

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Markov Decision Process . . . . .	5
2.1.1	Definition . . . . .	5
2.1.2	Bellman Equation . . . . .	10
2.1.3	Value Iteration vs Policy Iteration . . . . .	11
2.2	Solving MDP . . . . .	13
2.2.1	Dynamic Programming . . . . .	13
2.2.2	Reinforcement Learning . . . . .	14
2.2.3	Monte Carlo Methods . . . . .	17
2.2.4	Temporal Difference Learning . . . . .	18
2.2.5	Sarsa Algorithm : On-policy TD Control . . . . .	19
2.2.6	Q-learning Algorithm : Off-policy TD Control . . . . .	20
2.2.7	SARSA( $\lambda$ ) Algorithm . . . . .	21
<b>3</b>	<b>Black-Box Optimization</b>	<b>25</b>
3.0.1	Gaussian Processes . . . . .	25
3.0.2	Bayesian Optimization . . . . .	26
3.0.3	An RL Approach To Black-Box Optimization . . . . .	28
3.0.4	BURLAP Java Library . . . . .	34
3.0.5	Humans in the Loop . . . . .	36
<b>4</b>	<b>Experimental Setting</b>	<b>39</b>
4.1	Benchmark . . . . .	39
4.1.1	Test Functions . . . . .	39
4.1.2	Algorithm's Configuration . . . . .	46
4.2	Experiment on Humans . . . . .	49

<b>5 Results</b>	<b>51</b>
5.1 Sarsa( $\lambda$ )'s Performances on Test Functions . . . . .	52
5.1.1 Himmelblau' s Function . . . . .	52
5.1.2 Sphere Function . . . . .	55
5.1.3 Beale Function . . . . .	58
5.1.4 Styblinski-Tang Revised Function . . . . .	60
<b>6 Appendix A</b>	<b>65</b>

# Chapter 1

## Introduction

Nowadays every aspect of life is characterized by the need to make choices. Each choice is conditioned by a variable number of parameters. In many applicative domains (recommendation systems, medical analysis tools, real time game engines, speech recognizers...), this number could be very high. The possibility to manage in an efficient way such a number of parameters is guaranteed by a significant number of mathematical techniques increasingly performing and refined.

Every decisional problem can be translated into a mathematical model that represents the essence of the problem itself. A crucial step in formulating a model is the construction of the *objective function*. This requires developing a quantitative measure of performance relative to each of the decision maker's ultimate objectives identified while the problem was being defined. [HL01]

However it is not always possible to define the objective function in advance. It could be unknown or, at least, partially unknown at the time the problem is dealt with. In those cases the optimization process of the objective function takes the name of *black-box optimization*.

Mathematically, the problem of finding a global maximizer (or minimizer) of an unknown objective function  $f$  is formalized as:

$$x^* = \arg \max_{x \in \mathcal{X}} f(x) \quad (1.1)$$

where  $\mathcal{X}$  is some design space of interest. In global optimization,  $\mathcal{X}$  is often a compact subset of  $\mathbb{R}^d$  [SSW<sup>+</sup>16].

Bayesian optimization (BO) applied to Gaussian Processes (GPs) represents one of the best known black-box optimization techniques. It represents a powerful tool for the joint optimization of design choices that is gaining great popularity in recent years [AS08]. A prior belief over the possible objective functions is defined. After that, as data are observed, the model is sequentially refined via Bayesian posterior updating. The Bayesian posterior represents the updated belief -given data- on the likely objective function [SSW<sup>+</sup>16].

In this thesis, an innovative approach to black-box optimization based on the Reinforcement Learning (RL) technique is proposed. RL is the problem faced by a learner that must behaviour through trial-and-error interactions with a dynamic environment. It can be considered the problem of mapping situations to actions in order to maximize a numerical reward signal [KLM96].

In the last part of the current work, humans' function learning techniques are briefly analysed. Performances of human and artificial approaches are compared and results are deeply described.

# Chapter 2

## Background

### 2.1 Markov Decision Process

Markov Decision Processes (MDPs) are a classical formalization of sequential decision making under uncertainty, where actions influence not just immediate responses, but also subsequent situations. Hence an MDP can be described as a controlled *Markov chain*<sup>1</sup>, where the control is given at each step by the chosen action. In this chapter the structure of an MDP and many techniques to solve its are presented.

#### 2.1.1 Definition

Markov decision processes are defined as controlled stochastic processes satisfying the *Markov property*<sup>2</sup> and assigning reward values to state transitions [Put94]. Formally, they are described by the 5-tuple  $(S, A, T, p(), R())$  where:

- $S$  is the state space in which the process' evolution takes place;
- $A$  is the set of all possible actions which control the state dynamics;
- $T$  is the set of time steps where decisions need to be made;
- $p()$  denotes the state transition probability function;
- $r()$  provides the reward function defined on state transitions.

---

<sup>1</sup>A sequence of random variables  $X_0, X_1, \dots$  with values in a countable set  $S$  is a *Markov chain* if at any time  $n$ , the future states (or values)  $X_{n+1}, X_{n+2}, \dots$  depend on the history  $X_0, \dots, X_n$  only through the present state  $X_n$  [Kon09].

<sup>2</sup>A stochastic process has the *Markov property* if the probabilistic behaviour of the chain in the future depends only on its present value and discards its past behaviour.

**States** The set of environmental states  $S$  is defined as the finite set  $\{s_1, \dots, s_N\}$  where the size of the state space is  $N$ , i.e.  $|S| = N$  [WvO12]. Each state  $s \in S$  is a vector of attributes (*state variables*) that describes the current configuration of the system [NFK06]. More specifically, a state variable is the minimally dimensioned function of history that is necessary and sufficient to compute the "state of knowledge" [Pow07].

**Actions** The set of actions  $A$  is defined as the finite set  $\{a_1, \dots, a_K\}$  where the size of the action space is  $K$ , i.e.  $|A| = K$ . Actions can be used to control the system state. The set of actions that can be applied in some particular state  $s \in S$ , is denoted  $A(s)$ , where  $A(s) \subseteq A$ . In more structured representations the fact that some actions are not applicable in some states, is modelled by a precondition function :  $S \times A \rightarrow \text{true}, \text{false}$ , stating whether action  $a \in A$  is applicable in state  $s \in S$  [WvO12].

**Time Steps** The set of time steps  $T$  is defined as the finite set  $\{t_1, \dots, t_M\}$  where the size of the action space is  $M$ , i.e.  $|T| = M$ . Speaking about time a distinction between *epochs* and *episodes* has to be made. An episode consists of a fixed number of epochs. An epoch is the smallest time unit in an MDP.

**Transition Probability** The transition probabilities  $p()$  characterize the state dynamics of the system, i.e. indicate which states are likely to appear after the current state. For a given action  $a$ ,  $p(s'|s, a)$  represents the probability for the system to transit to state  $s'$  after undertaking action  $a$  in state  $s$ . This  $p()$  function is usually represented in a matrix form where  $P_a$  is the  $|S| \times |S|$  matrix containing elements  $\forall s, s', P_{a, s, s'} = p(s'|s, a)$ . Since each line of these matrices sums to one, the  $P_a$  is said to be stochastic. The  $p()$  probability distributions over the next state  $s'$  follow the fundamental property which gives their name to Markov Decision Processes. If one writes  $h_t = (s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t)$  the history of states and actions until time step  $t$ , then the probability of reaching state  $s_{t+1}$  consecutively to action  $a_t$  is only a function of  $a_t$  and  $s_t$ , and not of the entire history  $h_t$  [SB10]. This concept can be easily resumed through the following equation :

$$\forall h_t, a_t, s_{t+1} \quad P(s_{t+1}|h_t, a_t) = P(s_{t+1}|s_t, a_t) = p(s_{t+1}|s_t, a_t) \quad (2.1)$$

**Reward Function** The reward function specifies rewards for being in a state, or doing some action in a state. The state reward function is defined

as  $R : S \rightarrow \mathbb{R}$ , and it specifies the reward obtained in states. The reward function is an important part of the MDP that specifies implicitly the *goal* of learning. Thus, the reward function is used to give direction in which way the system, i.e. the MDP, should be controlled [WvO12]. It is critical that the reward set up truly indicate the goal that has to be accomplished. If the achievement of sub-goals is rewarded, then the agent might find a way to achieve them without achieving the real goal. Reward signal is the way of communicating the agent *what* goal to achieve, not *how* the goal has to be achieved [SB18].

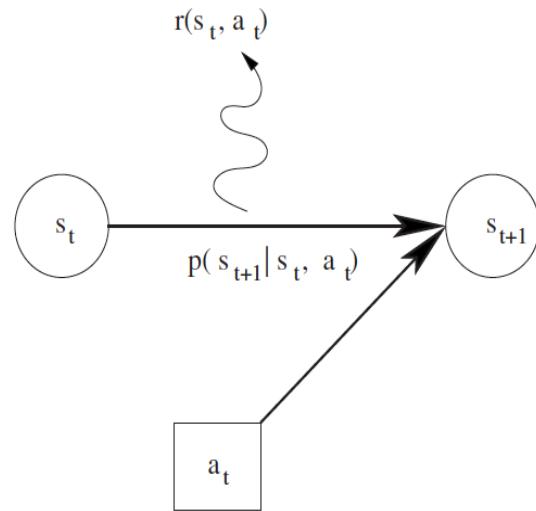


Figure 2.1: Markov decision process [SB10].

Markov decision processes allow to model the state evolution dynamics of a stochastic system when this system is controlled by an agent choosing and applying the actions  $a_t$  at every time step  $t$ . The procedure of choosing such actions is called an action policy, or strategy, and is written as  $\pi$  [SB10].

**Policy** Formally, given an MDP  $\langle S, A, p(), R() \rangle$ , a policy is a computable function that outputs for each state  $s \in S$  an action  $a \in A(s)$  [WvO12]. A policy can decide deterministically upon the action to apply or can define a probability distribution over the possible applicable actions. Then, a policy can be based on the whole history  $h_t$  (history-dependent policy) or can only consider the current state  $s_t$ . Thus, four main families of policies can be defined (2.1).

Policy $\pi_t$	Deterministic	Stochastic
Markov	$s_t \rightarrow a_t$	$a_t, s_t \rightarrow [0, 1]$
History-dependent	$h_t \rightarrow a_t$	$h_t, s_t \rightarrow [0, 1]$

Table 2.1: Different policy families for MDPs [SB10]

For a deterministic policy,  $\pi_t(s_t)$  or  $\pi_t(h_t)$  defines the chosen action  $a_t$ . For a stochastic policy,  $\pi_t(a, s_t)$  or  $\pi_t(a, h_t)$  represents the probability of selecting  $a \in A$  for  $a_t$  [SB10]. The sets so defined are included in each other, from the most general case of stochastic, history-dependent policies, to the very specific case of deterministic, Markov policies, as shown in figure 2.2.

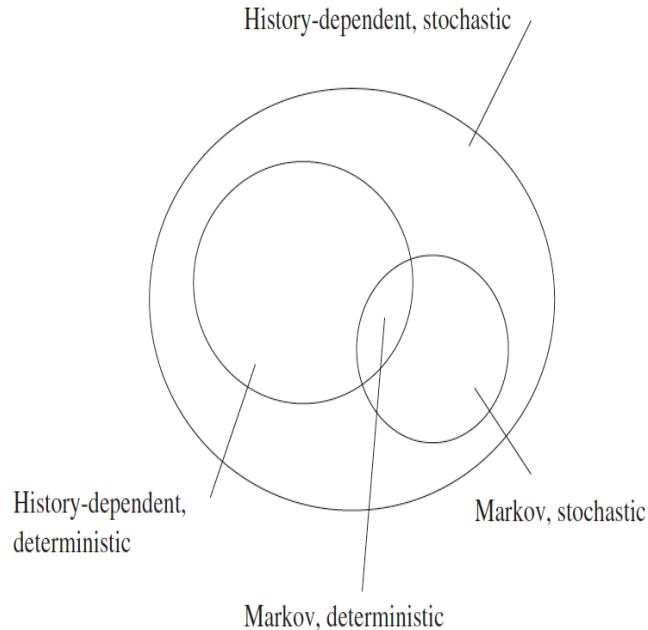
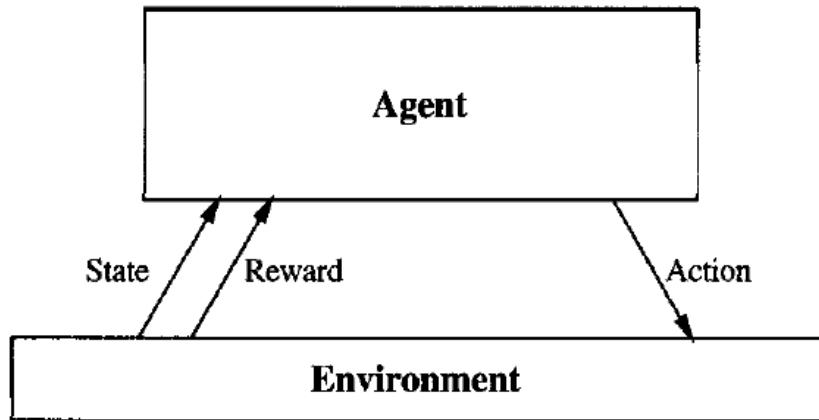


Figure 2.2: Relationship between the different sets of policies [SB10].

As shown in figure 2.3 each time an agent performs an action  $a_t$  in a state  $s_t$ , it receives a real-valued reward  $r_t$  that indicates the immediate value of this state-action transition. This produces a sequence of states  $s_i$ , actions  $a_i$ , and immediate rewards  $r_i$  as shown in the figure. The agent's task is to learn a control policy,  $\{\pi : S \rightarrow A\}$ , that maximizes the expected sum of these rewards, with future rewards discounted exponentially by their delay [Mit97].



$$s_0 \xrightarrow[r_0]{a_0} s_1 \xrightarrow[r_1]{a_1} s_2 \xrightarrow[r_2]{a_2} \dots$$

Figure 2.3: Policy application schema [SB18].

As previously said, the final goal of an optimal policy is to gather the greatest reward possible. In MDPs, there are several ways of taking into account the future in the current behaviour but there are basically three models of optimality which are sufficient to cover most of the approaches in the literature :

$$E\left[\sum_{t=0}^h r_t\right] \quad E\left[\sum_{t=0}^{\infty} \gamma^t r_t\right] \quad \lim_{h \rightarrow \infty} E\left[1/h \sum_{t=0}^h r_t\right]$$

Figure 2.4: Optimality : **a)** finite horizon, **b)** discounted, infinite horizon, **c)** average reward

**Discount Factor** The *finite horizon* model simply takes a finite horizon of length  $h$  and states that the agent should optimize its expected reward

over this horizon. In the *infinite-horizon model*, the long-run reward is taken into account, but the rewards that are received in the future are discounted according to how far away in time they will be received. A *discount factor*  $\gamma$ , with  $0 \leq \gamma \leq 1$  is used for this. Note that in this discounted case, rewards obtained later are discounted more than rewards obtained earlier. Additionally, the discount factor, ensures that, even with infinite horizon, the sum of the rewards obtained is finite. In episodic tasks, i.e. in tasks where the horizon is finite, the discount factor is not needed or can equivalently be set to 1. If  $\gamma = 0$  the agent is said to be myopic, which means that it is only concerned about immediate rewards. The last optimality model is *average-reward* model, maximizing the long-run *average-reward*. Sometimes this is called the *gain optimal* policy and in the limit, it is equal to the infinite-horizon discounted model [WvO12]

### 2.1.2 Bellman Equation

The concept of *value function* is the link between optimality criteria and policies. Most learning algorithms for MDPs compute optimal policies by learning value functions. The value of a state  $s$  under policy  $\pi$ , denoted  $V^\pi(s)$  is the expected return when starting in  $s$  and following  $\pi$  thereafter. Using the infinite-horizon, discounted model :

$$V^\pi(s) = E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s \right\} \quad (2.2)$$

A similar state-action value function :  $Q : S \times A \rightarrow \mathbb{R}$  can be defined as the expected return starting from state  $s$ , taking action  $a$  and thereafter following policy  $\pi$  :

$$Q^\pi(s, a) = E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s, a_t = a \right\} \quad (2.3)$$

For any policy  $\pi$  and any state  $s$  the expression 2.3 can recursively be defined in terms of a so-called *Bellman Equation* :

$$\begin{aligned} V^\pi(s) &= E_\pi \{r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = t\} \\ &= E_\pi \{r_t + \gamma V^\pi(s_{t+1}) | s_t = s\} \\ &\quad \sum_{s'} p(s, \pi(s), s') (R(s, a, s') + \gamma V^\pi(s')) \end{aligned} \quad (2.4)$$

It denotes that the expected value of state is defined in terms of immediate reward and values of possible next state weighted by the transition probabilities, and additionally a discount factor. Note that multiple policies can have the same value function, but for a given policy  $\pi$ ,  $V^\pi$  is unique. The goal for any MDP is to find a best policy, i.e. the policy that receives the most reward. This means maximizing the value function of equation 2.3 for all states  $s \in S$ . An optimal policy, denoted  $\pi^*$ , is such that  $V^{\pi^*}(s) \geq V^\pi(s)$  for all  $s \in S$  and all policies  $\pi$  :

$$V^*(s) = \max_{a \in A} \sum_{s'} p(s, \pi(s), s') (R(s, a, s') + \gamma V^\pi(s')) \quad (2.5)$$

This expression is called the *Bellman optimality equation*. It states that the value of a state under an optimal policy must be equal to the expected return for the best action in a state [WvO12].

### 2.1.3 Value Iteration vs Policy Iteration

Before studying how to solve MDPs let's clarify main differences between two key concepts that will be recursively find : *value iteration* and *policy iteration*.

Both value iteration and policy iteration compute the same thing (all optimal values), i.e. they work with Bellman updates. In value iteration one starts with a random value function and then he finds a new (improved) value function in an iterative process, until reaching the optimal value function. Value iteration computation of Bellman Equation is:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} p(s, a, s') [R(s, a, s') + \gamma V_k(s')] \quad (2.6)$$

Value iteration has three disadvantages: it is slow ( $O(S^2A)$  per iteration), the max at each state rarely changes, the policy often converges long before the values (figure 2.5).

Instead, policy iteration computations of Bellman Equation are:

- Evaluation : for a fixed current policy  $p$ , find values with policy evaluation:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} p(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')] \quad (2.7)$$

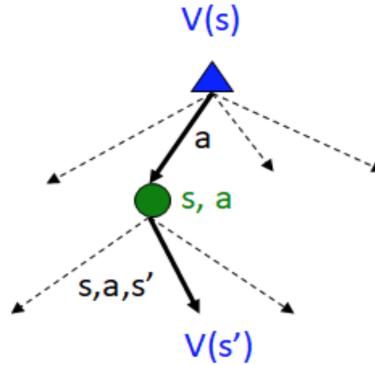


Figure 2.5: Value Iteration Schema

- Improvement : for fixed values, get a better policy using policy extraction :

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} p(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')] \quad (2.8)$$

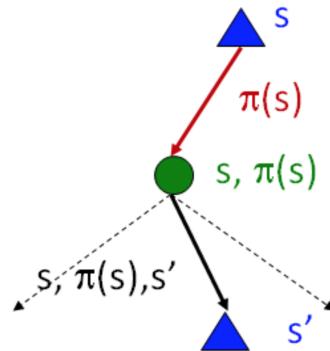


Figure 2.6: Policy Iteration Schema

Complexity of policy iteration is  $O(S^2)$ .

Summarizing what said until now it is possible to write:

In value iteration :

- Every iteration updates both the values and implicitly the policy;

- The policy is not tracked, but taking the max over actions implicitly recomputes it.

In policy iteration :

- Several passes that update utilities with fixed policy are made;
- After the policy is evaluated, a new policy is chosen;
- The new policy will be always better until the end.

## 2.2 Solving MDP

Now that the concepts of MDP, policy, optimality criteria and value function, have been formally defined, it is time to consider the question of how to solve an MDP computing an optimal policy  $\pi^*$ . Several dimensions exists along which algorithms have been developed for this purpose. The most important distinction is that between *model-based* and *model-free* algorithms [WvO12].

	Model-based algorithms	Model-free algorithms
General name	DP	RL
Basic assumption	A model of the MDP is known beforehand, and can be used to compute value functions and policies using the Bellman equation.	Rely on interaction with the environment. Because a model of the MDP is not known, the agent has to explore the MDP to obtain information.
Planning	Yes	Yes
Learning	No	Yes

Table 2.2: Main differences between model-based and model-free algorithms.

### 2.2.1 Dynamic Programming

The term Dynamic Programming (DP) refers to a class of algorithms that is able to compute optimal policies in the presence of a perfect model of the environment described as a Markov Decision Process. These algorithms are known as *planning algorithms*. In planning, the idea is that some description of a starting state or states, a goal state or states, and some set of possible actions that the agent can take are given. The objective is to search through a tree that is the sequences of actions that can be taken, and try to find a nice short plan. In so doing two of the possible planning algorithms are BFS algorithm and DFS algorithm.

The method of Dynamic Programming systematically records solutions for all sub-problems of increasing lengths. Using this programming paradigm the optimal policy is constructed through the step-by-step definition of optimal sub-policies. According to Bellman optimality principle, all optimal sub-policies of an optimal policy are optimal sub-policies. DP algorithms are obtained by turning Bellman equations into update rules for improving approximations of the desired value functions.

Studying DP it is possible to distinguish two different main methods. *Policy evaluation* refers to the typically iterative computation of the value functions for a given policy. *Policy improvement* refers to the computation of an improved policy given the value function for that policy. Either of these can be used to reliably compute optimal policies and value functions for finite MDPs given complete knowledge of the MDP.

Classical DP methods operate in sweeps through the state set, performing an *expected update* operation on each state. Each such operation updates the value of one state based on the values of all possible successor states and their probabilities of occurring. Expected updates are closely related to Bellman equations: they are little more than these equations turned into assignment statements. When the updates no longer result in any changes in value, convergence has occurred to values that satisfy the corresponding Bellman equation [SB18].

The assumption that a model is available will be hard to ensure for many applications, however DP algorithms are very relevant because they define fundamental computational mechanism which are also used when no model is available.

### 2.2.2 Reinforcement Learning

Dynamic Programming's methods compute optimal policies for an MDP assuming that a perfect model is available. Reinforcement Learning (or *approximate dynamic programming*, or *neuro-dynamic programming*) is primarily concerned with how to obtain an optimal policy when such a model is not available. RL adds to MDPs a focus on approximation and incomplete information, and the need for sampling and exploration. In contrast with DP's algorithms, model-free methods do not rely on the availability of priori known transition and reward models. The lack of the model generates a need to *sample* the MDP to gather statistical knowledge about this unknown model. Many model-free RL techniques exist that probe the environment by doing actions, thereby estimating the same kind of state value

and state-action value functions as model-based techniques [WvO12].

Roughly speaking Reinforcement Learning (RL) is the problem faced by a learner that must behaviour through trial-and-error interactions with a dynamic environment. It can be considered a problem of mapping situations to actions in order to maximize a numerical reward signal [KLM96].

In RL the learner must select an action to take in each time step: every choice done by the agent changes the environment in an unknown fashion and receives a reward which value is based on the consequences. The objective of the learner is to choose a sequence of actions based on observations of the current environment that maximizes cumulative reward or minimizes cumulative cost over all time steps [LM16]. The general class of algorithms that interact with the environment and update their estimates after each experience is called *online* RL.

---

**Algorithm 1:** A general algorithm for online RL [WvO12]

---

```

1 foreach episode do
2   s  $\in S$  is initialized as the starting state;
3   t := 0;
4   repeat
5     choose an action a  $\in A(s)$ ;
6     perform action a ;
7     observe the new state s' and received reward r update  $\tilde{T}, \tilde{R}, \tilde{Q}$ 
      and/or  $\tilde{V}$  using the experience  $\langle s, a, r, s' \rangle$ ;
8     s := s';
9   until s' is a goal state;
10 end
```

---

Studying RL, one of the main aspect is the distinction between *direct* and *indirect* Reinforcement Learning. The difference between these two approaches is shown in figure 2.3.

RL is different both from *supervised learning* and from *unsupervised learning*. It is not a sample of learning from a training set of labelled examples provided by a knowledgeable external supervisor (*supervised learning*) and it is not a sample of searching and finding structure hidden in a collection of unlabelled data (*unsupervised learning*). More specifically RL can be distinguished from other forms of learning based on the following char-

	Model-based/Indirect	Model-free/Direct
Basic Assumption	First of all it is necessary to learn the transition and reward model from interaction with the environment. After that, when the model is (approximately or sufficiently) correct, all the DP methods previously described can be applied.	It is possible to estimate values of functions, without even estimating the model of the MDP.

Table 2.3: RL's approaches classification.

acteristics :

- Reinforcement Learning deals with temporal sequences. In contrast with non-supervised learning problems where the order in which the examples are presented is not relevant, the choice of an action at a given time step will have consequences on the examples that are received at a subsequent time steps [SB10].
- In contrast with supervised learning, the environment does not tell the agent what would be the best possible action. Instead, the agent may just receive a scalar reward representing the value of its action and it must *explore* the possible alternative actions to determine whether its action was the best or not [SB10].

According to what just said, one of the challenges that arise in RL, and not in other kind of learning, is the trade-off between *exploration* and *exploitation*. To obtain an higher reward, an RL agent must prefer actions that it has tried in the past and found to be effective in producing reward. In order to discover such actions, it has to try actions that it has not selected before. The agent *exploits* what it has already experienced in order to obtain reward, but it has also to *explore* in order to eventually make better action selections in the future [SB18]. In other words *exploitation* consists of doing again actions which have proven fruitful in the past, whereas *exploration* consists of trying new actions, looking for a larger cumulated reward, but eventually leading to a worse performance. Dealing with the exploration/exploitation trade-off consists of determining how the agent should explore to get as fast as possible a policy that is optimal or close enough to the optimum. The most basic exploration strategy is the  $\epsilon$ -greedy policy, i.e. the learner takes its current best action with probability  $(1 - \epsilon)$  and a (randomly selected) other action with probability  $\epsilon$ . [SB10].

### 2.2.3 Monte Carlo Methods

Monte Carlo Methods (MC) represent a first instance of model-free RL's algorithms. They are one way solving the Reinforcement Learning problem based on averaging sample returns.

The Monte Carlo approach consists of performing a large number of trajectories from all states  $s$  in  $S$ , and estimating  $V(s)$  as an average of the cumulated rewards observed along these trajectories. In each trial, the agent records its transitions and rewards, and updates the estimates of the value of the encountered states according to a discounted reward scheme. The value of each state then converges to  $V^\pi(s)$  for each  $s$  if the agent follows policy  $\pi$ .

More formally, let  $(s_0, s_1, \dots, s_N)$  be a trajectory consistent with the policy  $\pi$  and the unknown transition function  $p()$ , and let  $(r_1, \dots, r_N)$  be the rewards observed along this trajectory. In MC method, the  $N$  values  $V(s_t)$ ,  $t = 0, \dots, N - 1$  are updated according to :

$$V(s_t) \leftarrow V(s_t) + \alpha(s_t)(r_{t+1} + r_{t+2} + \dots + r_N - V(s_t)) \quad (2.9)$$

with the learning rates  $\alpha(s_t)$  converging to 0 along the iterations [SB10] . MC algorithms treat the long-term reward as a random variable and take as its estimate the sampled mean. Because the sampling is dependent on the current policy  $\pi$ , only returns for actions suggested by  $\pi$  are evaluated. Thus, *exploration* is of key importance here, just as in other model-free methods. One way of ensuring enough exploration is to use exploring starts, i.e. each state-action pair has a non-zero probability of being selected as the initial pair.

A distinction can be made between *every-visit* MC, which averages over all visits of a state  $s \in S$  in all episodes, and *first-visit* MC, which averages over just the returns obtained from the first visit to a state  $s \in S$  for all episodes. Both variants will converge to  $V^\pi$  for the current policy  $\pi$  over time [WvO12] .

Studying RL methods it is necessary to distinguish between *on-policy methods* and *off-policy methods*. On-policy methods attempt to evaluate or improve the policy that is used to make decisions. Off-policy methods evaluate or improve a policy different from that used to generate the data. MC methods can be used for both on-policy and off-policy control, and the general pattern complies with the generalized policy iteration procedure.

	<b>On-policy methods</b>	<b>Off-policy methods</b>
<b>Basic assumption</b>	They attempt to evaluate or improve the policy that is used to make decisions.	No specific policy is required to be followed; the agent could even behave randomly and despite this it can still find the optimal policy.

Table 2.4: Difference between on-policy and off-policy methods.

### 2.2.4 Temporal Difference Learning

An important problem faced by RL is the so called *temporal credit assignment problem*. In model-free contexts it is difficult to assess the utility of some action, if the real effects of this particular action can only be perceived much later. One possibility is to wait until the "end" (e.g. of an episode) and punish or reward specific actions along the path taken. However, this will take a lot of memory and often, with ongoing tasks, it is not known beforehand whether, or when, there will be an "end". Instead, one can use similar mechanisms as in value iteration to adjust the estimated value of a state based on the immediate reward and the estimated (discounted) value of the next state. This is generally called *temporal difference learning* which is a general mechanism underlying the model-free methods [WvO12].

The basic temporal difference algorithm is  $TD(0)$ . This algorithm relies on a comparison between the actually received reward and the reward expected from the previous estimations.

If the estimated values  $V(s_t)$  and  $V(s_{t+1})$  in states  $s_t$  and  $s_{t+1}$  were exact, one would have

$$V(s_t) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots, \quad (2.10)$$

$$V(s_{t+1}) = r_{t+2} + \gamma r_{t+3} + \gamma^2 r_{t+4} + \dots \quad (2.11)$$

Thus one would have

$$V(s_t) = r_{t+1} + \gamma V(s_{t+1}). \quad (2.12)$$

The temporal difference error

$$\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t) \quad (2.13)$$

measures the error between the current estimation  $V(s_t)$  and the corrected estimation  $r_{t+1} + V(s_{t+1})$ . The temporal difference method consists of correcting this error little by little by modifying  $V()$  according to a Windrow-Hoff equation :

$$V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)] = V(s_t) + \alpha\delta_t. \quad (2.14)$$

Like Monte Carlo methods, TD ones perform estimation of action-value functions based on the experience of the agent and can do it without a model of the underlying MDP; however, they combine this estimation using local estimation propagation mechanisms coming from dynamic programming, resulting in their incremental properties. Thus TD methods, which are at the heart of most reinforcement learning algorithms, are characterized by this combination of estimation methods with local updates incremental properties [SB10].

### 2.2.5 Sarsa Algorithm : On-policy TD Control

Knowing the exact value of all states is not always enough to determine what to do. If the agent does not know which action results in reaching any particular state, i.e. if the agent does not have a model of the transition function, knowing  $V$  does not help it determine its policy. To solve this problem, Watkins [WD92] introduced the action-value function  $Q$ , whose knowledge is similar to the knowledge of  $V$  when  $p$  is known. The action-value function of a fixed policy  $\pi$  whose value function is  $V^\pi$  is :

$$\forall s \in S, a \in A, \quad Q^\pi(s, a) = R(s, a) + \gamma \sum_{s'} p(s'|s, a) V^\pi(s'). \quad (2.15)$$

The value of  $Q^\pi(s, a)$  is interpreted as the expected value when starting from  $s$ , executing  $a$  and then following the policy  $\pi$  afterwards. Having  $V^\pi(x) = Q^\pi(x, \pi(x))$ , the corresponding Bellman equation is :

$$\forall s \in S, a \in A \quad Q^*(s, a) = R(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_b Q^*(s', b) \quad (2.16)$$

Than :

$$\forall s \in S, \quad V^*(s) = \max_a Q^*(s, a), \quad \pi^*(s) = \arg \max_a Q^*(s, a). \quad (2.17)$$

The SARSA algorithm works on state-action pairs rather than on states. Its update equation is :

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]. \quad (2.18)$$

the information necessary to perform such an update is  $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$ , hence the name of the algorithm : SARSA. This algorithm suffers from one conceptual drawback: performing the updates as stated above implies knowing in advance what will be the next action  $a_{t+1}$  for any possible next state  $s_{t+1}$ . As a result, the learning process is tightly coupled to the current policy (the algorithm is called "on-policy") and this complicates the exploration process. As a result, proving the convergence of SARSA was more difficult than proving the convergence of "off-policy" algorithms such as Q-learning. Empirical studies often demonstrates the better performance of SARSA compared to Q-learning [SB10].

---

**Algorithm 2:** SARSA Algorithm : On-policy TD Control

---

```

1 Initialize  $Q(s, a)$  arbitrarily;
2 Repeat the next cycle until  $s$  is terminal;
3 foreach episode do
4   Initialize  $s$ ;
5   Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy);
6   foreach step of episode do
7     Take action  $a$ , observer  $r, s'$ ;
8     Choose  $a'$  from  $s'$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy);
9      $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$  ;
10     $s \leftarrow s'$  ;
11     $a \leftarrow a'$  ;
12  end
13 end

```

---

### 2.2.6 Q-learning Algorithm : Off-policy TD Control

The Q-Learning Algorithm can be seen as a simplification of the algorithm, given that it is no more necessary to determine the action at the next step to calculate updates. Its update equation is :

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (2.19)$$

---

**Algorithm 3:** Q-learning Algorithm [SB10]

---

```

1  $\alpha_t$  is the learning rate ;
2 Initialize ( $Q_0$ );
3 for  $t = 0$  ;  $t \leq T_{tot} - 1$  ;  $t++$  do
4    $s_t \leftarrow \text{ChooseState};$ 
5    $a_t \leftarrow \text{ChooseAction};$ 
6    $s_{t+1}, r_{t+1} \leftarrow \text{Simulate}(s_t, a_t);$ 
7   update  $Q_t$  ;
8   begin
9      $Q_{t+1} \leftarrow Q_t;$ 
10     $\delta_t \leftarrow r_{t+1} + \gamma \max_b Q_t(s_{t+1}, b) - Q_t(s_t, a_t)$  ;
11     $Q_{t+1}(s_t, a_t) \leftarrow Q_t(s_t, a_t) + \alpha_t(s_t, a_t)\delta_t$  ;
12  end
13 end
14 return  $Q_{Tot}$ 

```

---

The main difference between SARSA and Q-Learning lies in the definition of the error term. The  $Q(s_{t+1}, a_{t+1})$  term in the equation 2.18 is replaced by  $\max_a Q(s_{t+1}, a)$  in equation 2.19. Updates are based on instantaneously available information. In this algorithm, the  $T_{tot}$  parameter corresponds to the number of iterations. There is here one learning rate  $\alpha_t(s, a)$  for each state-action pair, it decreases at each visit of the corresponding pair. The **Simulate** function returns a new state and the corresponding reward according to the dynamics of the system. The choice of the current state and of the executed action is performed by functions **ChooseState** and **ChooseAction**. The **Initialize** function initializes the  $Q$  function with  $Q_0$ , which is often initialized with **null** values, whereas more adequate choices can highly improve the performance [SB10].

### 2.2.7 SARSA( $\lambda$ ) Algorithm

Previous algorithms only perform one update per time step in the state that the agent is visiting. This update process is particularly slow. Indeed, an agent deprived of any information on the structure of the value function needs at least  $n$  trials to propagate the immediate reward of a state to another state that is  $n$  transitions away. Before this propagation is achieved, if the initial values are **null**, the agent performs a random walk in the state space, which means that it needs an exponential number of steps as a function of  $n$  before reaching the reward "trail" (figure 2.7).

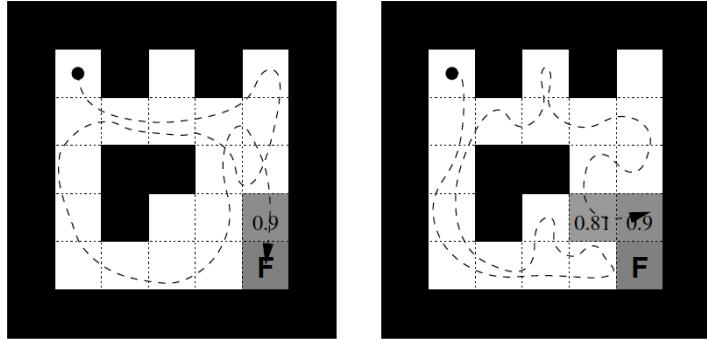


Figure 2.7: Q-learning : first and second trial. It can be seen that, all values being initially null, the propagation of non-null values does not starts until the agent finds the reward source for the first time, and only progress once per trial [SB10].

A naive way to solve the problem consists of using a memory of trajectory and to propagate all the information backwards along the performed transitions each time a reward is reached. Such a memory of performed transitions is called an "eligibility trace". A problem with this naive approach is that the required memory grows with length of trajectories, which is obviously not feasible in the infinite horizon context. In SARSA( $\lambda$ ) algorithm a more sophisticated approach that addresses the infinite horizon context is used.

As already saw, Q-learning integrates the temporal difference error idea. With the update rule of Q-learning :

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha_t \{r_{t+1} + \gamma V_t(s_{t+1}) - Q_t(s_t, a_t)\} \quad (2.20)$$

for transition  $(s_t, a_t, s_{t+1}, r_{t+1})$ , and in the case where action  $a_t$  executed in state  $s_t$  is the optima action for  $Q_t$ , then the error term is  $r_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)$ . For a generic parameter  $\lambda \in [0, 1]$  it can be generalized as shown in algorithm 4.

Here  $z_t(s, a)$  is the eligibility trace and *absorbing state* is the terminal state [SB10]. The eligibility trace represents a short-term memory vector,  $\mathbf{z}_t \in \mathbb{R}^d$ , that parallels the long-term weight vector  $\mathbf{w}_t \in \mathbb{R}^d$ . The rough idea is that when a component of  $\mathbf{w}_t$  participates in producing an estimated value, then the corresponding component of  $\mathbf{z}_t$  is bumped up and then begins to

---

**Algorithm 4:** SARSA ( $\lambda$ ) Algorithm [SB10]

---

```

1 /*  $\alpha$  is a learning rate */ ;
2 Initialize  $Q_0$  ;
3  $z_0 \leftarrow 0$  ;
4  $s_0 \leftarrow \text{ChooseState}$  ;
5  $a_0 \leftarrow \text{ChooseAction}$  ;
6  $t \leftarrow 0$  ;
7 while  $t \leq T_{tot} - 1$  do
8    $(s'_t, r_{t+1}) \leftarrow \text{Simulate}(s_t, a_t)$  ;
9    $a'_t \leftarrow \text{ChooseAction}$  ;
10  update  $Q_t$  and  $z_t$  ;
11  begin
12     $\delta_t \leftarrow r_{t+1} + \gamma Q_t(s'_t, a'_t) - Q_t(s_t, a_t)$  ;
13     $z_t(s_t, a_t) \leftarrow z_t(s_t, a_t) + 1$  ;
14    for  $s \in S, a \in A$  do
15       $Q_{t+1}(s, a) \leftarrow Q_t(s, a) + \alpha_t(s, a)z_t(s, a)\delta_t$  ;
16       $z_{t+1}(s, a) \leftarrow \gamma\lambda z_t(s, a)$ 
17    end
18    if  $s'_t$  non absorbing then
19       $s_{t+1} \leftarrow s'_t$  and  $a_{t+1} \leftarrow a'_t$  ;
20    end
21    else
22       $s_{t+1} \leftarrow \text{ChooseState}$  ;
23       $a_{t+1} \leftarrow \text{ChooseAction}$  ;
24    end
25  end
26 end
27 return  $Q_{Tot}$ 

```

---

fade away. Learning will then occur in that component of  $w_t$  if a non-zero TD error occurs before the trace falls back to zero. The trace-decay parameter  $\lambda \in [0, 1]$  determines the rates at which the trace falls [SB18]. Roughly speaking eligibility trace is a parameter used to control the “memory” of the algorithm, associated to a given state, enabling the assignment of more credit, i.e., higher trace, to the states and actions performed closer to the end of the episode.

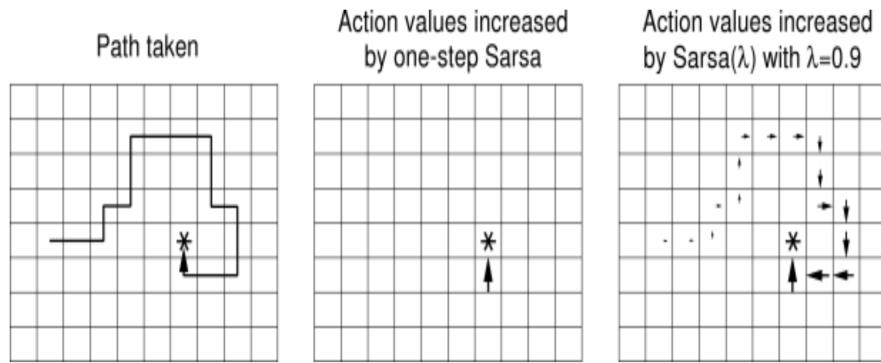


Figure 2.8: Grid-world example of the speed up of policy learning due to the use of eligibility traces.

The effect of the  $\lambda$  value is still poorly understood and tuning optimally its value for a given problem is still an open empirical problem [SB10].

## Chapter 3

# Black-Box Optimization

One of the most exciting challenges in optimization practice is the possibility to optimize in absence of an algebraic model of the system to be optimized. This kind of optimization is known as *black-box optimization* and it is the subject matter of this dissertation.

A black-box function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  is a function for which the analytic form is not known. Nowadays there are lots of mathematical models to succeed in optimize this kind of functions. One of the best known of these methods is the *Bayesian Optimization Model* (BO) applied to *Gaussian Processes* (GPs).

In this chapter the BO model will be analysed and then an innovative, Reinforcement Learning based approach to solve the same problem will be proposed.

### 3.0.1 Gaussian Processes

Before introducing BO it is necessary to describe what *Gaussian Processes* (GPs) are. GPs are an alternative approach to regression problems. GP is a *non-parametric* approach (there is not a priori knowledge of how many parameters will be useful for the regression) to find a distribution over the possible functions  $f(x)$  that are consistent with observed data. A GP is a generalization of the Gaussian probability distribution. Whereas a probability distribution describes random variables which are scalars or vectors (for multivariate distributions), a *stochastic process* governs the properties of functions.

GP is a convenient and powerful prior distribution on functions, which will be taken here to be of the form

$$f : \mathcal{X} \leftarrow \mathbb{R}. \quad (3.1)$$

The GP is defined by the property that any finite set of  $N$  points  $\{x_n \in \mathcal{X}\}_{n=1}^N$  induces a multivariate Gaussian distribution on  $\mathbb{R}^N$ . The  $n$ th of these points is taken to be the function value  $f(x_n)$  [SLA12]. The support and properties of the resulting distribution on functions are determined by a mean function

$$m : \mathcal{X} \leftarrow \mathbb{R} \quad (3.2)$$

and by a positive definite covariance function

$$k : \mathcal{X} \times \mathcal{X} \leftarrow \mathbb{R}. \quad (3.3)$$

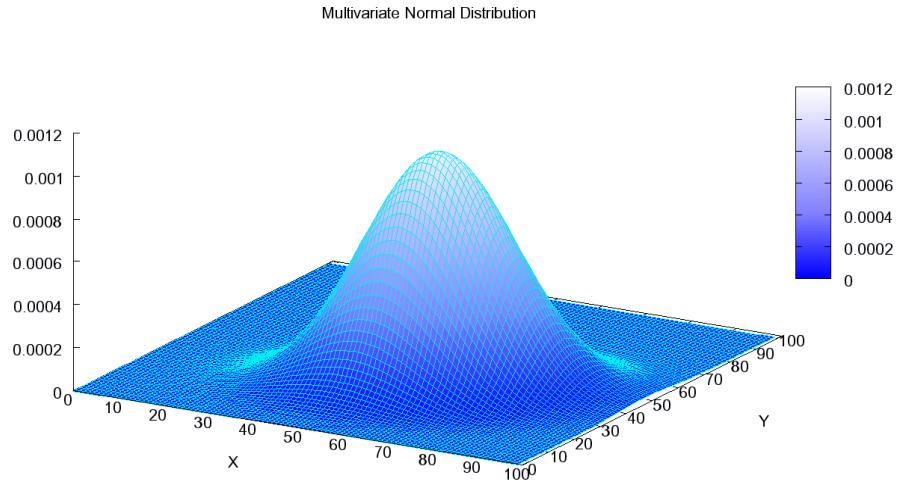


Figure 3.1: Multivariate Gaussian Distribution [MND].

### 3.0.2 Bayesian Optimization

Let's assume that the function  $f(x)$  is drawn from a GP prior and that our observations are of the form  $\{x_n \in \mathcal{X}\}_{n=1}^N$ , with  $y_n \sim \mathcal{N}(f(x_n), v)$  where  $v$  is

the variance of noise introduced into the function observations. This prior and those data induce posterior over functions; the acquisition function, which is denoted by

$$a : \mathcal{X} \leftarrow \mathbb{R}^+, \quad (3.4)$$

determines what point in  $\mathcal{X}$  should be evaluated next via a proxy optimization

$$x_{\text{next}} = \arg \max_x a(x), \quad (3.5)$$

where several different functions have been proposed. There are several popular choices of acquisition functions. Under the Gaussian process prior, those functions depend on the model solely through its predictive mean function  $\mu(x; \{x_n, y_n\})$  and predictive variance function  $\sigma^2(x; \{x_n, y_n\})$  [SLA12].

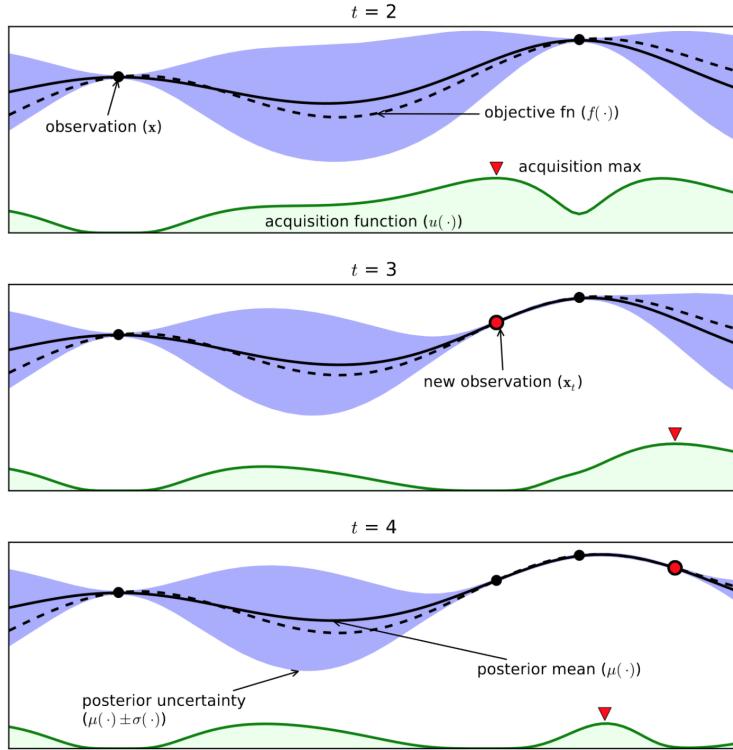


Figure 3.2: 2-d Bayesian Optimization Process Example [Bay]

According to what said above, it is possible to summarize the Bayesian Optimization's algorithm as follow:

---

**Algorithm 5:** Bayesian Optimization[SSW<sup>+</sup>16]

---

```

1 for  $n = 1, 2, \dots$  do
2   select new  $x_{n+1}$  by optimizing acquisition function  $\alpha$ 
      
$$x_{n+1} = \arg \max_x \alpha(x, \mathcal{D}_n)$$

      query objective function to obtain  $y_{n+1}$ 
      augment data  $\mathcal{D}_{n+1} = \{\mathcal{D}_n, (x_{n+1}, y_{n+1})\}$ 
      update statistical model
3 end
```

---

### 3.0.3 An RL Approach To Black-Box Optimization

As previously said, the aim of this thesis is to describe an innovative RL based approach to the black-box function optimization problem. The process to reach this goal has required about three months of work. A lot of attempts have been made. A full, detailed chronology of the development of the method presented in this chapter is described in **Appendix A**. Here the final product of all efforts is described.

Let's consider the following scenario :

- **Agent** : The Reinforcement Learning agent has to maximize a black-box bivariate function. The function is continuously defined over a specific domain. The function is mapped on a space of  $X \times Y$  pixels. The agent has to complete its job having exactly  $e$  epochs for each one of the  $E$  episodes. In each epoch the agent has a specific position in space described through the two coordinates  $(x, y)$ . Each time the agent makes an action, the angle between  $(x, y)$  and  $(x', y')$  and the value of the function  $f(x', y')$ , are computed.
- **State** : The state is represented by two lists: the first one contains the last two computed *angles* and the second one contains the correspondent last two *actions*.
- **Actions** : In each epoch the agent can make one of four different actions : *move north*, *move south*, *move east*, *move west*. Each time the

agent moves itself of  $p$  pixels in one of the previously described directions. The resultant effective movement over function is computed as shown in algorithm 6.

---

**Algorithm 6:** From pixels to real values

---

```

1 /* knowing pixel_X and pixel_Y */;
2 /* knowing pixel_X_Range and pixel_Y_Range */ ;
3 /* knowing the function */;
4
5 domain = function.getDomain() ;
6 x_Range = domain.max_X - domain.min_X ;
7 y_Range = domain.max_Y - domain.min_Y ;
8
9 x_Real = domain.min_X + (pixel_X * x_Range) / pixel_X_Range ;
10 y_Real = domain.min_Y + (pixel_Y * y_Range) / pixel_Y_Range ;
11
12 return x_Real, y_Real

```

---

- **Reward :** In the context just described, a real terminal state doesn't exist. Cause black-box functions are considered, it is not possible to specify two coordinates  $(x, y)$  and the corresponding value function  $z$  as maximum, that is, as the terminal state. For this reason every action of the agent is rewarded. The simulation ends after the  $E$  episodes are completed. In each explored state a  $\Delta$  computed as follow is defined:

$$\Delta = f(x, y) - \text{currentMax} \quad (3.6)$$

where  $f(x, y)$  is the value function computed in the *current state*. Reward at each epoch is equal to  $\Delta$ .

The nature of the movement the agent can make in each epoch can be of two different types : *linear movement* or *parametric movement*. A linear movement can be thought as a crow flies one. In this case the structure of the function is not really considered. On the other hand, a parametric movement can be thought as a real movement over the function curves. In this case the structure of the function influences the amount of the movement really done.

If the movement is linear, angles are computed as follow :

---

**Algorithm 7:** Angle computation in linear movement case.

---

```

1 /* knowing  $(x, y)$  */ ;
2 /* knowing  $(x', y')$  */ ;
3 /* movementAmount =  $M$  */ ;
4 /* currentMax =  $\max f(x_n, y_n)$  */ ;
5
6  $z = f(x, y)$  ;
7  $z' = f(x', y')$  ;
8
9  $\delta = ((x' - x), (y' - y))$  ;
10  $\alpha = \arctan\left(\frac{\delta}{\text{movementAmount}}\right)$ 
11
12  $\Delta = z' - \text{currentMax}$ 

```

---

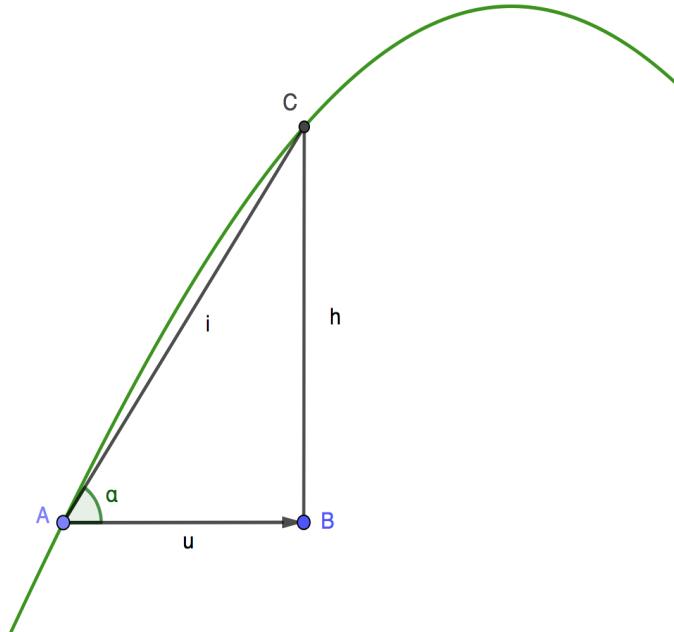


Figure 3.3: Linear movement computations.

This algorithm can be explained recalling some trigonometry. Let's suppose to be in a simple 2d-case like the one represented in figure 3.3.

Let's suppose that the starting point of the RL agent at epoch  $e_i$  is the point  $A$  with coordinates  $(x, y)$ . In the  $2d$ -case the agent can make only one of two actions for each epoch : *move on* and *go back*. Making a linear movement of the type *move on* the arriving point at epoch  $e_{i+1}$  is  $B$  with coordinates  $(x', y)$ . Knowing that the hidden objective function is  $f(x) = \sin(2x)$ , it is possible to compute  $f(x')$ . At this time point  $C$  with coordinates  $(x', f(x'))$  should be known.  $u$  is the **movement amount** and  $\overline{CB}$  equals to  $f(x') - f(x)$  is the  $\delta$ . From trigonometry

$$\tan \alpha = \frac{\delta}{\text{movementAmount}}, \quad (3.7)$$

so it is finally possible to compute the angle  $\alpha$  as follow :

$$\alpha = \arctan \frac{\delta}{\text{movementAmount}} \quad (3.8)$$

If the movement is parametric it is necessary to be more specific about the amount of movement over the function. This can be done using the following algorithm :

This algorithm can be easily understood looking at figure 3.4.

As previously done, let's suppose that the starting point of the RL agent at epoch  $e_i$  is the point  $A$  with coordinates  $(x, y)$ . In the  $2d$ -case the agent can make only one of two actions for each epoch : *move on* and *go back*. Making a linear movement of the type *move on* the arriving point at epoch  $e_{i+1}$  is  $B$  with coordinates  $(x', y)$ . Knowing that the hidden objective function is  $f(x) = \sin(2x)$ , it is now possible to compute  $f(x')$ . Point  $C$  with coordinates  $(x', f(x'))$  should now be known.  $u$  is the **movement amount** and  $\overline{CB}$  equals to  $f(x') - f(x)$  is  $\delta$ . Using trigonometry it is possible to compute  $\alpha$ . The goal is to know the real amount of movement done by the agent on the function. In order to do this the first thing to do is to know how much is  $\overline{AE}$ . From the first theorem of Euclid it is possible to know that *in a right-angled triangle, the square constructed on a cathetus is equivalent to the rectangle that has for dimensions the hypotenuse and the projection of the cathetus on the hypotenuse*. This means that in a right-angled triangle, the cathetus is proportional medium between the hypotenuse and its own projection on it. According to this it is possible to write the following proportion :

$$h : u = u : AD. \quad (3.9)$$

So :

---

**Algorithm 8:** Angle computation in parametric movement case.

---

```

1 /* knowing  $(x, y)$  */ ;
2 /* knowing  $(x', y')$  */ ;
3 /* movementAmount =  $M$  */ ;
4 /* currentMax =  $\max f(x_n, y_n)$  */ ;
5
6  $z = f(x, y)$  ;
7  $z' = f(x', y')$  ;
8
9  $\delta = ((x' - x), (y' - y))$  ;
10  $\alpha = \arctan\left(\frac{\delta}{\text{movementAmount}}\right)$  ;
11
12  $\text{hypotenuse} = \frac{\text{movementAmount}}{\cos \alpha}$  ;
13
14  $\text{projectionOnHypotenuse} = \frac{\text{movementAmount} * \text{movementAmount}}{\text{hypotenuse}}$  ;
15
16  $\text{realMovementAmount} = \text{projectionOnHypotenuse} * \cos \alpha$  ;
17
18  $\Delta = z' - \text{currentMax}$ ;
```

---

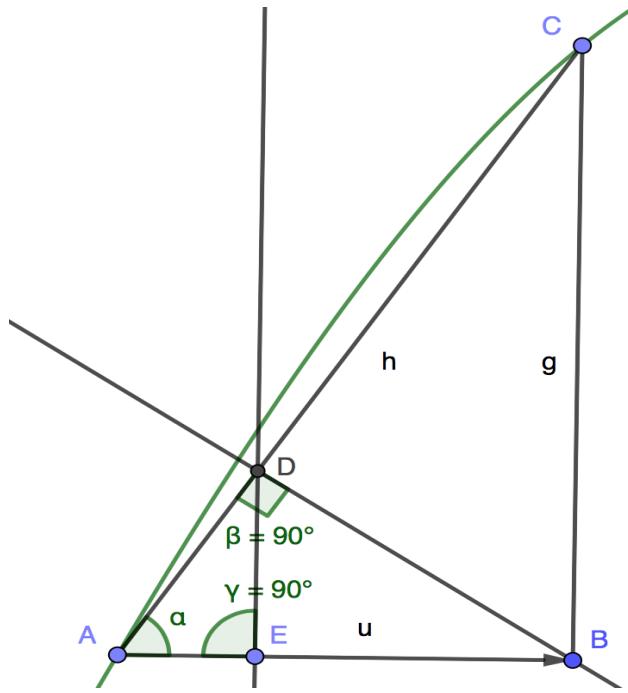


Figure 3.4: Parametric movement computations.

$$AD = \frac{u * u}{h}, \quad (3.10)$$

that is the same of :

$$\text{projectionOnHypotenuse} = \frac{\text{movementAmount} * \text{movementAmount}}{\text{hypotenuse}} \quad (3.11)$$

of algorithm 8. Now all elements useful to compute the *real movement amount* represented by segment  $\overline{AE}$  are known. It can be computed as follow:

$$\overline{AE} = \overline{DA} \cos \alpha. \quad (3.12)$$

Both parametric and linear approach are effective. Therefore the question is why select one method instead of the other one. The answer to this question is simple. Adopting the parametric approach the optimization process is longer but more accurate. This depends on the fact that the real amount of movement is correlated with the largeness of the angle. For this

reason, it is lesser at each time.

On the other hand using the parametric approach, the agent could be ideally trained on a set of specific functions and then it should make very good performances also on a generic, never previously explored, function using its knowledge about angles, actions and slope.

Adopting the linear approach the optimization process is slower but less accurate. This depends on the fact that the movement is as the crow flies. Using the linear approach the concept of function cannot be abstracted. This means that the agent has to be trained for every specific function in order to optimize its. Any change on a specific function, however small, invalidates the previous training.

### 3.0.4 BURLAP Java Library

In order to formalize and solve the problem previously described adopting an RL approach, BURLAP (Brown-UMBC Reinforcement Learning and Planning) Java library developed and maintained by James MacGlashan of Brown University, has been used.

BURLAP uses a highly flexible system for defining states and actions of nearly any kind of form, supporting discrete continuous, and relational domains. Planning and learning algorithms range from classic forward search planning to value function-based stochastic planning and learning algorithms [BUR].

In order to define customized MDPs, BURLAP offers a set of classes and interfaces (figure 3.5).

The main features offered by BURLAP are :

- **State** : implementing this interface the state variables of the MDP state space can be defined. An instance of this object will specify a single state from the state space [BUR] .
- **Action** : implementing this interface a possible action that the agent can select, can be defined. If the MDP action set is discrete and unparametrized, the provided concrete implementation `SimpleAction`, which defines an action entirely by a single string name, can be adopted [BUR].

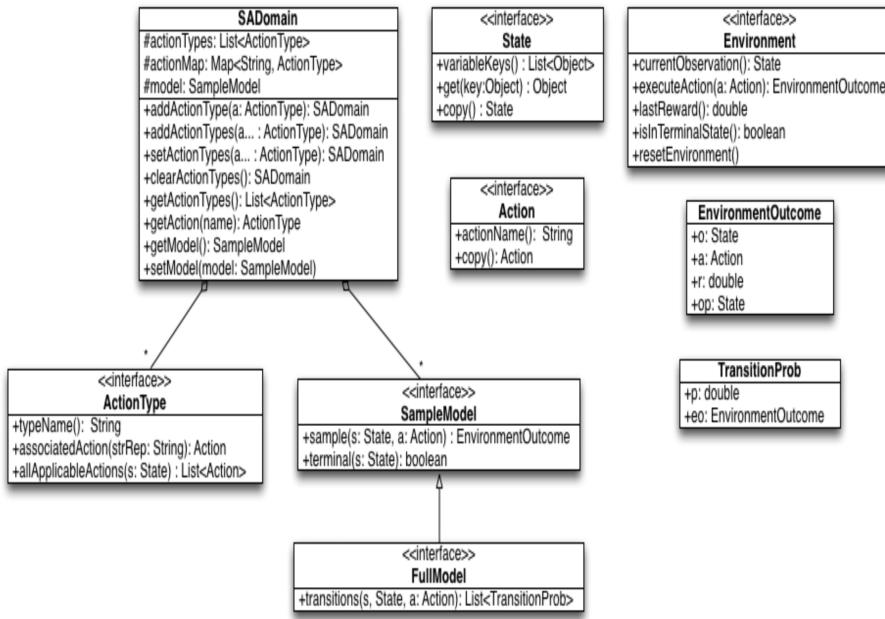


Figure 3.5: UML Diagram of the Java interfaces/classes for an MDP definition.

- **SampleModel** : implementing this interface the model of the MDP can be defined. This interface only requires to implement methods that can sample a transition: spit back out a possible next state and reward given a prior state and action taken [BUR].
- **Environment** : An MDP defines the nature of an environment, but ultimately, an agent will want to interact with an actual environment, either through learning or to execute a policy it computed from planning for the MDP. An environment has a specific state of the world that the agent can only modify by using the MDP actions. Implement this interface to provide an environment with which BURLAP agents can interact. If one defines a new MDP, then he will probably don't want to implement **Environment** himself and instead use the provided concrete **SimulatedEnvironment** class, which takes an **SADomain** with a **SampleModel**, and simulates an environment for it [BUR].

An extended definition of classes and interfaces can be found at [http:](http://)

//burlap.cs.brown.edu/doc/index.html.

In order to model the problem described above, starting from features offered by BURLAP, interfaces and abstract classes have been implemented and extended as shown in figure 3.6.

Because of the detailed explanation of choices made to model the problem given in the first section of this chapter and in **Appendix A**, the mere analysis of architectural definition is left to the reader. The full code can be found at (...).

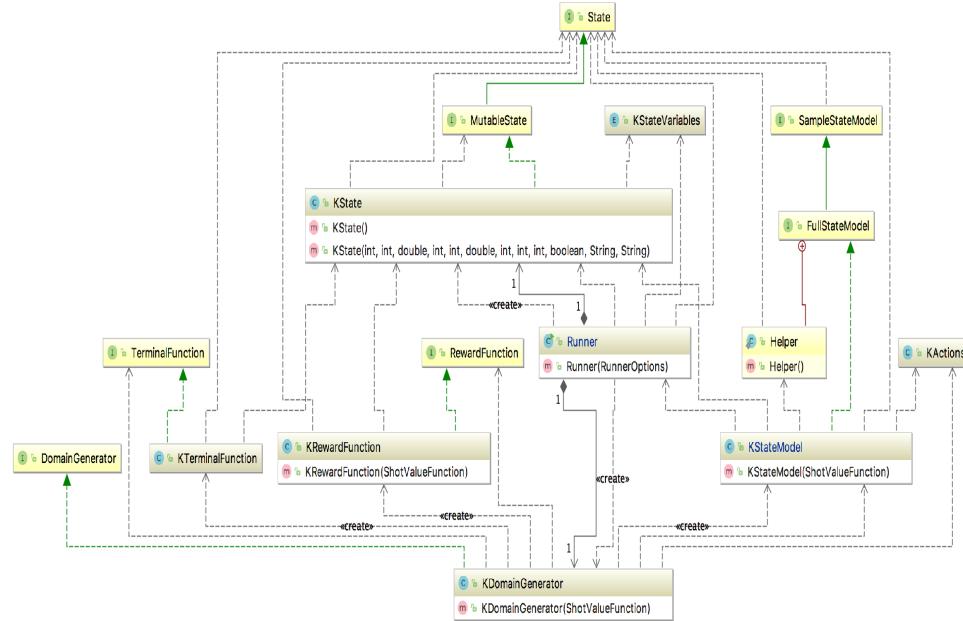


Figure 3.6: Class diagram of RL based optimization approach.

### 3.0.5 Humans in the Loop

Every time one gets into a rental car, he/she have to learn how hard to press the gas pedal for a given amount of acceleration. Solving this problem requires learning a relationship between two continuous variables. Over the past 50 years, several studies of *function learning* have shed light on how people come to understand continuous relationships (Carroll 1963; Brehmer

1971; 1974; Koh and Meyer 1991; Busemeyer et al. 1997; DeLosh et al. 1997; Kalish et al. 2004; McDaniel and Busemeyer 2005). It has become clear that people can learn and recall a wide variety of relationships, but demonstrate certain systematic biases that tell about the mental representations and implicit assumptions that humans employ when solving function learning problems. Several models have been developed to understand the cognitive mechanisms behind function learning. Those models tend to fall into two different theoretical camps. The first includes *rule-based* theories (e.g., Carroll, 1963, Brehmer, 1974, Koh and Meyer, 1991), which suggest that people learn an explicit function from a given family, such as polynomials (Carroll 1963; McDaniel and Busemeyer 2005) or power-law functions (Koh and Meyer 1991). This approach attributes rich representations to human learners, but has traditionally given limited treatment to how such representations could be acquired. A second approach includes *similarity-based* theories (e.g., DeLosh et al., 1997, Busemeyer et al., 1997), which focus on the idea that people learn by forming associations: if  $x$  is used to predict  $y$ , observations with similar  $x$  values should also have similar  $y$  values. Most recently, hybrids of these two approaches have been proposed (e.g., Kalish et al., 2004, McDaniel and Busemeyer, 2005), with an associative learning process that acts on explicitly represented functions. Almost all past research on computational models of function learning has been oriented towards understanding the psychological processes that underlie human performance, or the steps by which people update and deploy their mental representations of continuous relationships [LGWK15]. In this thesis, a different approach is taken, presenting a rational analysis of function learning.

It was seen that Gaussian Process regression is a powerful tool to model, explore, and exploit unknown functions. However, Gaussian Process regression might also be applied in a different, more psychological context, namely as a model of human cognition in general and function learning in particular. Recently, Lucas, Giffiths, Williams, and Kalish (2015) have proposed also use Gaussian process regression as a rational model of function learning that can explain various effects within the literature of human function learning. Within the exploration-exploitation context, Borji and Itti (2013) and Wu, Schultz, Speekenbrink, Nelson, and Meder (2017) showed that Gaussian Process-based optimization can explain how participants actively search for the best output when trying to optimize 2d functions [SSK17].

In the next few chapters, starting from those studies on 2d-function learning, and using the previously cited mathematical tools, an experiment

about 3d-function learning will be presented. After that, human learning techniques will be compared to RL agent's policy .

# Chapter 4

## Experimental Setting

The current chapter is divided into two sections. In the first one the benchmark of the experiment will be described. In the second one the same experiment applied to humans and its results will be discussed.

### 4.1 Benchmark

#### 4.1.1 Test Functions

The experiment described in this work is about maximize black-box functions adopting an RL based approach. The first important choice is about selecting suitable *test functions*. In applied mathematics, test functions, also known as *artificial landscapes*, are useful to evaluate characteristics of optimization algorithms. Four test functions have been chosen for this work:

- Himmelblau' s Function;
- Sphere Function;
- Beale Function;
- Styblinski-Tang' s Revised Function.

**Himmelblau' s Function** In mathematical optimization, Himmelblau' s function is a continuous, bivariate, multi-modal function. The original function is defined by:

$$f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2 \quad (4.1)$$

It has four local minima :

- $f(3.0, 2.0) = 0.0$ ;
- $f(-2.805118, 3.131312) = 0.0$ ;
- $f(-3.779310, -3.283186) = 0.0$ ;
- $f(3.584428, -1848126) = 0.0$ .

The function can be defined on any input domain but it is usually evaluated on  $x \in [-5, 5]$  and  $y \in [-5, 5]$ .

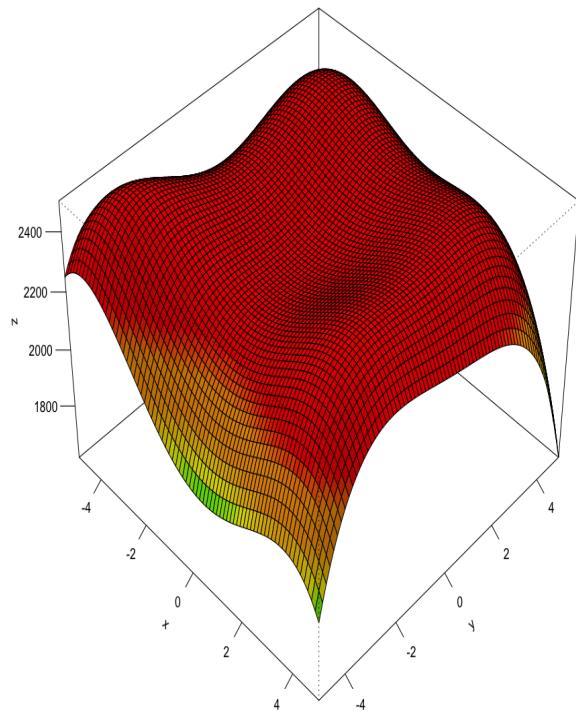


Figure 4.1: Customized Himmelblau' s Function.

Because of the aim to maximize the function has been inverted as follow

:

$$f(x, y) = -(x^2 + y - 11)^2 + (x + y^2 - 7)^2 \quad (4.2)$$

it has been picked up of 2500 units in order to have as less as possible negative values. So the final adopted function is :

$$f(x, y) = -(x^2 + y - 11)^2 + (x + y^2 - 7)^2 + 2500 \quad (4.3)$$

This function has its global maximum in  $f(x, y) = 2500$ .

In order to represent this customized version of Himmelblau' s Function using Java Graphical Environment, it has been mapped in a space of  $600 \times 600$  pixels and it has been properly rotated. The resulting contour plot is the one represented in figure 4.2.

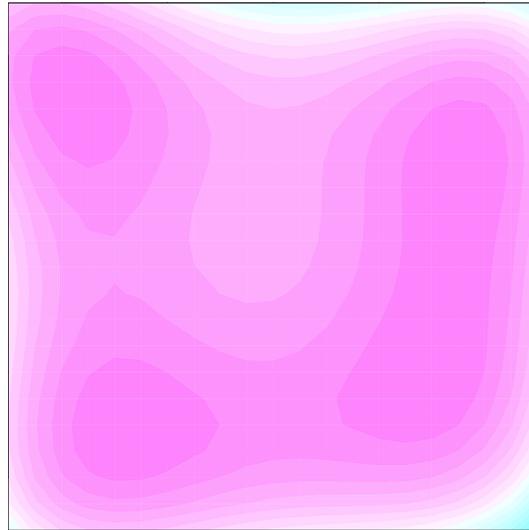


Figure 4.2: Contour plot of customized version of Himmelblau' s Function.

**Sphere Function** In mathematical optimization, Sphere function is a continuous, bivariate, multi-modal function.

The original function is defined by:

$$f(x, y) = (x^2 + y^2) \quad (4.4)$$

It has a global minimum in  $f(x, y) = 0$ .

The function can be defined on any input domain but it is usually evaluated on  $x \in [-10, 10]$  and  $y \in [-10, 10]$ .

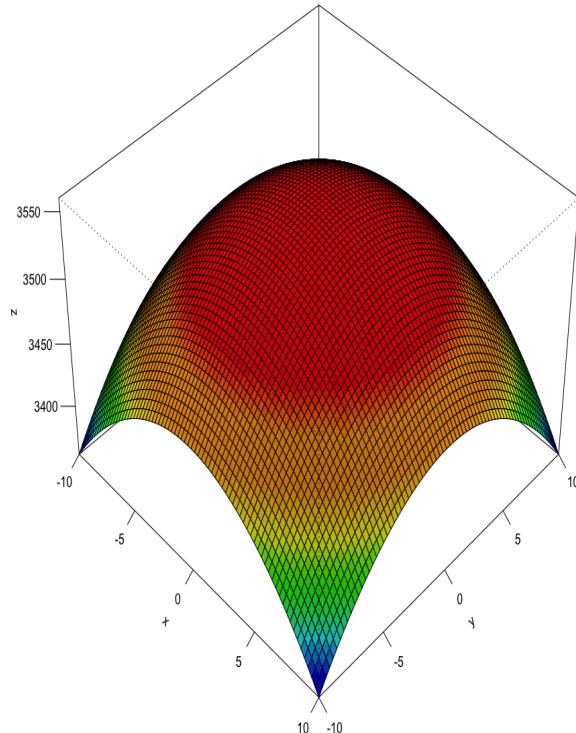


Figure 4.3: Customized Paraboloid of Revolution.

Because of the aim to maximize, the function has been inverted as follow :

$$f(x, y) = -(x^2 + y^2) \quad (4.5)$$

and it has been picked up of 3560 units in order to have as less as possible negative values. So the final adopted function is :

$$f(x, y) = -(x^2 + y^2) + 3560 \quad (4.6)$$

This function has its local maximum in  $f(x, y) = 3560$ .

In order to represent this customized version of Sphere function function using Java Graphical Environment, it has been mapped in a space of  $600 \times 600$  pixels and it has been properly rotated. The resulting contour plot is the one represented in figure 4.4

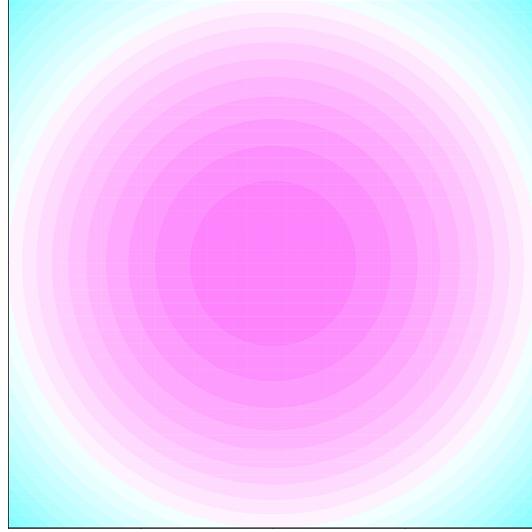


Figure 4.4: Contour plot of customized Parabolic Function.

**Beale Function** In mathematical optimization, Beale Function is a continuous, multi-modal function defined on a multi-dimensional space. The function is defined by:

$$f(x, y) = (1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2 \quad (4.7)$$

The function can be defined on any input domain but it is usually evaluated on  $x \in [-3, 3]$  and  $y \in [-3, 3]$ .

It has one global minimum in:  $f(x, y) = 0$ .

Because of the aim to maximize, it has been inverted as follow:

$$f(x, y) = -((1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2) \quad (4.8)$$

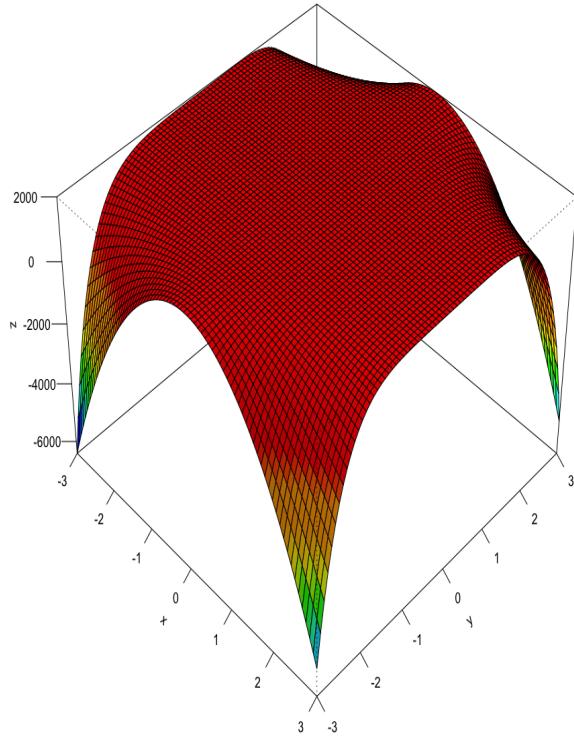


Figure 4.5: Customized Beale Function.

and it has been picked up of 2000 units in order to have as less as possible negative values. So the final adopted function is :

$$f(x, y) = -((1.5-x+xy)^2 + (2.25-x+xy^2)^2 + (2.625-x+xy^3)^2) + 2000 \quad (4.9)$$

This customized function has its global maximum in  $f(x, y) = 2000$ .

In order to represent this function using Java Graphical Environment it has been mapped in a space of  $600 \times 600$  pixels and it has been properly rotated. The resulting contour plot is the one represented in figure 4.6

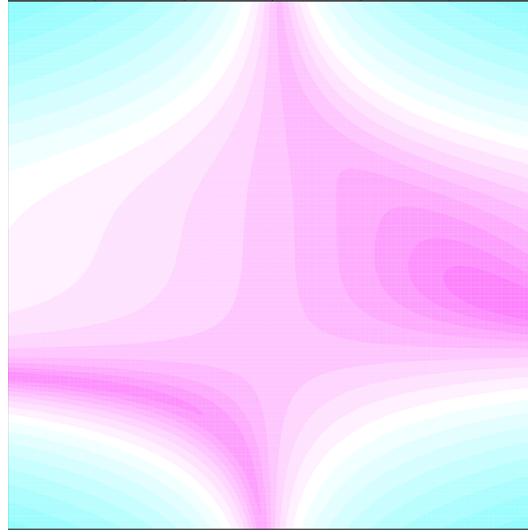


Figure 4.6: Contour plot of customized Beale Function.

**Styblinski-Tang Revised Function** In mathematical optimization, Styblinski-Tang Function is a continuous, multi-modal function defined on a multi-dimensional space. The original function is defined by :

$$f(x) = \frac{\sum_{i=1}^n x_i^4 - 16x_i^2 + 5x_i}{2} \quad (4.10)$$

In this thesis the bivariate version of the original function multiplied by two in order to make it less flatten, is considered:

$$f(x, y) = (x^4 - 16x^2 + 5x) + (y^4 - 16y^2 + 5y). \quad (4.11)$$

The function can be defined on any input domain but it is usually evaluated on  $x \in [-5, 5]$  and  $y \in [-5, 5]$ .

Because of the aim to maximize, it has been inverted as follow :

$$f(x, y) = -((x^4 - 16 * x^2 + 5 * x) + (y^4 - 16 * y^2 + 5 * y)) \quad (4.12)$$

and it has been picked up of 5000 units in order to have as less as possible negative values. So the final adopted function is:

$$f(x, y) = -((x^4 - 16 * x^2 + 5 * x) + (y^4 - 16 * y^2 + 5 * y)) + 5000 \quad (4.13)$$

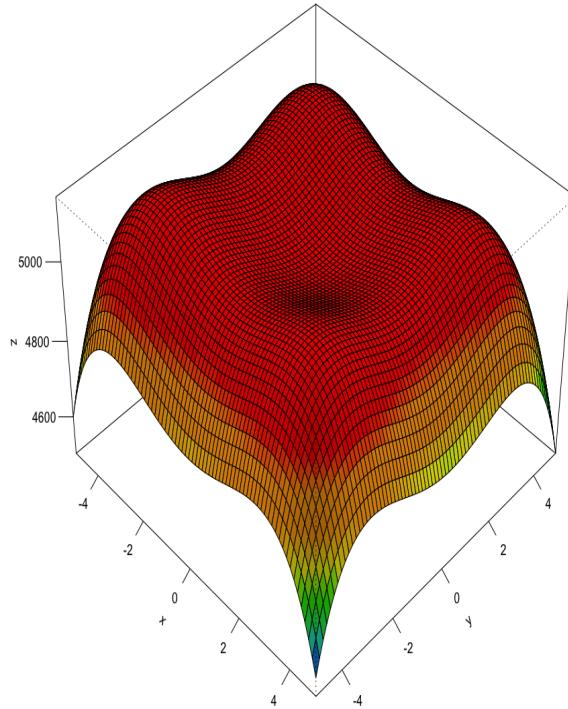


Figure 4.7: Customized Styblinski Function.

This customized function has its global maximum in  $f(x, y) = 5156.6638$ . In order to represent this function using Java Graphical Environment it has been mapped in a space of  $600 \times 600$  pixels and it has been properly rotated. The resulting contour plot is the one represented in figure 4.8.

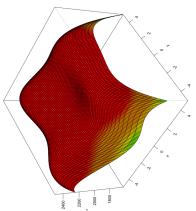
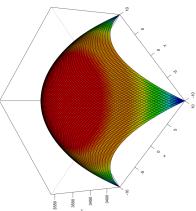
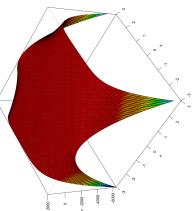
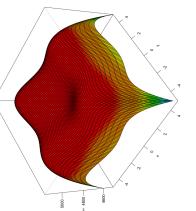
#### 4.1.2 Algorithm's Configuration

The RL algorithm employed is the SARSA( $\lambda$ ) one. It has only one configuration (summarized in table 4.2) and four different declinations (summarized in table 4.3).

In the proposed configuration a low  $\lambda$  value is employed. Typically it

4.1

Table 4.1: Objective Functions' Summary Table

Name	Formula	Domain	Optimal $f$	Sketch
Himmelblau	$f(x, y) = -((x^2 + y - 11)^2 + (x + y^2 - 7)^2) + 2500$	$x, y \in [-5; +5]$	2500	
Sphere	$f(x, y) = -(x^2 + y^2) + 3560$	$x, y \in [-10; +10]$	3560	
Beale	$f(x, y) = -((1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2) + 2000$	$x, y \in [-3; 3]$	1000	
Styblinski-Tang	$f(x, y) = -((x^4 - 16 * x^2 + 5 * x) + (y^4 - 16 * y^2 + 5 * y)) + 5000$	$x, y \in [-5; +5]$	5156.6638	

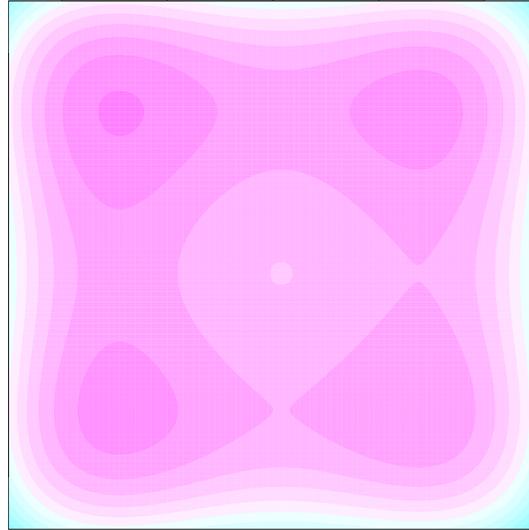


Figure 4.8: Contour plot of customized Styblinski Function.

	$\lambda$	$\epsilon$	qInit	Learning Rate	Discount Factor	Episodes	Epochs
Configuration 1	0.1	0.1	0	1	0.99	150	100

Table 4.2: Configuration Table

Declination	Type of search
Declination 1	Linear, not random search
Declination 2	Linear, random search
Declination 3	Parametric, not random search
Declination 4	Parametric, random search

Table 4.3: Declination table

should be used lower learning rate when a higher value of  $\lambda$  is employed. However, in this case, it has been decided to use the highest possible learning rate. The logical consequence to this choice is to decrease the  $\lambda$  value.

An  $\epsilon$  value equals to 0.1 means that 90% of selected actions depends on the history, whereas the remaining 10% represents a random choice having absolutely no relations with the previous history.

A discount factor set to 0.99 means the rewards will propagate as long as possible through time.

The main goal presented in this work is to compare human optimization performances with  $\text{SARSA}(\lambda)$  ones in a context of uncertainty. Assuming that each try has an high cost, the number of possible tries must be as low as possible assuring, at the same time, a reasonable training space. In order to satisfy these requirements, the number of episodes is set to 150. Each episode is composed by 100 epochs. Note that in addition to the 150 training episodes there is also a greedy one. In this last episode the  $\epsilon$  value and the *learning rate* parameter are set to 0. Instead, the number of epochs, the `qInit` parameter and the discount factor are left unchanged.

Each declination of  $\text{SARSA}(\lambda)$  corresponds to one combination of search techniques described in the previous chapter. In the implementation of random search, the BURLAP `RandomFactory` class has been used. `RandomFactory` allows to logically group various random generators. The `id` is set to 0 and the `seed` is set to 1.

## 4.2 Experiment on Humans

The test conducted on humans affects a cross-section of twenty-three people. In the selection process of test subjects the following three parameters with corresponding bounds were considered :

- **Sex**  $\in \{\text{Male}, \text{Female}\}$
- **Age**  $\in [4, 52]$
- **Calculus Knowledge**  $\in \{\text{Null}, \text{Basic}, \text{Advanced}\}$

Each one of the test subjects was placed in front of a personal computer. The following rules have been proposed :

- *You will be offered four different levels.*
- *The first three levels will need to be repeated three times.*
- *The last level can be executed only once.*
- *Look carefully at the legend before starting the experiment.*
- *At each click you will receive a score with a colour associated with it.*
- *The goal is to make as many points as possible with 15 clicks available for each game.*

- *No further explanation will be provided.*

After that a chromatic scale indicating how positive the reward was for each click was proposed.



Figure 4.9: Chromatic Scale.

A white screen (with a specific hidden function) of  $600 \times 600$  pixels with which to interact according to the rules previously listed was finally proposed.

# Chapter 5

## Results

This chapter is divided into two sections. In the first section results obtained from the experiment described in the first section of the previous chapter are analysed. For each function results obtained applying the four different declinations of **SARSA**( $\lambda$ ) algorithm are compared :

- Linear, not random search
- Linear, random search
- Parametric, not random search
- Parametric, random search

Three different plots for three different metrics are employed. The first one measures the *Euclidean distance*<sup>1</sup> between point  $(x_i, y_i)$  at epoch  $e_i$ ,  $\forall i \in n$  where  $n$  is the number of epochs of greedy episode, and point  $(x^*, y^*)$  which is the maximum. The second metric measures the *proximity in percentage*<sup>2</sup> between the current value function and the optimal one in the greedy

---

<sup>1</sup>In mathematics, the Euclidean distance is the "ordinary" straight-line distance between two points in an Euclidean space. It is computed as follow:

$$\text{Euclidean\_distance} = \sqrt{(x_{i+n} - x_i)^2 + (y_{i+n} - y_i)^2} \quad (5.1)$$

with  $n \geq 0$ .

Note that this performance metric is evaluated on the true function values but this information is not available for the optimization methods.

<sup>2</sup>The *proximity in percentage* between the current value function and the optimal one is computed as follow:

episode. The third one represents the *gap metric*<sup>3</sup>.

## 5.1 Sarsa( $\lambda$ )’s Performances on Test Functions

### 5.1.1 Himmelblau’ s Function

Plot (a) of figure 5.1 graphically represents performances of the four possible declinations of SARSA( $\lambda$ ) algorithm obtained considering the *Euclidean distance metric*.

Looking at the `linear, not random` declination’s performance, it can be possible to note a starting high Euclidean difference between the agent’s position and the global maximum. It decreases until epoch 5. Starting from this epoch, a recurrent pattern starts to be developed. This behaviour depends on the fact that once the agent has achieved a good position compared to the maximum, it has to continue to make movements until the terminal state is achieved. Each movement implies a reward. The agent selects a set of movements that maximizes the reward and systematically repeats them. The stability just described is also proved by the fact that the standard deviation of the current declination is the lowest of the set with a value of 25.3 (table 5.1).

`Linear, random` declination has a greater instability compared to the `linear, not random` one. In the first fifty epochs the *Euclidean distance* from the global maximum is constantly higher but then it starts to stabilize. Starting from the fiftieth epoch it develops a recurrent pattern and achieves best results. With an Euclidean distance mean of 71.6 pixels from the maximum, it has the best performance of the set (table 5.1).

---


$$\text{Proximity\_in\_percentage} = \frac{\text{currentValueFunction} \times 100}{\text{optimalValueFunction}} \quad (5.2)$$

Note that this performance metric is evaluated on the true function values but this information is not available to the optimization methods.

<sup>3</sup>This metric measures how effective each method is at finding the global maximum.

$$G_t = \frac{f(x^+, y^+) - f(x_1, y_1)}{f(x^*y^*) - f(x_1, y_1)} \quad (5.3)$$

Where  $x^+$  is the incumbent or best function sample found up to epoch  $e$ . The gap  $G_t$  will therefore be a number between 0 (indicating no improvement over the initial sample) and 1 (if the incumbent is the maximum). Note that this performance metric is evaluated on the true function values but this information is not available to the optimization methods. [HBdF11]

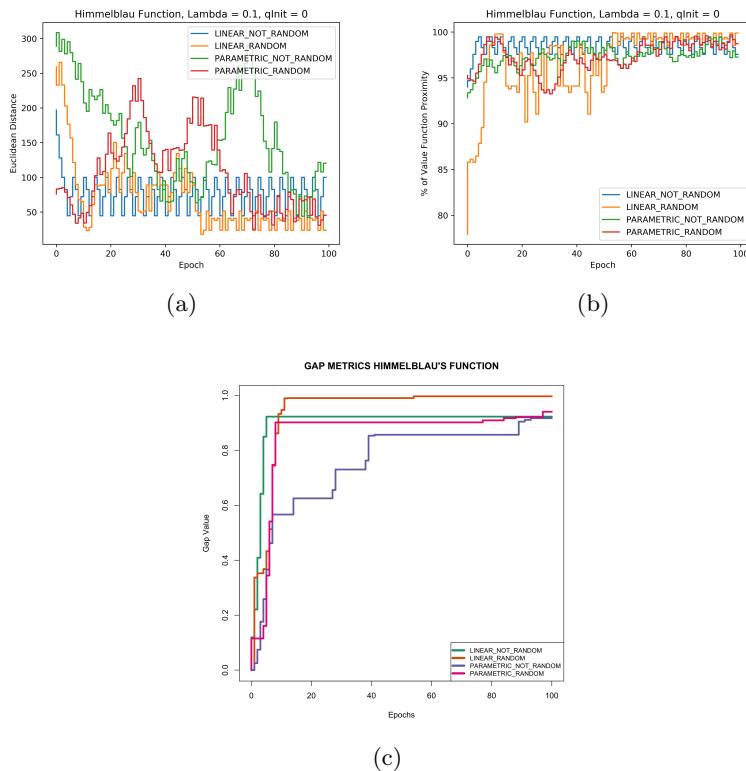


Figure 5.1: Himmelblau's Function.

The parametric, not random declination has worst performances of the declinations' set. It is unable to minimize the Euclidean distance from the global maximum and it is highly unstable. Looking at table 5.1 it is possible to note that the current declination as highest values in all statistic measures.

Finally, the parametric, random declination selects a starting point near to the global maximum but it is unable to maintain a low Euclidean distance from it until epoch seventy. Starting from this epoch it develops a recurrent pattern similar to one previously described. Despite of this it is possible to say that its performances are worse compared to the first two ones.

Slower convergence of linear declinations compared to parametric ones depends on the fact that in the first case, as already explained in chapter 3, greater movements are done at each epoch and a lower training time is

required to achieve the maximum. The amount of parametric movements are strictly conditioned by the function's shape. If this represents an obstacle to a fast convergence, however it is an incentive to an higher precision.

Himmelblau' s Function	Mean (pixels)	Standard deviation
Linear not Random	77.7	25.3
Linear Random	71.6	52.2
Parametric not Random	158.6	71.3
Parametric Random	105.5	56.0

Table 5.1: Euclidean Metric's Performances

Plot (b) of figure 5.1 graphically represents performances of the four possible declinations of  $\text{SARSA}(\lambda)$  algorithm obtained considering the *Percentage of Value Function's Proximity metric*.

All declinations, except for the `linear`, `random` one, start with an high level of value function's proximity to the maximum. It depends on the function's shape.

Once it has achieved a good enough proximal point to the maximum, the `linear`, `not random` declination maintains always the same vale function's proximity level developing a recurrent pattern. The average level of proximity achieved by this declination is the best of the set (table 5.2). The high stability of the current declination is proved by the lowest standard deviation of the set (table 5.2).

The `linear`, `random` declination starts from a point not very close to the maximum in terms of value function's value. Its great instability reveals an inadequate training space. It starts to stabilize starting from the fiftieth epoch. Its starting instability is proved by the highest standard deviation of the set (table 5.2).

The `parametric`, `not random` declination reveals a general stability. The growth of proximity is slower then the corresponding linear declination because of the reduced amount of movement depending on the function's shape.

Finally, the `parametric`, `random` declination, as the corresponding linear declination, shows an initial instability. In general performances are more stable than the linear corresponding ones because of the lower amount of movement done at each epoch. Starting from the sixty-fifth epoch the proximity slowly grows and stabilizes.

Himmelblau's Function	Mean (%)	Standard deviation
Linear not Random	98.5	0.96
Linear Random	96.8	4.0
Parametric not Random	97.4	1.2
Parametric Random	97.3	1.6

Table 5.2: Value Function's Proximity.

Plot (c) of figure 5.1 graphically represents performances of the four possible declinations of **SARSA**( $\lambda$ ) algorithm obtained considering the *gap metric*. According to this plot the most effective declination of **SARSA**( $\lambda$ ) algorithm in finding the maximum is the **linear, random** one. Looking at table 5.9 it is possible to note that **linear, random** declination achieves point  $(3.67, -1.57)$  with a value function of 2498.457 out of 2500.0.

### 5.1.2 Sphere Function

Before starting to analyse performances of different declinations of **SARSA**( $\lambda$ ) algorithm in maximizing the Sphere function, it is important to underline that the starting point for the current function in non-random declinations is  $(0, 0)$ . The reason for this choice is that  $f(300, 300)$  is itself the maximum of the function.

Plot (a) of figure 5.2 graphically represents performances of the four possible declinations of **SARSA**( $\lambda$ ) algorithm obtained considering the *Euclidean distance metric*. Looking at the **linear, not random** declination it is possible to note a constant, high Euclidean difference between the agent's position and the global maximum. This distance never decreases and still from the first epochs, a recurrent pattern starts to be developed. The inefficiency of the algorithm's declination depends on an inadequate training space. The relationship between different explored states and total possible states does not permit the agent to develop an efficient policy. Bad performance's results of the current declination are also certified by an high metric's average (table 5.3). On the other hand the great stability of this declination is underlined by the lowest standard deviation of the set (table 5.3).

**Linear, random** declination starts from a point relatively close to the maximum, but, one more time, the inadequate training space prevents its to develop a rally optimal policy. This declination is the less efficient one

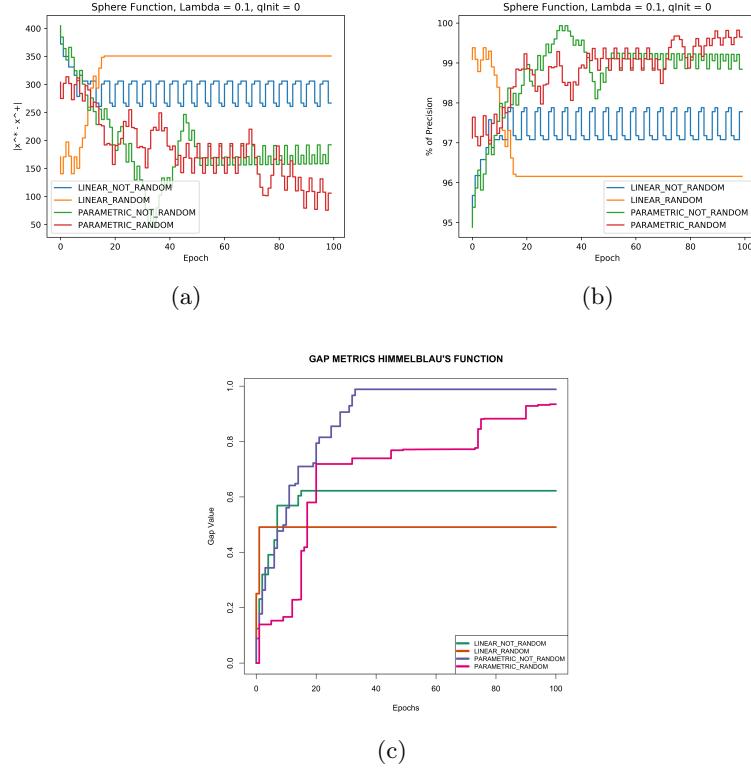


Figure 5.2: Sphere Function.

with an average Euclidean distance from the maximum of 328.9 pixels.

Better performances can be registered looking at parametric declinations of the algorithm. The **parametric, not random** declination starts with a physiological high Euclidean difference from the maximum. It rapidly decreases until the thirty-fifth epoch and then it has another soft increase. Starting from the fifty-fifth epoch a recurrent pattern starts to be developed. Despite of the fact that this declination is the most unstable one with a standard deviation of 67.9 (table 5.3), it has one of the lowest average Euclidean distance from the maximum.

The **parametric, random** declination of SARSA( $\lambda$ ) algorithm has the best performance of the set. After a starting, considerable instability it reaches an Euclidean difference of around fifty pixels. The good performance is also proved by the fact that the average Euclidean distance from the maximum is 185.3 pixels.

Sphere Function	Mean (pixels)	Standard deviation
Linear not Random	293.2	25.6
Linear Random	328.9	56.3
Parametric not Random	190.7	67.9
Parametric Random	185.3	58.0

Table 5.3: Euclidean Metric's Performances

Plot (b) of figure 5.2 graphically represents performances of the four possible declinations of SARSA( $\lambda$ ) algorithm obtained considering the *Percentage of Value Function's Proximity metric*.

The **linear, not random** declination starts from a value function's proximity of about 97.5% from the maximum. It is completely unable to improve this percentage. Already starting from the first epochs it develops a recurrent pattern. This behaviour allows its to have the lowest standard deviation of the declinations' set (table 5.4). Unfortunately the value function's proximity mean is the second lowest one.

The **linear, random** declination starts from a very high value function's proximity level. Because of an inadequate training space it decreases until the twentieth epoch. Starting from this epoch it stars to stabilize. Effects of this behaviour on the average value function's proximity to the maximum and on average standard deviation are relevant.

The **parametric, not random** declination is the most unstable one. It starts with a low value function's proximity level and slowly grows achieving peaks of about 100%. Starting from the fiftieth epoch it stabilizes developing a recurrent pattern. The initial high instability has effects on average value function's proximity variance making its the highest one (table 5.4).

The **parametric, random** declination of SARSA( $\lambda$ ) algorithm is the most efficient one. It starts with a value function's proximity level of about 97% and slowly grows first stabilizing around 99% and than growing more achieving peaks of 99.8%. Looking at table 5.4 it is possible to note that the average value function's proximity is the highest one with a level of 98.8%.

In this case it is possible to underline that parametric declinations do better than linear ones. In particular the random starting position allows the agent to develop a better policy. The main reason of the failure of non random declinations is the reduced training space compared to the distance between the starting point and the maximum.

Plot (c) of figure 5.1 graphically represents performances of the four

Sphere Function	Mean (%)	Standard deviation
Linear not Random	97.3	0.49
Linear Random	95.5	0.92
Parametric not Random	98.7	0.98
Parametric Random	98.8	0.72

Table 5.4: Value Function's Proximity.

possible declinations of  $\text{SARSA}(\lambda)$  algorithm obtained using the *gap metric*. According to this plot the most effective declination of  $\text{SARSA}(\lambda)$  algorithm in finding the maximum is the **parametric, not random** one. Looking at table 5.9 it is possible to note that **parametric, not random** declination achieves point  $(-1.033, -1.099)$  with a value function of 3557.722 out of 3560.0.

### 5.1.3 Beale Function

Giving an overview to figure 5.3 it is possible to see that Beale Function represents the case in which  $\text{SARSA}(\lambda)$  algorithm's declinations do worst.

Plot (a) of figure 5.3 graphically represents performances of the four possible declinations of  $\text{SARSA}(\lambda)$  algorithm obtained considering the *Euclidean distance metric*. Looking at the **linear, not random** declination it is possible to note a growing performance's worsening until tenth epoch. Starting from epoch ten algorithm's declination starts to stabilize developing a recurrent pattern. Using this declination, the average Euclidean distance from the maximum is 462.4 pixels.

**Linear, random** declination seems to be the one which performs better. Its average Euclidean distance from the maximum is 201.6 pixels and its standard deviation is the lowest of the set with a value of 15.4. This means that it has a relatively good stability. Looking at plot (a) what just said can be easily proved. The current declination starts from a point relatively near to the maximum and never departs from its developing a recurrent pattern.

Looking at **parametric, not random** declination of  $\text{SARSA}(\lambda)$  algorithm it is possible to note an high instability. In the very first epochs the distance between the starting point and the maximum slowly decreases, but starting from epoch thirty it starts to messily increase achieving a level near to 350 pixels. This high instability is reflected also in the high level of standard deviation (table 5.5).

**Parametric, random** declination combines instability and inaccuracy.

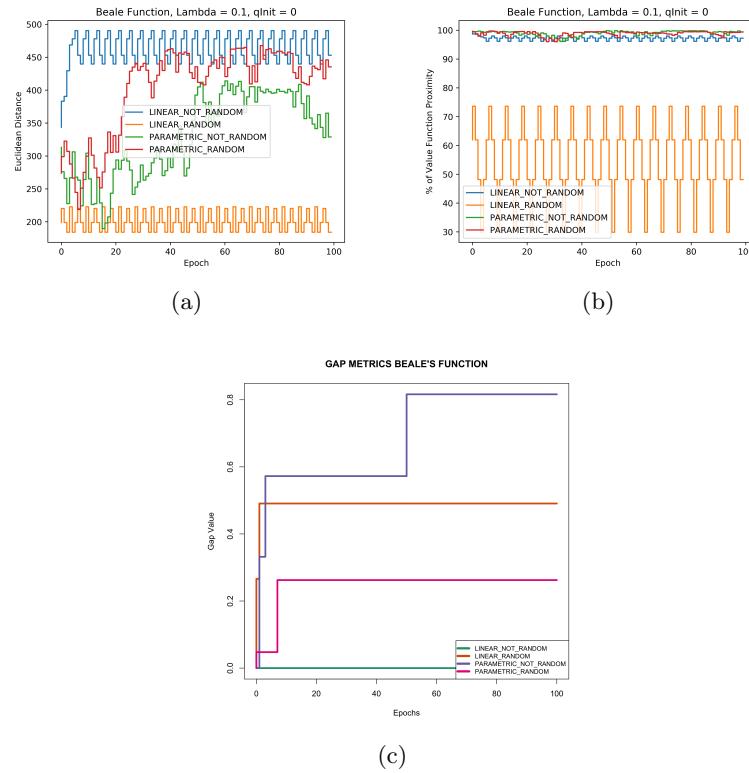


Figure 5.3: Beale Function.

It starts from a distance of about 300 pixels and messily increase its. Its performance is the worst one with an average Euclidean distance of 402.99 pixels and a standard deviation of 66.6 (table 5.5).

Generally speaking it is possible to affirm that those bad performances are the result of an inadequate training space and of a particular function's shape. In addition to this, as already explained, parametric declinations do worst than corresponding ones because of the reduced amount of movement at each epoch. The relatively better performance of `linear`, `random` declination depends on the possibility to select a starting point closer to the maximum.

Looking at plot (b) of figure 5.3 it is possible to note an inconsistency with what said until now. Only `linear`, `random` configuration makes an highly swinging performance varying between a percentage of value function'

Beale Function	Mean (pixels)	Standard deviation
Linear not Random	462.4	25.7
Linear Random	201.6	15.4
Parametric not Random	329.8	62.7
Parametric Random	402.99	66.6

Table 5.5: Euclidean Metric's Performances

s proximity of about 75% to a percentage of 30%. All other declinations make performances around 100% of value function's proximity. Even looking at table 5.6 this impression is proved by the high standard deviation of `linear`, `random` configuration and by the high average of value function's proximity of all other declinations. These apparently good performances depends on the function's shape that is very flat.

Beale Function	Mean (%)	Standard deviation
Linear not Random	97.25	0.70
Linear Random	54.26	13.89
Parametric not Random	99.06	0.81
Parametric Random	98.88	0.798

Table 5.6: Value Function's Proximity.

Deception can also be revealed by observing plot (c) of figure 5.3. The relatively best performing declination is the `parametric, not random` one which achieves an effectiveness in finding the maximum of about 0.8 out of 1.0. It is also proved by table 5.9. `Parametric, not random` declination achieves point  $(-0.77, 1.5599)$  with a value function of 1997.392 out of 2000. The worst performing one is the `linear, not random` declination with a constant *gap value* equals to 0.

#### 5.1.4 Styblinski-Tang Revised Function

Plot (a) of figure 5.4 graphically represents performances of the four possible declinations of  $\text{SARSA}(\lambda)$  algorithm obtained considering the *Euclidean distance metric*.

Looking at the `linear, not random` declination's performance, the Euclidean difference between the agent's position and the global maximum is about 200 pixels. Staring from the twentieth epoch it starts to stabilize

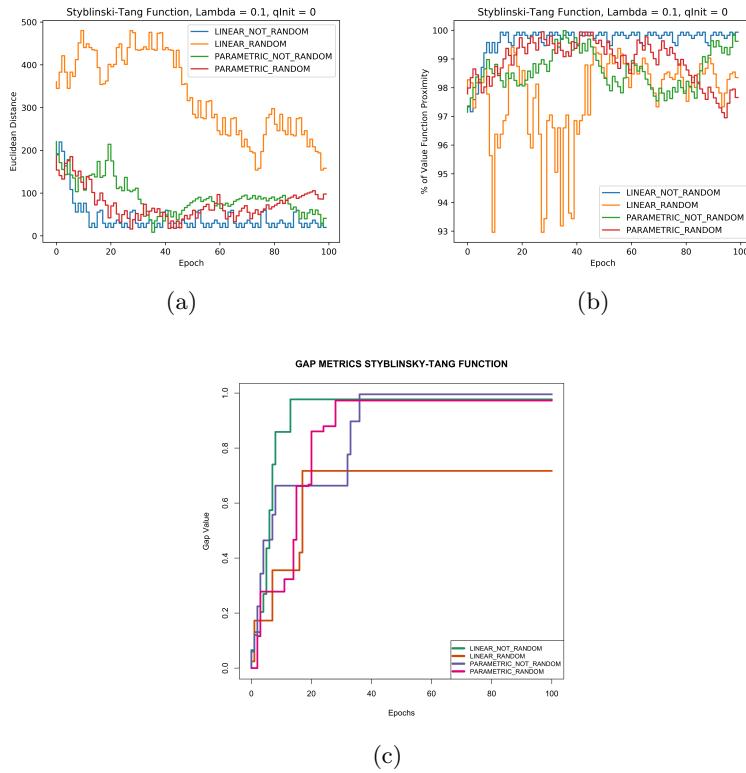


Figure 5.4: Styblinski Function.

around an Euclidean distance level from the maximum of about 40 pixels. Note that in less than twenty epochs the Euclidean distance from the maximum is about the 20% of the starting one. The current declination also reveals a great stability. With an average Euclidean distance of 43.63 pixels and a standard deviation of 41.05, the `linear, not random` declination is the best of the set (table 5.7).

A decisively worse performance is given by the `linear, random` declination. It starts with an Euclidean distance from the maximum of about 360 pixels and messily decreases until achieving a level of 180 pixels. The noise just described is also proved by the highest average Euclidean distance and by the highest standard deviation of the set (table 5.7).

Last two, parametric declinations give for the first time the possibility to really see behaviour of parametric declination. Starting from about the fiftieth epoch, `parametric, not random` movement starts to imperceptibly

increase. This lower amount of movement compared to the non parametric movements depends on the definition of the parametric movement itself. The same behaviour can be seen in `parametric`, `random` declination starting from epoch 75. This declination performs better than the `parametric`, `not random` one in the first epochs but it performs worse in the last ones. With an average Euclidean distance from the maximum of 75.26 pixels, and with an average standard deviation of 37.34, the `parametric`, `random` declinations is generally better than the `parametric`, `not random` one.

Styblinski-Tang Function	Mean (pixels)	Standard deviation
Linear not Random	43.63	41.05
Linear Random	330.16	95.39
Parametric not Random	91.79	42.24
Parametric Random	75.26	37.34

Table 5.7: Euclidean Metric's Performances

The performance just described is also proved by the percentage of value function's proximity to the maximum. With an average mean of 99.63% of precision, and a standard deviation of 0.57, the `linear`, `not random` declination is the best of the set. General good performances are also made by all other declinations except for the `linear`, `random` one. Looking at plot (b) of figure 5.4 it is possible to note an high instability of this last declination especially until epoch fifty. It then starts to stabilize varying between a proximity level of 97% and 99%.

Styblinski-Tang Function	Mean (%)	Standard deviation
Linear not Random	99.63	0.57
Linear Random	97.73	1.46
Parametric not Random	98.61	0.66
Parametric Random	98.90	0.74

Table 5.8: Value Function's Proximity.

Looking at the third plot it is possible to note that the best performance is the one of the `parametric`, `not random` declination. This result obviously depends on the fact that around the thirty-ninth epoch it reaches the maximum for then decisively distance itself from it. This behaviour is again due to the absence of an adequate training space. This lack prevents the agent from developing an optimal policy.

Table 5.9: Best Soundings

Name	Linear	Not Random	Linear Random	Parametric	Not Random	Parametric Random
Himmelblau	2486.938 (2.67, 1.34)		2498.457 (3.67, -1.57)	2485.972 (-2.1, 3.30)		2492.111 (-2.283, 3.234)
Sphere	3484.44 (-8.8, -0.67)		3537.986 (-3.24, 3.40)	3557.722 (-1.034, -1.099)		3553.626 (1.69, 1.88)
Beale	1985.797 (0, 0)		1472.184 (1.69, 2.269)	1997.392 (-0.77, 1.559)		1994.983 (1.309, -0.89)
Styblinski-Tang	5153.086 (-2.67, -2.67)		5126.192 (3, -2.97)	5155.97 (-2.77, -2.95)		5154.053 (-3.17, -2.99)



# Chapter 6

# Appendix A

One of the most amazing aspect of the work presented in this thesis is the manner in which the goal has been reached. It is important to clarify that during the first weeks, it was not easy to understand the real target of the work. This initial uncertainty led to implement some strange solutions. Even if some of these approaches were not useful for the goal described in the previous chapters of this thesis, although they were very efficient in different contexts. For this reason the full history of the current research activity has been here recorded. In so doing, ten different main steps have been identified.

**First Step** In the first attempt to reach the solution of the problem, the agent was allowed to take five different actions :

- Move North
- Move South
- Move West
- Move East
- Dig.

In so doing, the agent could take a *movement action* or a *dig action*. This choice was inspired by the possibility for humans to make movements for no specific reasons. It was later clear that this implementation was redundant. Humans move themselves for no apparently reason before decide to take a goal-directed action. A RL agent *decides* what action taking in

the next step immediately after have completed the previous action without needing *time to think*.

In this step the state was represented by three elements :

- Coordinate  $x$
- Coordinate  $y$
- Number of remaining *function soundings*

Unfortunately such a state representation prevents from the possibility to generalize the training. Given two different bivariate functions  $f(x, y)$  and  $g(x, y)$  and supposing to have trained the RL agent using function  $f(x, y)$ , when one decides to sound  $g(x, y)$  on point  $x$  and  $y$  already sounded in  $f(x, y)$ , the result could be different (figure 6.1).

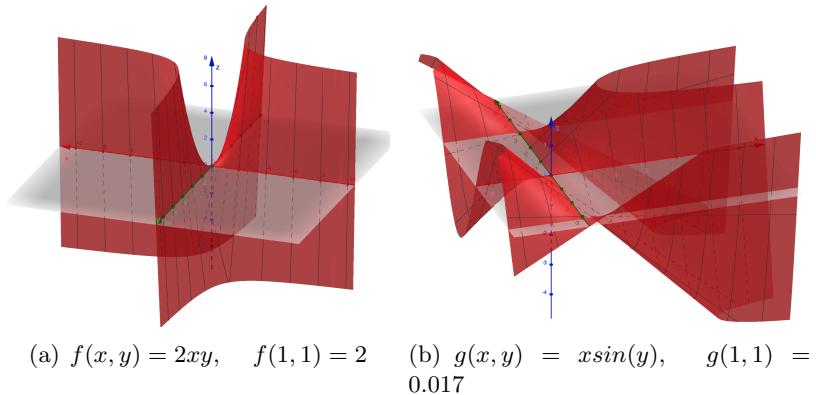


Figure 6.1: Different value functions for different functions.

In addition to this, it became soon clear that the information concerning the number of remaining function's soundings is redundant. It did not add any information to knowledge. In the implementation described above, the agent was in a terminal state when the number of remaining function soundings was equal to zero.

**Second Step** The second step was globally equal to the first one except for two fundamental aspects. In the current step the possibility for the agent to make movements of different sizes was introduced. This possibility allowed its to reach the goal with a lower training time.

The reward function was also redefined. It was equal to the percentage of proximity to the maximum of the function. It became soon clear that a so defined reward function was completely wrong. In a context of black-box optimization the maximum of the function is unknown and it was not possible to compute the percentage of proximity to the maximum.

**Third Step** In the third step an additionally change to the previous one was made. At the end of each episode, the agent was artificially forced to restart from the coordinates of the previous episode's best rewarded epoch. Results given by the combination of the last two approaches resulted to be very fruitful in the context well-known function optimization but not in the context described in this thesis. This choice forced the agent to take a specific action, breaking, in so doing, for at least one time the concept of environment formalized through a Markov Decision Process.

**Fourth Step** In the fourth step performances obtained through changes made in the second and third step were improved. This was possible imposing the agent to start each episode from the best rewarded epoch from all previous episodes.

**Fifth Step** The fifth step is one of the most important because of innovations introduced. The first change made regarded actions. In this step the possibility for the agent to directly choose the size of the movement was introduced. There were still four different possible *movement actions* :

- Move North
- Move South
- Move West
- Move East

but the agent could select to move itself in one direction of :

- *movement amount*;
- *movement amount*  $\times 2$ ;
- *movement amount*  $\times 3$ .

The unit of measurement of *movement amount* was the *pixel*. The basic amount of movement was 40 pixels. The real movement amount over function was computed through algorithm 6 in Chapter 3. It was proved that using this set of *action movements* the algorithm converged more quickly to the maximum of the function but the precision was lower.

In addition to this, the possibility for the agent to distinguish between *move action* and *dig action* was removed. Each time the agent made a movement it also sounded the function obtaining a reward. The reward was positive if its proximity to the global maximum was in percentage  $\geq 75\%$ , negative in other cases. As already explained, the disadvantage of the last choice is that in the context described in this thesis the objective function is unknown so it is not possible to compute the percentage of proximity to the maximum.

**Sixth Step** In this step the concept of *delta* was introduced. *delta* was computed as the difference between the value of the function sounded at time  $t$  and the maximum value sounded until time  $t$ .

In this step an higher control about *exploration* and *exploitation* was achieved. In order to do this a decreasing  $\epsilon$  factor between episodes was introduced. Different decreasing functions were tested and it was clear that the best one was the logarithmic one.

In this step the reward was still based on the concept of percentage with an engagement of the *delta*. To be more specific if the *delta* was  $\leq 0$  or the percentage was  $\leq 78\%$  the reward was negative, otherwise it was positive.

**Seventh Step** In this step the definition of the reward function was improved. The concept of percentage of proximity to maximum was removed and a reward based exclusively on the difference between the sounded value function at time  $t$  and the maximum until time  $t$  was introduced. This was a little but fundamental adjustment because it removed a relevant mistake from the project.

**Eighth Step** This step can be defined as an experimental step. Here the agent was made able to explore more. In order to do this the possibility to start each episode from a random point was introduced. Through this experiment it was proved that more exploration gave the possibility to prevent many limit cases but it did not always entail an improvement in algorithm performances in terms of time and precision.

**Ninth Step** The current one is the second revolutionary step. It introduces two innovative elements: the *parametric movement* and the *stop action*. As explained in chapter 3, the parametric movement allows to compute the *real agent's movement amount* keeping in consideration an approximation of function's shape.

Speaking about the second innovation introduced, it is important to say that when *stop action* was selected, the simulation of the current episode was interrupted. Consequences from introduction of this new action were very interesting. Given an inadequate training space, the *stop action* led to an immediate stop. This depended on the fact that rather than receiving more negative rewards before reaching a maximum better than the current one (i.e. the starting point itself), it computed that the optimal policy consisted in staying around the starting point itself, still moving according to a recurring pattern that minimizes the negative rewards.

**Tenth Step** This represents the last step before the final version of the project. In this version the state was revolutionised. It was now composed of the last three angles of improvements computed as described in chapter 3. Angles varied between  $-\pi$  and  $+\pi$  because also a worsening is here considered. Angles were rounded to the nearest multiple of two. The choices introduced in this step were unexpectedly not fruitful. The failure depended on two factors. Last three angles as the unique component of the state were not enough and the rounding cannot be chosen a priori. The correct choice of rounding should depend on the *Lipschitz Constant*.

Let's consider a single-variable function  $f(x)$  for  $x$  inside a domain  $D$ . The magnitude of the slope of  $f(x)$  for two points  $(x_1, f(x_1))$  and  $(x_2, f(x_2))$  is thus

$$\left| \frac{f(x_1) - f(x_2)}{x_1 - x_2} \right|. \quad (6.1)$$

A non-negative real number  $L$ , which is the smallest upper bound of the slope is called the *Lipschitz Constant* :

$$L = \left| \frac{f(x_1) - f(x_2)}{x_1 - x_2} \right| = \sup \left| \frac{df}{dx} \right|. \quad (6.2)$$

The Lipschitz Constant limits how fast the function can change. If there is no Lipschitz Constant, the function can change extremely fast without border, that is, the function becomes discontinuous[Lip].

For multivariate function the *Lipschitz Constant* is defined as:

$$L_{x_n} = \sup \left| \frac{\partial f}{\partial x_n} \right|. \quad (6.3)$$

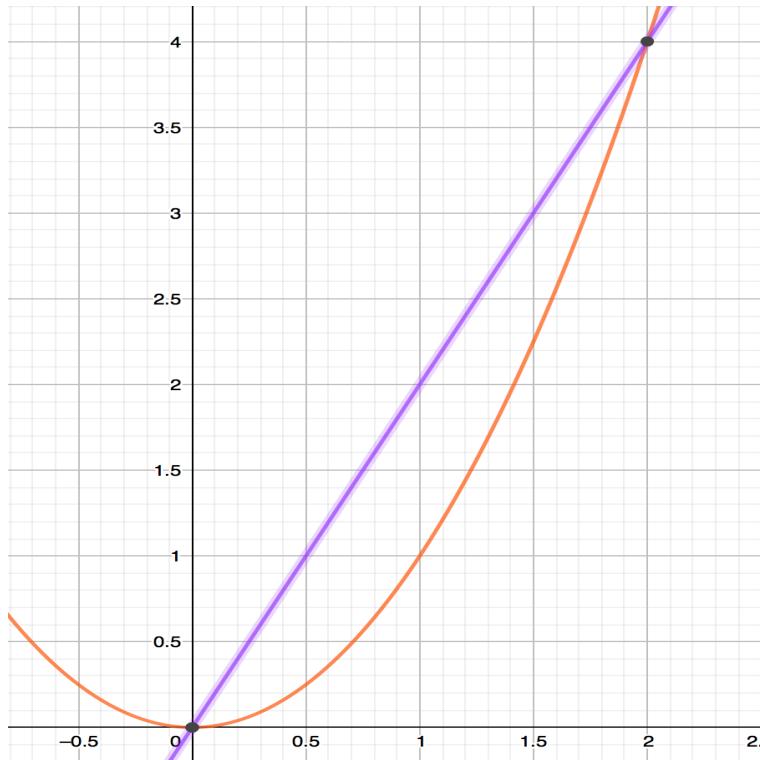


Figure 6.2: Let's consider a 2d-function  $f(x) = x^2$  defined in the domain  $x \in [0, 2]$ . If  $| \frac{df}{dx} f(x) | = 2x$ , then  $L = \sup | \frac{df}{dx} | = 2 \times 2 = 4$ .

# Bibliography

- [AS08] Ryan P. Adams and Oliver Stegle. Gaussian process product models for nonparametric nonstationarity. In *ICML*, 2008.
- [Bay] <https://towardsdatascience.com/shallow-understanding-on-bayesian-optimization-324b6c1f7083>. Accessed: 2018-05-20.
- [BUR] <http://burlap.cs.brown.edu/tutorials/bd/p2.html>. Accessed: 2018-05-20.
- [HBdF11] Matthew Hoffman, Eric Brochu, and Nando de Freitas. Portfolio allocation for bayesian optimization. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, UAI'11, pages 327–336, Arlington, Virginia, United States, 2011. AUAI Press.
- [HL01] Frederick S. Hillier and Gerald J. Lieberman. *Introduction to Operations Research*. McGraw-Hill, New York, NY, USA, seventh edition, 2001.
- [KLM96] Leslie Pack Kaelbling, Michael L. Littman, and Andrew W. Moore. Reinforcement learning: A survey. *CoRR*, cs.AI/9605103, 1996.
- [Kon09] Takis Konstantopoulos. Markov chains and random walks. 2009.
- [LGWK15] Christopher G. Lucas, Thomas L. Griffiths, Joseph J. Williams, and Michael L. Kalish. A rational model of function learning. *Psychonomic Bulletin & Review*, 22(5):1193–1215, Oct 2015.
- [Lip] [https://www.eee.hku.hk/~msang/Numerical\\_Lipschitz.pdf](https://www.eee.hku.hk/~msang/Numerical_Lipschitz.pdf). Accessed: 2018-06-05.

- [LM16] Ke Li and Jitendra Malik. Learning to optimize. *CoRR*, abs/1606.01885, 2016.
- [Mit97] Tom M. Mitchell. *Machine Learning*. 1997.
- [MND] [https://en.wikipedia.org/wiki/Multivariate\\_normal\\_distribution/media/File:Multivariate\\_Gaussian.png](https://en.wikipedia.org/wiki/Multivariate_normal_distribution/media/File:Multivariate_Gaussian.png). Accessed: 2018-05-20.
- [NFK06] Yuriy Nevmyvaka, Yi Feng, and Michael Kearns. Reinforcement learning for optimized trade execution. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, pages 673–680, New York, NY, USA, 2006. ACM.
- [Pow07] Warren B. Powell. *Approximate Dynamic Programming: Solving the Curses of Dimensionality (Wiley Series in Probability and Statistics)*. Wiley-Interscience, 2007.
- [Put94] M. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley and Sons, 1994.
- [SB10] Olivier Sigaud and Olivier Buffet. *Markov Decision Processes in Artificial Intelligence*. Wiley-IEEE Press, 2010.
- [SB18] R. S. Sutton and A. G. Barto. *Reinforcement Learning : An Introduction*. 2018.
- [SLA12] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 2951–2959. Curran Associates, Inc., 2012.
- [SSK17] Eric Schulz, Maarten Speekenbrink, and Andreas Krause. A tutorial on gaussian process regression: Modelling, exploring, and exploiting functions. *bioRxiv*, 2017.
- [SSW<sup>+</sup>16] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
- [WD92] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.

- [WvO12] M. Wiering and M. van Otterlo. *Reinforcement Learning: State-of-the-Art*. Adaptation, Learning, and Optimization. Springer Berlin Heidelberg, 2012.