

Reinforcement Learning Pseudo Code

Wouter Deketelaere

September 23, 2022

1 Agent Learning Algorithm

Algorithm 1 Agent Learning Algorithm

```
strategy  $\leftarrow$  LearningStrategy                                ▷ strategy = Learning Algorithm
environment  $\leftarrow$  Environment                                ▷ OpenAI Gym Environment

function LEARN(n_episodes)                                     ▷ maximum number of episodes
  while episode_count < n_episodes do
    episode  $\leftarrow$  Episode()                                   ▷ Create a new Episode
    state  $\leftarrow$  environment.state                             ▷ Where is the Agent?
    while episode not done do
      action  $\leftarrow$  strategy.next_action                       ▷ The Agent chooses its next action
      percept  $\leftarrow$  environment.step(action)                 ▷ The Agent ends up in new state
      episode.add(percept)                                       ▷ Add Percept to Episode
      strategy.learn(episode)                                    ▷ The Agent learns from Percepts in the Episode
      state = percept.next_state
    end while
  end while
end function
```

2 Learning Strategy

Algorithm 2 LearningStrategy

```
 $\lambda$                                 ▷ exponential decay rate used for exploration/exploitation
 $\varepsilon$                              ▷ (decaying) probability of selecting random action
 $\varepsilon_{max} = 1.0$                  ▷ exploration probability at start
 $\varepsilon_{min} = 0.01$                  ▷ minimum exploration probability

function LEARN(episode: Episode)                                ▷ abstract method
  Subclass implementation                                         ▷ here we insert the actual learning algorithm
  t  $\leftarrow$  t + 1                                                  ▷ increase episode time step
   $\tau \leftarrow \tau + 1$                                            ▷ increase overall time step
end function
```

3 Tabular Learner

Algorithm 3 TabularLearner implements LearningStrategy

α ▷ learning rate
 π ▷ 2D Numpy array holding the $\pi(s)$ values
 \mathbf{v} ▷ 1D Numpy array holding $v(s)$ -values
 \mathbf{q} ▷ 2D Numpy array holding all $q(s, a)$ -values

function INIT()
 $\pi \leftarrow$ uniform distributed policy table
 $\forall s : v(s) \leftarrow 0$
 $\forall s, a : q(s, a) \leftarrow 0$ ▷ $q(s, a)$ is the q-value of state s and action a
end function

function LEARN(episode: Episode)
 (SUBCLASS IMPLEMENTATION) ▷ subclasses insert their own code at this point
 EVALUATE() ▷ evaluate policy
 IMPROVE() ▷ improve policy
end function

function NEXTACTION() **return** action
 $a = \text{sample}(\pi(s))$ ▷ use Python's random.choice with p-parameter and add tie-breaking
return a
end function

function EVALUATE()
for each $s \in \mathcal{S}$ **do** ▷ $v(s)$ is equal to the maximum $q(s, a)$ -value in state s
 $v(s) \leftarrow \max_a \mathbf{q}(s)$ ▷ $\mathbf{q}(s)$ is a 1D Numpy array holding all $q(s, a)$ -values for state s
end for
end function

function IMPROVE
for each $s \in \mathcal{S}$ **do**
 $a_* = \arg \max_a \mathbf{q}(s)$ ▷ Numpys' argmax doesn't do tie-breaking, add this yourself

$$\pi(a|s) = \begin{cases} 1 - \varepsilon + \frac{\varepsilon}{|A(s)|}, & \text{if } a = a_* \\ \frac{\varepsilon}{|A(s)|}, & \text{if } a \neq a_* \end{cases}$$
 ▷ $|A(s)|$ is the number of actions in state s
end for
 $\varepsilon(\tau) = \varepsilon_{min} + (\varepsilon_{max} - \varepsilon_{min}) \cdot e^{-\lambda \cdot \tau}$ ▷ τ increases only after the episode ends
end function

4 Q-learning

Algorithm 4 Q-Learning implements TabularLearner

```
function LEARN(episode: Episode)
  p ← last Percept in Episode                                ▷  $s = p.s$  en  $a = p.a$ 
   $q(s, a) \leftarrow q(s, a) + \alpha \cdot [r(s, a) + \gamma \cdot \max_{a'} (q(s', a')) - q(s, a)]$     ▷ Update rule
  ▷  $\max_{a'}$  = maximum q-value over the next action  $a'$ 
  ▷  $r(s, a)$  = reward for taking action  $a$  in state  $s$ 
  SUPER( )                                                    ▷ now call EVALUATE and IMPROVE in the superclass

end function
```

5 N-step Q-learning

Algorithm 5 N-Step Q-Learning implements TabularLearner

```
N                                                                    ▷ N = number of steps

function LEARN(episode: Episode)

  if  $|\mathcal{E}| \geq N$  then                                          ▷ Do we have enough Percepts in the Episode  $\mathcal{E}$ ?
    for each  $p \in \mathcal{P}_N$  do                                          ▷  $P_N = N$  last percepts in Episode

       $q(s, a) \leftarrow q(s, a) - \alpha \cdot (q(s, a) - [r(s, a) + \gamma \cdot \max_{a'} (q(s', a'))])$     ▷ Update rule
      ▷  $s = p.s$  en  $a = p.a$ 

    end for
  end if
  SUPER( )                                                    ▷ now call EVALUATE and IMPROVE in the superclass
end function
```

6 Monte Carlo

Monte Carlo is equal to N-step Q-learning with $N = \text{length of the entire episode}$. By subclassing N-step Q-learning, you can reuse the code for N-step Q-learning.

7 Deep Q-learning (DQN)

Algorithm 6 Deep Q-Learning implements LearningStrategy

```

batch_size                                     ▷ batch size used by NNs

 $Q_1(\theta_1) \leftarrow$  Keras NN Model           ▷  $Q_1$  = Neural Network Model with weights  $\theta_1$ 
 $Q_2(\theta_2) \leftarrow$  Keras NN Model           ▷  $Q_2$  = Neural Network Model with weights  $\theta_2$ 

 $C$                                              ▷ Ring counter value for updating weights of  $Q_2$ 
function LEARN(episode: Episode)

    if  $|\mathcal{E}| \geq \text{batch\_size}$  then           ▷ Enough percepts in Episode  $\mathcal{E}$ ?
         $\mathcal{P} \leftarrow \mathcal{E}.\text{sample}()$        ▷ Sample random Percepts from Episode  $\mathcal{E}$ 
        LEARNFROMBATCH( $\mathcal{P}$ )
    end if
end function

function LEARNFROMBATCH( $\mathcal{P}$ )
     $\text{count} \leftarrow 0$ 

     $\mathcal{D} = \text{BUILDTRAININGSET}(\mathcal{P})$ 

    TRAINNETWORK( $\mathcal{D}$ )

     $\text{count} \leftarrow \text{count} + 1$ 
    if  $\text{count} \bmod C$  then                       ▷ Copy weights  $\theta_1$  from  $Q_1$  into  $\theta_2$  of  $Q_2$  every  $C$  times
         $\theta_2 = \theta_1$ 
    end if
end function

function BUILDTRAININGSET( $\mathcal{P}$ ) return  $\mathcal{D}$ 
     $\mathcal{D}$                                              ▷  $\mathcal{D}$  = empty trainingsset
    for each  $p \in \mathcal{P}$  do                         ▷  $P$  = sample of Percepts
         $\mathbf{q}(s) = Q_1(s; \theta_1)$                  ▷ Predict q-values (numpy array) for state s (numpy array)
         $q_* = \max(Q_2(s'; \theta_2))$              ▷  $q_*$  = max q-waarde according to  $Q_2$ 
         $q(s, a) = \begin{cases} r(s, a), & \text{if episode done} \\ r(s, a) + \gamma \cdot q_*, & \text{otherwise} \end{cases}$    ▷  $q(s, a) \neq \text{array}$ 
         $\mathcal{D}.\text{append}(s, q(s, a))$ 
    end for
    return  $\mathcal{D}$ 
end function

function TRAINNETWORK( $\mathcal{D}$ )
    for each  $(s, q(s, a)) \in \mathcal{D}$  do
         $\theta_1 = \theta_1 + \nabla_{\theta_1} Q_1(\theta_1)$    ▷ Train NN  $Q_1$  on  $(s, q(s, a))$  using fit-methode
    end for
end function

```

8 Double Deep Q-learning (DDQN)

Algorithm 7 Deep Q-Learning

```

batch_size                                     ▷ batch size voor de NNs

 $Q_1(\theta_1) \leftarrow$  Keras NN Model          ▷  $Q_1$  = Neural Network Model with weights  $\theta_1$ 
 $Q_2(\theta_2) \leftarrow$  Keras NN Model          ▷  $Q_2$  = Neural Network Model with weights  $\theta_2$ 

 $C$                                              ▷ Ring counter value for updating weights of  $Q_2$ 
function LEARN(episode: Episode)

    if  $|\mathcal{E}| \geq \text{batch\_size}$  then                ▷ Enough percepts in Episode  $\mathcal{E}$ ?
         $\mathcal{P} \leftarrow \mathcal{E}.\text{sample}()$             ▷ Sample random Percepts from Episode  $\mathcal{E}$ 
        LEARNFROMBATCH( $\mathcal{P}$ )
    end if
end function

function LEARNFROMBATCH( $\mathcal{P}$ )
     $count \leftarrow 0$ 

     $\mathcal{D} = \text{BUILDTRAININGSET}(\mathcal{P})$ 

    TRAINNETWORK( $\mathcal{D}$ )

     $count \leftarrow count + 1$ 
    if  $count \bmod C$  then                                ▷ Copy weights  $\theta_1$  from  $Q_1$  into  $\theta_2$  of  $Q_2$  every  $C$  times
         $\theta_2 = \theta_1$ 
    end if
end function

function BUILDTRAININGSET( $\mathcal{P}$ ) return  $\mathcal{D}$ 
     $\mathcal{D}$                                              ▷  $\mathcal{D}$  = empty trainingsset
    for each  $p \in \mathcal{P}$  do                                ▷  $P$  = sample of percepts
         $\mathbf{q}(s) = Q_1(s; \theta_1)$                     ▷ Predict q-values (numpy array) for state  $s$ 
         $a_* = \arg \max_a Q_1(s'; \theta_1)$                 ▷ Best action  $a_*$  as per  $Q_1$ 
         $q_* = Q_2(s'; \theta_2)[a_*]$                     ▷ q-value  $q_*$  as per  $Q_2$  for that best actie  $a_*$ 
         $q(s, a) = \begin{cases} r(s, a), & \text{if episode done} \\ r(s, a) + \gamma \cdot q_*, & \text{otherwise} \end{cases}$     ▷  $q(s, a) \neq \text{array}$ 
         $\mathcal{D}.\text{append}(s, q(s, a))$ 
    end for
    return  $\mathcal{D}$ 
end function

function TRAINNETWORK( $\mathcal{D}$ )
    for each  $(s, q(s, a)) \in \mathcal{D}$  do
         $\theta_1 = \theta_1 + \nabla_{\theta_1} Q_1(\theta_1)$     ▷ Train NN  $Q_1$  on  $(s, q(s, a))$  using fit-method
    end for
end function

```

9 Next action voor (D)DQN

Algorithm 8 Next Action (D)DQN

function NEXTACTION() **return** action

$$a = \begin{cases} \text{sample}(A(s)), & \text{random } j\varepsilon \\ \arg \max_a Q_1(s; \theta_1), & \text{otherwise} \end{cases}$$

▷ $Q_1(s; \theta_1)$ = predict q-values using Q_1

return a

end function
