

Reinforcement Learning.

①

Erdem KOYUNCU.

- Studies the question of how agents should take actions in an environment to maximize a cumulative reward.
 - Many applications, including:
 - * Robot control · Individual robots, swarm intelligence
 - * Games / Strategy ·
 - Chess, backgammon, Go, checkers, ...
 - Atari games (Space Invaders, etc ...)
 - Very complex games like Starcraft.
- In all examples RL agents much better than best humans or usually best classical algorithms (non-RL).

* Telecommunication systems:

How should the users of the spectrum use the available resources in an efficient manner?

* ...

ANY scenarios where a model of the environment is known & one can interact with it, but the objective is too complex to be solved analytically or with other approaches.

Mathematical formulation of RL relies on Markov Decision Processes. (MDP)

An mDP can be defined via

(3)

\mathcal{X} : a finite state space

\mathcal{A} : a finite action space.

r : a reward function.

Specifically $r(x, a, y)$ is the reward obtained for a transition from state $x \in \mathcal{X}$ to state $y \in \mathcal{X}$ due to action $a \in \mathcal{A}$.

$P_{a|x}(y)$: Probability of transitioning from state x to y given that one takes action a .

EX: 3×3 grid environment.

4

(3,1)	(3,2)	(3,3)
(2,1)	(2,2)	(2,3)
(1,1)	(2,2)	(1,3)

Person at initial location $(1,1)$, wishes to travel to target $(3,3)$.

State space $X = \{(1,1), \dots, (3,3)\}$

Actions $A = \{U, D, L, R\}$. terminal up, down, left, right state.

- A step fails with prob ϵ , successful with probability $1-\epsilon$
in which case the person cannot move from its current location.
- Goal: Reach $(3,3)$ with the min. number of actions taken.
- Rewards: A reward of -1 for every action. A large negative reward for every action that attempts to step outside. 0 reward after $(3,3)$.
This reward design tries to ensure:
 $E[\text{number of actions to reach } (3,3)]$
- sequence of actions that minimize
 $E[\text{number of actions to reach } (3,3)]$
- sequence of actions that maximize
 $E[\text{cumulative reward}]$.

Value of an initial state

(5)

$x_0 = x$ given controls $\{A_t\}$:

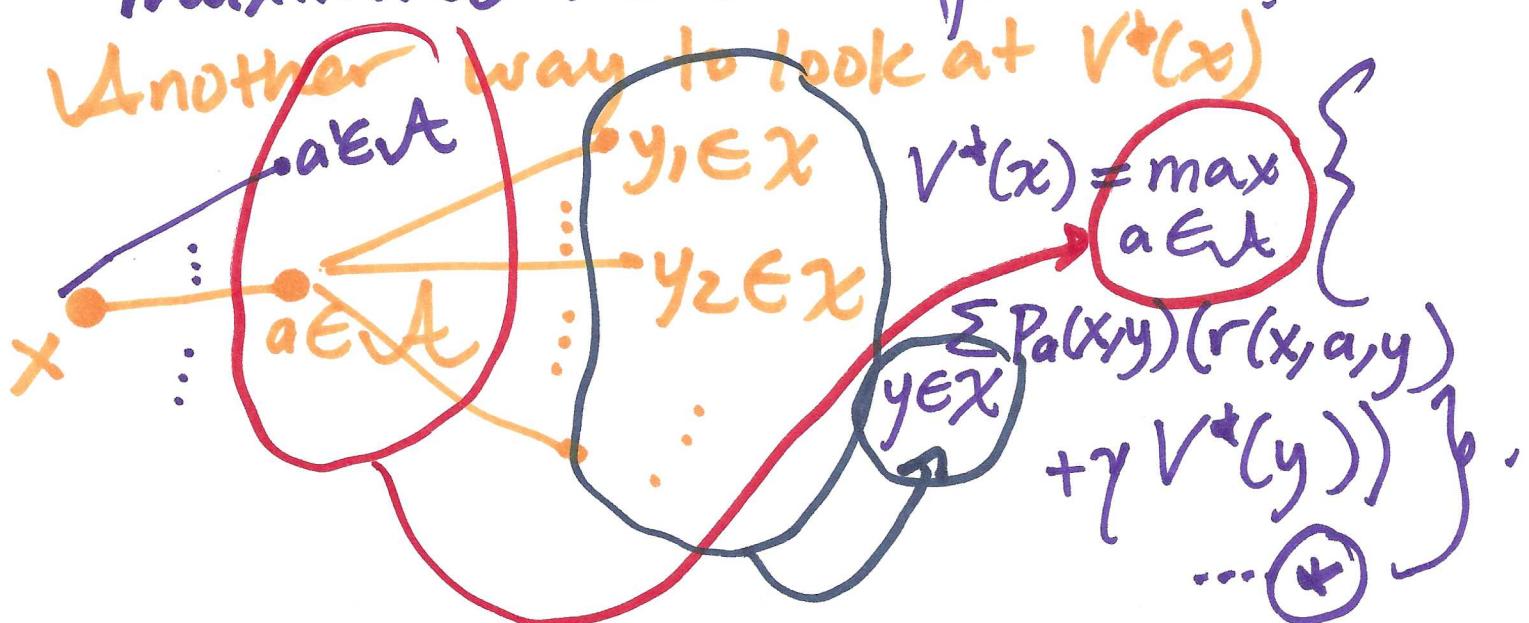
$$J(x, \{A_t\}) = E \left[\sum_{t=0}^{\infty} \gamma^t R(X_t, A_t) \mid X_0 = x \right]$$

Reward
Discount
(ensures value is finite).

Optimal Value Function:

$$V^*(x) = \max_{\{A_t\}} J(x, \{A_t\})$$

i.e. which sequence of (possibly random) actions should be taken so as to maximize the value function?



To understand how the equation is derived, note that the optimal value function is the best (maximum) expected cumulative reward for a given initial state x . (6)

Given initial state x , the expected reward of a policy that chooses action a all the time is

$$\sum_{y \in X} P_{a|X}(x,y) (r(x,a,y) + \gamma V^*(y)).$$

maximizing over a should give the optimal value function given by Θ^* .

We can now define:

$$\Theta^*(x,a) = \sum_{y \in X} P_{a|X}(x,y) (r(x,a,y) + \gamma \max_a \hat{\Theta}^*(y,a))$$

where

$$V^*(x) = \max_a \Theta^*(x,a)$$

If we can find this optimal Θ^* ,
 we have found the optimal policy
 to maximize the expected rewards.

The Θ learning algorithm can be
 used for this purpose. To intuitively
 understand how the algorithm is
 derived, note:

$$\Theta^*(x, a) = E_y \left[r(x, a, y) + \gamma \max_a \Theta^*(y, a) \right]$$

(expectation calculated over the
 next state).

Given x, a , and one $y \in \mathcal{X}$ we imagine
 updating $\Theta(x, a)$ to be close to ~~Θ^*~~ for the
 particular samples (x, a, y) we
 consider. This leads to:

$$\Theta(x, a) \leftarrow (1 - \alpha) \Theta(x, a) + \alpha (r(x, a, y) + \gamma \max_a \Theta(y, a))$$

This is the Θ learning algorithm.

8

In practice, the algorithm is run through many episodes of the agent starting from the initial state of the environment and exploring it.

- ① Initialize Θ table (e.g. randomly)
- ② for episodes = 1 to N:

 2.1 $x \leftarrow$ Initial State.

 2.2 for $t = 1$ to T :

 2.2.3 . With prob. ϵ select a random action a
(this is called exploration)

 Otherwise set $a = \arg\max_b \Theta(x, b)$
 to be the optimal action for state
 x (this is called exploitation).

 2.2.4 . Observe reward r and next state y_{t+1} .

 2.2.5 Update $\Theta(x, a) \leftarrow (1-\alpha)\Theta(x, a) + \alpha(r + \gamma \max_c \Theta(y, c))$.

 2.2.6 Set $x \leftarrow y$.

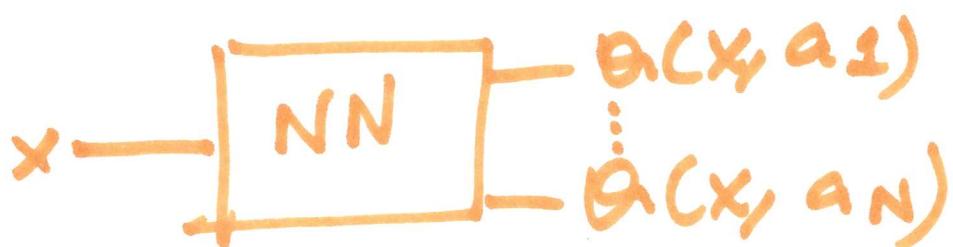
The algorithm is guaranteed to converge to the optimal policy (when $\epsilon=0$)
Setting $\epsilon=0$ in the beginning may allow faster convergence.

9

Deep Q Networks

In certain applications, the dimension/size of the state space may be impractically large so that ordinary Q-learning will not be applicable (e.g. atari games, chess, ...). One can then approximate $Q(x, a)$ via a neural network!

A typical architecture accepts the state x as the input and provides the Q values $Q(x, a)$ for every possible action $a \in A$.



The reason this design is favored is that one needs all the Q values $Q(x, a)$, $a \in A$ to determine the optimal action anyway.

How to train a deep Q-network?

10

In an optimal policy, we need:

$$Q^*(x, a) = \mathbb{E}_y \left(r(x, a, y) + \gamma \max_a Q^*(y, a) \right)$$

We can convert this to a loss function form that is more suitable to neural network training. Suppose we have a sequence of (state, action, reward, next state) vectors, i.e., (x_i, a_i, r_i, y_i) , $i=1 \dots K$. Then, we essentially wish to minimize:

$$\sum_{i=1}^K \left(Q(x_i, a_i) - \left(r_i + \gamma \max_a Q(y_i, a) \right) \right)^2$$

If we can minimize this, we can find the optimal policy (provided y_i is generated by the environment as a function of x_i, a_i).

→ The nice thing about this form is that it is almost in the form of a NN loss function, $Q(x_i, a_i)$ is the NN output and the rest of the terms is the desired output for the particular inputs...

One problem is that the desired output is not fixed, but will change as we change the network parameters. (11)

This results in instability during training. To fix this, the idea is to create two networks. One will be used to keep track of the target values

$$r_i + \gamma \max_a \Theta(y_i, a; W_t)$$

and will thus be called a target network. The other will be used to set & update the policy

$$\Theta(x_i, a_i; W_p)$$

and thus called the policy network.

Note during training that the networks will have different parameters.

Only the policy network will be optimized using the backpropagation algorithm. Every once in a while, we will copy the parameters of the policy network to the target network, updating $W_t \leftarrow W_p$.

Another issue is related to how the sequence $(x_i, a_i, r_i, y_i)_{i=1, \dots, K}$ is generated to train $D(x_i, a_i; W_p)$. Ideally, this sequence will be generated in a model-free fashion in a way similar to the Q learning algorithm. But then, there will be a lot of dependency between most elements of the sequence. In particular, $x_{i+1} = y_i$, for example. We have discussed that it is very hard to train neural networks over correlated data. The solution is to use experience replay:

Store experiences (x_i, a_i, r_i, y_i) to a buffer as the agents experience them. When training the policy network, sample a random batch of samples from this replay memory and use the batch to train W_p .

The DQN Algorithm.

13

- ① Initialize replay memory D to capacity M .
 - ② " policy network with weights W_p .
 - ③ " target " " " W_t .
 - ④ for episodes = 1 to N
 - ④.1 Initialize $x \leftarrow$ initial state.
 - ④.2 for $t=1$ to T .
 - ④.2.1 with prob ϵ . set $a =$ random action.
 - otherwise set $a = \arg\max_b Q(x, b; W_p)$
 - ④.2.2 Observe reward r , new state y .
 - ④.2.3 Store (x, r, a, y) in replay memory
 - ④.2.4 Sample random minibatch of transitions
 $(x_i, r_i, a_i, y_i), i=1, \dots, L$
from replay memory.
Let $z_i = r_i + \gamma \max_c Q(y_i, c; W_t)$
be the target values.
- Use the Backprop algorithm to minimize
- $$\sum_{i=1}^L (Q(x_i, a_i; W_p) - z_i)^2$$
- ④.2.5 Every C steps set $W_t \leftarrow W_p$.
- ④.2.6 Set $x \leftarrow y$.
- !! careful
 W_t , not W_p