

Introduction to Reinforcement Learning

Model-free

Raphael Cousin

Avril 19, 2019

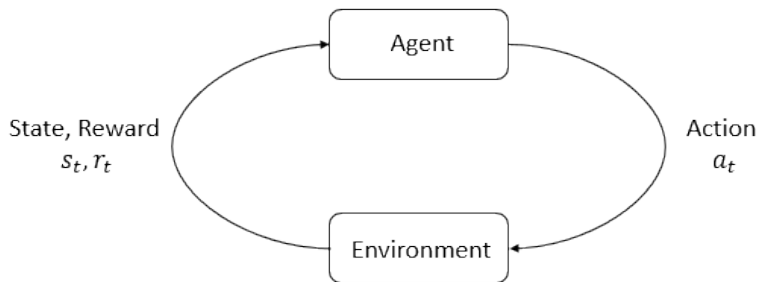
I.1) Introduction

- Supervised learning \rightarrow maps an input to an output based on example input-output pairs.
- Unsupervised learning \rightarrow models the probability density of inputs (using latent variables).
- Reinforcement learning (RL) \rightarrow maps actions to take in an environment to maximize cumulative reward.

I.2) Introduction

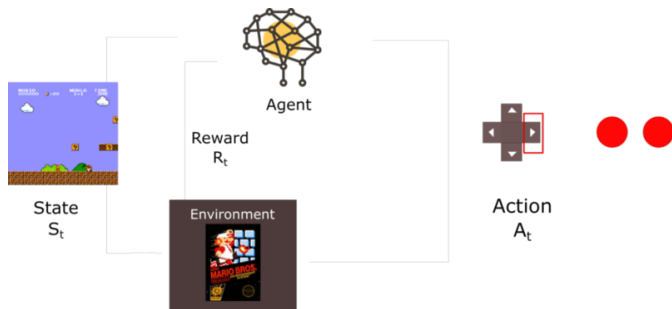
In RL problems:

- Unknown “correct” actions → no explicit supervision for a learning algorithm to try to mimic.
- Provide our algorithms only a reward function → indicates to the learning agent when it is doing well, and when it is doing poorly.



I.3) Introduction

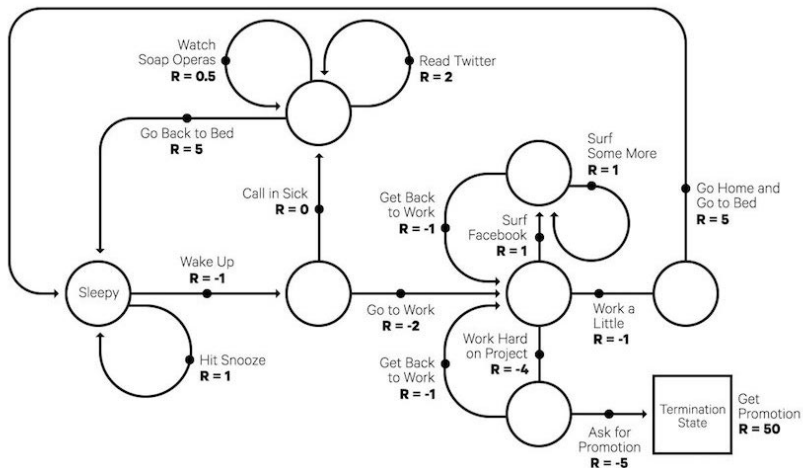
Applications: autonomous helicopter flight, robot legged locomotion, cell-phone network routing, marketing strategy selection, factory control, efficient web-page indexing, chess, video games.



I.4) Formalism: Markov Decision Process

- S set of all valid states
- A set of all valid actions
- $R : S \times A \times S \rightarrow \mathbb{R}$ reward function
 $r_t = R(s_t, a_t, s_{t+1})$
- $P : S \times A \rightarrow \mathcal{P}(S)$ transition probability function
 $P(s'|s, a)$ probability of transitioning into state s' if you are in state s and take action a
- ρ_0 starting state distribution

I.5) Example: Markov Decision Process



I.6) Policy and Trajectories

- Policy (agent brain): rule used to decide what actions to take
 $a_t \sim \pi(\cdot | s_t)$
- Trajectories: sequence of states and actions
 $\tau = (s_0, a_0, s_1, a_1, \dots)$
- States follow:
 $s_0 \sim \rho_0(\cdot)$
 $s_{t+1} \sim P(\cdot | s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0) = P(\cdot | s_t, a_t)$

I.7) Reward and Return

Keep in mind: the reward function R : $r_t = R(s_t, a_t, s_{t+1})$

- Finite-horizon undiscounted return:

$$R(\tau) = \sum_{t=0}^T r_t.$$

- Infinite-horizon discounted return:

$$R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$$

I.8) The RL Problem

- The probability of a T -step trajectory is:
$$P(\tau|\pi) = \rho_0(s_0) \prod_{t=0}^{T-1} P(s_{t+1}|s_t, a_t) \pi(a_t|s_t)$$
- The expected return:
$$J(\pi) = \int_{\tau} P(\tau|\pi) R(\tau) = E_{\tau \sim \pi}[R(\tau)]$$
- The optimization problem in RL:
$$\pi^* = \arg \max_{\pi} J(\pi)$$

I.9) Value Functions

- Value Function, \rightarrow expected return in state s , according to π :

$$V^\pi(s) = E_{\tau \sim \pi}[R(\tau) | s_0 = s]$$

- The Optimal Value Function \rightarrow expected return in state s , according to optimal policy:

$$V^*(s) = \max_{\pi} E_{\tau \sim \pi}[R(\tau) | s_0 = s]$$

I.10) Action Value Functions

- The Action-Value Function \rightarrow expected return taking action a in state s , and then according π :

$$Q^\pi(s, a) = E_{\tau \sim \pi}[R(\tau) | s_0 = s, a_0 = a]$$

- The Optimal Action-Value Function \rightarrow expected return taking action a in state s , and then according to optimal policy:

$$Q^*(s, a) = \max_{\pi} E_{\tau \sim \pi}[R(\tau) | s_0 = s, a_0 = a]$$

I.11) Bellman Equations

- **Idea :** The value of your starting point is the reward you expect to get from being there, plus the value of wherever you land next.
- $V^\pi(s) = E_{a \sim \pi, s' \sim P}[r(s, a) + \gamma V^\pi(s')]$
- $Q^\pi(s, a) = E_{s' \sim P}[r(s, a) + \gamma E_{a' \sim \pi}[Q^\pi(s', a')]]$

With $s' \sim P(\cdot|s, a)$, $a \sim \pi(\cdot|s)$; and $a' \sim \pi(\cdot|s')$.

I.12) Bellman Equations

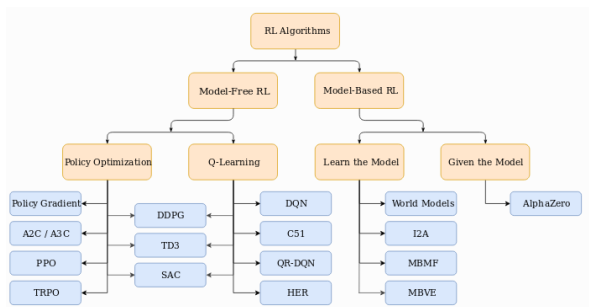
Bellman equations for the optimal value functions

- $V^*(s) = \max_a E_{s' \sim P}[r(s, a) + \gamma V^*(s')]$
- $Q^*(s, a) = E_{s' \sim P}[r(s, a) + \gamma \max_{a'} Q^*(s', a')]$

I.13) Advantage Functions

- **Idea :** It is not how good an action is in an absolute sense, but only how much better it is than others on average.
- The advantage function \rightarrow how much better to take action a in state s , over randomly ($\pi(\cdot|s)$), following π forever after.
- $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$.

I.14) Bestiary of RL Algorithms



- **Model free** : large data + low complexity
- **Model based** : less data + high complexity

II.1) Dynamic Programming

- Policy Evaluation: compute the state-value V_π for a given policy π :

$$V_{t+1}(s) = \mathbb{E}_\pi[r + \gamma V_t(s') | S_t = s] = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) (r + \gamma V_k(s')) \quad (1)$$

- Policy Improvement: generates a better policy $\pi' \geq \pi$ by acting greedily.

$$Q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s, A_t = a] = \sum_{s', r} P(s', r | s, a) (r + \gamma V_\pi(s')) \quad (2)$$

$$\pi'(s) = \arg \max_{a \in \mathcal{A}} Q_\pi(s, a)$$

- Policy Iteration: iterative procedure to improve the policy when combining policy evaluation and improvement.

$$\pi_0 \xrightarrow{\text{evaluation}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluation}} V_{\pi_1} \xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{evaluation}} \dots \xrightarrow{\text{improve}} \pi^* \quad (3)$$

This policy iteration process works and always converges to the optimality.

II.2) Monte-Carlo Methods

- First, let's recall $V(s) = \mathbb{E}[G_t | S_t = s]$. To compute the empirical return G_t , MC methods need to learn from complete episodes $S_1, A_1, R_2, \dots, S_T$ to compute $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$ and all the episodes must eventually terminate.
- The empirical value function for state s is: $V(s) = \frac{\sum_{t=1}^T \mathbb{1}[S_t=s] G_t}{\sum_{t=1}^T \mathbb{1}[S_t=s]}$.
- As well, the empirical action-value function is :
$$Q(s, a) = \frac{\sum_{t=1}^T \mathbb{1}[S_t=s, A_t=a] G_t}{\sum_{t=1}^T \mathbb{1}[S_t=s, A_t=a]}$$
- Improve the policy greedily with respect to the current value function: $\pi(s) = \arg \max_{a \in \mathcal{A}} Q(s, a)$
- Generate a new episode with the new policy π (i.e. using algorithms like ϵ -greedy helps us balance between exploitation and exploration.)
- Estimate Q using the new episode:
$$q_\pi(s, a) = \frac{\sum_{t=1}^T \left(\mathbb{1}[S_t=s, A_t=a] \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1} \right)}{\sum_{t=1}^T \mathbb{1}[S_t=s, A_t=a]}$$

II.3) Temporal-Difference Learning

TD learning can learn from incomplete episodes and hence we don't need to track the episode up to termination.

- Value Estimation: The key idea in TD learning is to update the value function $V(S_t)$ towards an estimated return $R_{t+1} + \gamma V(S_{t+1})$. To what extent we want to update the value function is controlled by the learning rate hyperparameter α :

$$V(S_t) \leftarrow (1 - \alpha)V(S_t) + \alpha G_t \quad (4)$$

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t)) \quad (5)$$

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)) \quad (6)$$

- Similarly, for action-value estimation:
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)).$

II.3) SARSA: On-Policy TD control

At time step t , we start from state S_t and pick action according to Q values, $A_t = \arg \max_{a \in \mathcal{A}} Q(S_t, a)$; ϵ -greedy is commonly applied.

With action A_t , we observe reward R_{t+1} and get into the next state S_{t+1}

Then pick the next action in the same way as in step 1

$$A_{t+1} = \arg \max_{a \in \mathcal{A}} Q(S_{t+1}, a)$$

Update the action-value function

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

$t = t + 1$ and repeat from step 1.

In each update of SARSA, we need to choose actions for two steps by following the current policy twice (in Step 1. 3.).

II.4) Q-Learning: Off-policy TD control

TD learning can learn from incomplete episodes and hence we don't need to track the episode up to termination.

- At time step t , we start from state S_t and pick action according to Q values, $A_t = \arg \max_{a \in \mathcal{A}} Q(S_t, a)$; ϵ -greedy is commonly applied. With action A_t , we observe reward R_{t+1} and get into the next state S_{t+1}

Update the action-value function

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_{a \in \mathcal{A}} Q(S_{t+1}, a) - Q(S_t, A_t))$$

$t = t + 1$ and repeat from step 1.

II.5) Deep Q-Network

Theoretically, we can memorize $Q^*(.)$ for all state-action pairs in Q-learning, like in a gigantic table. However, it quickly becomes computationally infeasible when the state and action space are large. Thus people use functions (i.e. a machine learning model) to approximate Q values and this is called function approximation. For example, if we use a function with parameter θ to calculate Q values, we can label Q value function as $Q(s,a;\theta)$.

II.6) Deep Q-Network improvement

Unfortunately Q-learning may suffer from instability and divergence when combined with an nonlinear Q-value function approximation and bootstrapping (See Problems 2).

Deep Q-Network (“DQN”; Mnih et al. 2015) aims to greatly improve and stabilize the training procedure of Q-learning by two innovative mechanisms:

- Experience Replay: All the episode steps $e_t = (S_t, A_t, R_t, S_{t+1})$ are stored in one replay memory $D_t = \{e_1, \dots, e_t\}$
- *PeriodicallyUpdatedTarget* :
- The loss function looks like this:
$$\mathcal{L}(\theta) = \mathbb{E}_{(s,a,r,s') \sim U(D)} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right)^2 \right].$$
where $U(D)$ is a uniform distribution over the replay memory D ; θ^- is the parameters of the frozen target Q-network.
- clip the error term to be between $[-1, 1]$.

III.1) Policy Optimization

- Parametrization of policy, π_θ .
- We aim to maximize the expected return $J(\pi_\theta) = E_{\tau \sim \pi_\theta}[R(\tau)]$.
- Gradient ascent:
$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_\theta)|_{\theta_k}.$$
- We can proof that:
$$\nabla_\theta J(\pi_\theta) = E_{\tau \sim \pi_\theta}[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t) R(\tau)]$$

III.2) policy gradient proof

- $P(\tau|\theta) = \rho_0(s_0) \prod_{t=0}^T P(s_{t+1}|s_t, a_t) \pi_\theta(a_t|s_t)$
- The Log-Derivative Trick. $\nabla_\theta P(\tau|\theta) = P(\tau|\theta) \nabla_\theta \log P(\tau|\theta)$.
- Log-Probability of a Trajectory. $\log P(\tau|\theta) = \log \rho_0(s_0) + \sum_{t=0}^T \left(\log P(s_{t+1}|s_t, a_t) + \log \pi_\theta(a_t|s_t) \right)$.
- The environment has no dependence on θ
 $\rho_0(s_0)$, $P(s_{t+1}|s_t, a_t)$, and $R(\tau)$ are zero.
- Grad-Log-Prob of a Trajectory.
$$\begin{aligned} & \nabla_\theta \log P(\tau|\theta) \\ &= \nabla_\theta \log \rho_0(s_0) + \sum_{t=0}^T \left(\nabla_\theta \log P(s_{t+1}|s_t, a_t) + \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \\ &= \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t). \end{aligned}$$

III.3) policy gradient proof

- $\nabla_{\theta} J(\pi_{\theta}) = \nabla_{\theta} E_{\tau \sim \pi_{\theta}}[R(\tau)]$
= $\nabla_{\theta} \int_{\tau} P(\tau|\theta) R(\tau)$ Expand expectation
= $\int_{\tau} \nabla_{\theta} P(\tau|\theta) R(\tau)$ Bring gradient under integral
= $\int_{\tau} P(\tau|\theta) \nabla_{\theta} \log P(\tau|\theta) R(\tau)$ Log-derivative trick
= $E_{\tau \sim \pi_{\theta}}[\nabla_{\theta} \log P(\tau|\theta) R(\tau)]$ Return to expectation form
- $\nabla_{\theta} J(\pi_{\theta}) =$
 $E_{\tau \sim \pi_{\theta}}[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) R(\tau)]$ Expression for grad-log-prob

III.4) Policy gradient estimation

- Set of trajectories $\mathcal{D} = \{\tau_i\}_{i=1,\dots,N}$ from the policy π_θ
- $\hat{g} = \frac{1}{|\mathcal{D}|} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t) R(\tau),$

III.5) REINFORCE On-Policy Policy Gradient

- ➊ Initialize Θ at random
- ➋ Generate one episode $\tau S_1, A_1, R_2, S_2, A_2, \dots, S_T$
- ➌ For $t = 1, 2, \dots, T$
 - ➊ Estimate the the return R_t since the time step t .
 - ➋ Apply gradient ascent $\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla \ln \pi(A_t | S_t, \theta)$
- Set of trajectories $\mathcal{D} = \{\tau_i\}_{i=1, \dots, N}$ from the policy π_θ
- $\hat{g} = \frac{1}{|\mathcal{D}|} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) R(\tau),$

III.6) Policy gradient improvement

Don't Let the Past Distract You

- $\nabla_{\theta} J(\pi_{\theta}) = E_{\tau \sim \pi_{\theta}} [\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau)]$.
- Taking a step with this gradient pushes up the log-probabilities of each action in proportion to $R(\tau)$, the sum of all rewards ever obtained. But this doesn't make much sense.
- Agents should really only reinforce actions on the basis of their consequences. Rewards obtained before taking an action have no bearing on how good that action was: only rewards that come after.

III.7) Useful lemma

- We want to proof:
 $E_{x \sim P_\theta}[\nabla_\theta \log P_\theta(x)] = 0.$
- Recall $\int_x P_\theta(x) = 1.$
- Thus
 $\nabla_\theta \int_x P_\theta(x) = \nabla_\theta 1 = 0.$
- Thus
$$\begin{aligned} 0 &= \nabla_\theta \int_x P_\theta(x) \\ &= \int_x \nabla_\theta P_\theta(x) \\ &= \int_x P_\theta(x) \nabla_\theta \log P_\theta(x) \\ 0 &= E_{x \sim P_\theta}[\nabla_\theta \log P_\theta(x)]. \end{aligned}$$

III.8) Don't Let the Past Distract You

- $\nabla_{\theta} J(\pi_{\theta}) = E_{\tau \sim \pi_{\theta}} [\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1})].$
- In this form, actions are only reinforced based on rewards obtained after they are taken.
- $\hat{R}_t \doteq \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}),$

III.9) Don't Let the Past Distract You - proof

- $$\begin{aligned}\nabla_{\theta} J(\pi_{\theta}) &= E_{\tau \sim \pi_{\theta}} [\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau)] \\ &\quad \sum_{t=0}^T \sum_{t'=0}^T E_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(s_{t'}, a_{t'}, s_{t'+1})] \\ &= \\ &\quad \sum_{t=0}^T \sum_{t'=0}^T E_{\tau \sim \pi_{\theta}} [E_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(s_{t'}, a_{t'}, s_{t'+1})] | s_0, \dots, s_t]\end{aligned}$$
- For the case of $t' \leq t$ (the reward comes before the action being reinforced), this term is zero, thanks to lemma.

III.10) Baselines in Policy

- For any function b which only depends on state,
- $E_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0$.
- Add or subtract any number of terms like this from our expression, don't change expectation.
- $\nabla_\theta J(\pi_\theta) = E_{\tau \sim \pi_\theta} [\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) - b(s_t) \right)]$.

III.11) Baselines in Policy

- Any function b used in this way is called a baseline.
- Interesting choice is $V^\pi(s_t)$.
- Empirically, $b(s_t) = V^\pi(s_t)$ reduce variance. Conceptually, it encodes the intuition that if an agent gets what it expected, it should “feel” neutral about it.

III.12) Approximate $V^\pi(s_t)$

- Basic method for learning $V_\phi \rightarrow$ minimize a mean-squared-error objective:
- $\phi_k = \arg \min_{\phi} E_{s_t, \hat{R}_t \sim \pi_k} [(V_\phi(s_t) - \hat{R}_t)^2],$
- π_k is the policy at epoch k. Using gradient descent, from previous ϕ_{k-1} .

III.13) Gradient estimation summary

What we have seen is that the policy gradient has the general form

$$\nabla_{\theta} J(\pi_{\theta}) = E_{\tau \sim \pi_{\theta}} [\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \Phi_t],$$

where Φ_t could be any of

- $\Phi_t = R(\tau)$
- $\Phi_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1})$
- $\Phi_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) - b(s_t)$ All of these choices lead to the same expected value for the policy gradient, despite having different variances.

III.14) Other gradient estimation

- Action-Value Function.

$$\Phi_t = Q^{\pi_\theta}(s_t, a_t)$$

- The Advantage Function. Recall $A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t)$, describes how much better or worse it is than other actions on average (relative to the current policy).

$$\Phi_t = A^{\pi_\theta}(s_t, a_t)$$

III.15) Estimation with Advantage Function - proof

- We want to proof:

$$\nabla_{\theta} J(\pi_{\theta}) = E_{\tau \sim \pi_{\theta}} [\sum_{t=0}^T \left(\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) Q^{\pi_{\theta}}(s_t, a_t)]$$

- We know:

$$\begin{aligned} \nabla_{\theta} J(\pi_{\theta}) &= E_{\tau \sim \pi_{\theta}} [\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{R}_t] \\ &= \sum_{t=0}^T E_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{R}_t] \end{aligned}$$

- Define $\tau_{:t} = (s_0, a_0, \dots, s_t, a_t)$ as the trajectory up to time t , and τ_t as the remainder of the trajectory after that.

$$\nabla_{\theta} J(\pi_{\theta}) = \sum_{t=0}^T E_{\tau_{:t} \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) E_{\tau_t \sim \pi_{\theta}} [\hat{R}_t | \tau_{:t}]]$$

- The grad-log-prob is constant with respect to the inner expectation (because it depends on s_t and a_t , which the inner expectation conditions on as fixed in $\tau_{:t}$), so it can be pulled out, leaving:

$$\nabla_{\theta} J(\pi_{\theta}) = \sum_{t=0}^T E_{\tau_{:t} \sim \pi_{\theta}} [E_{\tau_t \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{R}_t | \tau_{:t}]]$$

- Markov chain gives us

$$E_{\tau_t \sim \pi_{\theta}} [\hat{R}_t | \tau_{:t}] = E_{\tau_t \sim \pi_{\theta}} [\hat{R}_t | s_t, a_t]$$

- Which is the definition of $Q^{\pi_{\theta}}(s_t, a_t)$

III.16) Actor critic On-Policy Policy Gradient

If the value function is learned in addition to the policy, we would get Actor-Critic algorithm.

Critic: updates value function parameters w and depending on the algorithm it could be action-value $Q(a|s;w)$ or state-value $V(s;w)$.

Actor: updates policy parameters θ , in the direction suggested by the critic, $\nabla_{\theta} Q(a|s;w)$.

III.17) Off-Policy Policy Gradient

TODO

IV.1) Other approach

It exists a large panel of algorithms, build with different intuition and purpose. But as deep learning they are few mathematical proofs.

We're still going to name a few:

- For Model-Bases, the most famous is AlphaGo

*rlsilver2017mastering, but it exists many other I2Arl, acaniere2017imc
Free algorithms, for example proximal Policy optimization rl_schulman*

- hybrid method where the value function is learned in addition to the policy, called the Actor-Critic algorithm, asynchronous advantage actor-critic a3c, *baeizadeh2016reinforcement, Soft Actor – Critics sac, harnoja2018soft*

IV.2) Other approach

The evaluation of these techniques is mainly empirical. The most common evaluation was on Atari

rl_kaiser2019modelgames.Theenvironmentalmostalreadyexist,itwasjust

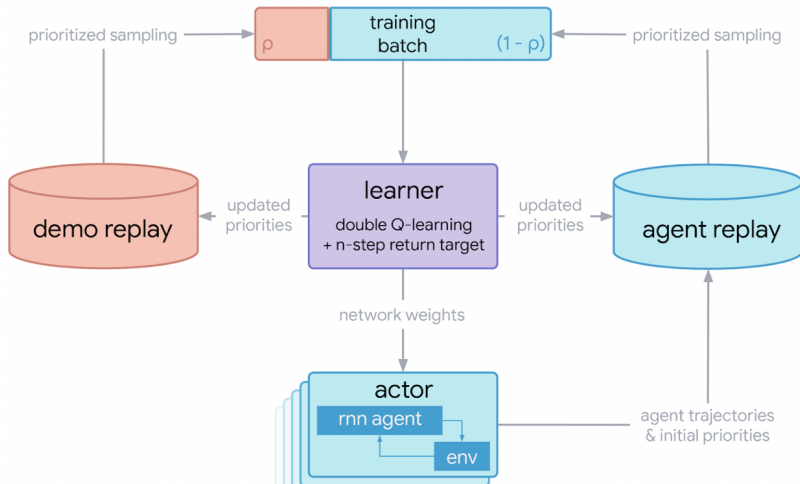
For a lot of games, a simple *DQN* algorithm was able to be much more efficient than human in the games. See ??

But for other it was not good at all. If we take a look at the games with the worst results. They are 2 types of games:

- the games with very sparse reward, for example platform games that implies a long suit of specific actions before obtain a reward. In this kind of games, the random exploration has no chance to obtain a reward (even after million of games), so their are no matter to learn, the agent is stuck to play randomly for ever
- the gamers with lot of negative reward, for example a game that many actions bring a loss, it the kind of game the agent learn to do nothing to avoid bad action

IV.3) rl with demonstration

An idea to help the learning process is based on human demonstrations. We create trajectories made by human and we add this trajectories to the memory of the agent. Then we let the agent overtake the master.



This approach is not always practicable. Furthermore the fact of learning has to be not only measured by the score achieved but also the complexity and energy consumption. Other methods try to improve the fact of learning:

- PILCO

PILCO_{levine2018reinforcementisbasedonBayesianinference}(andGdrivenalgorithms_{rl_pathak2017curiosity,rl_burda2018largecompute}

- Behaviour Suite for Reinforcement Learning

bsuit_{osband2019behaviourisacollectionofcarefully} – designedexperimentssthatinvestigatecorecapabilitiesofareinforceme

- To collect clear, informative and scalable problems that capture key issues in the design of efficient and general learning algorithms.
- To study agent behavior through their performance on these shared benchmarks.

Exploration

Credit Assignment

References

TODO