

# A Brief Survey of Reinforcement Learning

Short walk-through on building learning agents

---

Giancarlo Frison

December 11, 2018

# Table of contents i

---

1. Introduction
2. Genetic Programming
3. Multi-armed Bandit
4. Markov Decision Process
5. Neural Networks in RL

## Table of contents ii

---

6. Actor-Critic

7. Deep Deterministic Policy Gradient

8. Asgard

# Introduction

---

# What is Reinforcement Learning



[9]

REINFORCEMENT LEARNING is an area of machine learning concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.

# Challenges in RL

RL vision is about creating systems that are capable of learning how to adapt in the real world, solely based on trial-and-error. Challenges:

## Reward signal

The optimal policy must be inferred by interacting with the environment. The only learning signal the agent receives is the reward.

## Credit assignment problem

Agents must deal with long-range time dependencies: Often the consequences of an action only materialise after many transitions of the environment [8].

## Non stationary environments

Especially in multi-agent environments, in the same circumstances of state  $s_t$  and action  $a_t$  outcomes might be different.

## Difference from other optimizations

---

There are major differences within:

- Supervised learning

## Difference from other optimizations

---

There are major differences within:

- Supervised learning
- Mathematical optimization

## Difference from other optimizations

---

There are major differences within:

- Supervised learning
- Mathematical optimization
- **Genetic programming**

# Genetic Programming

---

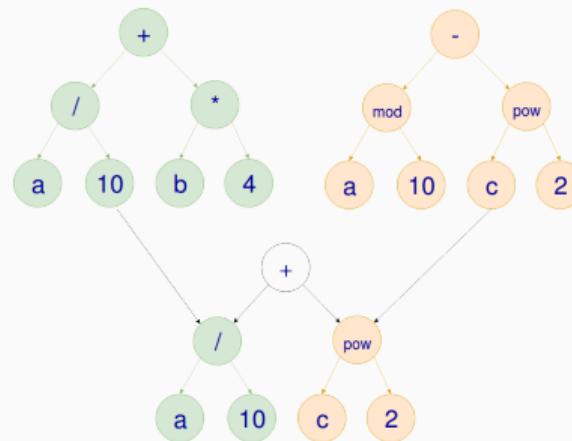
# Genetic Programming

---

GP is a technique whereby computer programs are encoded as a set of genes that are then modified using an evolutionary algorithm.

- A number of chromosomes are randomly created.
- Each chromosome is evaluated through a fitness function.
- Best ones are selected, the others are disposed.
- Chromosomes could be breded among the selected for a new generation.
- Offsprings are randomly mutated.
- Repeat until the score threshold is reached [3].

# Crossover



**Figure 1:** The “breeding” is called crossover. The chromosome (in this case an AST), is merged between two individuals for searching a better function.

## Strengths on many local minima

GP may overtake gradient-based algorithms when the solution space has many local minima

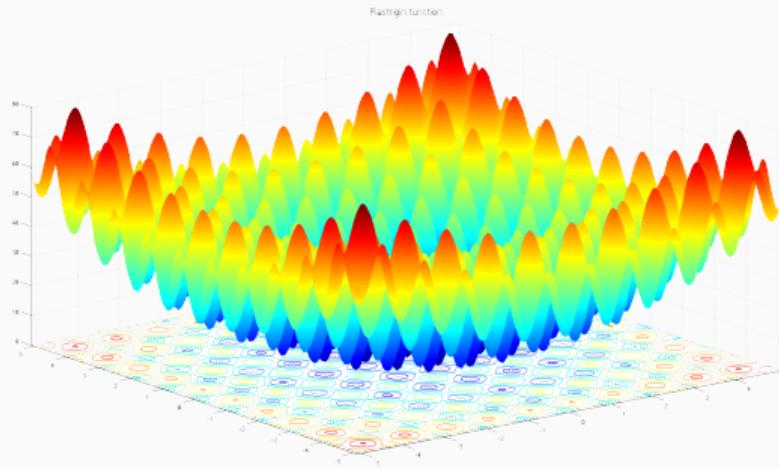


Figure 2: Rastrigin function

## Strengths on no-differentiable functions

$f(x)$  is DIFFERENTIABLE when it has always a finite derivative along the domain.

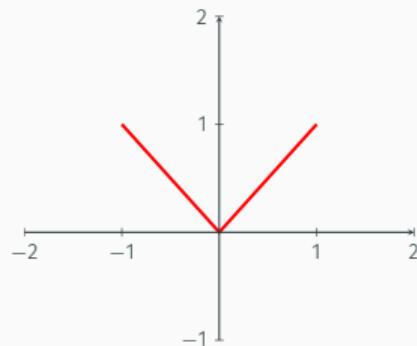


Figure 3:  $f(x) = |x|$  has no derivative at  $x = 0$ .

GP is tollerant to *latent* no-differentiable functions.

## Multi-armed Bandit

---

# Multi-armed Bandit



The multi-armed bandit problem has been the subject of decades of intense study in statistics, operations research, A/B testing, computer science, and economics [12]

# $Q$ Value's Action

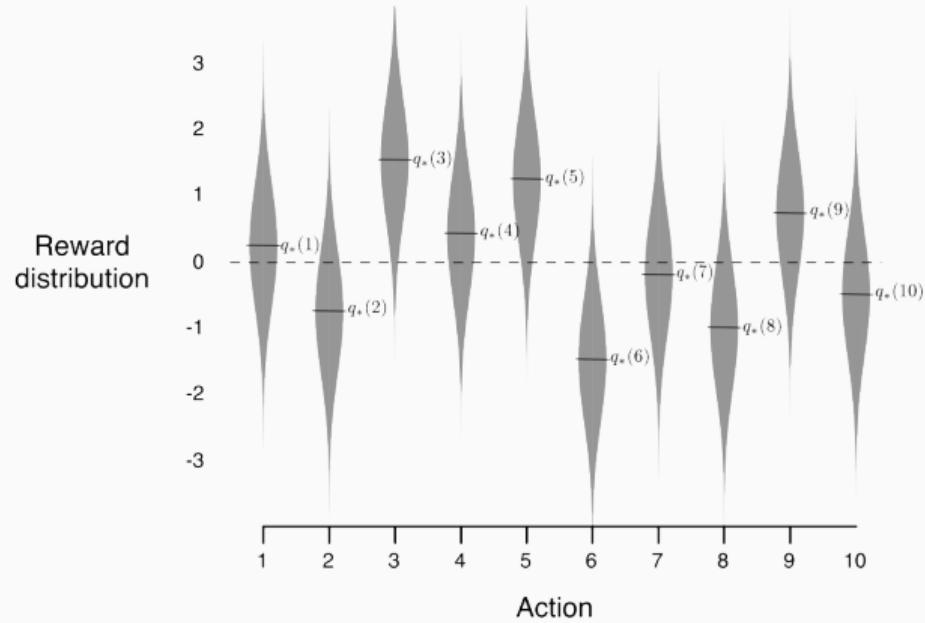


Figure 4: Obtained measures after repeated pullings with 10 arms [8].

## *Q* Value's Action

$Q_n$  is the estimated value of its action after  $n$  selections.

$$Q_{n+1} = \frac{R_1 + R_2 + \dots + R_n}{n}$$

A more scalable formula, updates the average with incremental and small constant:

$$Q_{n+1} = Q_n + \frac{1}{n}(R_n - Q_n)$$

General expression of the bandit algorithm, recurrent pattern in RL (*Target* could be considered the reward  $R$  by now).

$$\text{NewEstimate} = \text{OldEstimate} + \text{StepSize}(\text{Target} - \text{OldEstimate})$$

# Gambler's Dilemma

When pulled, an arm produces a random payout drawn independently of the past. Because the distribution of payouts corresponding to each arm is not listed, the player can learn it only by experimenting.

## Exploitation

Earn more money by exploiting arms that yielded high payouts in the past.

## Exploration

Exploring alternative arms may return higher payouts in the future.

There is a *tradeoff* between EXPLORATION and REGRET.

## Markov Decision Process

---

## Definition of MDP



Figure 5: 2015 DARPA Robotics Challenge [6]

Despite as in bandits, MDP formalizes the decision making (Policy  $\pi$ ) in sequential steps, aggregated in Episodes.

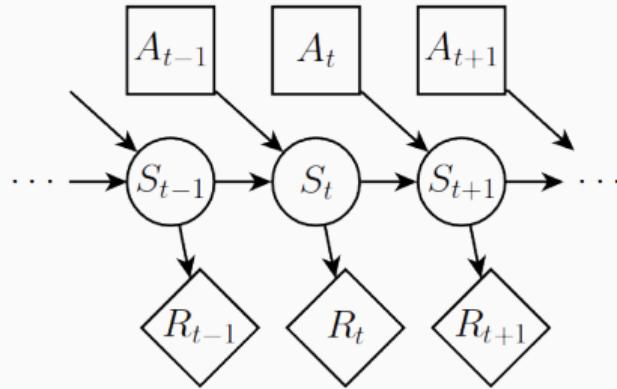


Figure 6: Model representation of MDP [5]

MDP strives to find the best  $\pi$  to all possible states. In Markov processes, the selected action depends **only** on current state.

## Bellman equation

In a multi-step process the total return  $G$  is the *discounted* sum of rewards, from the begin state to the terminal state:

$$G_\pi = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \gamma^n r_n$$

BELLMAN EQUATION: The value of a state  $V_s$  is the *averaged* reward obtained in that state plus the sum of all values following  $\pi$  thereafter.

The goal of RL is to find the *optimal*  $\pi$  with highest  $G$

$$\pi = \operatorname{argmax}_s r + \gamma v_\pi(s)$$

# Grid World

☒	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	☒

- $A = \text{up, down, right, left}$
- Terminal states are the flagged boxes.
- $R_t = -1$  for all states.

The problem is to define the best  $\pi$ . Value function is computed by iterative policy evaluation.

# Iteration 1

Calculated  $V_1$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

Policy  $\pi_1$

FLAG	←	?	?
↑	?	?	?
?	?	?	↓
?	?	→	FLAG

## Iteration 2

Calculated  $V_2$

0	-1.7	-2	-2
-1.7	-2	-2	-2
-2	-2	-2	-1.7
-2	-2	-1.7	0

Policy  $\pi_2$

	←	←	?
↑	↔	?	↓
↑	?	↶	↓
?	→	→	

## Iteration 3

Calculated  $V_3$

0	-2.4	-2.9	-3
-2.4	-2.9	-3	-2.9
-2.9	-3	-2.9	-2.4
-3	-2.9	-2.4	0

Policy  $\pi_3$

	←	←	↖
↑	↑↖	↖	↓
↑	↖	↖	↓
↖	→	→	

## Model-free learning

The learning comes only through raw experiences (*episodes*) when there is no prior knowledge of the environment.

## Bootstrapping

Value estimations are based partly on other learned estimations.

Bootstrapping helps to increase the stability of estimations.

$$V(s_t) = V(s_t) + \alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

## Neural Networks in RL

---

## Beyond the Gridworld

In the Gridworld every state value  $v_t$  is stored in a table. The approach lacks scalability and is inherently limited to fairly low-dimension problems.



Figure 7: The state  $v_t$  might be a frame of a videogame [11]

# Deep Reinforcement Learning

---

Deep learning enables RL to scale to decision-making problems that were previously intractable, i.e., settings with high-dimensional state and action spaces. [1]

- **Universal function approximator** Multilayer neural networks can approximate any continuous function to any degree of accuracy [4].

# Deep Reinforcement Learning

---

Deep learning enables RL to scale to decision-making problems that were previously intractable, i.e., settings with high-dimensional state and action spaces. [1]

- **Universal function approximator** Multilayer neural networks can approximate any continuous function to any degree of accuracy [4].
- **Representation learning** Automatically discover the representations needed for feature detection or classification from raw data, also known as FEATURE ENGINEERING [2].

Neural networks can learn the approximated value of  $V_s$  or  $Q(s, a)$  for any given STATE/ACTION SPACE.

## Discrete space

When the ACTION SPACE is limited by a set of options, (e.g. *turn right, left*) the output layer could be a set of neurons as the number of possible actions.

A1	A2	A3	A4
0.21	0.27	<b>0.30</b>	0.21

**Figure 8:** Output layer with softmax

The choosed action is the one with the biggest  $Q_t$  (A3).

## Exploration in discrete space

---

The balance between EXPLORATION and EXPLOITATION is one of the most challenging task in RL. The most common is the  $\epsilon$ -greedy:

$$\pi(s) = \begin{cases} \text{random action,} & \text{if } \xi < \epsilon. \\ \text{argmax}Q(s, a), & \text{otherwise.} \end{cases} \quad (1)$$

With  $\epsilon$ -GREEDY, at each time step, the agent selects a random action with a fixed probability,  $0 \leq \epsilon \leq 1$ , instead of selecting the learned optimal action.

# Exploration in discrete space

$\epsilon$ -GREEDY must be tuned and it is fixed during the training. In contrast, VALUE-DIFFERENCE BASED EXPLORATION adapts the exploration-probability based on fluctuations of confidence [10].

A1	A2	A3	A4
0.21	0.27	<b>0.30</b>	0.21

**Figure 9:** Low confidence, flat likelihood

A1	A2	A3	A4
0.21	0.07	<b>0.70</b>	0.01

**Figure 10:** High confidence, sharp likelihood

Applications of RL require continuous state spaces defined by means of continuous variables *e.g. position, velocity, torque.*

## Single Output

Like in DDPG algorithm. The NN yields only one output.

## Probability distribution

The network yields a distribution. For convenience usually it is a NORMAL DISTRIBUTION defined by 2 outputs:  $\mu$  and  $\sigma$ .

NN comes into play by approximating state-action value pairs  $Q_t(s, a)$ .

## Double network

Other the main network, a sibling network (TARGET) is used in training for adjusting values  $Q$ , the second network is periodically aligned the the main network.

## Priority experience replay

Past state/action/reward tuples are kept also for future iterations. Tuples with higher error in evaluation are retained longer.

► First disruptive advance in DEEP RL in 2015 by DeepMind [7].

## Actor-Critic

---

## Policy-based algorithm

---

ACTOR-CRITIC methods lay in the family of POLICY-BASED algorithms. Differently by VALUE-BASED ones - like DQN - the main network (ACTOR) learns directly the action  $a_\pi(s)$  to return for a given state  $s$ .

☛ Therefore there is no  $\text{argmax } Q_\pi(s, A)$  among all  $Q_\pi(s)$  for a given state

# Actor-Critic

---

ACTOR-CRITIC defines an architecture based on 2 parts:



## Actor

Main network which infers the action for a given state  
 $(s_t \rightarrow a)$



## Critic

Secondary network used only for training, evaluates the ACTOR output returning the *averaged* value of state.

$$s_t \rightarrow V(s)$$

## Advantage function for incremental learning

The ADVANTAGE technique is used to give a evaluation of the last action *compared* to the average of other actions.

## Advantage function for incremental learning

The ADVANTAGE technique is used to give a evaluation of the last action *compared* to the average of other actions.

It answers the following question:

Is last action more rewarding than others in such scenario?

## Advantage function for incremental learning

The ADVANTAGE technique is used to give a evaluation of the last action *compared* to the average of other actions.

It answers the following question:

Is last action more rewarding than others in such scenario?

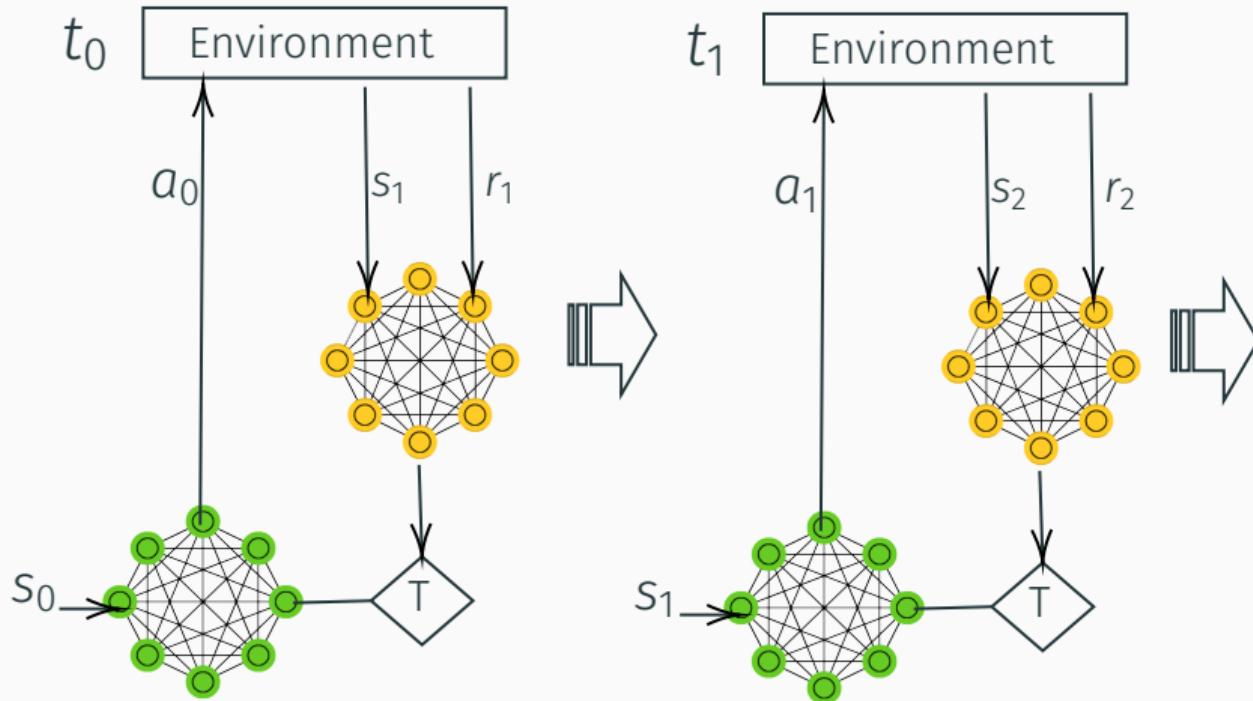
Yes

thumb up Action's likelihood is encouraged. Gradients are pushed *toward* this action.

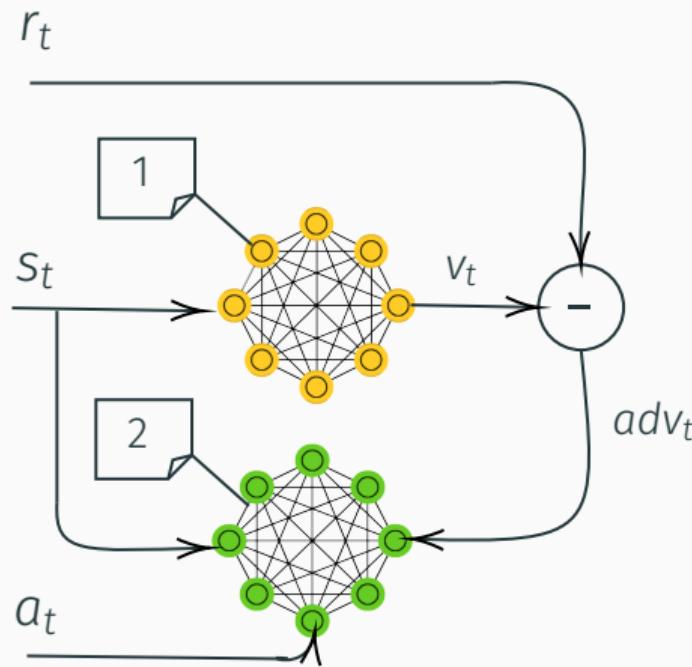
No

thumb down Action's likelihood is discouraged. Gradients are pushed *away* from this action.

# Training



## Actor training



- 1 Critic returns the averaged value of the state.
- 2 Actor is trained considering  $s_t$ ,  $a_t$  and  $adv_t = r_t - v_t$

## Deep Deterministic Policy Gradient

---

Especially effective in *continuous space*. Derived from ACTOR-CRITIC and DQN it uses 4 networks.

## **Value/action networks**

Likewise AC there is a distinction between value and action networks.

## **Duplication of AC networks**

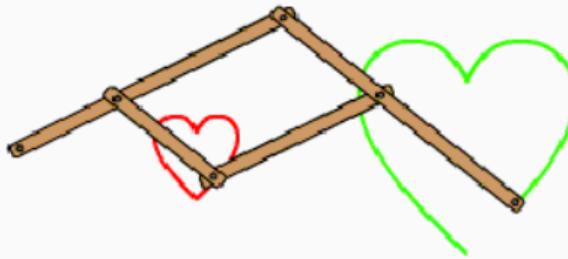
For stability those networks are duplicated.

## **Replay buffer**

Likewise DQN an history of past experiences is reused.

## Gradient propagation

Despite others tecniques, DDPG improves the actual  $\pi$  by applying the gradient obtained during the training iteration, from the value network to the output network.



**Figure 12:** Like with a pantograph, the gradient  $V_s$  is applied to the  $\pi$  network for a similar effect, without knowing the value of  $\pi$ .

# Asgard

---

# Asgard utility framework

Java utility for training model-free agents by means of RL. It uses DL4J for neural networks and brings following algorithms:

## DQN

Deep-Q with PRIORITIZED MEMORY REPLAY and VALUE-DIFFERENCE BASED EXPLORATION.

## A3C

Asynchronous advantage actor critic.



## Easy API

```
1 var output1 = init.step(stateVector1);
2 var action1 = output1.intValue();
3 //perform action1 in the environment
4 // get reward and new state
5 var output2 = output1.next(reward).step(stateVector2);
6 //....
7 output2.done(finalReward);
```

asgard1.java

ASGARD provides extensible interfaces for reactive flow patterns.

## Agent producer

Creates training (markovian or not) agents.

## Environment

Represent the environment for training e.g. OpenAI Gym.

## Episodes collector

Gathers a batch of episodes and then execute the training iteration.

## Threshold guards

When some metrics (e.g. averaged reward) are met, the training stops.

## Player

Agents are ready for inference and model can be stored.

# Demo

-  K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath.  
**A brief survey of deep reinforcement learning.**  
*arXiv preprint arXiv:1708.05866*, 2017.
-  Y. Bengio, A. Courville, and P. Vincent.  
**Representation learning: A review and new perspectives.**  
*IEEE transactions on pattern analysis and machine intelligence*,  
35(8):1798–1828, 2013.
-  Giancarlo Frison.  
**First Steps on Evolutionary Systems.**  
2018.  
link <https://labs.cx.sap.com/2018/07/02/first-steps-on-evolutionary-systems/>.

## References ii

---

-  K. Hornik, M. Stinchcombe, and H. White.  
**Multilayer feedforward networks are universal approximators.**  
*Neural networks*, 2(5):359–366, 1989.
-  J. Zico Kolter.  
**Markov Decision Process.**  
<http://www.cs.cmu.edu/afs/cs/academic/class/15780-s16/www/slides/mdps.pdf>.
-  John Williams.  
**2015 darpa robotics challenge.**  
Public domain.

## References iii

---

-  V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al.  
**Human-level control through deep reinforcement learning.**  
*Nature*, 518(7540):529, 2015.
-  P. Montague.  
**Reinforcement Learning: An Introduction**, by Sutton, R.S. and Barto, A.G.  
*Trends in Cognitive Sciences*, 3(9):360, 1999.
-  Omanomanoman.  
CC BY-SA 4.0, <https://commons.wikimedia.org/>.

-  M. Tokic and G. Palm.  
**Value-difference based exploration: adaptive control between epsilon-greedy and softmax.**  
In *Annual Conference on Artificial Intelligence*, pages 335–346. Springer, 2011.
-  Valve Corporation.  
A screenshot from the video game Dota 2.
-  Yamaguchi???.  
CC-BY-SA-3.0 (<http://creativecommons.org/licenses/by-sa/3.0/>)].