

# Balloma Video Game playing Reinforcement Learning Agent.

Luis Rojas Aguilera  
Udacity  
Capstone Project Proposal

## CONTENTS

<b>I</b>	<b>Introduction</b>	<b>1</b>
<b>II</b>	<b>Deep Reinforcement Learning</b>	<b>1</b>
<b>III</b>	<b>Balloma video game</b>	<b>2</b>
<b>IV</b>	<b>Environment</b>	<b>2</b>
<b>V</b>	<b>Reward function</b>	<b>2</b>
	<b>References</b>	<b>3</b>

## LIST OF FIGURES

1	Balloma splash screen. . . . .	1
2	Balloma’s scene sample. Elements of scene are: a) The ball (blue), b) Floating elements (green), c) Score records (red), d) Target position (yellow) . . . . .	2

## LIST OF TABLES

# Balloma Video Game playing Reinforcement Learning Agent.

**Abstract**—This work presents a capstone project proposal to achieve a Nanodegree in Machine Learning from Udacity. The objective of this work is to construct a video game playing robot, through the use of Deep Reinforcement Learning techniques. Also it should provide a methodology for test automation on game development process. This document provide descriptions on: *a) problem to be solved, b) objectives, c) context, d) tools, e) metrics and f) techniques.*

## I. INTRODUCTION

Balloma is a single-player, android video game, developed by Black River Studios<sup>1</sup> at brasilian SIDIA<sup>2</sup>. During the game development process, after any new features are integrated, regression and functional tests must be executed in order to validate if new functionalities are properly working and side effects caused by new code have not appeared.



Fig. 1. Balloma splash screen.

Currently such tests are executed manually by humans, thus as functionalities' stack increases also does testing complexity and workload. In order to aim testing process this proposal presents a Proof Of Concept method for test automation using Deep Reinforcement Learning techniques. Next sections presents further details on Balloma video game, the RL framework and how to use it for game automatic controlling.

## II. DEEP REINFORCEMENT LEARNING

Reinforcement Learning (RL) is a field of Machine Learning that studies autonomous interactions among software agents and environments, and the way an agent can enhance its behavioral rules (policy) in order to maximize the reward it obtains from an environment. That is, at time step  $t$ , agent takes action  $a_t$  on a given environment  $E$  with transition dynamics

$p(s_{t+1}|s_t, a_t)$  and current state  $s_t$  obtaining a reward  $r_t$  from a reward function  $r(s_t, a_t)$  and transforming current state into  $s_{t+1}$ .

Currently RL proposes two main approaches to represent an agent behavioral rules: a) Value-based and b) Policy-based algorithms.

In value-based algorithms it is used a function called action-value ( $q_*(s, a)$ ) that contains the expected cumulative reward of performing a given action while in a given environment's state. Such function can be used by the agent to decide which action to take in a given state by selecting from  $q$  the action that maximizes reward with a given probability  $\epsilon$  that allows a desired exploration behavior so as it is applicable for non-deterministic environment dynamics. This approach is constrained to discrete action and state spaces.

Policy-based algorithms directly maps states and actions by a approximated policy function that represents the underlying stochastic distribution presented by environment's dynamics. While in the Value-based approach an heuristic should be defined to provide exploration behavior and avoid conditional randomness influenced by agent greediness, in Policy-based such heuristic is expressed in the policy function and is not necessarily fixed for each possible action-space setup. This approach is applicable to continuous action and state spaces.

For both approaches, agent behavior is guided by an optimal policy  $\pi_*(s, a, \theta)$  evaluated by objective function  $J(\theta) = \mathbb{E}[R(\tau)]$  of policy parameters  $\theta$ , based on state-action-reward expectations on trajectory  $\tau = S_0, A_0, R_1, S_1, \dots$  drawn from a probability distribution  $p(s', r|s, a)$ . In such case is convenient using the **Bellman Expectation Equation** [5] to represent expected return of taking an action  $a_t$  while in state  $s_t$  and following a policy  $\pi$ :

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} [r(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1} \sim \pi} [Q^\pi(s_{t+1}, a_{t+1})]] \quad (1)$$

In this proposal I intend to apply a hybrid approach for RL called Deep Deterministic Policy Gradients (DDPG) [1]. It follows the actor-critic [2] algorithm DPG [3] in which the agent's behavior is represented by function  $\mu(s|\theta^\mu)$ , called the actor, that is evaluated by the critic, a non-linear action-value function  $Q(s, a)$  parametrized by  $\theta^Q$ , adjusted by minimizing the loss:

$$L(\theta^Q) = \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \mu, r_t \sim E} [(Q(s_t, a_t|\theta^Q) - y_t)^2] \quad (2)$$

where

$$y_t = r(s_t, a_t) + \gamma Q(s_{t+1}, \mu(s_{t+1})|\theta^Q) \quad (3)$$

<sup>1</sup><http://blackriverstudios.net/>

<sup>2</sup><https://www.sidia.com>

Actor function  $\mu(s|\theta^\mu)$  specifies the agent's policy by deterministically mapping states to a specific action. It is adjusted to maximize the expected reward from a start distribution  $J = \mathbb{E}_{r_i, s_i \sim E, a_i \sim \pi}[R_1]$ . That is attained by updating parameters  $\theta^\mu$  following the policy's performance gradients, expressed as:

$$\begin{aligned} \nabla_{\theta^\mu} &\approx \mathbb{E}_{s_t \sim \rho^\beta} [\nabla_{\theta^\mu} Q(s, a|\theta^Q)|_{s=s_t, a=\mu(s_t|\theta^\mu)}] \\ &= \mathbb{E}_{s_t \sim \rho^\beta} [\nabla_a Q(s, a|\theta^Q)|_{s=s_t, a=\mu(s_t)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s=s_t}] \end{aligned} \quad (4)$$

Both actor and critic functions are expressed by neuronal networks, updated following an approximation to a supervised learning approach. For that matter, target networks  $\mu'(s|\theta^{\mu'})$  and  $Q'(s, a|\theta^{Q'})$  are created, which are soft copies of actor and critic networks respectively. That soft copy is made from updated parameters as follows:  $\theta' \leftarrow \tau\theta + (1-\tau)\theta'$  with  $\tau \ll 1$

In order to provide exploration capabilities, during training, a different policy than that of the agent is used, obtained by adding noise sampled from a noise process  $\eta$  to the actor function. Such exploration policy is expressed as follows:

$$\mu'(s_t) = \mu(s_t|\theta_t^\mu) + \eta \quad (5)$$

Training follows an episodic process similar to Q-Learning. In each episode a set of agent-environment interactions, represented as a transition tuple  $(s_t, a_t, r_t, s_{t+1})$  are stored in a replay buffer  $R$  [4]. A set of tuples are sampled from  $R$  and used to apply the updates previously explained in this section. This iterative process can be truncated based on optimization or time constraining conditions are met.

### III. BALLOMA VIDEO GAME

Balloma is a single-player video game that runs on Android devices. It is composed of several scenes that initially appears locked and gets unlocked as the user completes unlocked scene's objective.

Each scene contains a constrained little terrain world with elements arranged and boundaries that if crossed makes the scene end. The scene's main element is a ball the player should control and carry to a remarked target position in the scene. To control such ball the player should touch and swipe the device's screen inputting a direction and speed he wants the ball to take. Once the ball is positioned at target, scene ends, a score is given to user and next scene is unlocked.



Fig. 2. Balloma's scene sample. Elements of scene are: a) The ball (blue), b) Floating elements (green), c) Score records (red), d) Target position (yellow)

While the user carries the ball to target he can make the ball go trough floating elements in the scene which makes the score increase. Also a timer is included in the scene which influences the final score, lower scene's completion running times provides higher rewards.

Figure 2 presents the game's first scene which appears unlocked by default. In Figure 2 the ball is bounded by a blue box. The blue diamonds bounded with a green box appears floating in the scene, if reached by the ball makes the score (upper left red box) increase. Target position is remarked inside a yellow box. Also there is a timer (lower left corner red box) that influences the final score.

### IV. ENVIRONMENT

An interface is implemented to provide a RL suitable environment by which the agent can interact with the game to construct the optimized policy. Since the game runs in Android, the environment should provide a controlling interface with an android physical or virtual device. Such functionality can be attained through the Android Debug Bridge (adb<sup>3</sup>).

In this proposal the states are represented by raw scene frames of 240x240. It means at each time step  $t$  the environment is at state  $s_t \in \mathbb{R}^3$  formatted as a tensor of shape  $(240, 240, 3)$ . A similar approach was presented in [1] and [4], however scenes in those works are 2D spaces and in this case it is 3D. Frames are cropped to exclude score indicators and focus on more important pixels.

On each time step the agent can chose to take an action  $a_t \in \mathbb{R}^4$ . Actions are presented as a tuple  $a_t = (x_0, x_1, y_0, y_1, \nu)$  containing swipe attributes: a)  $x_0$  starting swipe  $x$  coordinate, b)  $x_1$  final swipe  $x$  coordinate, c)  $y_0$  starting swipe  $y$  coordinate, d)  $y_1$  final swipe  $y$  coordinate and e) swipe velocity  $\nu$ .

Each episode starts with the first swipe inputted by the agent. If the balls falls or get to target position environment is set into a terminal state. Also it will be constrained in time to avoid long episodes.

### V. REWARD FUNCTION

At each time step, after the agent choses an action  $a_t$  and such action is applied to environment, a reward  $r_t$  is observed. Such reward is computed taking into account two coefficients: a) the ratio of gathered floating diamonds over all such elements in the scene ( $\omega$ ) and b) episode time passed since it started ( $\varphi$ ). Then reward is calculated as:

$$r(s_t, a_t) = \frac{\omega_t}{\varphi_t} \quad (6)$$

Such coefficients appears as elements of the scene at a fixed frame position. It should be cropped and processed so as to convert such pixels into a suitable representation of its digits. In order to extract the score on each time step from those scene elements, in this work it will be trained a digit recognizer on the well known MNIST handwritten digit database [6].

<sup>3</sup><https://developer.android.com/studio/command-line/adb>

## REFERENCES

- [1] Lillicrap, Timothy P., et al. "Continuous control with deep reinforcement learning." arXiv preprint arXiv:1509.02971 (2015).
- [2] Konda, Vijay R., and John N. Tsitsiklis. "Actor-critic algorithms." *Advances in neural information processing systems*. 2000.
- [3] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. 2014. Deterministic policy gradient algorithms. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32 (ICML'14)*, Eric P. Xing and Tony Jebara (Eds.), Vol. 32. JMLR.org I-387-I-395.
- [4] Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David, Rusu, Andrei A, Veness, Joel, Bellemare, Marc G, Graves, Alex, Riedmiller, Martin, Fidjeland, Andreas K, Ostrovski, Georg, et al. Humanlevel control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [5] (2011) Bellman Equation. In: Sammut C., Webb G.I. (eds) *Encyclopedia of Machine Learning*. Springer, Boston, MA
- [6] LeCun, Y. & Cortes, C. (2010), 'MNIST handwritten digit database'.