Policy : behavior f.

Value function :

Model : subjective representation  
of environment

;

policy 1<sup>st</sup>  $a = \pi(s)$

$$2^{nd} \quad \pi(a|s) = P[A_t = a | S_t = s]$$

value function      prediction of future reward  
用于 evaluate how state  
用于 select action

基于 state, 不同 action.  $\Rightarrow$  future reward

$$V_{\pi}(s) = E_{\pi}(R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots | S_t = s)$$

slope:  $\gamma^n$  衰减  
 discount

Model : prediction of environment

Transition  $P$  : next state

Rewards  $R$  : immediate reward

$$P_{ss'}^a = P(S_{t+1} = s' | A_t = a, S_t = s)$$

$$R_{ss'}^a = P(R_{t+1} | S_t = s, A_t = a)$$

3. 分类

Value-based

value  $f \rightarrow$  policy

Policy-based  
policy

actor critic 均有

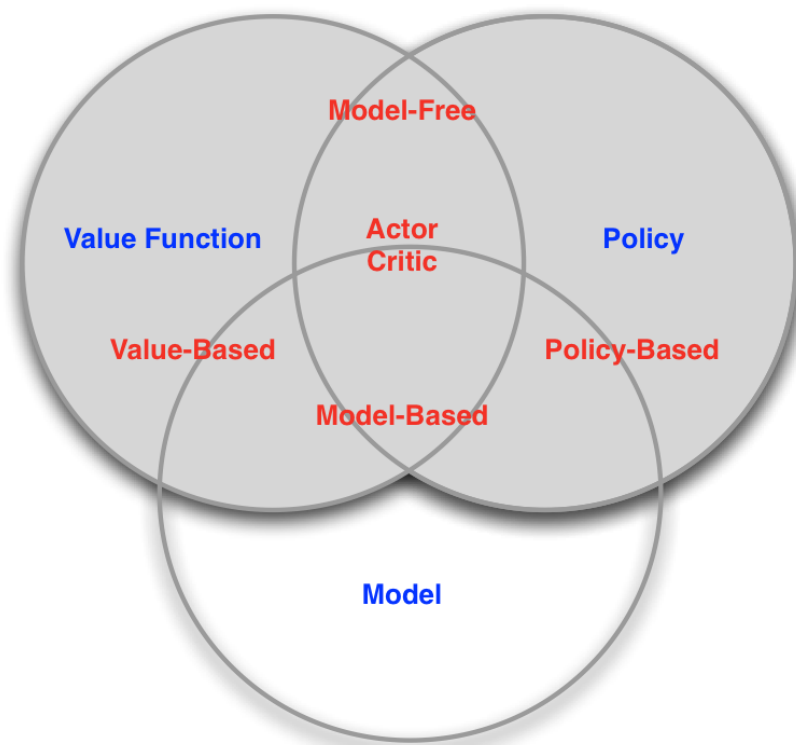
model free

无 model

只根据 Policy  
or value f 选择  
action 选择

model based

有 model



4.

Problems

# 1. Sequential decision making

## ① RL problem

- environment 环境 unknown

- 5 environment

## ② planning

- improve policy

- know environment

- 5 model (internal

model: perfect interaction)

- improve policy

## 2. 对未知

Exploration

vs

Exploitation

try more  
info

use known  
info  $\rightarrow$  reward

### 3. 对 action

predict

vs

control

policy  $\rightarrow$  future  
predict value f

寻找 best policy  
-

Model-free  
不理环境

(envi)

Q learning

Sarsa

Policy gradient

Model-Based  
理环境

(envir) → (model)

建模 →

Q learning

Sarsa

Policy gradient

+ imagination.

Policy Based

概率

Policy gradient

Value-Based

价值

Q learning

Sarsa

↙ ↘  
actor critic

Monte-Carlo update

Temporal -



回合更新

单步更新

On-Policy

自己学 自己采样

Sarsa

Sarsa( $\lambda$ )

Off-Policy

自己

| 观察

其他人

Q

Deep Q