

Reward

- The goal of agent is maximize Reward
- All goals is described as maximization of expected cumulative Reward
- Cumulative Reward at step t

$$G_t = R_{t+1} + R_{t+2} + \dots$$

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

R_t : Reward at step t

- Problem:** Rewards that come soon (in the begin) are more probable to happen, since they are more predictable than long term.
- As we need to **discount** the reward as we go further from the current step.

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{k+t+1} \quad \text{where } \gamma \in [0, 1]$$

$$= \gamma^0 \cdot R_{t+1} + \gamma^1 \cdot R_{t+2} + \gamma^2 \cdot R_{t+3} \dots$$

- As the timestep increase the future rewards get less and less.

Episodic or Continuous Task

- Episodic have a start and a end point (terminal state)
- This create an episode
- Example:** In Mario Bros, an episode begins at the launch of a new Mario; and ends when you've killed or reach the end of the level.
- Continuous are tasks that continues forever (no terminal state).
- In this case the Agent has to learn how to choose the best action and simultaneously interacts with the environment
- there's not start/end point
- Example:** Automated stock trading

Monte Carlo

- Monte Carlo: collect rewards at the end of the episode and then calculate the maximum expected future reward.
- When an episode ends (terminal state) the agent looks at the total cumulative reward to see how well it did.
- In this approach, rewards are only received at ends of the game
- then we start a new game with the added knowledge.

$$V(S_t) = V(S_t) + \alpha [G_t - V(S_t)]$$

Learning Rate α Discounted Cumulative Rewards
 Maximum expected future reward starting from the state S_t (Updated Value) Actual estimation of maximum expected future reward starting at S_t (Previous Value)

Example:



If we take the maze environment:

- We always start at the same starting point.
- We terminate the episode if the cat eats us or if we move > 20 steps.
- At the end of the episode, we have a list of State, Actions, Rewards, and New States.
- The agent will sum the total rewards G_t (to see how well it did).
- It will then update $V(st)$ based on the formula above.
- Then start a new game with this new knowledge.

Temporal Difference TD

- Estimates the Expected Reward at each step, not waiting until the end of the episode to update the Maximum Expected Future Reward.
- It will update its value estimation V for the non-terminal state S_t occurring at the experience
- this method updates the Value Function after any individual steps.

$$V(S_t) = V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

↑
Current Value

Reward at $t+1$ ←
↓
Discounted Value on
the next step

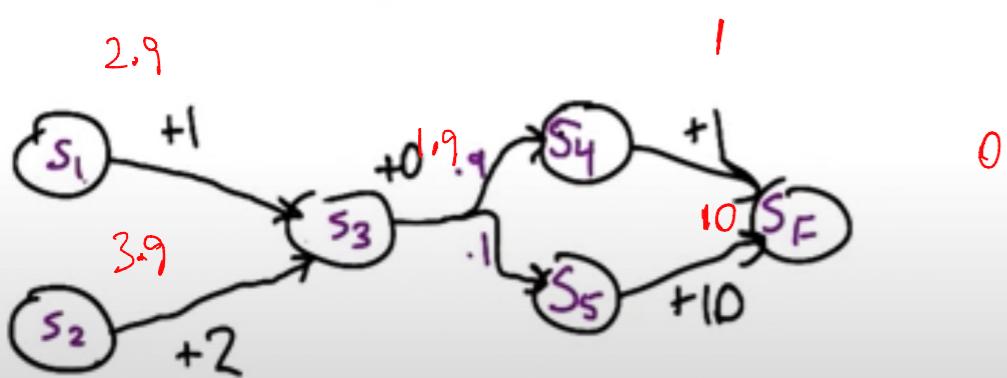
- At time $t+1$ they immediately form a TD using the observed Reward R_{t+1} and the current estimated $V(S_{t+1})$

Exploration / Exploitation

- Exploration: find more information about the environment
- Exploitation: exploit known information to Maximize Reward
- In the beginning the Agent must Explore to get knowledge of the environment.
- As times goes, the Agent starts to exploit its own knowledge to maximize the Reward.

Example

Learn $V(s) = \begin{cases} 0, & \text{if } s = s_F \\ E[r + \gamma V(s')], & \text{otherwise} \end{cases}$



• What's the value of each state? $\gamma=1.0$

1) Start backward, as we know the value of $s_F = 0$

$$V(s) = E[r + \gamma V(s)]$$

$$V(s_4) = E[1 + 1 \cdot V(s_F)]$$

$$1 + 1 \cdot 0 = 1$$

$$V(s_5) = E[1 + 1 \cdot V(s_F)]$$

$$V(s_5) = 10 + 0 = 10$$

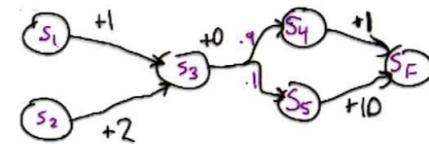
$$V(s_3) = E[0 + 1, \{ V(s_4) * 0.9 + V(s_5) * 0.1 \}]$$

$$0 + 1 \cdot \{ 1 + 1.0 \} = 1.9$$

Estimating from Data

$$\gamma=1$$

1. $s_1 \xrightarrow{+1} s_3 \xrightarrow{+0} s_4 \xrightarrow{+1} s_F$
2. $s_1 \xrightarrow{+1} s_3 \xrightarrow{+0} s_5 \xrightarrow{+0} s_F$
3. $s_1 \xrightarrow{+1} s_3 \xrightarrow{+0} s_4 \xrightarrow{+1} s_F$
4. $s_1 \xrightarrow{+1} s_3 \xrightarrow{+0} s_4 \xrightarrow{-1} s_F$
5. $s_2 \xrightarrow{+2} s_3 \xrightarrow{+0} s_5 \xrightarrow{+0} s_F$



QUIZ

What is an appropriate estimate for $V(s_1)$ after 3 episodes? 4?



• Now we want to calculate the $V(s_i)$ given those 5 samples

1) $s_1 \rightarrow s_3 \rightarrow s_4 \rightarrow s_F$

2

2) $s_1 \rightarrow s_3 \rightarrow s_5 \rightarrow s_F$

11

3) $s_1 - s_3 - s_4 - s_F$

2

4) $s_1 - s_3 - s_4 - s_F$

2

$$s_F = 0 \uparrow$$

$$s_4 = 1 + 1 \cdot 0 = 1$$

$$s_3 = 0 + 1 \cdot s_4 = 1$$

$$s_1 = 1 + 1 \cdot s_3 = 2$$

$$s_5 = 10 + 1 \cdot s_F = 10$$

$$s_3 = 0 + 1 \cdot s_5 = 10$$

$$s_1 = 1 + 1 \cdot s_3 = 11$$

$$s_4 = 1 + 1 \cdot s_F = 1$$

$$s_3 = 0 + 1 \cdot s_4 = 1$$

$$s_1 = 1 + 1 \cdot s_3 = 2$$

Estimation = Average

$$\text{up to } 3^{\text{rd}} \text{ episode: } (2+11+2)/3 = 5$$

$$4^{\text{th}} \text{ episode: } (2+11+2+2)/4 = 4.25$$

Three Approaches of RL

Value Based

- The goal is to optimize the value function $V(s)$

The $V(s)$ tell us the maximum expected future reward the agent will get at each state

The $V(s)$ of each state is the total amount of reward an agent can expect to accumulate over the future, starting at that state.

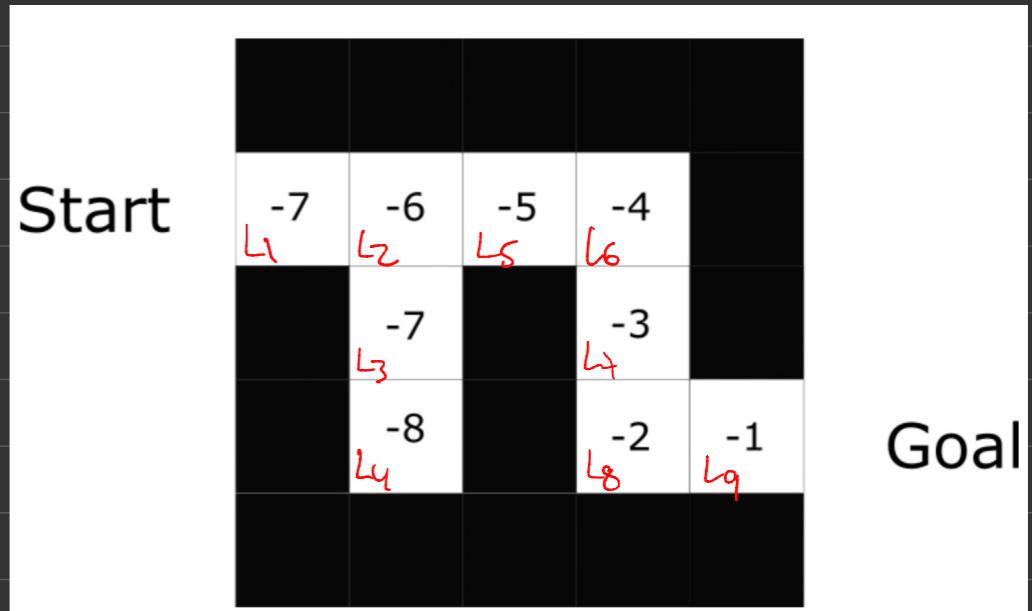
$$V_\pi(s) = E_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

↓
Discounted Reward
given that State s

expected ↴

- The Agent will use this value function to select which state to choose at each time step. The agent takes the biggest value.

- In this example, at each time step, the agent will take the biggest value: $\pi, -6, -3, -1$.



L1: at this location only can choose L2
L2: check the $V(s)$ of $[L_1, L_3, L_5]$. Choose the higher

Policy Based

- The goal is to directly optimize the policy function $\pi(s)$ without using a value function

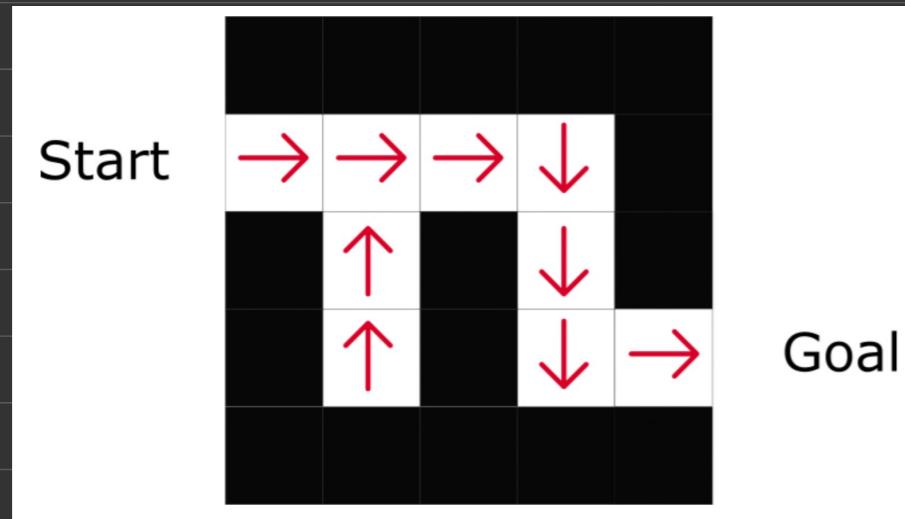
$$a = \pi(s) \quad \text{action} = \text{policy(state)}$$

- We learn the policy. This let us map each state to the best corresponding action

- Deterministic: A policy at a given state will always return the same action

- Stochastic: output a prob. distribution over actions

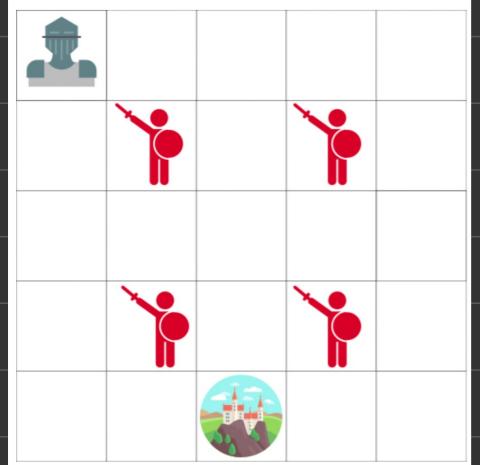
$$\pi'(a|s) = P[A_t=a | S_t=s]$$



- In this example the policy directly indicate the best action to take for each steps.

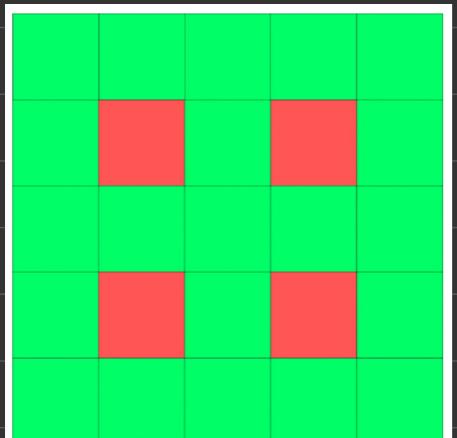
Q-Learning

The Problem:



- Move one tile each time
- lose -1 at each step (help to learn faster)
- If touch the enemy lose -100 and the episode finish
- If touch the castle, win 100 and episode finish

Strategy 1

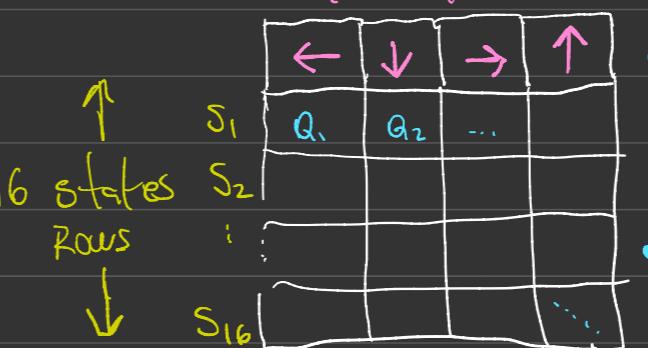


The same map, but colored in to show which tiles are safe to visit.

- Color each tile as green "safe" and red "enemy"
- Tell agent to move only in green tiles.
- Question:
 - how to know the best path?
 - how avoid infinite loop between adjacent green tiles?

Strategy 2: Q-Table

- Create a table where we'll calculate the Max Expected future Reward for each action at the state.
- We transform this grid into a table where:
 - Rows: Represent all states
 - Cols: Represent the actions
- In the castle example we have 4 actions ($\uparrow \leftarrow \downarrow \rightarrow$) and 16 states.



- each Q will represent the Quality or the Max Expected Future Reward for that given state and action
- that means each value will be the reward I'll get if I take that action at that state.
- Example: If I am at state s_1 and take the action \downarrow I'll receive the Q_2 as the reward.

• How to calculate each Q-value of the table?

Using Q-Learn Algorithm



0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

0 are impossible moves (if you're in top left hand corner you can't go left or up!)

- the algorithm take two inputs: (state, action)
- Return the expected future reward of that action when you are at that state.

$$Q(s_t, a_t) = E[R_{t+1} + \gamma^1 R_{t+2} + \gamma^2 R_{t+3} \dots | s_t, a_t]$$

Q-Value for that state given that action
Expected discounted cumulative reward given that state and action

- We can see this Q-function as a table lookup where it finds the row "state" and scan for the highest value in the col "Action"

To construct the values it is necessary the Agent explore the environment. As the explore is being it starts updating the values and after that it starts to exploit the knowledge acquired.

- To control the explore / exploit we use epsilon-greedy

```
if np.random.random() < epsilon:
    choose a random action # EXPLORE
else:
    choose a action in memory # EXPLOIT
reduce value of epsilon
```

epsilon [0, 1]

1. Initialize Q-values ($Q(s, a)$) arbitrarily for all state-action pairs.
2. For life or until learning is stopped...
3. Choose an action (a) in the current world state (s) based on current Q-value estimates ($Q(s, \cdot)$).
4. Take the action (a) and observe the outcome state (s') and reward (r).
5. Update $Q(s, a) := Q(s, a) + \alpha [r + \gamma \max_a Q(s', a') - Q(s, a)]$

Example:

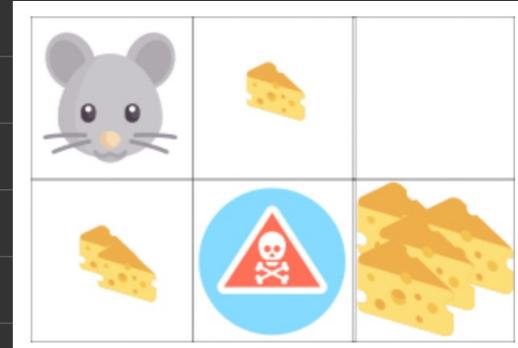
- One cheese +1
- two cheese +2
- Big pile +10 (end episode)
- Rat poison -10 (end episode)

$$LR = 0.1$$

$$\gamma = 0.9$$

	\leftarrow	\rightarrow	\uparrow	\downarrow
Start	0	0	0	0
Small cheese	0	0	0	0
Nothing	0	0	0	0
2 small cheese	0	0	0	0
Death	0	0	0	0
Big cheese	0	0	0	0

We move at random (for instance, right)



1. We are in start state and take a random action \rightarrow
2. We found a piece of cheese, so Reward +1
3. Update the Q-Value of being at start and doing action \rightarrow . $Q(\text{start}, \rightarrow)$

$$NewQ(s, a) = Q(s, a) + \alpha [R(s, a) + \gamma \max Q'(s', a') - Q(s, a)]$$

\ New Q value for that state and that action
\ Current Q value
\ Reward for taking that action at that state
\ Learning Rate
\ Discount rate
\ Maximum expected future reward given the new s' and all possible actions at that new state

0 0.1 1.0 0.9 0 0

$$NewQ(\text{start}, \rightarrow) = Q(\text{start}, \rightarrow) + \alpha [\Delta Q(\text{start}, \rightarrow)]$$

$$\Delta Q(\text{start}, \rightarrow) = R(\text{start}, \rightarrow) + \gamma \max Q'(1\text{cheese}, \rightarrow) - Q(\text{start}, \rightarrow)$$

$$= 1.0 + 0.9 * \max [Q(1\text{cheese}, \leftarrow), Q(1\text{cheese}, \downarrow), Q(1\text{cheese}, \rightarrow)] - Q(\text{start}, \rightarrow)$$

$$= 1.0 + 0.9 * 0.0 - 0 = 1$$



$$\text{New } Q(\text{start}, \rightarrow) = 0 + 0.1 \times 1.0 = 0.1$$

	\leftarrow	\rightarrow	\uparrow	\downarrow
Start	0	0.1	0	0
Small cheese	0	0	0	0
Nothing	0	0	0	0
2 small cheese	0	0	0	0
Death	0	0	0	0
Big cheese	0	0	0	0

The updated Q-table

```

def decrement_epsilon(self):
    self.epsilon = self.epsilon * self.eps_dec if self.epsilon > self.eps_end else self.eps_end

def learn(self, state, action, reward, state_):
    # get all possible actions given the new state_ the agent is now
    actions = np.array([self.Q[(state_, a)] for a in range(self.n_actions)]) # all actions given t

    # select the action that get the highest value
    a_max = np.argmax(actions)

    # calculate the discounted reward
    discounted_reward = reward + self.gamma * self.Q[(state_, a_max)] - self.Q[(state, action)]

    # update the previous state value with the new value
    self.Q[(state, action)] += self.lr * discounted_reward

    self.decrement_epsilon()

```

A recap...

- Q-learning is a value-based Reinforcement Learning algorithm that is used to find the optimal action-selection policy using a q function.
- It evaluates which action to take based on an action-value function that determines the value of being in a certain state and taking a certain action at that state.
- Goal: maximize the value function Q (expected future reward given a state and action).
- Q table helps us to find the best action for each state.
- To maximize the expected reward by selecting the best of all possible actions.
- The Q come from quality of a certain action in a certain state.
- Function $Q(\text{state}, \text{action}) \rightarrow$ returns expected future reward of that action at that state.
- This function can be estimated using Q-learning, which iteratively updates $Q(s, a)$ using the Bellman Equation
- Before we explore the environment: Q table gives the same arbitrary fixed value
→ but as we explore the environment → Q gives us a better and better approximation.

Recap: (Tabular) Q-Learning

Algorithm:

```

Start with  $Q_0(s, a)$  for all  $s, a$ .
Get initial state  $s$ 
For  $k = 1, 2, \dots$  till convergence
    Sample action  $a$ , get next state  $s'$ 
    If  $s'$  is terminal:
        target =  $R(s, a, s')$ 
    Sample new initial state  $s'$ 
    else:
        target =  $R(s, a, s') + \gamma \max Q_k(s', a')$ 
         $Q_{k+1}(s, a) \leftarrow (1 - \alpha)Q_k(s, a) + \alpha [\text{target}]$ 
         $s \leftarrow s'$ 

```

Deep Q-Learn

- High level idea: make Q-learning look like supervised learning
- Replace the Q-table by a function approximation that produces the Q-value
- This enable to interact with env. that has many states, that would be impossible to store in a table format
- **Q-learning:** take the current state and do a table lookup to (rows) and get the max value (col) that maxi the reward. This max will be the action the agent must perform.
- **DQN:** Instead a table, take the current state and use a Neural Net (function Approximation) to predict what would be the next best action to take. The table is replaced by the Neural Net weights.

Problems

1) What would be the target value so that we can calculate loss and backprop to update weights?

2) How avoid overfitting, that is, if the agent keep going to one direction, the network will "learn" to do just this action

Recap: Approximate Q-Learning

Algorithm:

Start with $Q_0(s, a)$ for all s, a .

Get initial state s

For $k = 1, 2, \dots$ till convergence

 Sample action a , get next state s'

 If s' is terminal:

 target = $R(s, a, s')$

 Sample new initial state s'

 else:

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

$\theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta} \mathbb{E}_{s' \sim P(s'|s, a)} [(Q_{\theta}(s, a) - \text{target}(s'))^2] |_{\theta=\theta_k}$

$s \leftarrow s'$

Chasing a nonstationary target!

Updates are correlated within a trajectory!

- As the same network is producing the predictions and targets, there will be too much instability, as in a supervised learning the targets are fixed given the same inputs.
- In RL, the targets are being updated all time as the agent keep exploring the environment. This ends up as the N.N chasing its own.
- To solve this, use two identical N.N.
 - 1^o used to predict and being updated the weights
 - 2^o other used to produce the targets, so that the targets are a little fixed. This network do not backprop.
- After some timesteps C , copy the weights from 1 to 2, so that the targets get better through time and the predictions of 1 get better.

- Use a replay memory with a capacity N to store

(s_t , a_t , r_t , s_{t+1})

the agent is in state s_0 , perform the action a_0 , get reward r_0 and move to new state s_1

DQN Training Algorithm

Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights θ

Initialize target action-value function \hat{Q} with weights $\theta^- = \theta$

For episode = 1, M do

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

 For $t = 1, T$ do

 With probability ϵ select a random action a_t
 otherwise select $a_t = \text{argmax}_a Q(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in D

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from D

 Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters θ

 Every C steps reset $\hat{Q} = Q$

 End For

End For

