

MARKOV DECISION PROCESS

- creem un 'factor de decizie' \rightarrow Agent - Agent care intervine cu mediul incognitor
- interacțiunile se desfășoară pas cu pas
- la fiecare pas agentul are o reprezentare a stării mediului
- agentul selectează o acțiune iar mediul face tranziția către o nouă stare, iar agentul primește o recompensă

Componentele unui MDP

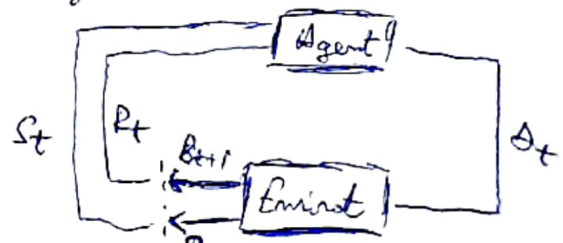
- mediul
- agentul
- toate stările posibile ale mediului ~~actuale~~
- toate acțiunile pe care le poate lua agentul
- toate recompensele pe care agentul le poate lua în funcție de acțiunile alese

Scopul : să maximizăm recompensa totală pe care o primim în funcție de acțiunile alese

Notatii :

- $S \rightarrow$ stări
- $A \rightarrow$ acțiuni
- $R \rightarrow$ reward
- $t \rightarrow$ timestep
- $S_t^{e^s} \rightarrow$ starea primită de agent la momentul t
- $A_t \rightarrow$ acțiunea selectată la momentul t pornind de la S_t
- $(S_t, A_t) \rightarrow$ face tranziția la $S_{t+1}^{e^s}$, și primește un reward R_{t+1}
- procesul primirii unei recompense poate fi notat ca o funcție $f : f(S_t, A_t) = R_{t+1}$
- traiectoria reprezentată procesul regitral :

$S_0, A_0, R_1, S_1, A_1, R_2, \dots$



Sequl MDP is to maximize it's cumulative reward. So, we will define Expected Return of the rewards at a given time

- we can define $G_t = R_{t+1} + R_{t+2} + \dots + R_T$, T - the final step

- the agent goal is to maximize G

- but G_t could be infinite so we need to modify it:

- ~~discounted~~ the agent goal is to maximize the discounted reward

- discount rate $= \gamma \in [0, 1]$

- $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+1+k}$

- $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+1+k}$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

- also: the agent will care more about the immediate reward rather than the future ones because they will be heavily discounted

What's the probability that an agent will select a specific action from a specific state?



Policies (π)

"How good" is a specific action or a specific state for the agent?



Value Functions

Policy: a function that maps a given state to probabilities of selecting each possible action from that state

• an agent follows a policy

• e.g. if an agent follows a policy π at time t , then

$\pi(a|s)$ is the probability that $A_t = a$ if $S_t = s$.

This means that, at time t , under policy π , the probability of taking action a in state s is $\pi(a|s)$

• also: for each state $s \in S$, π is a probability distribution over $a \in A(s)$

Value Function: function of states, or of state, action pairs, that estimates how good it is for an agent to be in a given state, or how good is the agent to perform a given action in a given state

• Value function $\xrightarrow{\text{given in terms of}}$ expected return \longrightarrow the way the agent acts $\xrightarrow{\text{is that it's}} \pi$ policy

• $\{ \text{states} \}$ or $\{ \text{state, action} \}$



state-value function



action-value function

State-Value Function: V_{π}

- how good any given state is for an agent following policy π
- in other words, it gives the value of a state under π

$$V_{\pi}(s) = E[G_t | S_t = s] = E\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right]$$

Action-Value Function: Q_{π}

- how good it is for the agent to take any given action from any given state while following policy π
- in other words, it gives the value of an action under π

$$Q_{\pi}(s, a) = E[G_t | S_t = s, A_t = a] = E\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

- $Q_{\pi}(s, a)$ - is referred as Q-function
- the output of $Q_{\pi}(s, a)$ - is referred as Q-value

Optimal Policies.

- a policy π is better than a policy π' if the expected return of π is \geq the expected return for π' for all states
 $\pi \succ \pi' \Leftrightarrow V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s \in S$
- a policy that is better or at least the same as all policies is called "optimal policy"

Optimal Value-state function

$$V_*(s) = \max_{\pi} V_{\pi}(s) \quad \forall s \in S$$

Optimal Action-Value function

$$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a) \quad \forall s \in S, a \in A(s)$$

Bellman Optimality Equation for Q_* : It must satisfy the eqn.

$$Q_*(s, a) = E[R_{t+1} + \gamma \max_{a'} Q_*(s', a')]$$

Q-learning

- it's a reinforcement learning technique used for learning the optimal policy and a Markov Decision Process
- it's objective is to find a policy that is optimal in the sense that the expected return over all successive timesteps is the maximum achievable
- in other words the goal of Q-learning is to find the optimal policy by learning the optimal Q values for each (state, action) pair

How Does It Work?

- iteratively updates the Q values for each state, action pair using the Bellman Equation until the Q value converges to the optimal
- Ex. • At the start of the game the learner has no idea of how good any given action is from any given state, it is not aware of anything besides the current state of the environment
- therefore the Q values for each (state, action) pair will all be initialized with 0 since the learner knows nothing about environment at the start
- through the game the Q values will be iteratively updated using Value Iteration
- we will be making use of a table, called Q-Table to store the Q-values for each (state, action) pair

End Ex

- Exploration vs Exploitation
- to get this balanced we use an ϵ -greedy strategy

ϵ -Greedy Strategy

- helps us balance the exploration and exploitation
- we have an exploration rate ϵ that we initially set to 1
- this ϵ is the probability that our agent will explore the environment rather than exploit it
- ϵ decreases over time

Update Q-Value

$$Q_*(s, a) = E \left[R_{t+1} + \gamma \max_{a'} Q_*(s', a') \right]$$

$$\text{loss} = Q_*(s, a) - Q(s, a), \text{ we try to minimize this}$$

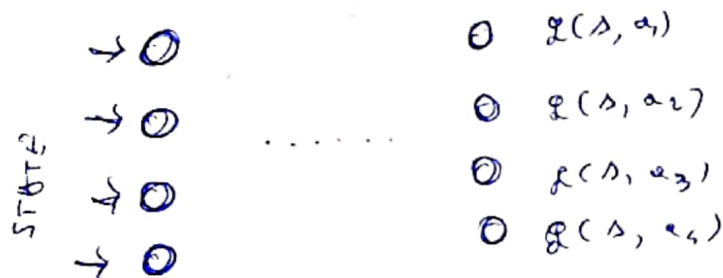
$$\text{loss} = E \left[R_{t+1} + \gamma \max_{a'} Q_*(s', a') \right] - E \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right]$$

- $\lambda \rightarrow$ learning rate

$$Q^{\text{new}}(s, a) = (1 - \lambda) \underbrace{Q(s, a)}_{\text{old Value}} + \lambda \left(\overbrace{R_{t+1} + \gamma \max_{a'} Q(s', a')}^{\text{Next Value}} \right)$$

DEEP Q-LEARNING / DEEP Q-NETWORKS

- instead of using value iteration to directly compute the Q values and find the optimal Q function, we instead use a function approximator to estimate the optimal Q function. What could it be? ... neural networks
- we will be using a deep neural network to estimate the Q values for each (state, action) pair in a given environment, and in turn the n.n. will approximate the optimal Q value
- the action combining Q learning with a deep n.n. is called Deep Q-learning
- a deep n.n. that approximates a Q function is called a Deep Q-Network / DQN
- n.n. receives states as an input and the network estimates Q values for each action that can be taken from that state



- it's objective is to approximate the optimal Q function, and the optimal Q will satisfy the Bellman Equation
- the loss of the network is calculated comparing the predicted Q values to the targeted Q values ~~from~~, and the objective is to minimize this loss
- in a ~~context~~ environment, we'll use images as our input to the network
- because a single frame does not contain sufficient information (eg. where the ball is moving), we will be using more a stack of frames which will represent a single input

Experience Replay / Replay Memory

- utilized in the training process of a Q -learning
- we store agent experiences at each timestep in the dataset called Replay Memory
- at time t the agent experiences $x_t = (S_t, a_t, R_{t+1}, S_{t+1})$
- all of the agent experiences at each timestep over all episodes replayed by the agent are stored in the Replay Memory
- in practice will usually see the replay memory set to some finite size like $N (= 10^6 \text{ for ex})$ and will store the last N experiences
- this replay memory dataset is actually what will randomly be sampling from to train the network
- the act of ~~storing~~^{pairing} experience and sampling from the Replay Memory that stores this experiences is actually what Experience Replay is
- why use Replay Memory? and not train sequentially?
 - to break the correlation between consecutive samples

How it works:

- initialize replay memory capacity
- initialize the network with random weights
- for each episode
 - initialize the starting state
 - for each timestep
 - select an action (via exploration or exploitation)
 - execute selected action in an emulator
 - observe reward and next state
 - store experience in replay memory
 - sample random batch from replay memory
 - process states from batch
 - pass batch of preprocessed states to policy network
 - calculate loss between output Q -values and target Q -values
 - requires a recall pass to the network for the next state
 - gradient descent updates weights in the policy network to minimize loss

Problems when Using Single Network

- when we update the weights, we also update the target ~~same~~ values ~~and~~ ~~value~~ because they are calculated using the same weights
- our q values will be updated with each iteration to move closer to the target q values, but the target q value will also be moving in the same direction
- rather than doing a second pass to the policy network to calculate the targeted q values, we instead obtain the targeted q values from a totally separate network called The Target Network
- this network is a clone of the policy network, it's weights are frozen with the original policy network's weights and we update the weights in the target network to the policy network's new weights every certain amount of timesteps
- this removes much of the instability introduced by only using one network to calculate both the q values as well as the target q values
- these values won't stay fixed the entire time, after x amount of time steps, will update the weights in the target network with the weight from our policy network which will then update the target q values with respect to what it's learned over those next timesteps; this will cause the policy network to start to approximate the updated targets

The new steps will look like:

- initialise replay memory exactly
- initialise the Policy network with random weights
- * • clone the policy network, and call it the Target network
- for each episode
 - initialise the starting state
 - for each time step
 - select an action via exploration or exploitation
 - execute selected action in an emulator
 - observe reward and next state
 - store experience in replay memory
 - sample random batch from replay memory
 - preprocess states from batch
 - pass batch of preprocessed states to Policy network
 - calculate loss between output Q-Values and Target Q-Values
 - * • require a pass to the Target network for the real state
 - gradient descent updates weight in the Policy network to minimize loss
 - * • after τ time steps, weights in the Target network are updated to the weights in the policy network