

# Key Concepts in RL

Agent : Agent takes actions

Environment : World in which agent exists & operates

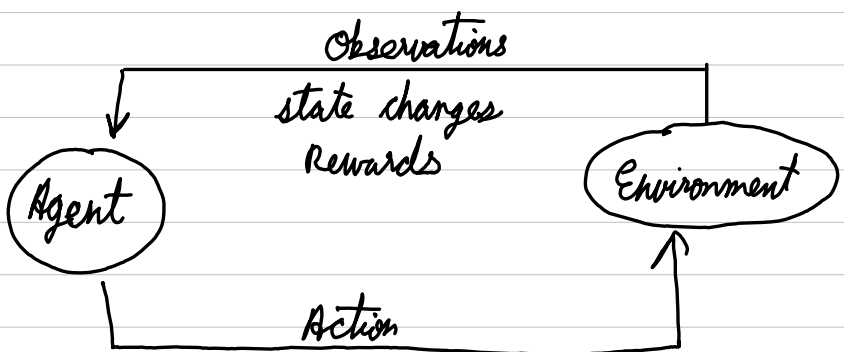
Actions : A move that the agent can make

Action space  $A$  : set of possible actions an agent can make in the environment  
Action space can be discrete or continuous

Observations : Input from the environment

state : situation which the agent perceives

Reward : Feedback that measures the success or failure of the agent's action



Policy : Which action to do in which state ( $\pi$ ) [agent's policy]

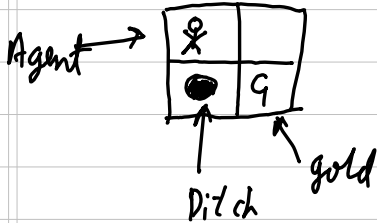
Value function : How good the situation is.  
Prediction of future reward

Model : Predicts what environment will do next  
eg atari game - prediction of opponent

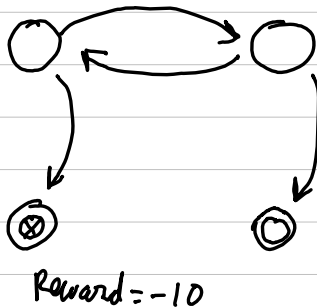
Return : Total reward obtained by the agent

# Example of environment

consider a simple environment

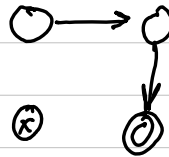


goal of agent is to obtain gold without falling in ditch

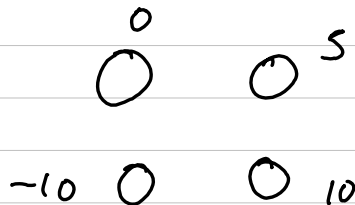


Possible actions

Optimal policy



Optimal value function



## Total Reward

Agent wants to maximize the reward

At any time  $t$ , the agent has a total reward

$$G_t = \sum_{i=t}^{\infty} \gamma^i r_i = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

↑  
Total reward at time  $t$   
also called  $G_t$

↑  
reward returned at time  $i$

↑  
discount factor

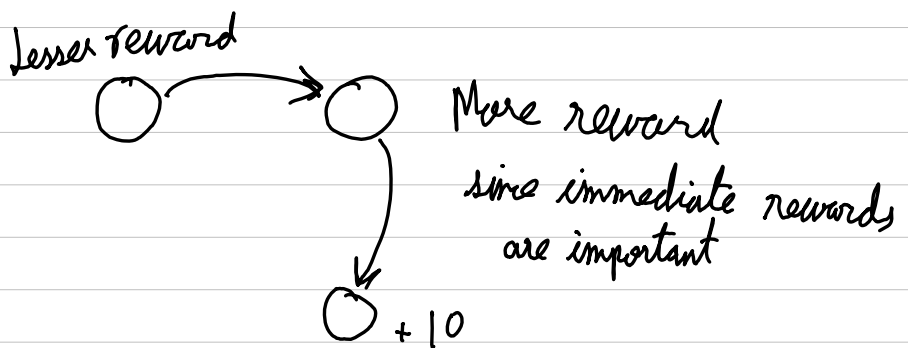
In order to give more importance to immediate rewards, a  $\gamma$  factor is added  
( $0 \leq \gamma \leq 1$ )

↑  
discounting factor

$G_t$  is the total discounted reward

$\gamma$  close to 0 leads to myopic evaluation

$\gamma$  close to 1 leads to far sighted evaluation



Intermediate rewards can be written as

$$R(s, a, s')$$

↑  
current state

↑  
action taken

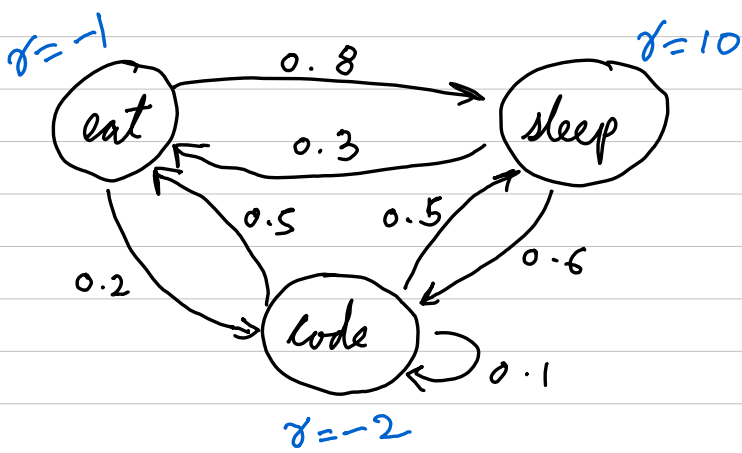
↑  
Next state

Rewards can be assigned to states (Markov reward process) <sup>as in</sup>  
or to actions (Markov decision process) <sup>as in</sup>

Total discounted reward is also known as return

## Markov Reward Process

Consider a simple Markov chain



We need the agent to spend maximum time in state sleep.

Hence we put the reward in state sleep.

We put negative rewards in other states

In order to prevent the agent from going round and round in intermediate states, negative weights are used.

Now we can perform a random walk and calculate the reward

$$\begin{aligned} \text{eg } \text{eat} - \text{eat} - \text{sleep} &= -2 + -2 + 10 = 6 \\ \text{eat} - \text{sleep} - \text{code} - \text{code} &= -2 + 10 + -2 = 6 \end{aligned}$$

A Markov reward process is a Markov chain with values.

$$\text{Tuple} \rightarrow \langle S, P, R, \gamma \rangle$$

$S \rightarrow$  states

$P \rightarrow$  state transition  $p$

$\gamma \rightarrow$  discount factor

$R \rightarrow$  Rewards

In Markov reward process, the model is entirely random walk.

We are just interested in the rewards ( $V$ ) of the state

there is only  $P$ , no  $\pi$

used to compute the value function.

## Value function in MRP

Since MRP is a random walk, we don't have any control over where the agent will go.

Where agent will go is decided by transition probability  $P$ .

$P(s|s')$  define the likelihood that the agent goes from state  $s$  to  $s'$

$$V(s) = R(s) + \gamma \sum P(s'|s) V(s')$$

Goodness of state = Reward obtained + Weighted of prob sum of Goodness of next possible states

$V$  depends upon  $V$  itself. So how to calculate it? It can be done by two methods

- ① Eigenvalues  $\rightarrow$  costly
- ② D.P.  $\rightarrow$  Better

# Dynamic Programming

Dynamic — sequential or temporal component

Programming — optimizing a linear program.

Method of solving complex problem by breaking into small steps.

Solve subproblems  $\rightarrow$  combine solutions

Output from previous state is retained to make the next step.

Optimization of policy is done. Random initial policy is updated to improve the agent.

In order to be solved by dynamic programming the problem must satisfy 2 properties —

1) Optimal substructure

Problem can be decomposed.

2) Overlapping subproblems

Subproblems recur many times

subproblems can be cached & reused

MDPs satisfy both properties

Bellman equations gives recursive decomposition  
Value function stores and reuses values

D.P. is a model based technique requiring knowledge about  $P$  &  $R$  in the environment<sup>K</sup>

# MRP Bellman equations

We know that for MRP

$$V(s) = R(s) + \gamma \sum P(s/s') V(s')$$

In Matrix format

$$\begin{bmatrix} V(1) \\ V(2) \\ V(3) \\ \vdots \end{bmatrix} = \begin{bmatrix} R(1) \\ R(2) \\ R(3) \\ \vdots \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} V(1) \\ V(2) \\ V(3) \\ \vdots \end{bmatrix}$$

$$V = R + \gamma P V \quad \text{--- (1)}$$

$$(I - \gamma P) V = R$$

$$\therefore V = (I - \gamma P)^{-1} R$$

Solving via eigenvectors for the Bellman equations — time complexity  $O(n^3)$

Only possible to solve for small no of states

for others use Monte Carlo, dynamic and iterative methods

## DP for MRP

In DP, we use the past values to find the next values and converge at the optimal values.

$$V_{k+1}(s) = \sum_{s' \in S} P(s'|s) [R(s, s') + \gamma V_k(s')]$$

In Matrix form

$$V_{k+1} = R + \gamma P V_k$$

D.P. assumes that full knowledge of MRP is present

Iterative application of Bellman expectation backup

$$V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V_{\pi}$$

At each iteration :

for all states  $s \in S$

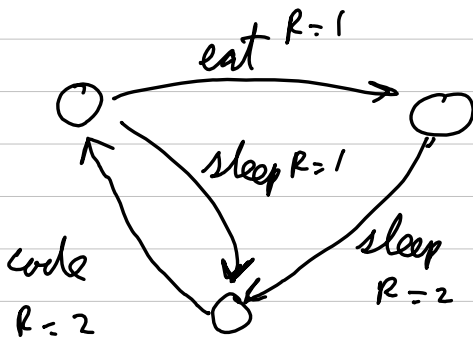
Update  $V_{k+1}(s)$  from  $V_k(s')$

( $s'$  is successor state of  $s$ )



# Markov Decision Process

Instead of rewarding the states, we will reward actions



Here decision, policy comes into picture

Tuple  $\langle S, \underline{A}, P, R, \gamma \rangle$

$A \rightarrow$  Actions.

It incorporates actions & policies.

The objective is to calculate the optimal policy that maximizes expected return

Where agent will go is decided by  $\pi$  &  $P$

$\pi \rightarrow$  Policy : where agent wants to go  
ie which action to select

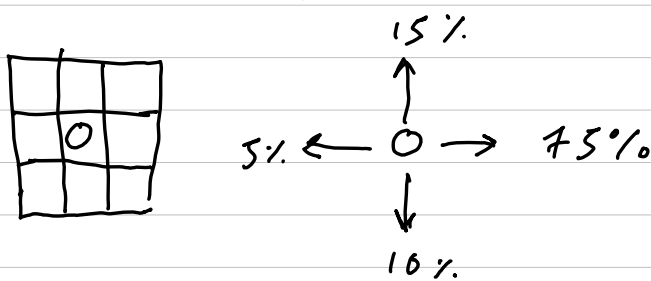
$P \rightarrow$  transition probability : where the environment takes the agent.

eg  $\pi$  may be choose action [up]

But for [up] action  $P$  for reaching state 2 may be 70% while 30% of times it may fall into ditch.

## Policy (MPP)

Policy is a strategy used by agent



It can be thought of as probability to choose action.

The policy decides what actions the agent must take in order to reach the goal

Initially the policy can be random and then improved iteratively

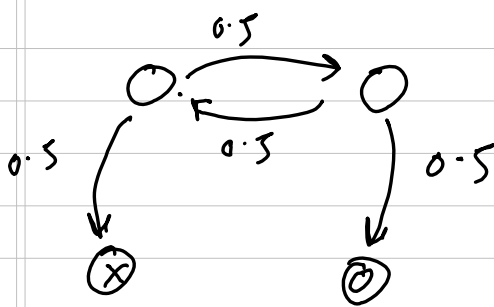
There are two types of policies

Deterministic  $\rightarrow$  always consider 1 action  
eg up  $\pi = [\text{up}]$

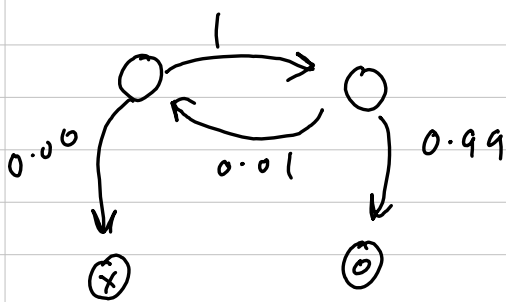
Stochastic  $\rightarrow$  consider probability of transitions

$$\pi = [0.15, 0.75, 0.1, 0.05]$$

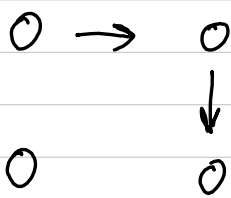
eg a random initial policy may be



But optimized policy may be



for deterministic



Deterministic policy may make moves predictable

# V function in MDP

Used to evaluate goodness or badness of a state  
Prediction of future reward

$$V_{\pi}(s) = E_{\pi} \{ G_t \mid s_t = s \}$$

Value w.r.t. Policy  $\pi$  = Expected total rewards to be obtained in the state

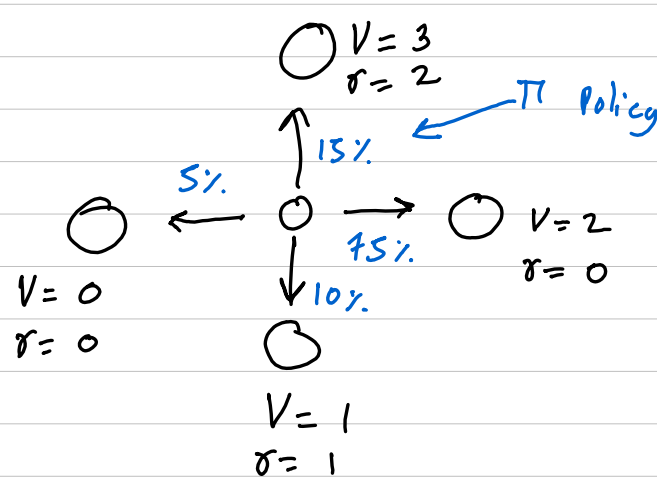
If  $\pi$  is stochastic then  $V_{\pi}$  is weighted probability of transitions.

$$V = \sum_{\text{for all actions } a} \left[ \text{Probability of choosing action } a \text{ from } \pi \times \left[ \text{Reward obtained Now} + \text{Expected reward of next state from } a \right] \right]$$

The expected reward of next state obtained from  $a$  is nothing but  $V$  of Next state

If we want to weigh  $\gamma$ , we can multiply  $V$  by  $\gamma$

for now assuming that agent can go wherever it wishes  
ie. No wind  
 $p(s, a, s') = 1$   
eg board games



at  $\gamma=1$

$$V = 0.15 \times (3+2) + 0.75 \times (2+0) + 0.1 \times (1+1) + 0$$

$V$  needs policy  $\pi$  because  $\pi$  will decide the behaviour of agent in the environment

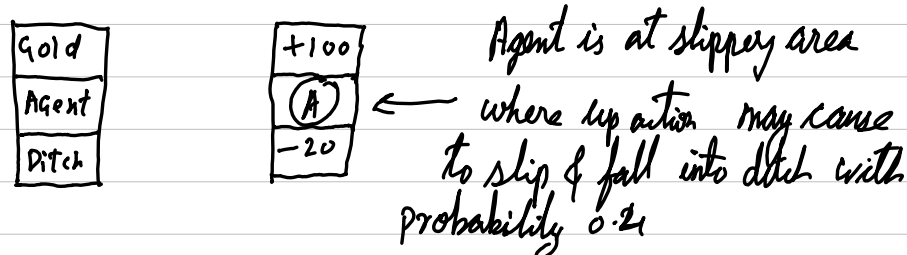
If policy  $\pi$  is deterministic then there will be a single action

$$V_{\pi} = \text{Reward obtained} + \text{Expected reward of next state}$$

for MDP,  $V$  always dependant on  $\pi$

$V$  is measure how good the states are w.r.t.  $\pi$ . Environment  $(P, R)$  also comes into picture

for MRP,  $V$  only depends on environment



for MRP,  $\uparrow 0.3$   $\downarrow 0.7$  then  $V = 100 \times 0.3 - 20 \times 0.7$

Since the agent has high chance of slipping, we say state is bad

for MDP, let  $P$  up:  $\uparrow 0.6$   $\downarrow 0.4$  down:  $\downarrow 1$

If  $\pi$   $\uparrow 0.9$   $\downarrow 0.1$  then  $V$  is still good

$$V = 0.9 \times (0.6 \times 100 - 20 \times 0.4) + 0.1 \times (-20) = 44.8$$

But if  $\pi$  itself is bad

$\pi$ :  $\uparrow 0.1$   $\downarrow 0.9$  Then there is no hope that agent will reach gold

$$V = 0.1 \times 0.6 \times 100 - 0.1 \times 0.4 \times 20 - 0.9 \times 1 \times 20 = -12.8$$

even if the environment permits success 60% times, agent has not learnt properly and takes success only  $0.1 \times 0.6 = 0.06 = 6\%$  of times

So  $V$  of state is still less, even lesser than random policy.

There is no use of the gold being there since agent cannot exploit it

so from view of  $\pi$ ,  $V$  is still bad.

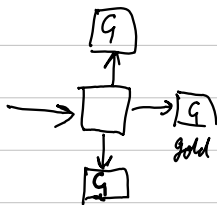
In MDP,  $V$  depends on the policy used and hence denoted by  $V_{\pi}$

That is, goodness or badness of state depends on  $\pi$  & environment both

$\gamma$  value (discount factor)

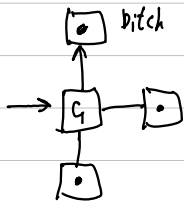
state is "good" if it gives me high rewards on reaching it

AND  
The next states that can be reached from that state are also "good"



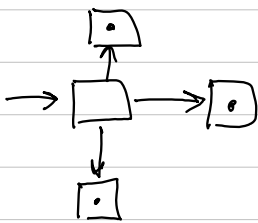
$R$  is low but  
Next states are good

(A)



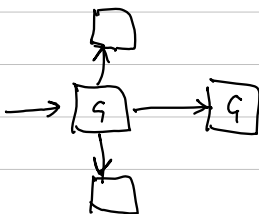
$R$  is high but  
Next states are bad

(B)



$R$ , Next states bad

(C)



$R$ , Next states  
good

(D)

If  $\gamma$  is reduced, then more focus will be given to immediate rewards  
at high  $\gamma$ , future rewards prioritized

$\gamma = 0$  (B), (D) Greedy

$0 < \gamma < 1$  (D) Near sight

$\gamma = 1$  (D), (A) far sight

$\gamma > 1$  (A) future rewards

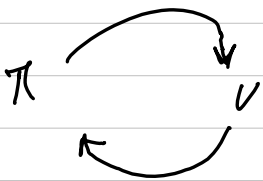
# Policy Iteration

Goal of RL is to find the policy  $\pi$  &  $V$ .

Policy can be chosen to maximize the expected rewards. or going to the best states always.

$V$  needs  $\pi$  to evaluate the goodness

$\pi$  needs  $V$  to find best path



chicken & egg problem  
expectation maximization

Without chicken egg, the equation can be directly solved, but is infeasible.

In RL, EM is called as  $\rightarrow$

① Policy evaluation (find  $V$ )  $\pi \rightarrow V$

② Policy improvement (find  $\pi$ )  $V \rightarrow \pi$

This is called as

- ① Policy iteration - Randomly select policy  $\pi$   
find value function  $V$   
Improve policy ....



$V^*$  &  $\pi^*$  are optimized value function & policy  $\pi$

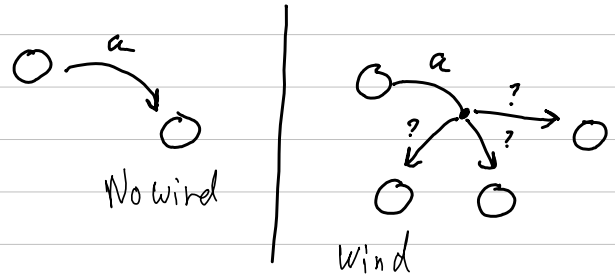
$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{M} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{M} \dots \pi_* \rightarrow V_*$$

This method is the dynamic programming method.

# Transition probabilities

This is used when taking an action does not guarantee going to a state

eg windy travel.



$$P(s' | s, a)$$

$\uparrow$  Next state       $\uparrow$  current state       $\nwarrow$  action

eg  $P(X | Y, a) = 0.8$  means that there is 80% chance of going from state  $X$  to state  $Y$  after taking action  $a$ .

so while finding  $V, \pi$ , these transition probabilities must also be taken into account

$$V = \underbrace{\left[ \text{Probability of choosing action } a \text{ from } \pi \right]}_{\text{all actions}} \times \underbrace{\left[ \text{Probability of that action going to state } s' \right]}_{\text{all states from } a} \times \left[ \text{Reward} + V_{\text{next}} \right]$$

# Mathematical Notation of Policy iteration

Policy evaluation  $\rightarrow$  finding  $V$  wrt  $\pi$

Probability that on choosing an action  $a$ , from states agent lands on state  $s'$ . These are known beforehand from how the environment functions

Transition Probabilities

discount factor

goodness of next state  $s'$

Reward obtained from transition  $s \rightarrow s'$  by action  $a$

set of all states

$\propto$  Next state  $s'$

set of all possible actions  $a$

Goodness of current state

$$\pi V(s) = \sum_{a \in A} \pi(a/s)$$

$$\sum_{s' \in S} P(s'/s, a) \times [R(s, a, s') + \gamma V^\pi(s')]$$

Probability weights that agent chooses action  $a$  while in state  $s$  (strategies of agent)

stochastic policy

## Mathematical Notation of Policy iteration

Policy evaluation  $\rightarrow$  finding  $V$  wrt  $\pi$

for deterministic policy, there will always be a single action

$\pi(a|s)$  will not be a probability list but rather a single action.

Hence we don't need to take sum of all but rather only one action  $\pi(s)$

$$V_{\pi} = \sum_{s' \in S} p[s', s, \pi(s)] \cdot (R(s, \pi(s), s') + \gamma V^{\pi}(s'))$$

Sum only once here



## Q function

Q function captures the expected total future reward an agent in state "s", can receive by executing a certain action "a"

$$Q(s_t, a_t) = E(G_t \mid s_t, a_t)$$

↑      ↑      ↑      ↑      ↑  
Q function   state   action   expectation function   total reward

Q function takes in a state and a possible action and tells how much total reward will be obtained.

Choose the action that maximizes future reward

Same as V function, but V function is goodness for a state while Q function is goodness for a state action pair.

Also known as action value function.

$$Q^\pi(s|a) = \sum_{s' \in S} P(s'|s, a) \left[ R(s, a, s') + \gamma V^\pi(s') \right]$$

## V Q $\pi$ P explanation

V  $\rightarrow$  How good a state is

Q  $\rightarrow$  How good an action (from a state is)

$\pi \rightarrow$  Probability that agent will take an action

P  $\rightarrow$  Probability of going to a state from an action

If we look at form of equations, something like

$$V = \sum \pi \sum P \times (\gamma + \gamma V')$$

represents Weighted sum (wrt  $\pi$ ) of expected reward of actions

$$Q = \sum P \times (\gamma + \gamma V')$$

represents expected reward of an action

which is Weighted sum (wrt P) of expected reward of next state

$$\text{expected reward of a state} = \gamma + \gamma V'$$

$$= \underset{\text{reward}}{\text{actual}} + \underset{\text{state}}{\text{goodness of}}$$

$\gamma(s, a, s')$  depends on current state  $s$   
action taken  $a$   
next state  $s'$

or can be simply reward of  $(s, a)$  or  $(s')$

Note that we are dealing with two probabilities

P — environment  
 $\pi$  — agent policy

Hence both are weighted.

In some cases P may be fixed as 1

## Mathematical Notation of Policy iteration

Policy improvement  $\rightarrow \pi$  from  $V^\pi$

get the action that greedily selects the action that maximizes the expected reward

for deterministic policy

$$\pi'(s) = \max_{a \in A} \sum_{s' \in S} P(s'|s, a) \left[ R(s, a, s') + \gamma V_\pi(s') \right]$$

for stochastic policy we need to consider the expected return from the action and state

$$\pi(s|a) \rightarrow \sum_{s' \in S} P(s'|s, a) \left[ R(s, a, s') + \gamma V_\pi(s') \right]$$

stochastic  $\rightarrow$  sum of return for action  $a$   
 $\pi(s, a)$

Deterministic  $\rightarrow$  sum of return for all actions  
 $\pi(s)$

Choose best action

By using  $Q$  function, we can write

$$\pi(s|a) \rightarrow Q^\pi[s, a]$$

However there is a problem that this is not normalized.

If we want to normalize it then we can use softmax function (Boltzmann function) to normalize it

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix} \rightarrow \frac{e^{z_i}}{\sum e^{z_j}} \rightarrow \begin{matrix} \text{Prob} \\ \begin{bmatrix} 0.02 \\ 0.90 \\ 0.03 \\ 0.01 \\ 0.02 \end{bmatrix} \end{matrix}$$

This is done to ensure that the output values are in range  $(0, 1)$  while summing up to 1

$$\text{Here } \pi'(a|s) = \frac{e^{Q^\pi(s, a)}}{\sum_{b \in A} e^{Q^\pi(s, b)}}$$

## Value Iteration

Instead of going EM on  $V$  &  $\pi$ , why not combine it in single step?



calculate optimal value of  $V$  first and then go to  $\pi$

This is called

② Value iteration

$$V(s) = \max_{a \in A} \sum_{s' \in S} P(s'|s, a) (R(s, a, s') + V(s'))$$

Optimal value of state is the maximum expected reward obtainable by choosing the best action and then following the optimal policy thereafter.

$$\text{Goodness of current state} = \text{Reward Now} + \text{Goodness of Next best state}$$

or more correctly

$$\text{Goodness of current state} = \text{Expected reward obtained by choosing best action}$$

↓

action which gives best expected reward.

At last the policy can be found out

$$\pi^*(s) = \max_{a \in A} \sum_{s' \in S} P(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$$

Here  $\pi^*(s)$  will be deterministic policy

If there is a tie, only then it will be stochastic

## Policy Iteration

Random Policy

$\pi \vee \pi \vee \pi \vee \dots \pi^* \vee^*$

Less computation

Faster

Few iterations

Prediction + Control Problems  
 $(\pi, v)$   $(\pi, v, a)$

## Value Iteration

Random Value function

$v \vee v \vee v \vee \dots v^* \pi^*$

More computation

Slower

Many iterations

Control problems

Both are

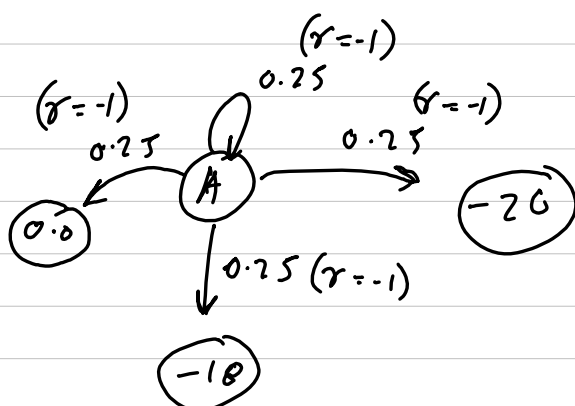
- ① Guaranteed to converge
- ② Variations of Bellman updates
- ③ One step look ahead

Bellman equations can be solved directly through eigenvalues, but it is a very time consuming method  $O(n^3)$  hence dynamic programming methods are used

Q1) Consider the grid world problem where agent follows equiprobable random policy in deterministic transitions ( $P=1$ ) transitions that leave the grid are treated as unchanged.  $r = -1$  for all transitions. Use Bellman equations to find  $V$  of state A.

$$\begin{array}{ccc} 0.0 & \boxed{A} & -20 \\ -14 & -18 & -20 \end{array} \left. \vphantom{\begin{array}{ccc} 0.0 & \boxed{A} & -20 \\ -14 & -18 & -20 \end{array}} \right\} \begin{array}{l} V \text{ table with} \\ V(A) \text{ missing} \end{array}$$

→ since  $P=1$  &  $\pi = [0.25, 0.25, 0.25, 0.25]$



$$V(A) = \sum_S \pi(A, a, s') \sum_S P(S, a, s') (r(S, a, s') + V(s'))$$

since  $P = 0$  for all states except  $s$

$$V(A) = \sum_S \pi(A, a, s') (r(A, a, s') + V(s'))$$

since  $r = -1$  for all,

$$\begin{aligned} V(A) = & 0.25 \times (-1 + 0) \\ & + 0.25 \times (-1 + V(A)) \\ & + 0.25 \times (-1 + -20) \\ & + 0.25 \times (-1 + -18) \end{aligned}$$

$$V(A) = -10.5 + 0.25 \times V(A)$$

$$+ 0.75 V(A) = -10.5$$

$$\therefore V(A) = -14$$

## DP for Value iteration MDP

We can similarly denote the equations in matrix form

$$V^\pi = (I - \gamma P^\pi)^{-1} R^\pi$$

$P^\pi \rightarrow$  state transition probability matrix under policy  $\pi$

$$P_{ss'}^\pi = \sum_{a \in A} \pi_{s,a} P_{s,a,s'}$$

$R^\pi$  is the reward vector under policy  $\pi$

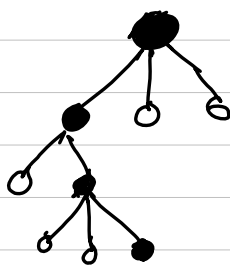
$$R_s^\pi = \sum_{a \in A} \pi_{s,a} \sum_{s' \in S} P_{s,a,s'} R_{sas'}$$

So similarly, for MDP, Dynamic programming can be used to solve the Bellman equations by step by step  $V_1, \pi_1$

$V_2, \pi_2$   
 $\vdots$   
so on



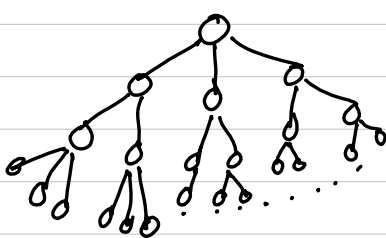
Dynamic programming solves MDPs through systemic algorithms like policy iteration and value iteration



} explores all options of a visited state

Repeat  $s'$   $\left\{ \begin{array}{l} \text{Visit state } s, \\ \text{search all possible next states from } s, \\ \text{choose state } s' \text{ that gives best result} \end{array} \right.$

It is different from exhaustive search



} explore all options of all states

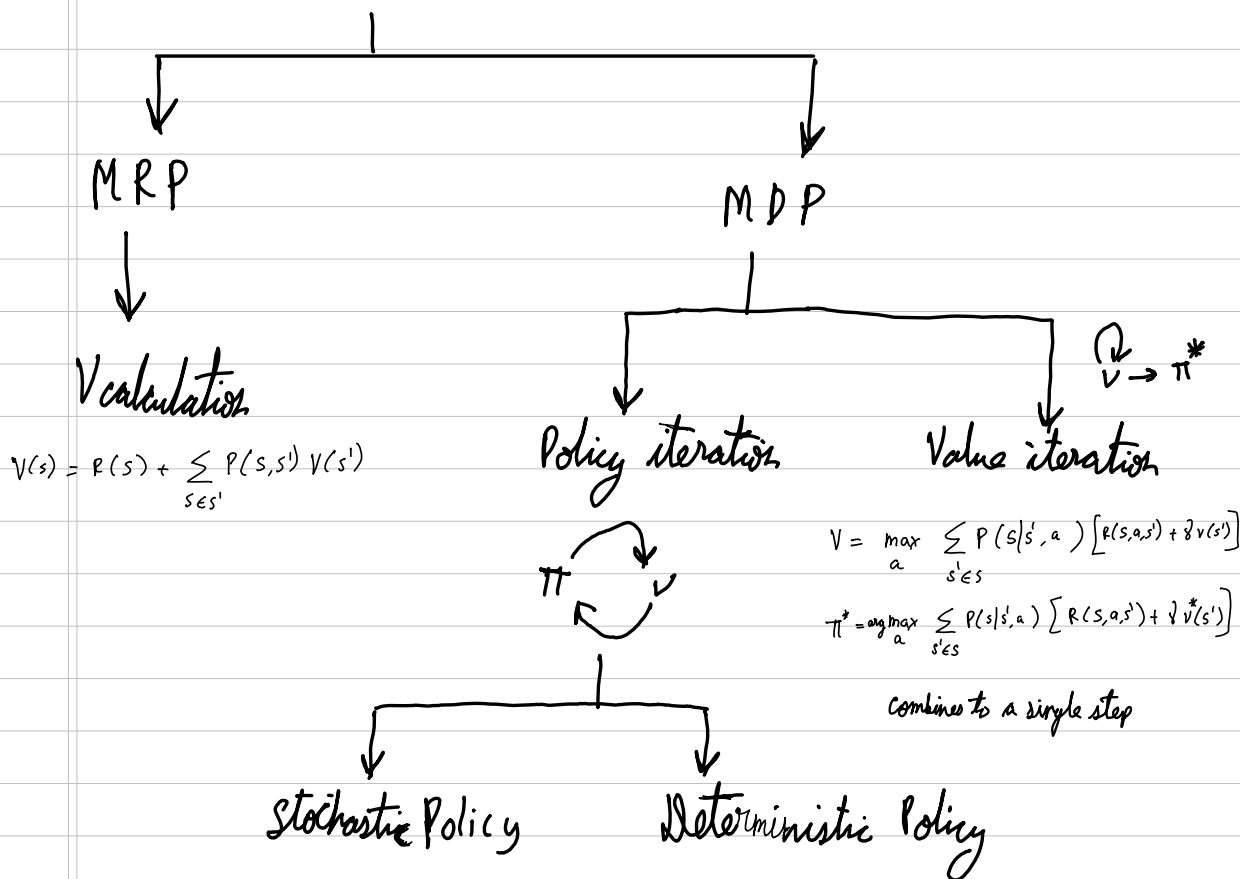
Dynamic programming requires knowledge of the environments' transition probabilities  $P(s, a, s')$

Most important equation  $\rightarrow$

$$\text{Goodness of state} = \text{Prob of going to next state} \times \left( \text{Reward of Next state} + \text{Goodness of Next state} \right)$$



# Bellman equations



$$\vec{V}(s) = \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P(s|s', a) [R(s, a, s') + \gamma V(s')]$$

$$V(s) = \sum_{s' \in S} P(s|s', \pi(s)) [R(s, \pi(s), s') + \gamma V(s')]$$

$$Q(s, a) = \sum_{s' \in S} P(s|s', a) [R(s, a, s') + \gamma V(s')]$$

$$\vec{V}(s) = \arg \max_a \sum_{s' \in S} P(s|s', a) [R(s, a, s') + \gamma \vec{V}(s')]$$

$\pi(s)$  returns action

$$\pi(s, a) = \frac{e^{Q(s, a)}}{\sum_{b \in A} e^{Q(s, b)}}$$

$\pi(a)$  returns probability matrix  
 $\pi(s, a)$  is a probability