

Reinforcement Learning

Changing features?

- Can an ML agent change features?
 - => Yes, provided it interacts with the environment that produces data
- Can features change every now and then by itself?
 - => When many entities participate and interact with environment (AI agents/Humans/Nature)
- What applications need this scenario?
 - => Where action is needed.
 - => Where environment / features are dynamically changing

Other Investigative Problems:

Sequence of suitable actions to be taken

- Play a game
- Explore a new city
- Take a heritage walk

Any previous label? Experience based? Reward based?

Reinforcement Q Learning

Real Life Applications of RL

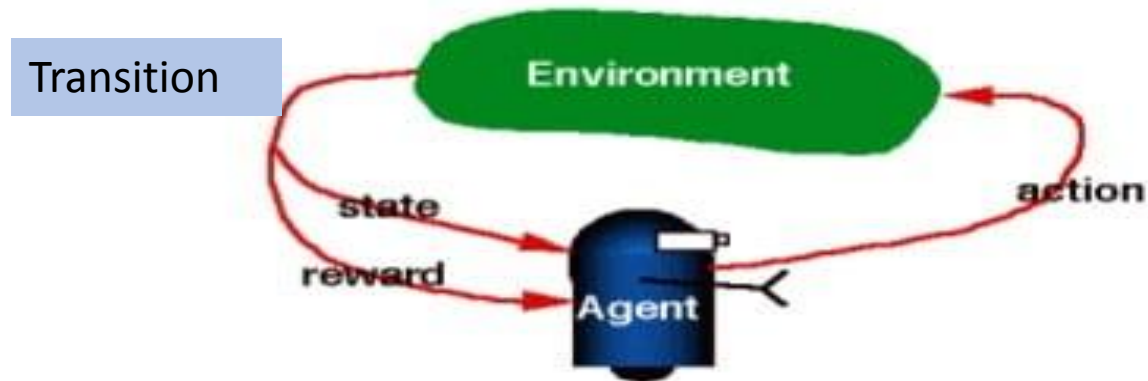
- Autonomous car.....What features does it change?
- Trading does time series or regression analysis do this? What features change?
- Making decisions at real time using multimedia data A swarm of robots operating
 - In manufacturing
 - In dangerous and unknown areas
- NLP – Text summarization, machine translation, Question Answering, chat-box, dialogue generation
 - supervised DL models to predict words
 - RL through rewards when to look for more words / where to look for important words
 - Need to define reward in terms of linguistic quality parameters such as cohesiveness, understandability etc.
- Dynamic Healthcare – diagnosis / treatment / drugs manufacture/ medical policies
- Auction bidding

What is Reinforcement Learning?

- Agent interacts with Environment
- RL \square Learning from interaction with an environment
- Environment state \square feature-vector
- Long term goal \square maximize rewards
- Agent must be able to partially/fully sense the environment state
- Take actions \square change environment state
- Tradeoff between exploration and exploitation

RL is learning from interaction

RL is learning from interaction



S T A R

Markov Decision Process (MDP)

- Set of states S – actually features vector
- Set of actions A
- Policy is mapping from States to possible actions $P: S \rightarrow A$
- State transition probabilities $p(s' \mid s, a)$.
- Immediate Reward Function $R: S \times A \rightarrow$ set of real numbers or discrete
- Discount factor γ in $[0, 1]$ – priority given to future experience
- Finite MDP if both S and A are finite

Lodhi Garden visit –Immediate rewards

	S ₁ Tomb, hungry, tired	S ₂ Garden, hungry, fresh	S ₃ Lake, filled, tired	S ₄ Eatery, filled, fresh	S ₅ Gazebo, hungry, fresh	S ₆ Out, filled, tired
S ₁ Tomb, hungry, tired	-10	50	70	100	50	0
S ₂ Garden, hungry, fresh	20	40	100	100	20	0
S ₃ Lake, filled, tired	0	70	20	0	0	50
S ₄ Eatery, filled, fresh	100	0	0	0	0	0
S ₅ Gazebo, hungry, fresh	0	0	100	100	0	0
S ₆ Out, filled, tired	0	0	60	0	0	100

Q learning basics

- Q is a quality or utility or Value Function -> helps assess decisions
- Maximizes sum of
 - Immediate reward
 - Projected future reward
- Recursive in nature
- Q must be updated with experience
- Initial Q : all zeros

Q learning algorithm

For each (s,a) , initialize $\hat{Q}(s,a) = 0$

Observe Current State

Do Forever:

1. Select action a and execute it
2. Receive Immediate Reward r
3. Observe new state s'
4. Update Utility Table for $\hat{Q}(s,a)$ as:

$$\hat{Q}(s,a) = r + \gamma \times \text{Max}_{a'}(s', a')$$

5. Enter New State

Learning Rate

$$\begin{aligned} Q[\text{state}, \text{action}] = & \\ & Q[\text{state}, \text{action}] + \\ & \text{learning Rate} * \\ & (\text{reward} + \text{Discount Factor} * \max_i \{Q[\text{new_state}, S_i]\} \\ & \quad - Q[\text{state}, \text{action}]) \end{aligned}$$

Initial Utilities – all zeros. Initial State S_3 ,
Discount factor $\gamma = 0.8$

	S_1 Tomb, hungry, tired	S_2 Garden, hungry, fresh	S_3 Lake, filled, tired	S_4 Eatery, filled, fresh	S_5 Gazebo, hungry, fresh	S_6 Out, filled, tired
S_1 Tomb, hungry, tired	0	0	0	0	0	0
S_2 Garden, hungry, fresh	0	0	0	0	0	0
S_3 Lake, filled, tired	0	0	0	0	0	0
S_4 Eatery, filled, fresh	0	0	0	0	0	0
S_5 Gazebo, hungry, fresh	0	0	0	0	0	0
S_6 Out, filled, tired	0	0	0	0	0	0

Learning by interacting – Episode 1

- Let first action randomly be $P(S_3 | S_2) \rightarrow$ Go to Garden.
- Immediate reward = 70
- $Q \rightarrow$ See next state $S_2 =$ Discount factor $\gamma * \text{Max}\{0,0,0,0,0,0\} = 0$
- Thus new $Q(S_3, S_2) = 70$, due to only instant reward

Updated Utility

	S ₁ Tomb, hungry, tired	S ₂ Garden, hungry, fresh	S ₃ Lake, filled, tired	S ₄ Eatery, filled, fresh	S ₅ Gazebo, hungry, fresh	S ₆ Out, filled, tired
S ₁ Tomb, hungry, tired	0	0	0	0	0	0
S ₂ Garden, hungry, fresh	0	0	0	0	0	0
S₃ Lake, filled, tired	0	70	0	0	0	0
S ₄ Eatery, filled, fresh	0	0	0	0	0	0
S ₅ Gazebo, hungry, fresh	0	0	0	0	0	0
S ₆ Out, filled, tired	0	0	0	0	0	0

Episode 2: Initial state = S_6 Outside

- Randomly move to S_3
- Reward = 60.
- $Q = 0.8 * \max\{70, 0\} = 56$
- Updated $Q(S_6, S_3) = 60 + 56 = 116$
- Repeat episodes to reach convergence

Updated Utility

	S ₁ Tomb, hungry, tired	S ₂ Garden, hungry, fresh	S ₃ Lake, filled, tired	S ₄ Eatery, filled, fresh	S ₅ Gazebo, hungry, fresh	S ₆ Out, filled, tired
S ₁ Tomb, hungry, tired	0	0	0	0	0	0
S ₂ Garden, hungry, fresh	0	0	0	0	0	0
S ₃ Lake, filled, tired	0	70	0	0	0	0
S ₄ Eatery, filled, fresh	0	0	0	0	0	0
S ₅ Gazebo, hungry, fresh	0	0	0	0	0	0
S ₆ Out, filled, tired	0	0	116	0	0	0

Deep Reinforcement Learning?

- A deep neural network is used to develop either a policy or a value function
- Deep neural networks require lots of real/simulated interaction with the environment to learn
- Lots of trials/interactions is possible in simulated environments
- We can easily parallelise the trials/interaction in simulated environments
- We cannot do this with robotics (no simulations) because action execution takes time, accidents/failures are expensive and there are safety concerns

Model-free versus Model-based

- A model of the environment allows inferences to be made about how the environment will behave
- Example: Given a state and an action to be taken while in that state, the model could ***predict the next state and the next reward***
- Models are used for planning, which means deciding on a course of action by considering possible future situations before they are experienced
- Model-based methods use models and planning. Think of this as modelling the dynamics $p(s' \mid s, a)$
- Model-free methods learn exclusively from trial-and-error (i.e. no modelling of the environment)