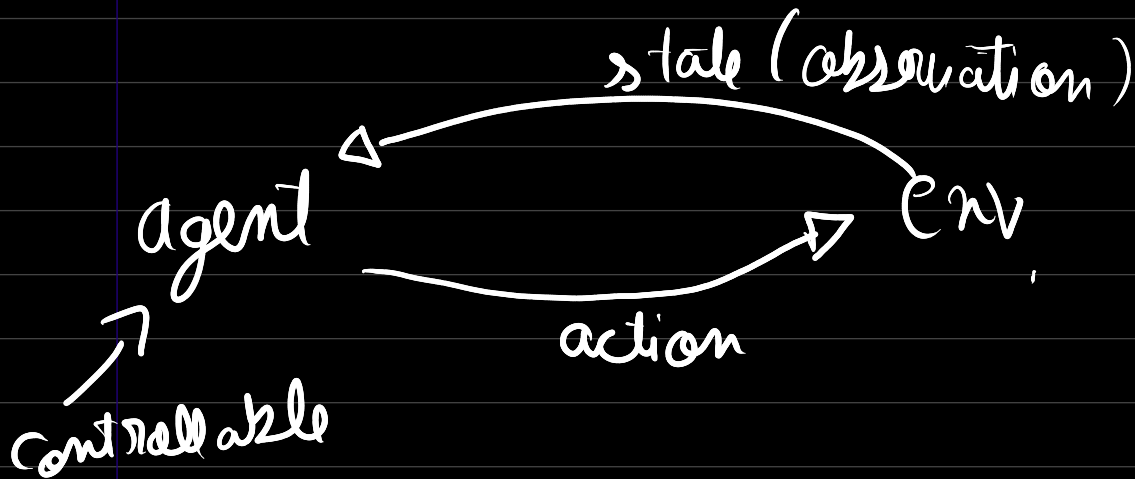


Reinforcement learning

Chap 1 : Markov Decision Process

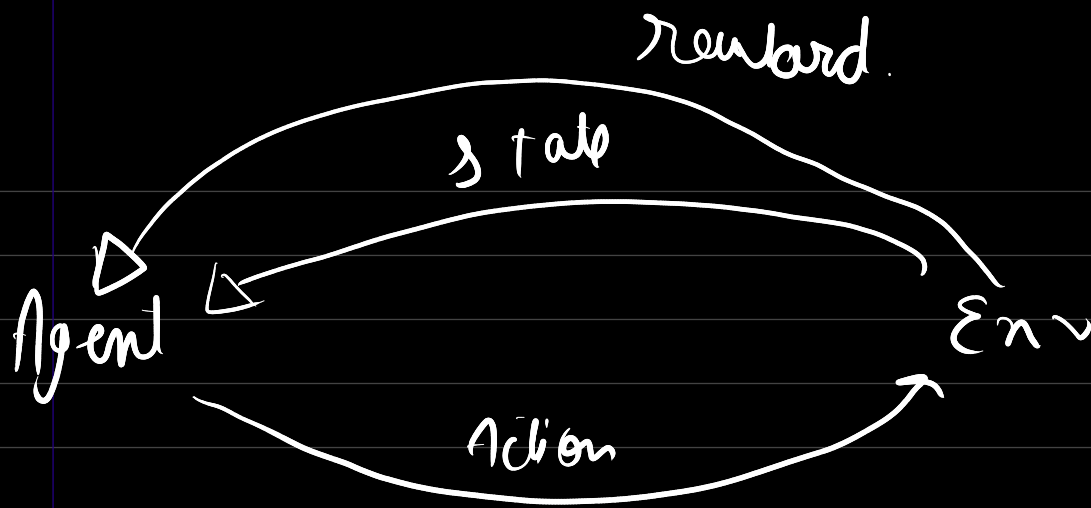


state - position
 - velocity
 - wind color
 - temperature

action - motor speed
 - direction

e. Driving car

program	agent	env	action	state
A	brain	limbs	neuron	embodiment
B	human	car controls	limb movement	elements of car
C	car	road	car movement	road situation



MDP Markov Decision Process

$s-a-r$, $s-a-r$

└─┬─┘
single
episode

$s_0, a_0, r_1, s_1, a_1, r_2, \dots$
 ↑ ↑ ↑
 state action reward

○ Markov principle:

state only dependent
or immediate prev state

(go to recent state

state can depend on more prev states

policy $(\pi) \rightarrow \pi: \mathcal{S} \rightarrow \mathcal{A}$
 takes state \rightarrow gives action

$\pi(a|s)$ = probability of action.

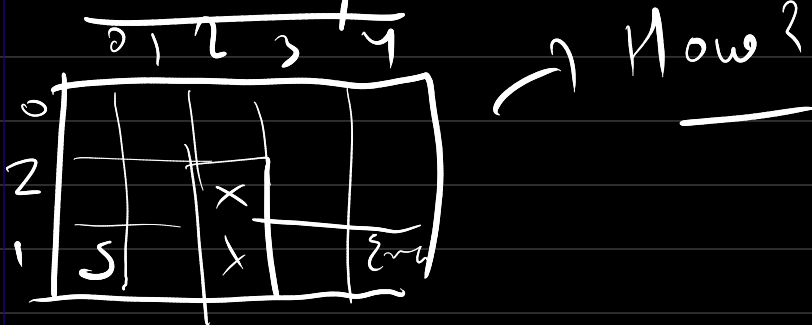
return $G_t \rightarrow$ reward $\times \gamma$

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

γ decay factor

Chap 2

\rightarrow Grid example:



State : 2 nums - x, y coords

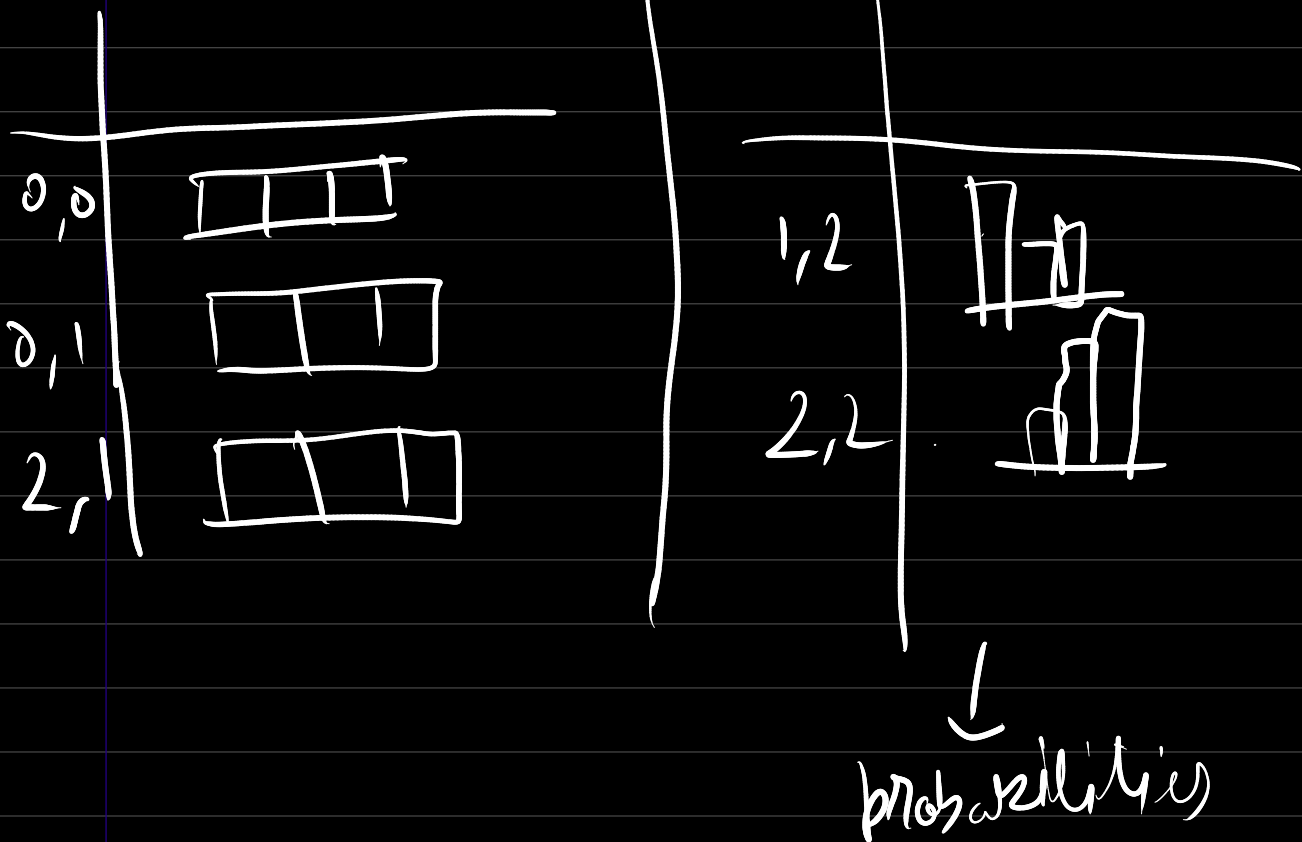
action : 1 num - 0/1/2/3 // up/down/left/right

reward : 0 if in target cell
 -1 otherwise

On } Model
off }

↳ no access to world model

2 policies



→ sampling → collect → Sample Trajectory

$\gamma = 0.8$ eg.

State	action	Next State	Reward	Return
				-9.33
3.0	1	3.1	-1	-1.8
3.1	2	4.1	-1	-1
4.1	1	4.2	0	0

$$\gamma = 0.8$$

It's successful one, when the agent reaches target

if it doesn't return is more negative.

o Policy gradient method : τ rate of change change of policy

State-value $V_{\pi}(s)$

action-value $Q_{\pi}(s,a)$

$$\underbrace{V_{\pi}(s) \quad Q_{\pi}(s|a)}_{\text{not as probability}}$$

- Policy gradient method
- Policy gradient method with value fn
- Just use action-value fn

we develop Q function \rightarrow that evaluates policy

Q_{π}

\downarrow
Better action

0	1	0.1
1	-1	0.5
2	0	2.3
3	0	5.3

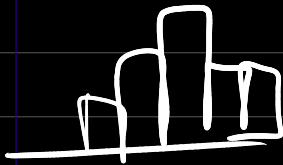
$Q \rightarrow Q_{\pi}$ Best case scenario

(-) maintaining balance between exploration and exploitation

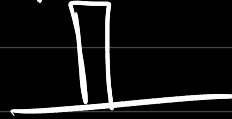
exploration



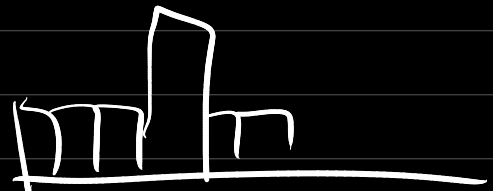
policy
gradient
balance



exploitation



(greedy balance



0.9

Random + High
probability



0.9 \rightarrow 0.1

Monte
Carlo method

Chapter 3 : Temporal
difference.

Monte Carlo \rightarrow works on episode basis
 \swarrow \searrow
good bad

→ need to wait for ep to end to evaluate

if long ep → problem

which action good?

... bad?

→ impact of individual action ← Credit assignment problem

→ SARSA
Choosing
Best

Choosing
 $\max(\text{action})$
all

[How much better or worse relates to new state]

[highest estimated]

On policy

behavior policy = target policy

off policy

\neq , , ,

no learning

{ sample efficient }

Off policy's better

Chapter - 4 : Deep Q Network

Discrete states \rightarrow continuous states \checkmark

Discrete actions \rightarrow continuous actions \times

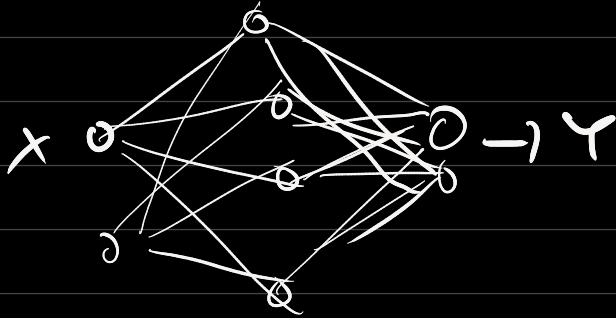
Value based $f_n \rightarrow$ DQN

Swap table with $n \times n$

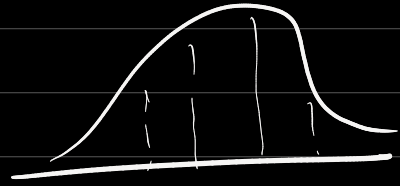
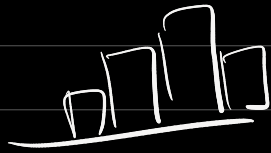
$f_{nn} \rightarrow f_n$ approximator.

Policy Gradient: Chap 5

⇒



(, track of dist of probability



↓ loss fn

inc probability of action acc to how good

it is

→ Probabilities adjust acc. to softmax

→ RL can be reliable.

Too much episodes

