

REINFORCEMENT LEARNING

Types of Learning

Supervised learning

- Given: training data + desired outputs (labels)

Unsupervised learning

- Given: training data (without labels)

Semi-supervised learning

- Given: training data + some labels

Reinforcement learning

- Given: observations and occasional rewards as the agent performs sequential actions in an environment

How Reinforcement Learning is Different

The aim of RL is to make sequential decisions in an environment:

- Driving a car
- Cooking
- Playing a videogame
- Controlling a power plant
- Treating a trauma patient
- Displaying online ads to a user

How to learn to do these things?

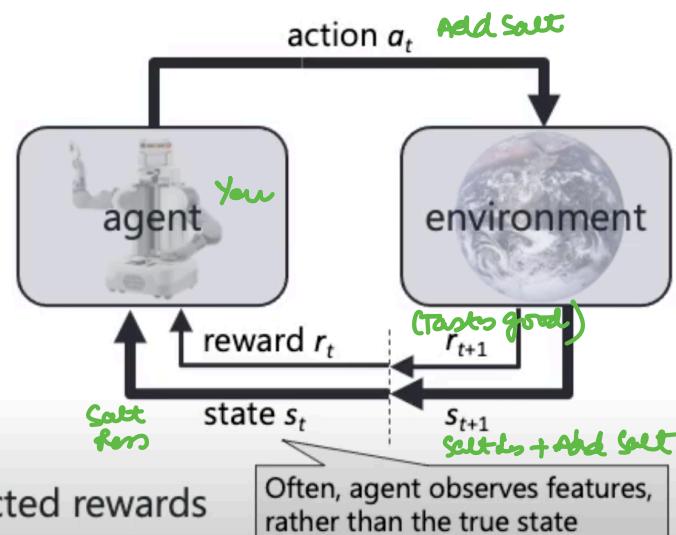
- RL assumes only occasional feedback, such as a tasty meal, or a car crash, or points in a video game.
- RL aims to use this feedback to learn through trial and error, as cleverly as possible. → to minimize the no. of trials.

no dataset in advance.
You receive your data as agent performs sequential actions in environment.
Reward /
actions good action bad

Reinforcement Learning

Main idea:

- Agent receives observations (state $s_t \in S$) and feedback (reward r_t) from the world
- Agent takes action $a_t \in A$
- Agent receives updated state s_{t+1} and reward r_{t+1}
- Agent's goal is to maximize expected rewards



Goal of RL is to learn a policy $\pi(s): S \rightarrow A$ for acting in the environment

How Reinforcement Learning is Different

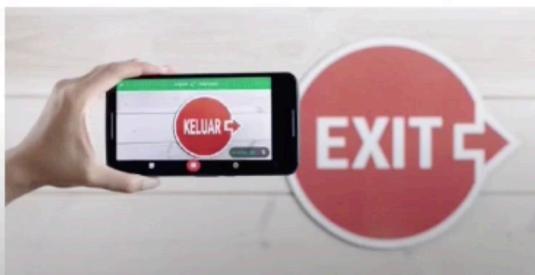
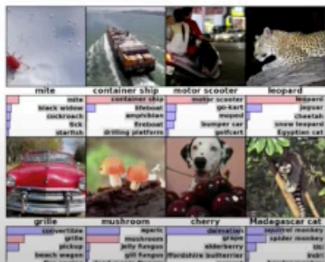
Characteristics of RL Problems:

- No supervision, only (occasional) rewards as feedback
- Sequential decision making \Rightarrow Data is generated as sequences (not i.i.d)
- Training data is generated by the learner's own behavior

→ independent & identical distribution.

When do we not need to worry about sequential decision making?

When your system is making a single isolated decision that does not affect future decision, e.g. classification, regression

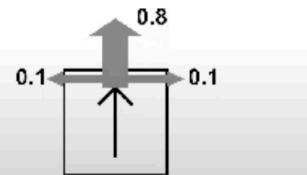
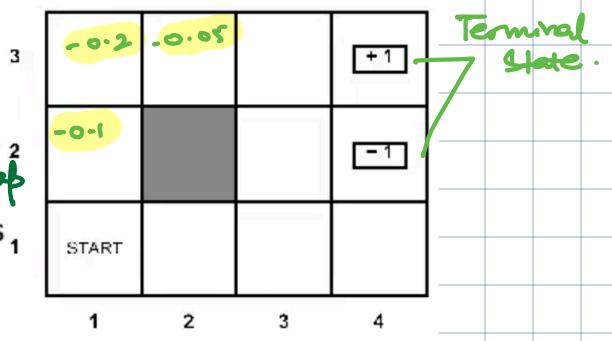


Reinforcement Learning Examples:
Robotics, Autonomous Driving, Mario

MARKOV DECISION PROCESS

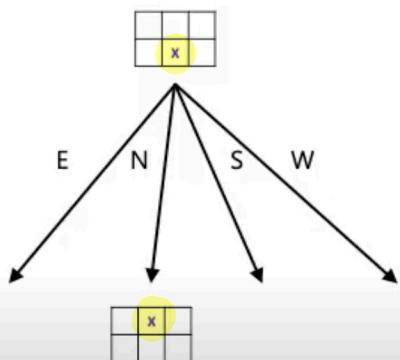
Grid World

- Agent operates in a grid with solid and open cells
- Each timestep, the agent receives a small negative "living" reward \rightarrow penalty sum: complete task asap
- There are two bigger magnitude rewards at terminal states, that end an episode
- The agent can move North, East, South, West
 - The agent remains where it is if it tries to move into a solid cell or outside the world
 - The chosen action succeeds 80% of the time (for an open cell)
 - 10% of the time, the agent ends up 90° off
 - Another 10% of the time, the agent ends up -90° off
 - For example, an agent surrounded by open cells and moving North will end up in the northern cell 80% of the time, in the eastern cell 10% of the time, and in the western cell 10% of the time
- Goal: maximize sum of rewards (i.e., maximize expected utility)

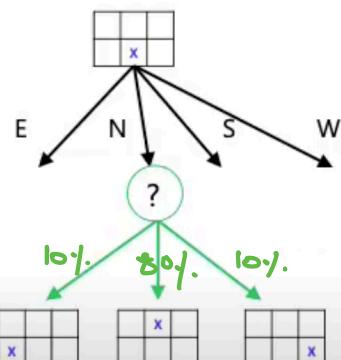


State Transition Diagram

Deterministic Grid World



(Probability) Stochastic Grid World



Markov Decision Process

An MDP (S, A, P, R) is defined by:

- Set of states $s \in S$ $S = \{s_0, s_1, s_2\}$
- Set of actions $a \in A$ $A = \{a_0, a_1\}$
- Transition function $P(s' | s, a)$
 - Probability $P(s' | s, a)$ that a from s leads to s'
 - Also called the "dynamics model" or just the "model"
- Reward function $R(s, a, s')$ or $R(s)$

we typically don't know $P(\cdot)$ and $R(\cdot)$

We don't know state transition diagrams in advance, we have to learn it through trial and error. In the process of learning through trial and error we will encounter samples from the transition diagram.

why the name Markov?

The Markov Property: given the present, the future and the past are independent

- i.e., Everything you need to know about the past is included in the present state

- For MDPs, the Markov property means that:

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1} = a_{t-1}, \dots, S_0 = s_0) = P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

Independent of past states and actions

Markov Assumption?

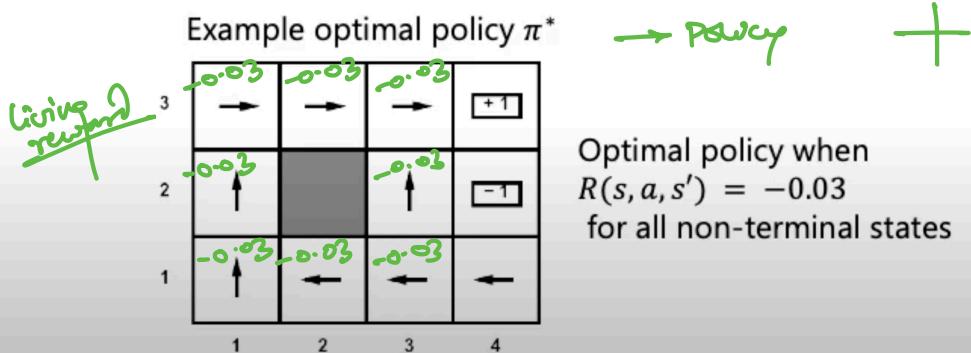
Solving the MDP

To solve an MDP (S, A, P, R) means to find the optimal policy $\pi^*(s): S \rightarrow A$

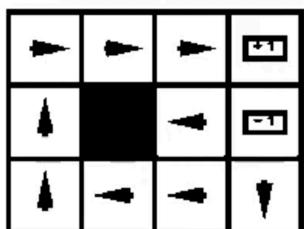
- “Optimal” \Rightarrow Following π^* maximizes total reward/utility (on average)

In RL, P and R are assumed unknown. But let's first assume that we do know the whole MDP

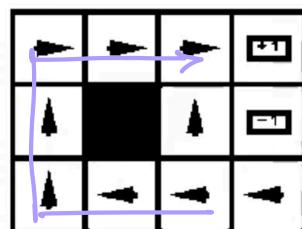
Policy is a mapping from state to Action



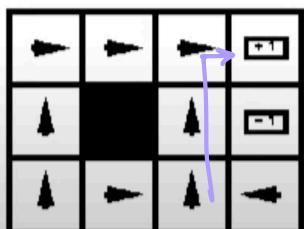
Example Optimal Policies



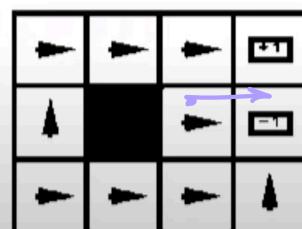
$$R(s) = -0.01$$



$$R(s) = -0.03$$

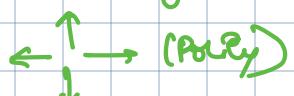


$$R(s) = -0.4$$



$$R(s) = -2.0$$

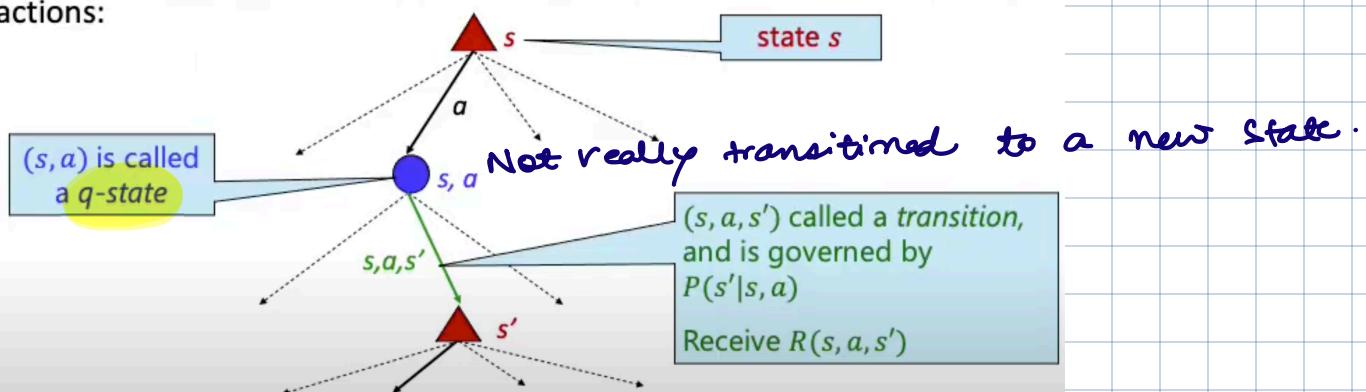
If your reward changes then policy will also change



Optimal policy is sensitive to reward fxn.

MDP Search Trees

- Each MDP state has an associated tree of future outcomes from various actions:

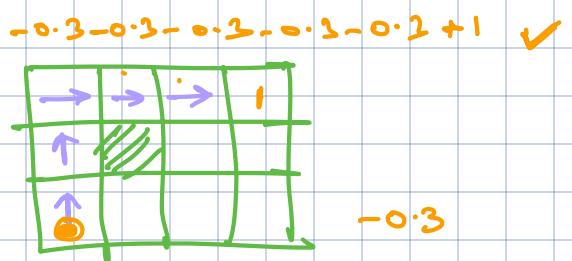


Defining Utility

- At each step, agent chooses an action to maximize expected rewards

So, we must consider utility over sequences of rewards
 $[r_t, r_{t+1}, r_{t+2}, \dots, r_\infty]$ → upper bound

Problem: infinite sequences yield infinite rewards



$$\text{Utility} / \text{Return} = \gamma^0(-0.3) + \gamma^1(-0.3) + \gamma^2(0.3) + \gamma^3(-0.3) + \gamma^4(-0.3) + \dots$$

$\gamma \rightarrow \text{discounted reward}$

Discounted Rewards

Idea: (uncertain) future rewards are worth exponentially less than current rewards

Future rewards are discounted by $0 < \gamma < 1$:

$$\text{Utility } U([r_1, \dots, r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_{t+1} \rightarrow \gamma^0 r_1 + \gamma^1 r_2 + \gamma^2 r_3 + \dots$$

$\gamma = 1$ would correspond to sum of all rewards (called additive utility)

γ controls the horizon of focus; smaller γ means a shorter horizon – more of a focus on the present

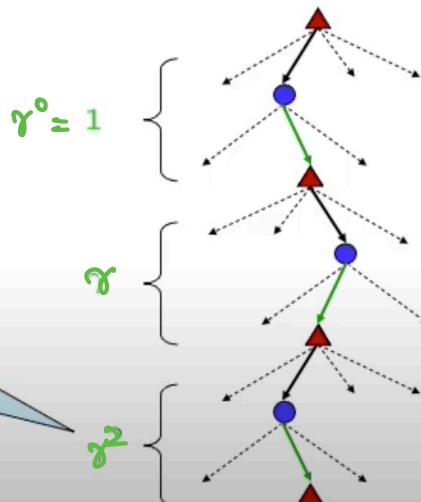
Discounted Rewards

Idea: (uncertain) future rewards are worth exponentially less than current rewards

Future rewards are discounted by $0 < \gamma < 1$:

$$U([r_1, \dots, r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_{t+1}$$

Future rewards matter less to the decision than more recent rewards

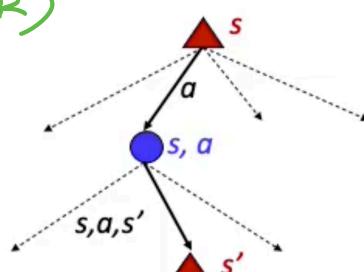


Keep diluting the contribution of future transitions

Recap: Defining MDPs

- Markov decision processes: (S, A, P, R)

- States S (and start state s_0)
- Actions A
- Transitions $P(s'|s, a)$
- Rewards $R(s, a, s')$ (and discount γ)



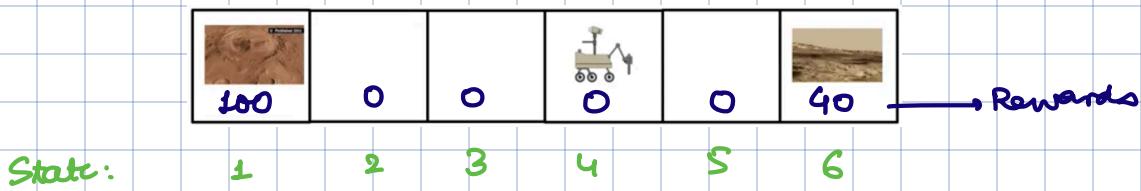
- MDP quantities so far:

- Policy = Choice of action for each state $\pi : S \rightarrow A$
- Utility (or return) = sum of discounted rewards

Mars Rover Example



Terminal States = 1, 6
(After reaching terminal state you stop)



Action ← left → right

NO discounted reward

State 4 3 2 1 : 4 5 6 : 4 5 4 3 2 1 :	$0 + 0 + 0 + 100 = 100$ $0 + 0 + 40 = 40$ $0 + 0 + 0 + 0 + 0 = 0$	Rewards 100 40 0
--	---	---------------------------

RETURN (UTILITY) (Discounted Reward)

$$\begin{aligned}
 4 & 3 & 2 & 1 : (0.9)^0 0 + (0.9)^1 0 + (0.9)^2 0 + (0.9)^3 100 \\
 \gamma = 0.9 & \quad \quad \quad = 0 + 0 + 0 + 0.729 \times 100 = 72.9
 \end{aligned}$$

Return: $R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$ (until terminal state)

Rewards

$$\gamma = 0.5$$

$$\begin{aligned}
 4 & 3 & 2 & 1 : 0 + (0.5)^0 0 + (0.5)^1 0 + (0.5)^2 100 \\
 & \quad \quad \quad : 12.5
 \end{aligned}$$

Example of Return

Return depends on Reward
Reward depends on Action

} Return depends on action you take

100	0	0	0	0	40
1	2	3	4	5	6

$$\gamma = 0.5$$

$v(s)$

100	50	25	12.5	6.25	40
1	2	3	4	5	6

$$(0.5)^0 \cdot 100 + (0.5)^1 \cdot (100)$$

$$0 + (0.5)^0 + (0.5)^1 \cdot 100$$

} Policy:
every state
← action

$$0 + (0.5)^0 + (0.5)^1 \cdot 100$$

v_s

100	2.5	5	10	20	40
1	2	3	4	5	6

$$0 + (0.5)^0 + (0.5)^1 \cdot 40$$

} Policy:
every state
→ action

r_s

100	50	25	12.5	20	40
1	2	3	4	5	6

Optimal Policy:
action that would give me maximum return

Policy
state → action
 s a

$$\begin{aligned} \pi(2) &= \leftarrow \\ \pi(3) &= \leftarrow \\ \pi(4) &= \leftarrow \\ \pi(5) &= \rightarrow \end{aligned}$$

} π^*

100	←	←	→	→	40
1	2	3	4	5	6

Policy:

Reach to the nearest terminal state in minimum moves.

State Value function & Action Value function

Solving MDPs: State Value Functions of Policies

Given MDP (S, A, P, R, γ) :

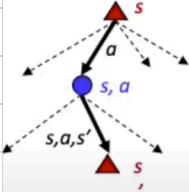
\hookrightarrow fixⁿ of state and of policy that you are following

Value of a state s under policy π :

$V^\pi(s) = \text{expected utility starting in } s \text{ and acting according to } \pi$

$$V^\pi(s) = \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r_{t+1} | S_0 = s \right)$$

Sequence of rewards generated by following π



Optimal value of a state s :

$V^*(s) = \text{expected utility starting in } s \text{ and acting optimally}$

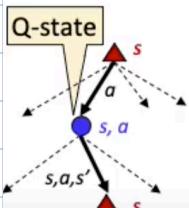
$$V^*(s) = V^{\pi^*}(s) = \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r_{t+1} | S_0 = s \right)$$

Rewards generated by following π^*

$$\underline{E(x)} = \sum x p(x)$$

Solving MDPs: Action Value Functions of Policies

- It is also helpful to define action-value functions



Q-value of taking action a in state s then following policy π :

$Q^\pi(s, a) = \text{expected utility taking } a \text{ in } s \text{ and then following } \pi$

$$Q^\pi(s, a) = \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r_{t+1} | S_0 = s, A_0 = a \right)$$

Optimal Q-value of taking action a in state s : $Q^*(s, a) = Q^{\pi^*}(s, a)$

$$\underline{E(x)} = \sum x p(x)$$

π^* can be greedily determined from Q^* :

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s, a)$$

$\left. \begin{array}{l} \text{Action that} \\ \text{maximizes } Q^*(s, a) \\ \text{will be the action} \\ \text{output by optimal} \\ \text{policy.} \end{array} \right\}$

$$\gamma^0 R_1 + \gamma^1 R_2 + \gamma^2 R_3 + \dots$$

State Value function
Action Value function

$Q(s, a) =$

- ① start from state s
- ② take action a
- ③ Behave optimally after that.

Return that you get?

						Possible Action Rewards
100	50	25	12.5	6.25	3.125	0
100	0	0	0	0	0	40
1	2	3	4	5	6	

$$Q(2, \rightarrow) = 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 \\ = 12.5$$

							all $Q(s, \rightarrow)$ $Q(s, \leftarrow)$
							Action value functions
100	50	25	12.5	6.25	3.125	0	40
1	2	3	4	5	6		

$Q(2, \leftarrow)$ $Q(2, \rightarrow)$
 $Q(3, \leftarrow)$ $Q(3, \rightarrow)$

[for all states
[for all actions $\leftarrow \rightarrow$
calculate return $Q(s, a)$]

100	50	25	12.5	6.25	3.125	0
1	2	3	4	5	6	

* $v^*(2)$ $v^*(3)$ $v^*(4)$ $v^*(5)$

State value function?

More Power Example:
Deterministic

Showing Calculations

		$Q(2, \leftarrow)$		$Q(2, \rightarrow)$			
		100	100	12.5	25	6.25	12.5
		100	0	0	0	0	40
1	2	3	4	5	6		

Terminal State:

$$Q(s, a) = R(s)$$

$$Q(2, \rightarrow) = R(2) + 0.5 \max_{a'} Q(3, a')$$

$$Q(2, \rightarrow) = 0 + 0.5 \max(25, 6.25)$$

$$Q(2, \rightarrow) = 0 + 0.5 \times 25 = 12.5$$

$$Q(4, \leftarrow) = R(4) + 0.5 \max_{a'} Q(3, a')$$

$$= 0 + 0.5 \max(25, 6.25)$$

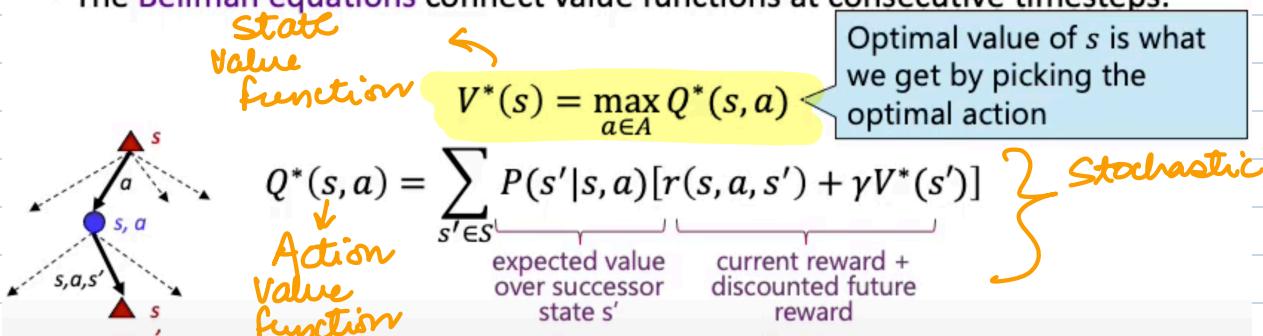
$$= 12.5$$

Bellman Equations

$$E(x) = \sum p(x)$$

Solving MDPs: Bellman Equations

- The Bellman equations connect value functions at consecutive timesteps:



- We can combine these together to get:

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

$Q(s, a) \rightarrow$ start from state s
 act in a
 behave optimally } Return .

s : current state
 a : current action

s' : state you get after taking action a
 a' : action that you take at state s'

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

Bellman Eqn

Deterministic

VALUE AND POLICY ITERATION

Solving MDPs with known P and R :
 Value and Policy Iteration.

VALUE ITERATION

Solving MDPs: Value Iteration

- Bellman Equation gives us a recursive definition of the optimal value:

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

Key Idea: solve iteratively via dynamic programming

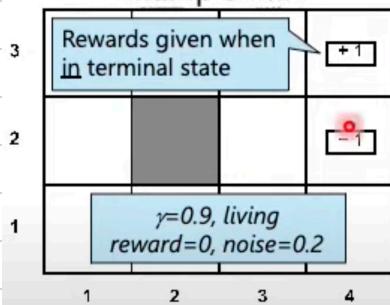
- Start with $V_0(s) = 0$ for all states s
- Iterate the Bellman update until convergence:

$$V_{i+1}(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a)[R(s, a, s') + \gamma V_i(s')]$$

Example: Value Iteration

Bellman Update Rule: $V_{i+1}(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a)[R(s, a, s') + \gamma V_i(s')]$

Example MDP



V ₀			
3	0	0	0
2	0		0
1	0	0	0
	1	2	3

Start with $V_0(s) = 0$

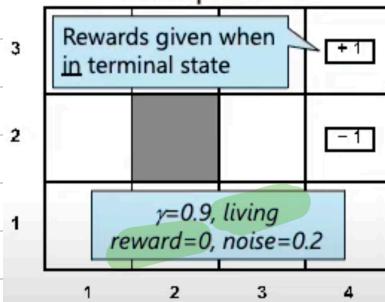
V ₁			
3	0	0	+1
2	0		-1
1	0	0	0
	1	2	3

Example: Value Iteration

Bellman Update Rule: $V_{i+1}(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a)[R(s, a, s') + \gamma V_i(s')]$

backward:

Example MDP



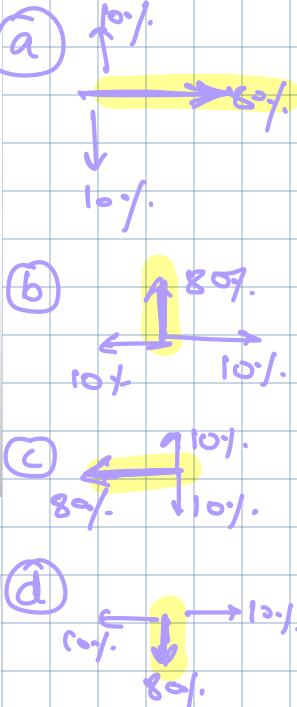
Optimal action will be "right"

V ₁			
3	0	0	+1
2	0		-1
1	0	0	0
	1	2	3

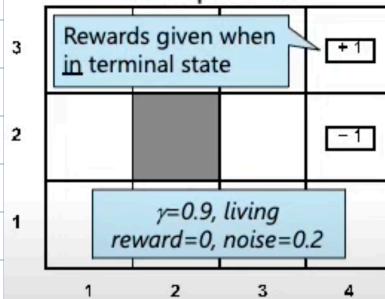
V ₂			
3	0	0	0.72
2	0		-1
1	0	0	0
	1	2	3

Example: Value Iteration

Bellman Update Rule: $V_{i+1}(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a)[R(s, a, s') + \gamma V_i(s')]$



Example MDP



V ₂			
3	0	0	0.72
2	0		-1
1	0	0	0
	1	2	3

V ₃			
3	0	0.52	0.78
2	0		-1
1	0	0	0
	1	2	3

- Information propagates outward from terminal states
- Eventually all states have correct value estimates

POLICY ITERATION

Value f^{π} is associated with some particular policy
and

optimal value f^{*} is associated with optimal policy.

In value iteration, we were dealing with optimal value f^{*} but now we will compute value f^{π} for some random policy.

$\pi: S \rightarrow A$

Policy Evaluation

- How do we calculate the V's for a fixed policy? Compute the value function corresponding to a policy
- Idea: Bellman updates for arbitrary policy:

$$V_0^\pi(s) = 0$$

$$V_{i+1}^\pi(s) \leftarrow \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_i^\pi(s')]$$

(value iteration update rule)

$$V_{i+1}(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V_i(s')]$$

No maximization here, we are just applying the action prescribed by the policy we are trying to evaluate.

Policy Iteration: An Alternative to Value Iteration

Repeat steps until convergence:

1. **Policy evaluation:** keep current policy π fixed, find value function $V^{\pi_k}(\cdot)$

Finding the value f^{π_k} for policy

- Iterate simplified Bellman update until values converge:

$$V_{i+1}^{\pi_k}(s) \leftarrow \sum_{s'} P(s'|s, \pi_k(s)) [R(s, \pi_k(s), s') + \gamma V_i^{\pi_k}(s')]$$

Chooses actions according to π

2. **Policy improvement:** find the best action for $V^{\pi_k}(\cdot)$ via one-step lookahead

$$\pi_{k+1}(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^{\pi_k}(s')]$$

find the best action corresponding to the value f^{π_k} you currently updated.

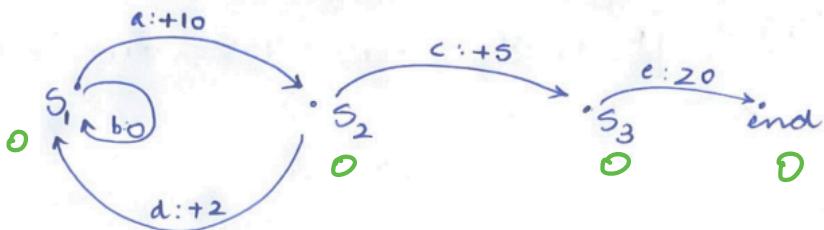
Finding the policy that would be optimal under the value f^{π_k} .

If you do one iteration \rightarrow Value Iteration

A car can operate in three modes: Electric Mode (S1), Hybrid Mode (S2), and Gasoline Mode (S3), each representing a state in the Markov Decision Processes (MDP). The car's system must decide which mode to switch to, aiming to optimize fuel efficiency and battery usage. Actions and their rewards are defined for each state: from S1, switch to S2 (+10) or stay in S1 (0); from S2, switch to S3 (+5) or back to S1 (+2); from S3, end the trip (+20). With a discount factor of 0.5 and initial state values of 0, perform one iteration of value iteration, calculating updated values for each state. [4] [CO1, CO3]

Q5:

a)



$$\gamma = 0.5$$

$$v_0(s_1) = 0$$

$$v_0(s_2) = 0$$

$$v_0(s_3) = 0$$

$$v_{i+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma v_i(s')]$$

from transitions

$$P(s'|s, a) = 1 \text{ for } \forall \text{ actions}$$

$v_1(s_1)$

action a:

$$s' = s_2 \quad 10 + 0.5 \times 0 = 10$$

action b:

$$s' = s_1 \quad 0 + 0.5 \times 0 = 0$$

$$v_1(s_1) = \max(10, 0) = 10 \rightarrow \text{action a}.$$

$v_1(s_2)$

action c:

$$s' = s_3 \quad 5 + 0.5 \times 0 = 5$$

action d:

$$s' = s_1 \quad 2 + 0.5 \times 0 = 2$$

$$v_1(s_2) = \max(5, 2) = 5 \rightarrow \text{action c}$$

$v_1(s_3)$

action e:

$$s' = \text{end} \quad 20 + 0.5 \times 0 = 20$$

Temporal Difference and Q-Learning

We have seen how to solve MDP when P and R are known - by using value and Policy Iteration.

Temporal Differencing (TD)

Policy Evaluation

- Start with $V_0(s) = 0$
- Iterate until convergence: $V_{i+1}^{\pi_k}(s) \leftarrow \sum_{s'} P(s'|s, \pi_k(s)) [R(s, \pi_k(s), s') + \gamma V_i^{\pi_k}(s')]$

How to extend this to when the functions $P(s'|s, a)$ and $R(s, a, s')$ are unknown and only revealed gradually through experience?

Every time you take action a from state s , you get a sample from the unknown $P(s'|s, a)$ and the corresponding reward $R(s, a, s')$

Bellman eqⁿ
written as update

By taking action in environment we get more info about these 2 unknown quantities.

Temporal Differencing (TD)

Policy Evaluation

- Start with $V_0(s) = 0$
- Iterate until convergence: $V_{i+1}^{\pi_k}(s) \leftarrow \sum_{s'} P(s'|s, \pi_k(s)) [R(s, \pi_k(s), s') + \gamma V_i^{\pi_k}(s')]$

Idea: Treat the single sample you get as representative of the distribution, and apply an *incremental* update to reduce the "Bellman error":

Sample of $V(s)$: $sample = R(s, \pi(s), s') + \gamma V_i^{\pi}(s')$

TD update: $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(sample - V^{\pi}(s))$

Doing this for each sample = computing the running average over samples

Expectation of
under P .

[] under P .

Everytime we take some action we get to see one sample.

we treat single sample as a proxy for entire distribution

TD treats every single sample you encounter as representative of distribution and you don't have to wait for large no. of interactions in environment to improve your

Consider an autonomous robot in a 4x4 warehouse grid, tasked with delivering items to a specific area. Each grid cell represents a different warehouse area. 'G' is the delivery area. The robot starts from area 1. The robot consumes energy (-0.2 reward) for each move and gains a significant energy saving (+5 reward) upon successful delivery to 'G'. The robot's route is: 1 → 2 → 3 → 4 → 8 → 12 → 11 → 15 → G. Using Temporal Difference Learning (TD(0)) with a learning rate (α) of 0.1 and a discount factor (γ) of 1, determine the updated efficiency value of area 1 after one delivery, starting from 0.

[4] |CO2|

-0.2	1	2	3	4
5	6	7	8	
9	10	11	12	
13	14	15	G	G

Fig. 1

b)
$$V^{\pi}(s_t) = V^{\pi}(s_t) + \alpha \left[R(s_t, \pi(s_t), s_{t+1}) + \gamma V^{\pi}(s_{t+1}) - V^{\pi}(s_t) \right]$$

$\alpha = 0.1$ $\gamma = 1$ $R = -0.2$

$s_1 \rightarrow s_2: V(1) = 0 + 0.1 (-0.2 + 1 \times 0 - 0) = \underline{-0.02}$

$s_2 \rightarrow s_3: V(2) = 0 + 0.1 (-0.2 + 1 \times 0 - 0) = \underline{-0.02}$

$s_3, s_4, \dots, s_{15} \quad ? \quad V(s) = \underline{-0.02}$

$s_{15} \rightarrow G: V(15) = 0 + 0.1 (5 + 1 \times 0 - 0) = \underline{0.5}$

Q Learning → Action Value fn.

Beyond Policy Evaluation: Learning $Q^*(s, a)$

Recall Bellman equation for optimal Q*:

$$Q^*(s, a) = \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')] \quad \text{Bellman used for Q.}$$

The corresponding Q-value iteration equation (analogous to the state value iteration) would be:

$$Q_{i+1}(s) \leftarrow \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_i(s', a')]$$

Again, this requires access to $P(s'|s, a)$ and $R(s, a, s')$, of which we only have samples from experience.

Apply the TD trick to this?

Q Learning

$$\text{Q-value iteration: } Q_{i+1}(s) \leftarrow \sum_{s' \in S} P(s'|s, a) \left[R(s, a, s') + \gamma \max_{a'} Q_i(s', a') \right]$$

The TD Version: Treat the single sample you get as representative of the distribution, and apply an *incremental* update to reduce the "Bellman error":

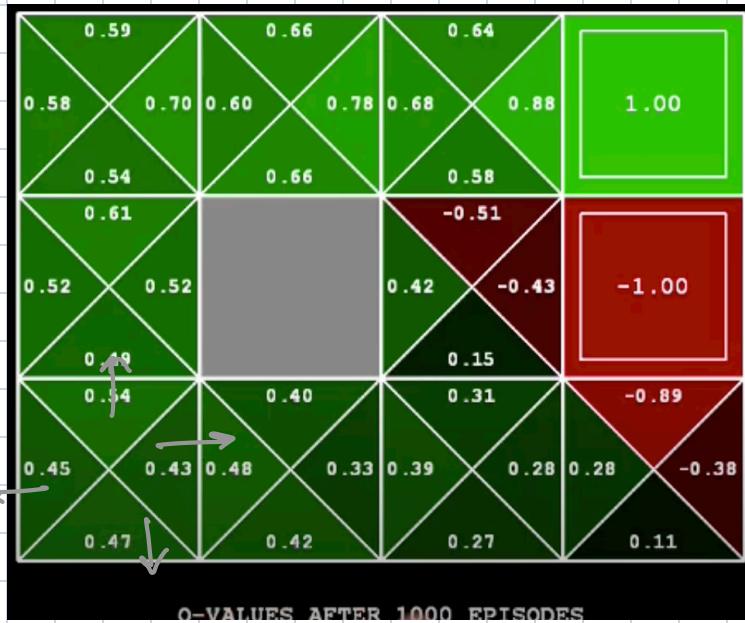
- Execute a single action a from state s and observe s' and R :

$$\text{sample} = R(s, a, s') + \gamma \max_{a'} Q^*(s', a')$$

- So, the incremental TD update is:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \underbrace{\left(\text{sample} + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)}_{\text{Bellman error}}$$

This is called "Q-Learning".



Resources :

→ Slides

https://www.youtube.com/playlist?list=PLYgyoWurxA_8ePNUuTLDtMvzyf-YW7im2

<http://incompleteideas.net/book/ebook/> → Book

<https://www.coursera.org/learn/unsupervised-learning-recommenders-reinforcement-learning>

↓
Machine Courses
Course.