



Università degli
Studi di Pavia

Deep Learning

11 – Deep Reinforcement Learning

Marco Piastra & Andrea Pedrini(*)

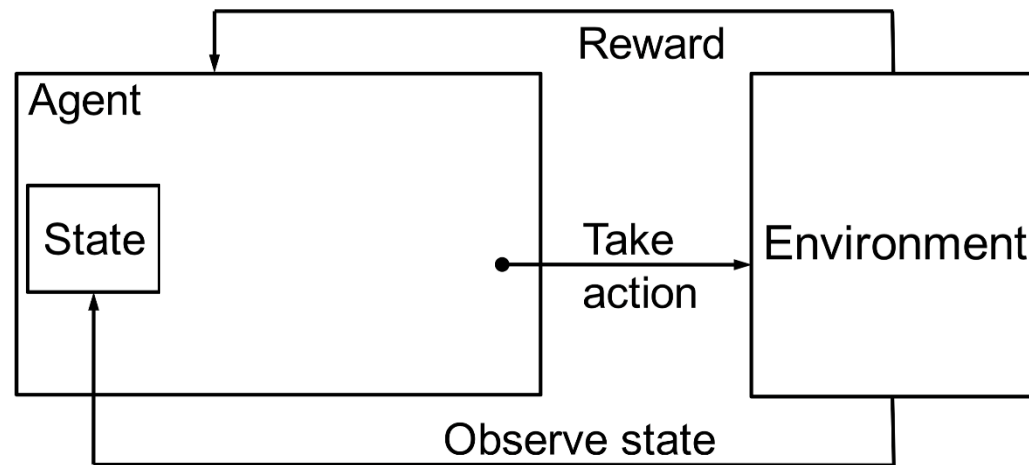
(*) Dipartimento di Matematica F. Casorati

This presentation can be downloaded at:
<http://vision.unipv.it/DL>

Basics (Intuition)

Deep Reinforcement Learning (DRL)

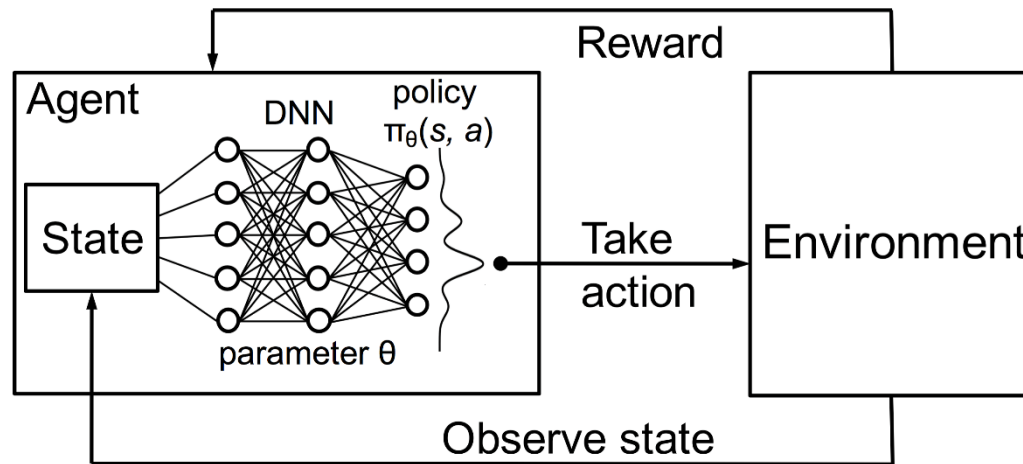
■ Reinforcement Learning



Deep Reinforcement Learning (DRL)

■ Deep Reinforcement Learning

Using a deep neural network as the approximator $\hat{Q}(s, a)$



The optimal policy is learnt incrementally by using a deep neural network

Q-Learning

▪ **Q-Learning Algorithm**

Initialize $\hat{Q}(s, a)$ at random, put the agent in a random state s

Repeat:

- 1) Select the action $\operatorname{argmax}_a \hat{Q}(s, a)$ with probability $(1 - \epsilon)$ otherwise, select a at random
- 2) The agent is now in state s' and has received the reward r
- 3) Update $\hat{Q}(s, a)$ by

$$\Delta \hat{Q}(s, a) = \alpha [r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)]$$

\mathcal{A}

Deep Reinforcement Learning

▪ Q-Learning Algorithm

Initialize $\hat{Q}(s, a)$ at random, put the agent in a random state s
Repeat:

- 1) Select the action $\operatorname{argmax}_a \hat{Q}(s, a)$ with probability $(1 - \epsilon)$ otherwise, select a at random
- 2) The agent is now in state s' and has received the reward r
- 3) Update $\hat{Q}(s, a)$ by

$$\Delta \hat{Q}(s, a) = \alpha [r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)]$$

Fundamental Idea:

use a deep neural network to learn the approximator $\hat{Q}(s, a)$
from the examples collected while **exploring - exploiting**

Deep Reinforcement Learning

▪ Q-Learning Algorithm

Initialize $\hat{Q}(s, a)$ at random, put the agent in a random state s

Repeat:

- 1) Select the action $\operatorname{argmax}_a \hat{Q}(s, a)$ with probability $(1 - \epsilon)$ otherwise, select a at random
- 2) The agent is now in state s' and has received the reward r
- 3) Update $\hat{Q}(s, a)$ by

$$\Delta \hat{Q}(s, a) = \alpha [r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)]$$

CAREFUL

maximizing $\hat{Q}(s, a)$ when this is a deep neural network may be non-trivial...

Reinforcement Learning Reformulation

Reinforcement Learning Reformulation

Trajectory

$$\tau := \langle (s_t, a_t) \rangle_{t=0}^T$$

i.e., a sequence of states and actions.

It can be either finite or infinite, depending on T

Reward

Reward function:

$$r_t := r(s_t, a_t, s_{t+1})$$

Depending on the application, it can be simplified:

$$r_t := r(s_t, a_t), \quad r_t := r(s_t)$$

Return

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

we will use these forms from now on, for brevity

It is discounted when $\gamma < 1$ or undiscounted, when $\gamma = 1$ (and the trajectory is finite)

Reinforcement Learning Reformulation

Probability of a trajectory

$$P(\tau|\pi) := \underbrace{P(s_0)}_{\text{probability of initial states}} \prod_{t=0}^{T-1} \underbrace{P(s_{t+1}|s_t, a_t)}_{\text{transition probability (i.e. the 'model')}} \pi(a_t|s_t)$$

Expected return of a policy

$$J(\pi) := \int_{\tau \sim \pi} P(\tau|\pi) R(\tau) = \mathbb{E}_{\tau \sim \pi} [R(\tau)]$$

where $\tau \sim \pi$ is the space of all the trajectories such that $a_t = \pi(s_t)$

*assuming that the policy is deterministic
(but see after)*

Central RL Problem

$$\pi^* := \operatorname{argmax}_{\pi} J(\pi)$$

i.e. finding the policy with the highest expected return

Reinforcement Learning Reformulation

Value Function (of a policy)

$$V^\pi(s) := \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s]$$

Action-Value function (of a policy)

$$Q^\pi(s, a) := \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s, a_0 = a]$$

Optimal Value Function

$$V^*(s) := \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s]$$

Optimal Action-Value Function

$$Q^*(s, a) := \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s, a_0 = a]$$

Reinforcement Learning Reformulation

■ Connecting Value and Action-Value Functions

$$V^\pi(s) = \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)]$$

$$V^*(s) = \max_a [Q^*(s, a)]$$

■ Optimal Policy

$$a^*(s) = \operatorname{argmax}_a [Q^*(s, a)]$$

■ Advantage Function

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$

It tells how advantageous (or disadvantageous) is a particular action w.r.t. what is prescribed by the policy

DQN Algorithm

Deep Q-Learning

■ Playing Atari with Deep Reinforcement Learning

[2013, V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, M. Riedmiller, <http://arxiv.org/abs/1312.5602>, see also <http://www.nature.com/nature/journal/v518/n7540/full/nature14236.html>]

A software system only

Runs on virtually any Linux-based system, it contains optional provisions for GPU

It's open source

<https://github.com/kuz/DeepMind-Atari-Deep-Q-Learner>

Sophisticated machine-learning techniques

Uses deep reinforcement learning
in combination with convolutional neural networks (CNN)

Same configuration, multiple games

Same configuration applied to arcade games
It learned to play 7 (2013) or 49 (2015) different games

It is autonomous

It learns by itself, it receives no human expertise as input
In many cases, it outperforms human players



(from GitHub)

Deep Q-Learning

- **DQN Algorithm** [<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>]

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Policy Gradient

Policy Gradient

Parametric Policy

A generic policy that depends on parameters θ

$$\pi_{\theta}$$

For instance, in the **DQN Algorithm**, the **Action-Value Function** is approximator is a Deep Neural Network

$$\hat{Q}(s, a; \theta)$$

Policy Gradient Ascent

At each iteration, improve parameters by:

$$\theta^{(k+1)} = \theta^{(k)} + \eta \nabla_{\theta} J(\pi_{\theta}) |_{\theta^{(k)}}$$

easier said than done ...

Policy Gradient

1) Probability of a trajectory, given a parametric policy

$$P(\tau|\pi_\theta) := P(s_0) \prod_{t=0}^{T-1} P(s_{t+1}|s_t, a_t) \pi_\theta(a_t|s_t)$$

2) Log-Derivative

By applying the chain rule:

$$\nabla_\theta \log P(\tau|\pi_\theta) = \frac{1}{P(\tau|\pi_\theta)} \nabla_\theta P(\tau|\pi_\theta)$$

It follows:

$$\nabla_\theta P(\tau|\pi_\theta) = P(\tau|\pi_\theta) \nabla_\theta \log P(\tau|\pi_\theta)$$

Policy Gradient

3) Log-Probability

$$\log P(\tau|\pi_\theta) := \log P(s_0) + \sum_{t=0}^{T-1} [\log P(s_{t+1}|s_t, a_t) + \log \pi_\theta(a_t|s_t)]$$

these terms do NOT depend on θ

4) Gradient of the Log-Probability

$$\nabla_\theta \log P(\tau|\pi_\theta) := \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t)$$

5) Expected return

$$J(\pi_\theta) := \int_{\tau \sim \pi_\theta} P(\tau|\pi_\theta) R(\tau) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$$

Policy Gradient

Basic Policy Gradient

$$\begin{aligned}\nabla_{\theta} J(\pi_{\theta}) &= \int_{\tau \sim \pi_{\theta}} \nabla_{\theta} P(\tau | \pi_{\theta}) R(\tau) \\ &= \int_{\tau \sim \pi_{\theta}} P(\tau | \pi_{\theta}) \nabla_{\theta} \log P(\tau | \pi_{\theta}) R(\tau) \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log P(\tau | \pi_{\theta}) R(\tau)] \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau) \right]\end{aligned}$$

this terms do NOT depend on θ

This last term is an expectation: it can be estimated with a sample mean

Policy Gradient

Basic Policy Gradient

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau) \right]$$

$$\hat{g} := \frac{1}{|\mathcal{D}|} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau)$$

estimated gradient (mean)

dataset: a sample of actual trajectories

Policy Gradient

Basic Policy Gradient

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau) \right]$$

an entire trajectory? even in the past?

More precisely:

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right]$$

*reward from t onward
(‘reward-to-go’)*

Simple Policy Gradient

▪ Pseudo-Algorithm

Initialize the weights θ of a DNN $\hat{Q}(s, a; \theta)$ at random

Repeat:

1) For M episodes

Start in initial state s_0

For t from 0 to T

play by $a_t \sim \pi_\theta(a|s_t)$

How can we 'sample a policy' in practice?

Collect the episode trajectory $\tau = \langle (s_t, a_t) \rangle_{t=0}^T$ and store it in \mathcal{D}

2) Sample a random minibatch $\mathcal{B} = \{(s_i, a_i)\}$ from \mathcal{D}

$$\Delta\theta = \eta \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{D}} \nabla_\theta \log \pi_\theta(a_t|s_t) R(\tau)$$

Sampling a Policy

Problem

Sampling actions from a stochastic policy

$$a_t \sim \pi_\theta(a|s_t)$$

Intended meaning:

$$\pi_\theta(a_t|s_t) \propto \hat{Q}(a_t, s_t; \theta)$$

the probability of each action should be proportional to the expected return

Discrete Action Space

Consider $\hat{Q}(a_t, s_t; \theta)$ as the **logit** of a softmax

$$\pi_\theta(a_t|s_t) := \frac{\exp(\hat{Q}(a_t, s_t; \theta))}{\sum_{a \in \mathcal{A}(s_t)} \exp(\hat{Q}(a, s_t; \theta))}$$

and sample accordingly

all possible actions in state s_t

The Continuous Case is a bit more complex ...

Actor-Critic

An Aside:

EGLP Lemma. Suppose that P_θ is a parameterized probability distribution over a random variable, x . Then:

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0.$$

Proof

Recall that all probability distributions are *normalized*:

$$\int_x P_\theta(x) = 1.$$

Take the gradient of both sides of the normalization condition:

$$\nabla_\theta \int_x P_\theta(x) = \nabla_\theta 1 = 0.$$

Use the log derivative trick to get:

$$\begin{aligned} 0 &= \nabla_\theta \int_x P_\theta(x) \\ &= \int_x \nabla_\theta P_\theta(x) \\ &= \int_x P_\theta(x) \nabla_\theta \log P_\theta(x) \\ \therefore 0 &= \mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)]. \end{aligned}$$

Actor-Critic

Policy Gradient

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right]$$

Due to the EGLP lemma:

$$\mathbb{E}_{a_t \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t)] = 0$$

for any function $b(s_t)$ that depends on s_t only

Policy Gradient with Baseline

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\left(\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) - b(s_t) \right) \right]$$

baseline ↘

We can subtract term-wise any function $b(s_t)$ without altering the expectation

Actor-Critic

Actor-Critic

(typical formulation)

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\left(\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) - V^{\pi}(s_t) \right) \right]$$

Note that:

$$\begin{aligned} \left(\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) - V^{\pi}(s_t) &= (r(s_t, a_t) + V^{\pi}(s_{t+1})) - V^{\pi}(s_t) \\ &= Q^{\pi}(s_t, a_t) - V^{\pi}(s_t) \\ &= A^{\pi}(s_t, a_t) \end{aligned}$$

it's the **advantage function**

Actor-Critic

Actor-Critic

(typical formulation)

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A^{\pi}(s_t, a_t) \right]$$

'Actor' 'Critic'

In practice, $V^{\pi}(s_t)$ is estimated via $\hat{V}(s; \phi)$
namely, another DNN with parameters ϕ from which

$$\hat{A}(s_t, a_t) := \left(r(s_t, a_t) + \hat{V}(s_{t+1}; \phi) \right) - \hat{V}(s_t; \phi)$$

What are the advantages? "It reduces variance"

Intuitively $\hat{Q}(s, a; \theta)$ depends also on how the action space is explored
whereas $\hat{V}(s_t; \phi)$ depends only on actual rewards $r(s_t, a_t)$

Actor-Critic

▪ Pseudo-Algorithm

Initialize the weights θ, ϕ of two DNNs $\pi_\theta(a|s), \hat{V}(s; \phi)$ at random

Repeat:

1) For M episodes

Start in initial state s_0

For t from 0 to T

play by $a_t \sim \pi_\theta(a|s_t)$

Collect all episode **transitions** $\tau_r := \langle (s_t, a_t, r_t, s_{t+1}) \rangle_{t=0}^T$ and store them in \mathcal{D}

2) For a random minibatch $\mathcal{B} = \{(s_i, a_i, r_i, s_{i+1})\}$ from \mathcal{D}

Evaluate

$$\hat{A}(s_i, a_i) = \left(r_i + \hat{V}(s_{i+1}, \phi) \right) - \hat{V}(s_i, \phi)$$

Update weights

$$\Delta\phi = -\eta_\phi \nabla_\phi \left(\hat{A}(s_i, a_i) \right)^2$$

$$\Delta\theta = \eta_\theta \nabla_\theta J(\pi_\theta) = \eta_\theta \nabla_\theta \log \pi_\theta(a_i|s_i) \hat{A}(s_i, a_i)$$

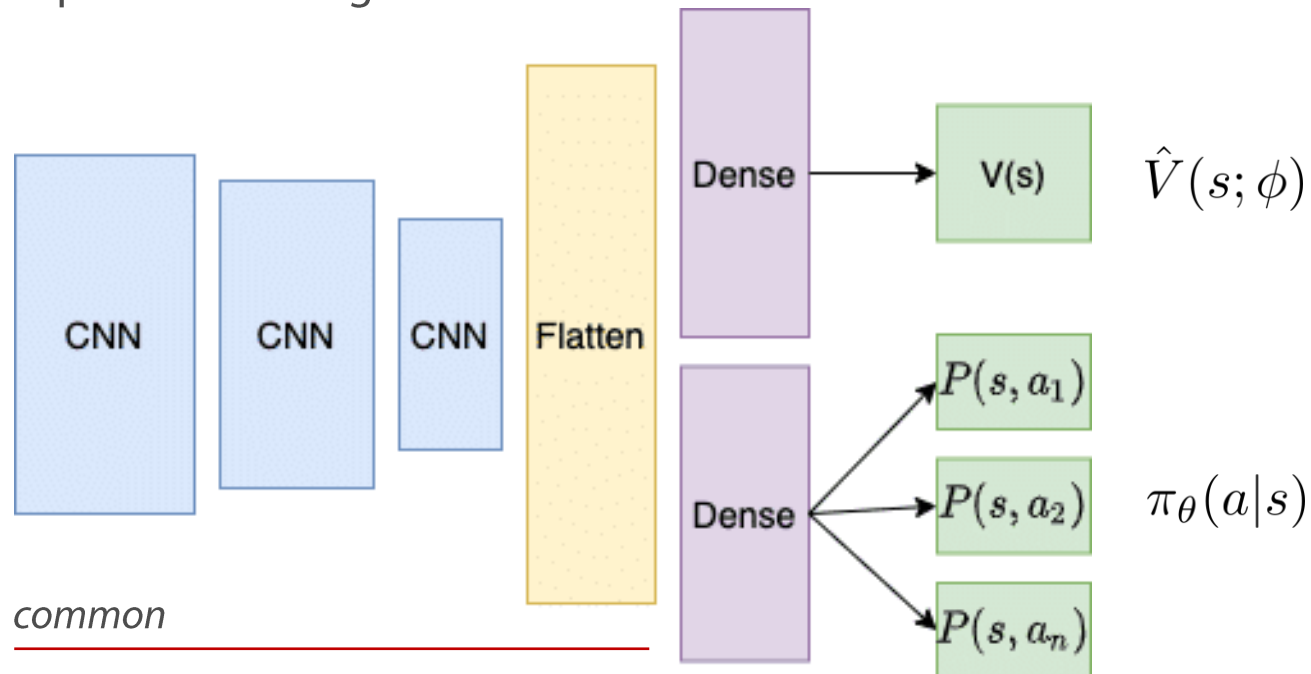
Actor-Critic

Network Architecture

A bifurcated structure which includes:

- A common part
- A V-head
- A π -head

It follows that part of the weights are *shared*



Normalized Advantage Function (NAF)

Normalized Advantage Function (NAF) algorithm

S. Gu, T. P. Lillicrap, I. Sutskever, S. Levine.
**Continuous deep Q-learning
with model-based acceleration**, 2016

Algorithm 1.2 NAF algorithm for continuous Q-learning

Randomly initialize $\tilde{Q}(s, a | \theta_{\text{PRED}}^Q)$ $\theta^Q := (\theta^\mu, \theta^P, \theta^V)$

Initialize the target network with $\theta_{\text{TAR}}^Q \leftarrow \theta_{\text{PRED}}^Q$

Initialize replay buffer $R \leftarrow 0$

for each episode do:

 Initialize random process \mathcal{N} for action exploration

$s_0 \leftarrow \text{Environment}(\text{reset})$

for $t = 0$ to T do:

$a_t \leftarrow \mu(s_t | \theta_{\text{PRED}}^\mu) + \mathcal{N}_t$

$r_t \leftarrow r(s_t, a_t)$

$s_{t+1} \leftarrow \text{Environment}(s_t, a_t)$

$RB \leftarrow RB \cup \{(s_t, a_t, r_t, s_{t+1})\}$ store transition in the replay buffer

 Sample at random and normalize the mini batch MB

for each sample $i = (s_i, a_i, r_i, s_{i+1})$ in m

$y_i = r_i + \gamma \tilde{V}(s_{i+1} | \theta_{\text{TAR}}^V)$

 Compute gradients

$\frac{\partial}{\partial \theta^Q} \left(y_i - Q(s_i, a_i | \theta_{\text{PRED}}^Q) \right)^2$ (Loss function $L(\theta^Q)$)

$\theta_{\text{PRED}}^Q \leftarrow \theta_{\text{PRED}}^Q - \eta \left(\frac{\partial}{\partial \theta^Q} L(\theta^Q) \right)$

$\theta_{\text{TAR}}^Q \leftarrow \tau \theta_{\text{PRED}}^Q + (1 + \tau) \theta_{\text{TAR}}^Q$

end for

end for

end for

Normalized Advantage Function (NAF) algorithm

Algorithm Highlights

- a deep neural network for $\hat{Q}(s, a)$

Algorithm 1.2 NAF algorithm for continuous Q-learning

Randomly initialize $\tilde{Q}(s, a | \theta_{\text{PRED}}^Q)$ $\theta^Q := (\theta^\mu, \theta^P, \theta^V)$

Initialize the target network with $\theta_{\text{TAR}}^Q \leftarrow \theta_{\text{PRED}}^Q$

Initialize replay buffer $R \leftarrow 0$

for each episode do:

 Initialize random process \mathcal{N} for action exploration

$s_0 \leftarrow \text{Environment}(\text{reset})$

for $t = 0$ to T do:

$a_t \leftarrow \mu(s_t | \theta_{\text{PRED}}^\mu) + \mathcal{N}_t$

$r_t \leftarrow r(s_t, a_t)$

$s_{t+1} \leftarrow \text{Environment}(s_t, a_t)$

$RB \leftarrow RB \cup \{(s_t, a_t, r_t, s_{t+1})\}$ store transition in the replay buffer

 Sample at random and normalize the mini batch MB

for each sample $i = (s_i, a_i, r_i, s_{i+1})$ in m

$y_i = r_i + \gamma \tilde{V}(s_{i+1} | \theta_{\text{TAR}}^V)$

 Compute gradients

$\frac{\partial}{\partial \theta^Q} \left(y_i - Q(s_i, a_i | \theta_{\text{PRED}}^Q) \right)^2$ (Loss function $L(\theta^Q)$)

$\theta_{\text{PRED}}^Q \leftarrow \theta_{\text{PRED}}^Q - \eta \left(\frac{\partial}{\partial \theta^Q} L(\theta^Q) \right)$

$\theta_{\text{TAR}}^Q \leftarrow \tau \theta_{\text{PRED}}^Q + (1 + \tau) \theta_{\text{TAR}}^Q$

end for

end for

end for

Normalized Advantage Function (NAF) algorithm

Algorithm Highlights

- a deep neural network for $\hat{Q}(s, a)$
- two deep networks:
one TARget, which is the objective
and one PREDictor for transient
approximations

Algorithm 1.2 NAF algorithm for continuous Q-learning

Randomly initialize $\tilde{Q}(s, a | \theta_{\text{PRED}}^Q)$ $\theta^Q := (\theta^\mu, \theta^P, \theta^V)$

Initialize the target network with $\theta_{\text{TAR}}^Q \leftarrow \theta_{\text{PRED}}^Q$

~~Initialize replay buffer $R \leftarrow \emptyset$~~

for each episode do:

 Initialize random process \mathcal{N} for action exploration

$s_0 \leftarrow \text{Environment}(\text{reset})$

 for $t = 0$ to T do:

$a_t \leftarrow \mu(s_t | \theta_{\text{PRED}}^\mu) + \mathcal{N}_t$

$r_t \leftarrow r(s_t, a_t)$

$s_{t+1} \leftarrow \text{Environment}(s_t, a_t)$

$RB \leftarrow RB \cup \{(s_t, a_t, r_t, s_{t+1})\}$ store transition in the replay buffer

 Sample at random and normalize the mini batch MB

 for each sample $i = (s_i, a_i, r_i, s_{i+1})$ in m

$y_i = r_i + \gamma \tilde{V}(s_{i+1} | \theta_{\text{TAR}}^V)$

 Compute gradients

$\frac{\partial}{\partial \theta^Q} \left(y_i - Q(s_i, a_i | \theta_{\text{PRED}}^Q) \right)^2$ (Loss function $L(\theta^Q)$)

$\theta_{\text{PRED}}^Q \leftarrow \theta_{\text{PRED}}^Q - \eta \left(\frac{\partial}{\partial \theta^Q} L(\theta^Q) \right)$

$\theta_{\text{TAR}}^Q \leftarrow \tau \theta_{\text{PRED}}^Q + (1 + \tau) \theta_{\text{TAR}}^Q$

 end for

end for

end for

Normalized Advantage Function (NAF) algorithm

Algorithm Highlights

- a deep neural network for $\hat{Q}(s, a)$
- two deep networks:
one TARget, which is the objective
and one PREDictor for transient
approximations
- careful tensorial formulation —————
to avoid the argmax step (see after)

Algorithm 1.2 NAF algorithm for continuous Q-learning

Randomly initialize $\tilde{Q}(s, a | \theta_{\text{PRED}}^Q)$ $\theta^Q := (\theta^\mu, \theta^P, \theta^V)$
Initialize the target network with $\theta_{\text{TAR}}^Q \leftarrow \theta_{\text{PRED}}^Q$
Initialize replay buffer $R \leftarrow 0$
for each episode do:
 Initialize random process \mathcal{N} for action exploration
 $s_0 \leftarrow \text{Environment}(\text{reset})$
 for $t = 0$ to T do:
 $a_t \leftarrow \mu(s_t | \theta_{\text{PRED}}^\mu) + \mathcal{N}_t$
 $r_t \leftarrow r(s_t, a_t)$
 $s_{t+1} \leftarrow \text{Environment}(s_t, a_t)$
 $RB \leftarrow RB \cup \{(s_t, a_t, r_t, s_{t+1})\}$ store transition in the replay buffer
 Sample at random and normalize the mini batch MB
 for each sample $i = (s_i, a_i, r_i, s_{i+1})$ in m
 $y_i = r_i + \gamma \tilde{V}(s_{i+1} | \theta_{\text{TAR}}^V)$
 Compute gradients
 $\frac{\partial}{\partial \theta^Q} \left(y_i - Q(s_i, a_i | \theta_{\text{PRED}}^Q) \right)^2$ (Loss function $L(\theta^Q)$)
 $\theta_{\text{PRED}}^Q \leftarrow \theta_{\text{PRED}}^Q - \eta \left(\frac{\partial}{\partial \theta^Q} L(\theta^Q) \right)$
 $\theta_{\text{TAR}}^Q \leftarrow \tau \theta_{\text{PRED}}^Q + (1 + \tau) \theta_{\text{TAR}}^Q$
 end for
 end for
end for

Normalized Advantage Function (NAF) algorithm

Algorithm Highlights

- a deep neural network for $\hat{Q}(s, a)$
- two deep networks:
one TARget, which is the objective
and one PREDictor for transient
approximations
- careful tensorial formulation
to avoid the argmax step (see after)
- noise based on a stochastic process
(i.e. a random walk, see later)
forcing **exploration**

Algorithm 1.2 NAF algorithm for continuous Q-learning

Randomly initialize $\tilde{Q}(s, a | \theta_{\text{PRED}}^Q)$ $\theta^Q := (\theta^\mu, \theta^P, \theta^V)$

Initialize the target network with $\theta_{\text{TAR}}^Q \leftarrow \theta_{\text{PRED}}^Q$

Initialize replay buffer $R \leftarrow 0$

for each episode do:

 Initialize random process \mathcal{N} for action exploration

$s_0 \leftarrow \text{Environment}(\text{reset})$

for $t = 0$ to T do:

$a_t \leftarrow \mu(s_t | \theta_{\text{PRED}}^\mu) + \mathcal{N}_t$

$r_t \leftarrow r(s_t, a_t)$

$s_{t+1} \leftarrow \text{Environment}(s_t, a_t)$

$RB \leftarrow RB \cup \{(s_t, a_t, r_t, s_{t+1})\}$ store transition in the replay buffer

 Sample at random and normalize the mini batch MB

for each sample $i = (s_i, a_i, r_i, s_{i+1})$ in m

$y_i = r_i + \gamma \tilde{V}(s_{i+1} | \theta_{\text{TAR}}^V)$

 Compute gradients

$\frac{\partial}{\partial \theta^Q} \left(y_i - Q(s_i, a_i | \theta_{\text{PRED}}^Q) \right)^2$ (Loss function $L(\theta^Q)$)

$\theta_{\text{PRED}}^Q \leftarrow \theta_{\text{PRED}}^Q - \eta \left(\frac{\partial}{\partial \theta^Q} L(\theta^Q) \right)$

$\theta_{\text{TAR}}^Q \leftarrow \tau \theta_{\text{PRED}}^Q + (1 + \tau) \theta_{\text{TAR}}^Q$

end for

end for

end for

Normalized Advantage Function (NAF) algorithm

Algorithm Highlights

- a deep neural network for $\hat{Q}(s, a)$
- two deep networks:
one TARget, which is the objective
and one PREDictor for transient
approximations
- careful tensorial formulation
to avoid the argmax step (see after)
- noise based on a stochastic process
(i.e. a random walk, see later)
forcing **exploration**
- **replay buffer** with random extraction
of **mini-batches** to avoid temporal
correlation arising from sequential
exploration

Algorithm 1.2 NAF algorithm for continuous Q-learning

Randomly initialize $\tilde{Q}(s, a | \theta_{\text{PRED}}^Q)$ $\theta^Q := (\theta^\mu, \theta^P, \theta^V)$

Initialize the target network with $\theta_{\text{TAR}}^Q \leftarrow \theta_{\text{PRED}}^Q$

Initialize replay buffer $R \leftarrow 0$

for each episode do:

 Initialize random process \mathcal{N} for action exploration

$s_0 \leftarrow \text{Environment}(\text{reset})$

for $t = 0$ to T do:

$a_t \leftarrow \mu(s_t | \theta_{\text{PRED}}^\mu) + \mathcal{N}_t$

$r_t \leftarrow r(s_t, a_t)$

$s_{t+1} \leftarrow \text{Environment}(s_t, a_t)$

$RB \leftarrow RB \cup \{(s_t, a_t, r_t, s_{t+1})\}$ store transition in the replay buffer

 Sample at random and normalize the mini batch MB

for each sample $i = (s_i, a_i, r_i, s_{i+1})$ in m

$y_i = r_i + \gamma \tilde{V}(s_{i+1} | \theta_{\text{TAR}}^V)$

 Compute gradients

$\frac{\partial}{\partial \theta^Q} \left(y_i - Q(s_i, a_i | \theta_{\text{PRED}}^Q) \right)^2$ (Loss function $L(\theta^Q)$)

$\theta_{\text{PRED}}^Q \leftarrow \theta_{\text{PRED}}^Q - \eta \left(\frac{\partial}{\partial \theta^Q} L(\theta^Q) \right)$

$\theta_{\text{TAR}}^Q \leftarrow \tau \theta_{\text{PRED}}^Q + (1 + \tau) \theta_{\text{TAR}}^Q$

end for

end for

end for

Normalized Advantage Function (NAF) algorithm

Algorithm Highlights

- a deep neural network for $\hat{Q}(s, a)$
- two deep networks:
one TARget, which is the objective
and one PREDictor for transient
approximations
- careful tensorial formulation
to avoid the argmax step (see after)
- noise based on a stochastic process
(i.e. a random walk, see later)
forcing **exploration**
- **replay buffer** with random extraction
of **mini-batches** to avoid temporal
correlation arising from sequential
exploration
- **Can cope with continuous \mathcal{A} and \mathcal{S}**

Algorithm 1.2 NAF algorithm for continuous Q-learning

Randomly initialize $\tilde{Q}(s, a | \theta_{\text{PRED}}^Q)$ $\theta^Q := (\theta^\mu, \theta^P, \theta^V)$
Initialize the target network with $\theta_{\text{TAR}}^Q \leftarrow \theta_{\text{PRED}}^Q$
Initialize replay buffer $R \leftarrow 0$
for each episode do:
 Initialize random process \mathcal{N} for action exploration
 $s_0 \leftarrow \text{Environment}(\text{reset})$
 for $t = 0$ to T do:
 $a_t \leftarrow \mu(s_t | \theta_{\text{PRED}}^\mu) + \mathcal{N}_t$
 $r_t \leftarrow r(s_t, a_t)$
 $s_{t+1} \leftarrow \text{Environment}(s_t, a_t)$
 $RB \leftarrow RB \cup \{(s_t, a_t, r_t, s_{t+1})\}$ store transition in the replay buffer
 Sample at random and normalize the mini batch MB
 for each sample $i = (s_i, a_i, r_i, s_{i+1})$ in m
 $y_i = r_i + \gamma \tilde{V}(s_{i+1} | \theta_{\text{TAR}}^V)$
 Compute gradients
 $\frac{\partial}{\partial \theta^Q} \left(y_i - Q(s_i, a_i | \theta_{\text{PRED}}^Q) \right)^2$ (Loss function $L(\theta^Q)$)
 $\theta_{\text{PRED}}^Q \leftarrow \theta_{\text{PRED}}^Q - \eta \left(\frac{\partial}{\partial \theta^Q} L(\theta^Q) \right)$
 $\theta_{\text{TAR}}^Q \leftarrow \tau \theta_{\text{PRED}}^Q + (1 + \tau) \theta_{\text{TAR}}^Q$
 end for
 end for
end for

Normalized Advantage Function (NAF) algorithm

▪ A special approximator

NOTE: all functions here are **continuous**
and of **vector** parameters

From the definition of the **Advantage Function**

$$A^\pi(s, a) := Q^\pi(s, a) - V^\pi(s)$$

The NAF approximator becomes:

$$\hat{Q}(s, a) := \hat{A}(s, a; \theta) - \hat{V}(s; \phi)$$

Define:

$$\hat{A}(s, a; \theta) = \frac{1}{2}(a - \mu(s; \theta_\mu))^T P(s; \theta_P)(a - \mu(s; \theta_\mu))$$

μ, P are Deep Neural Networks

Then the solution to

this is a quadratic form

$$\frac{\partial}{\partial a} \hat{Q}(s, a) = 0 \quad \iff \quad \frac{\partial}{\partial a} \hat{A}(s, a; \theta) = 0$$

is

$$a^* = \mu(s; \theta_\mu)$$