## Elements common to all
### control tasks

**Environment**



state $(s_t)$

action $(a_t)$

Reward $(r_t)$

agent
$\pi(s)$

aim

# Examples of Environment

Discuss St, at, rt, stp1 in the context of standard environment like chess, pacman, robot, etc

MDP :

$$(S, A, R, P)$$

MDP has no memory

$$P\left(S_{t+1} \mid S_t = s\right) = P\left(S_{t+1} \mid S_t = s_t, S_{t-1} = s_{t-1}, \dots\right)$$

Episodic Vs Continuing   MDPs ] examples ??

Reward : $R_t$   (instant Gratification)

Return : $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$
(cumulative long term) gain

$\gamma$ : discount factor

Q. Why bring in the discount factor?
[incentive to work faster]

— what happens if $\gamma = 0$ ?
— what happens if $\gamma = 1$ ?

# Policy

Central Problem in RL projects

Learn to get policy $\pi$

$\pi(s) \longrightarrow$ take a specific action at state $s$

$\pi(a|s) \longrightarrow$ probability of taking action $a$ given state $s$

Aim: find optimum policy $\pi_*$

## State Value

$$v_\pi(s) = E\left[G_t \mid S_t = s\right]$$

expected return starting from a particular state and following a policy $\pi$

$$= E\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots \mid S_t = s\right]$$

q-value of a state-action pair:

$$q_\pi(s,a) = E\left[G_t \mid S_t = s, A_t = a\right]$$

$$= \text{expected return starting from state } s,$$
taking action $a$ and following policy $\pi$
from there on.

# Bellman Equations

$$v_\pi(s) = E\left[ G_t \mid S_t = s \right]$$

$$= E\left[ R_{t+1} + \gamma G_{t+1} \mid S_t = s \right]$$

$$= \sum_a \pi(a|s) \sum_{s',r} P(s',r \mid s,a) \left[ r + \gamma v_\pi(s') \right]$$

Discuss Bellman Equation for v

$$q_\pi(s,a) = E\left[ G_t \mid S_t = s, A_t = a \right]$$

$$= E\left[ R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a \right]$$

$$= \sum_{s',r} P(s',r \mid s,a) \left[ r + \gamma \sum_{a'} \pi(a'|s') q_\pi(s',a') \right]$$

Discuss Bellman Equation for q

Optimal Policy $\pi_*$ depends on knowing optimal value

$$\pi_*(s) = \underset{a}{\arg\max} \sum_{r,s'} P(s',r \mid s,a) \left[ r + \gamma \, v_*(s) \right]$$

$$= \underset{a}{\arg\max} \; q_*(s,a)$$

Optimal values depends on optimal policy

$$v_*(s) = E_{\pi_*}\left[ G_t \mid S_t = s \right]$$

$$q_*(s,a) = E_{\pi_*}\left[ G_t \mid S_t = s, A_t = a \right]$$

DILEMMA
chicken & egg

# Bellman Optimality Equation

30 March 2023    11:30

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) \left[ r + \gamma\, v_\pi(s') \right]$$

] Value equation

$$v_*(s) = \max_a \sum_{s',r} P(s',r|s,a) \left[ r + \gamma\, v_*(s') \right]$$
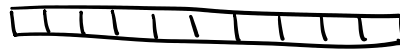
] equation for optimal $v$

$$q_\pi(s,a) = \sum_{s',r} P(s',r|s,a) \left[ r + \gamma \sum_{a'} \pi(a|s)\, q_\pi(s',a') \right]$$

} eqⁿ for $q$

$$q_*(s,a) = \sum_{s',r} P(s',r|s,a) \left[ r + \gamma \max_{a'} q_*(s',a') \right]$$

} eqⁿ for optimal $q_*$

# Simple Grid

Environment : 10 positions

$\rightarrow$ agent starts at a random position.

$\rightarrow$ at each position, agent can move left or right.

$\rightarrow$ one location is a target (fixed but not told to agent)

$\rightarrow$ reward = 1 if next state is the target,
  -1 else if next state is outside
  0 all else

## Exercise: Code this environment

- constructor function
- reset function
- step function
- view function

# Optimal Policy = ??

Recall the optimality equations

$$v_*(s) = \max_a \sum_{s',r} P(s',r|s,a) \left[ r + \gamma v_*(s') \right]$$

$$q_*(s,a) = \sum_{s',r} P(s',r|s,a) \left[ r + \gamma \max_{a'} q_*(s',a') \right]$$

These Equations are used as update rules :-

start with random $V(s)$ $\forall s$

$$V(s) = \max_a \sum_{s'r} P(s',r|s,a) \left[ r + \gamma V(s') \right]$$

This is very similar to the Jacobi & Gauss Siedel methods of solving equation

Limitation: $P(s',r|s,a)$ needs to be known
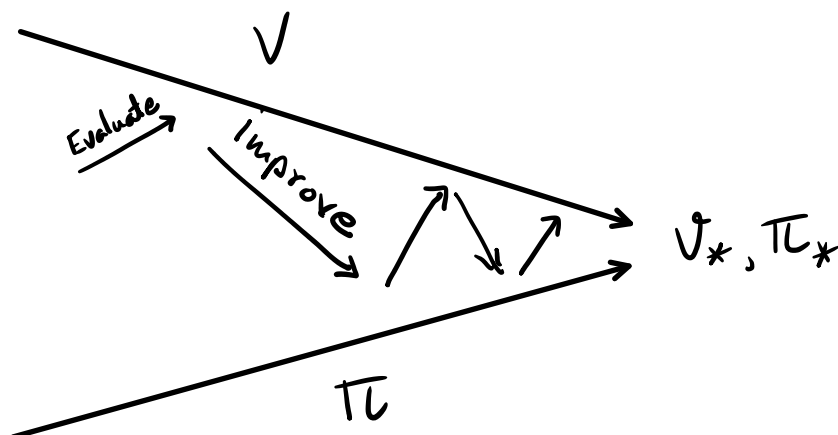
Coding exercise: Program this

# Value Iteration

30 March 2023

## Value **Iteration**

---
**Algorithm 2** Value Iteration

---
1:  **Input:** $\theta > 0$ tolerance parameter, $\gamma$ discount factor
2:  Initialize $V(s)$ arbitrarily, with $V(terminal) = 0$
3:  **repeat**
4:      $\Delta \leftarrow 0$
5:      **for** $s \in S$ **do**
6:          $v \leftarrow V(s)$
7:          $V(s) \leftarrow \max_{a \in A(s)} \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$
8:          $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
9:      **end for**
10: **until** $\Delta > \theta$
11: **Output:** $\pi$: greedy policy w.r.t. $V(s)$

---

A process that alternatively improves the policy



**Algorithm 2** Policy Iteration

1: **Input:** $\theta > 0$ tolerance parameter, $\gamma$ discount factor
2: Initialize $V(s)$ and $\pi(a|s)$ arbitrarily
3: **while** policy-stable $= false$ **do**
4:
5:     Policy Evaluation:
6:     **while** $\Delta > \theta$ **do**
7:         $\Delta \leftarrow 0$
8:         **for** $s \in S$ **do**
9:             $v \leftarrow V(s)$
10:            $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$
11:            $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
12:        **end for**
13:    **end while**
14:
15:    Policy Improvement:
16:    policy-stable $= true$
17:    **for** $s \in S$ **do**
18:        old-action $\leftarrow \pi(s)$
19:        $\pi(s) \leftarrow \arg\max_{a \in A(s)} \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$
20:        **if** old-action $\neq \pi(s)$ **then**
21:            policy-stable $\leftarrow false$
22:        **end if**
23:    **end for**
24:
25: **end while**
26: **Output:** Optimal policy $\pi(a|s)$ and state values $V(s)$