

# AI - Sequential Decision problems

①

3				+1
2				-1
1	Start			
	1	2	3	4

} sequential decision problem of a  $4 \times 3$  environment.

- An agent has to start from Start & reach +1/-1.
- Probability of moving up is 0.8 & that of moving left/right is 0.1
- On reaching goal, the reward is +1/-1 & terminate all other states have -0.04 as reward.
- Available actions are Up, Down, Left & Right
- The env. is assumed to be fully observable & also not deterministic.
- A specification of the outcome probabilities is called. For each action in each possible state is called a transition model
- $T(s, a, s')$  denote probability of reaching from states  $s$  to  $s'$  with action 'a'. Can be considered as a 3-D table.
- The probability is assumed to be Markovian - the probability depends only the immediately previous state only & not on the earlier states.
- The utility  $U^n$  with depend on sequence of states - an environmental history - rather than a single state.
- If the utility  $U^n$  is calculated as sum of the rewards received, then to reach +1 state, the utility will be 0.6 . after 10 steps.
- The specification of a sequential decision problem for a fully observable environment with a Markovian transitional model & additive rewards is called a Markovian decision process (MDP). Its components are

Initial State:  $S_0$

Transition Model:  $T(s, a, s')$

Reward Function:  $R(s) / R(s, a, s')$

## Solution / Policy

(2)

- A fixed action sequence can't be sol<sup>n</sup> as an agent might end in a non-goal state.
- Sol<sup>n</sup> must specify what is to be done for every state that agent might reach.
- A sol<sup>n</sup> of this kind is called a policy  $\pi$ .

The Policy of state  $s$  is denoted as  $\pi(s)$

- An agent with complete policy always knows what to do next.
  - The stochastic nature of the env. will lead to diff. environmental history each time policy is run.
  - The quality of a policy is  $\therefore$  measured by expected utility of the possible env. histories generated by that policy.
  - An optimal policy is a policy that generates highest expected utility. denoted as  $\pi^*$
  - Given  $\pi^*$  an agent consults its current percept to obtain the current state, & then executing action  $\pi^*(s)$ .
  - A policy that is guaranteed to reach terminal state is called proper policy; can use  $\lambda=1$  (additive rewards)
- Optimality in sequence decision problem

- To find possible choices for performance measures - choices for the utility  $f^h$  on env. histories, represented as  $U_h([s_0, s_1, \dots, s_n])$ .
- Finite horizon - There is a fixed time  $N$  after which nothing matters.

$$\text{Thus } U_h([s_0, s_1, \dots, s_{N+k}]) = U_h([s_0, s_1, \dots, s_N]) \text{ for } k > 0.$$

- In finite horizon, the optimal action in a given state could change over time, it is non-stationary. In stationary optimal policy, optimal action depends upon current state.



- Policies for infinite horizon are simpler.  
Compare to finite horizon.
- Infinite horizon necessarily means that there are no fixed deadlines & not infinite state sequences.
- The utility of state sequence can be calculated using multiattribute utility theory where each state  $s_i$  is viewed as an attribute of  $[s_0, s_1, \dots, s_n]$
- 2 ways to assign utilities to sequence
  - + Additive Rewards:  $U_h([s_0, s_1, \dots, s_n]) = R(s_0) + R(s_1) + \dots$
  - + Discounted Rewards:  $U_h([s_0, s_1, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$

where  $\gamma$  is a discount factor betw<sup>n</sup> 0 & 1. It desc<sup>ri</sup>bes the preference of the agent for current rewards over future rewards - 0 being used for no preference for rewards in distant future. A discount factor of  $\gamma$  is equivalent to interest rate of  $(1/\gamma) - 1$

- A value of policy is the expected sum of discounted rewards obtained where the expectation is taken over all possible state sequences that could occur, given the policy is executed. An optimal policy  $\pi^*$  is

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right]$$

### Utilities of State

- utility of states is defined in terms of the utility of state sequences.
- Roughly, the utility of a state is the utility of the state sequence that might follow it.

- If  $s_t$  be the state an agent reaches after executing  $\pi$  for 't' steps, then we have

$$U_{\pi}^{\pi}(s) = E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, s_0=s \right]$$

- The true utility of the state 's'  $\therefore$  is  $(U_{\pi^*}(s) =)$  expected sum of discounted rewards on executing the optimal policy.

- The utility  $U(s)$  allows to select action using Maximum Expected Utility - choose the action that maximises the expected utility of sub-state  $\pi^*(s) = \operatorname{argmax}_a \left( \sum_{s'} T(s, a, s') U(s') \right)$ .

- The utility of state is the immediate reward for that state plus the expected discounted utility of the next state, assuming the agent chooses the optimal action.

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s') \quad \text{Bellman's Equation}$$

- Value Iteration algo -  
+ Iterate to find optimal utility

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U_i(s')$$

Policy Iteration:

- Alternates betw 2 steps

+ Policy Evaluation: given  $\pi_i$ , calculate  $U_i = U^{\pi_i}$

+ Policy Improvement: Calculate MVEU policy  $\pi_{i+1}$ , using one step look ahead.

$$U_i(s) = R(s) + \gamma \sum_{s'} P(s' | s, \pi_i(s)) U_i(s')$$

+ Modified policy Iteration: Not all / Simplified step-s

$$U_{i+1} \leftarrow R(s) + \gamma \sum_{s'} P(s' | s, \pi_i(s)) U_i(s')$$

## Policy Iteration

- The policy iteration algorithm iterates bet<sup>w</sup> 2 steps
  1. Policy Evaluation: Given a policy  $\pi_i$ , calculate  $U_i$  for  $U_i$  the utility of each state if  $\pi_i$  ~~had~~ <sup>were</sup> to be executed
  2. Policy Improvement: Calculate a new maximum expected utility (MEU) policy  $\pi_{i+1}$ , using one step look ahead based on  $U_i$ .

The algo. terminates when policy improvement step yields no improvement in the utilities.

- Now, ~~then~~  $U_i$  is a fixed point in Bellman's eqn &  $\pi_i$  ~~is~~ must be the optimal policy.
- Implementing PI routine is much better than solving Bellman's Eqn.
- An  $i^{\text{th}}$  iteration specifies action  $\pi_i(s)$  for  $\pi_i$  in 's'.
- Leads to simplified version of Bellman's eqn.

$$U_i(s) = R(s) + \gamma \sum_{s'} T(s, \pi_i(s), s') U_i(s')$$

The eqn is also linear. For 'n' states there will be 'n' linear eqns which can be solved in  $O(n^3)$  by standard linear algebra methods.

- For small state space, policy evaluation using exact sol<sup>n</sup> methods is often the most efficient approach.
- For large state space, the no. of iterations can be reduced to give reasonably good approximation of the utilities. The simplified Bellman's eqn becomes.

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} T(s, \pi_i(s), s') U_i(s')$$

This is called modified policy iteration & much more efficient than standard policy iteration or value iteration algos.

- If instead of updating ~~of~~ all states, only certain states' utility or policy could be updated. This is asynchronous policy updation. Allows to design better heuristic algorithms.