

# REINFORCEMENT LEARNING (David Silver, DeepMind) UCL

- How does it compare to other machine learning paradigms
- No supervisor, just reward signal
  - Feedback delayed
  - Time matters (sequential)
  - Agent's actions affect subsequent observations

## • rewards

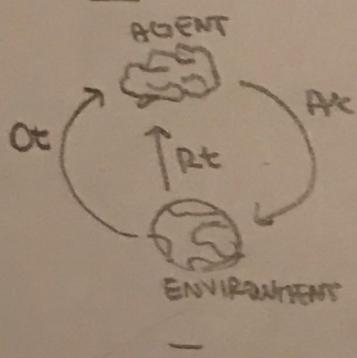
A reward  $R_t$  is a scalar feedback signal indicating how well it is doing at step  $t$ . The goal is to maximise cumulative reward in RL.

## • Agent hypothesis

All goals can be described by the maximisation of expected cumulative reward

## • sequential decision making

The goal is then to select actions to maximise the total future reward. Actions can have consequences which delay the reward. It may be better to sacrifice immediate reward to gain more long term reward.



At each step  $t$  the agent

- executes action  $A_t$
- receives observation  $O_t$
- receives scalar reward  $R_t$

The environment

- Receives action  $A_t$
- Emits observation  $O_{t+1}$
- Emits scalar reward  $R_{t+1}$

↑ increases at each env. step

The history is the sequence of observations, actions, rewards up to step  $t$   $H_t = O_1, R_1, A_1 \dots O_{t-2}, R_t, A_t$

What happens next depends on the history. We define a state  $s_t$   $S_t = f(H_t)$  MDP shows the policy  $\pi$  at the agent state  $s_t^a$

environment state  $s_t^e$  is the environment's private representation used to pick the next reward and observation. This is usually not visible to the agent

agent state  $s_t^a$  is the agent's internal representation used to pick the next action

A Motion state contains all useful info from the history

A state is motion iff  $P[S_{t+1} | S_t] = P[S_{t+1} | S_1 \dots S_t]$

An environment is fully observable when  $O_t = s_t^e = s_t^a$  the agent directly observes the environment state

It is partially observable when  $s_t^a \neq s_t^e$  the agent indirectly observes the environment state

↑ POMDP

## Inside on RL agent

### Policies

A policy is the map from state to action

deterministic policy  $a = \pi(s)$

stochastic  $\sim \pi(a|s) = P[A_t = a | S_t = s]$

### Value Function

A value function is a prediction of future rewards

value function  $V^\pi = E_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$   
given policy  $\pi$  discount factor  $\gamma$ , discounts future rewards  $V_m$  that are relevant

### Model

A model predicts what the environment will do next

next state  $P_{s,s'}^\pi = P[S_{t+1} = s' | S_t = s, A_t = a]$   
given  $s$   
action  $a$   
and state  $s$   $R_s^\pi = E[R_{t+1} | S_t = s, A_t = a]$   
next reward given action  $a$  and state  $s$

### Types of RL agents

Value based  
- just value  
function used

Policy based  
- just policy

Agent Critic  
- uses both value function  
and policy

### Model Free

- Policy or value function  
but no model

Model based  
- Policy and/or value function,  
and model

### Problems in sequential decision making

#### - Reinforcement learning

the environment is initially unknown and the agent improves its policy by interacting with the environment

#### - Planning

A model of the environment is known and the agent performs computations on its model (without doing any action) in order to improve its policy

### In RL

The agent should discover a good policy by exploring without losing too much reward along the way

Exploration, exploit prior info to maximise reward  
Exploration, find more information about the environment (both important)

Prediction, evaluate the future given a policy  
Control, optimise the future by finding the best policy

## Markov decision processes

they formally describe an environment, where the environment is fully observable. Almost all RL problems can be formalized as MDPs

[Talked about Markov property already]

For a Markov state  $s$  and a successor state  $s'$ , the state transition probability is defined as

$$P_{ss'} = P[S_{t+1} = s' | S_t = s]$$

The state transition matrix defines transition probabilities for every transition from any state to another

$$P = \begin{matrix} & \text{to} \\ \text{from} & \left[ \begin{matrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{matrix} \right] \end{matrix} \quad \begin{matrix} \text{each row} \\ \text{sums to 1} \end{matrix}$$

A Markov Process is a sequence of random states  $s_1, s_2, \dots$  with the Markov Property. (also called Markov chain)

It is defined as a tuple  $\langle S, P \rangle$

Finite set of states Transition matrix

A Markov Reward Process is a Markov chain with values.

Defined as a tuple  $\langle S, P, R, \gamma \rangle$

$$\text{Reward } R_s = E[R_{t+1} | S_t = s]$$

The return is the total discounted reward from time step  $t$

Given a  $\rightarrow$  Markov Chain  $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$   $\rightarrow$  avoids infinite cycles and uncomputability about the future because immediate rewards are better than delayed rewards (e.g. money)

The state value function  $v(s)$  of an MRP is the expected return starting from state  $s$

$$v(s) = E[G_t | S_t = s]$$

Bellman Equation for MRPs - decompose value function into  $R_{t+2}$  and  $v(S_{t+2})$

$$v(s) = E[G_t | S_t = s]$$

$$\text{immediate reward } = E[R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots | S_t = s]$$

$$= E[R_{t+2} + \gamma G_{t+2} | S_t = s]$$

$$= E[R_{t+2} + \gamma v(S_{t+2}) | S_t = s]$$

$$v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$$

A Markov decision process is a Markov reward process with decisions (all states are MDPs here as well)

It is a tuple  $\langle S, A, P, R, \gamma \rangle$  where  $A$  is a set of finite actions

A policy  $\pi$  is a distribution over actions given state

$$\pi(a|s) = P[A_t = a | S_t = s]$$

It defines the agent's actions and depends on the current state

The state value function  $v_\pi(s)$  is the expected return starting from state  $s$  and following policy  $\pi$   $v_\pi(s) = E[G_t | S_t = s]$

The action value function  $q_\pi(s, a)$  is the expected return starting from state  $s$ , taking action  $a$ , and the following policy  $\pi$

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

for  
MDPs

## Bellman expectation equation

$V_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s]$  For state value function  
 $q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, a_{t+1}) | s_t = s, a_t = a]$  For action value function  
 how good it is to be in state  $s$  and take action  $a$

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a) \quad q_{\pi}(s) = R_s^{\pi} + \gamma \sum_{s' \in S} P_{ss'}^{\pi} V_{\pi}(s')$$

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) (R_s^{\pi} + \gamma \sum_{s' \in S} P_{ss'}^{\pi} V_{\pi}(s'))$$

$$q_{\pi}(s, a) = R_s^{\pi} + \gamma \sum_{s' \in S} P_{ss'}^{\pi} \sum_{a' \in A} \pi(a'|s) q_{\pi}(s', a')$$

### Optimal value function

- The optimal state-value function  $v_*(s)$  is the maximum value function over all policies.  $v_*(s) = \max_{\pi} V_{\pi}(s)$
  - The optimal action-value function  $q_*(s, a)$  is the maximum action value function over all policies.  $q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$
- The optimal value function specifies the best possible performance in the MDP. The MDP is "solved".

For any MDP, there exists a policy  $\pi_*$  that is better than or equal to all other policies  $\pi_* \geq \pi, \forall \pi$ , given the partial ordering  $\pi \geq \pi'$  if  $V_{\pi}(s) \geq V_{\pi'}(s) \forall s$

$$\pi \geq (\pi(s)) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a \in A} q_{\pi}(s, a) \\ 0 & \text{otherwise} \end{cases}$$

There's always an optimal policy for any MDP

Bellman optimality equation (often known as Bellman equation)

$$v_*(s) = \max_a q_*(s, a)$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s')$$

$$v_*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s')$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a'} q_*(s', a')$$

There are no closed form solutions and we usually use iterative method to solve the equations.

### POMDPs

A partially observable Markov decision process is an MDP with hidden states. It is a hidden Markov model with actions.

A POMDP is a tuple  $\langle S, A, O, P, R, z, \gamma \rangle$

where  $O$  is a finite set of observations and  $z$  is an observation function

A belief state  $b(h)$  is a probability distribution over states conditioned on the history  $h$

$$b(h) = (P[S_t = s_1 | H_t = h], \dots, P[S_t = s_n | H_t = h])$$

The history  $H_t$  satisfies the Markov property

The belief state  $b(H_t)$  satisfies the Markov property

## Planning by Dynamic Programming

what is DP?      Dynamic sequential or temporal component of the problem  
Programming optimising a program i.e policy

It's a method for solving complex problems by breaking them down into subproblems. Then solve the subproblems and combine the solutions to find the solution to the original problem.

Problems must have two properties:

- Optimal substructure. The principle of optimality applies (the optimal solution to the subproblem tells you how to get the optimal solution to the overall problem), and the optimal solution can be decomposed into subproblems.
- Overlapping subproblems. The subproblems repeat many times and solutions can be cached and reused.

For MDP the Bellman equation gives recursive decomposition and the value function stores and reuses solution.

Both properties are satisfied

(transition and rewards)

DP assumes the full knowledge of the MDP and it's used for planning. Two special cases for planning: prediction & control

- Prediction. We're given an MDP  $\langle S, A, P, R, \gamma \rangle$  and  $\pi$  (or a  $\pi^*$ ) and we get  $V\pi$  (how good  $\pi$  is)
- control. We get the MDP  $\langle S, A, P, R \rangle$  and we get  $V\pi$  and  $\pi$ .

## Iter. Policy evaluation

Evaluate a given policy  $\pi$  using the Bellman expectation equation

$\rightarrow V_1 \rightarrow V_2 \rightarrow \dots V_T \leftarrow$  after many iterations we get the true value function  
arbitrary initial one bootstraps value function using Bellman eq.

We'll use synchronous backups

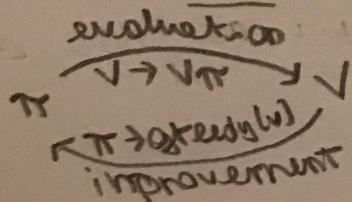
- For each iteration  $t+1$ , for all states  $s \in S$ , update  $V_{t+1}(s)$  from  $V_t(s)$  successor state from  $s$
- $$V_{t+1}(s) = \sum_a \pi(a|s) \left( p_{ss}^a + \gamma \sum_{s' \in S} p_{ss'}^a V_t(s') \right)$$

## How to improve a policy (policy iteration)

Given a policy  $\pi$ , evaluate it <sup>(1)</sup> (as seen before) and act greedily with respect to  $V\pi$  to improve it <sup>(2)</sup>

$$(1) V_\pi(s) = E[R_{t+1} + \gamma R_{t+2} + \dots | s_t = s]$$

$$(2) \pi' = \text{greedy}(V\pi) \leftarrow \text{choose the action with the lowest } V\pi(s) \text{ for all } s$$



Keep iterating until  $V\pi$  stops improving  
at that point you converged to  
 $\pi^*$  and  $V^*$  (optimal  $\pi$  and  $V$ )

## Policy improvement

Consider a deterministic policy  $\pi = \pi(s)$ , we can improve the policy by acting greedily  $\pi^*(s) = \arg\max_{a \in A} q_\pi(s, a)$

$q_{\pi^*} \geq q_\pi$  every time we act greedily

If improvement stops  $q_\pi(s, \pi^*(s)) = \max_{a \in A} q_\pi(s, a) = q_{\pi^*}(s, \pi(s)) = v_{\pi^*}(s)$   
 then the Bellman optimality equation  $v_{\pi^*}(s) = \max_{a \in A} q_{\pi^*}(s, a)$  has been satisfied (that means  $v_{\pi^*}$  is optimal)  
 therefore  $v_{\pi^*}(s) = v^*(s)$  for all  $s \in S$

Instead of converging, we can stop early

- While condition for stopping is met
- Stop after  $k$  iterations of policy evaluation (does not converge with  $k=1$ )

## Value iteration

### Principle of Optimality

Any optimal policy has an optimal first action  $a^*$  followed by an optimal policy for successor state  $s'$ .

Theorem (principle of optimality)

A policy  $\pi(s|s')$  achieves optimal value from  $s$ ,  $v_\pi(s) = v^*(s)$  iff for any state  $s'$  reachable from  $s$ ,  $\pi$  achieves the optimal value from state  $s'$ ,  $v_\pi(s') = v^*(s')$

### Deterministic value iteration

If we know the solution to subproblem  $v^*(s')$ , then solution  $v^*(s)$  can be found by one-step lookahead

$$v^*(s) \leftarrow \max_{a \in A} R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v^*(s')$$

The idea of value iteration is to apply these updates iteratively  
 intuition: start with final rewards and work backwards  
 still works with loops stochastic MDPs

## Value iteration

Problem: find optimal policy  $\pi$

Solution: iterative application of Bellman optimality backup

$$v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v^*$$

Using synchronous backups at each iteration  $k+1$  for all states  $s \in S$ , update  $v_{k+1}(s)$  from  $v_k(s')$

Unlike policy iteration, there's no explicit policy as intermediate value functions may not correspond to any policy

$$v_{k+1} \leftarrow \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_k(s') \right)$$

$$v_{k+1} = \max_{a \in A} R_s^a + \gamma P^a v_k$$

## Synchronous DP algorithms

Prediction	Bellman expectation equation	Iterative policy evaluation
Control	Bellman expectation equation + Greedy policy improvement	Policy Iteration
Control	Bellman optimality equation	Value Iteration

- Algorithms are based on state-value function  $V_P(s)$  or  $V_t(s)$
- Complexity  $O(mn^2)$  per iteration.  $m$  actions  $n$  states
- Could also apply to action value function  $q_P(s, a)$  or  $q_t(s, a)$
- Complexity  $O(m^2 n^2)$  per iteration

## Asynchronous DP

DP methods described ~~use~~ asynchronous backups (ie all states are backed up in parallel)

Asynchronous DP backs up states individually, in any order.

For each state, apply the appropriate backup.  
Can significantly reduce computation and it is guaranteed to converge if all states continue to be selected.

### In-place DP

Synchronous value iteration stores two copies of value function

$$\text{for all } s \in S \quad V_{\text{new}}(s) = \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\text{old}}(s') \right)$$

In-place value iteration only stores one copy of value function

$$\text{for all } s \in S \quad V(s) = \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V(s') \right)$$

### Prioritised sweeping

Use magnitude of Bellman error to guide state selection  
eg.  $\left| \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V(s') \right) - V(s) \right|$

- Backup the state with the largest remaining Bellman error
- Update Bellman error of affected states after each backup
- Requires knowledge of reverse dynamics (preferably stored)
- Can be implemented efficiently by maintaining a priority queue

### Real-time programming

- Use only states that are relevant to agent
- Use agent's experience to guide the selection of states
- After each time-step  $s_t, a_t, r_{t+1}$
- Backup the state  $s_t$

$$V(s_t) = \max_{a \in A} \left( R_{st}^a + \gamma \sum_{s' \in S} P_{st}^a V(s') \right)$$