

# Paper Note

**Remosy**

A report submitted for the course  
COMP8755 AND Individual Computing Project  
Supervised by: Dr. Penny Kyburz  
The Australian National University

May 2019

Except where otherwise indicated, this report is my own original work.

Remosy  
15 May 2019

---

# Contents

---

<b>1</b>	<b>Atari Gym:Ice Hockey</b>	<b>1</b>
1.1	Gym . . . . .	1
1.2	Ice Hockey . . . . .	1
<b>2</b>	<b>Learning from Demonstration</b>	<b>3</b>
2.1	Playing hard exploration games by watching YouTube . . . . .	3
2.2	Player Experience Extraction from Gameplay Video . . . . .	3
2.3	Learning Montezuma’s Revenge from a Single Demonstration . . . . .	3
2.4	Deep Q-learning from Demonstration . . . . .	3
<b>3</b>	<b>State Embeddings</b>	<b>5</b>
3.1	Dynamic image encoder . . . . .	5
3.2	Stack of difference of frame video-clip encoder . . . . .	5
3.3	Grad-CAM ++ . . . . .	5
3.4	VCG . . . . .	6
<b>4</b>	<b>RL</b>	<b>7</b>
4.1	Definitions . . . . .	7
4.2	Off-Policy . . . . .	8
4.3	On-Policy . . . . .	8
4.4	Model Based . . . . .	8
4.5	Model Free . . . . .	8
4.6	DQN . . . . .	8
4.7	A3C . . . . .	8
<b>5</b>	<b>IRL</b>	<b>9</b>
5.1	Apprenticeship Learning . . . . .	9
5.2	Bayesian Inverse Reinforcement Learning . . . . .	9
5.3	Maximum Entropy Reinforcement Learning . . . . .	9
5.4	Generative Adversarial Imitation Learning . . . . .	9
5.5	Project Scope . . . . .	9
5.6	Report Outline . . . . .	9
<b>6</b>	<b>Policy Optimisation</b>	<b>11</b>
6.1	Dynamic image encoder . . . . .	11
6.2	Stack of difference of frame video-clip encoder . . . . .	11
	<b>Bibliography</b>	<b>13</b>



---

# Atari Gym:Ice Hockey

---

## 1.1 Gym

Open AI published a reinforcement learning toolkit, gym. It includes the Arcade learning environment(Bellemare et al., 2012) ran on Stella Atari emulator. Therefore, user can train AI agents for a large number of Atari games from the gym library.

## 1.2 Ice Hockey

Ice Hockey, an Atari video game released in 1981. It has two game modes: single and 2-player multiplayer. Ice Hockey has two teams, one team wear yellow cloth and green pants, another team wear blue cloth and red pants. Each team has two roles, goalie and offense. In this game, the action control is taken by either goalie or offense of a team, who is closer to the hockey puck. The whole time for a game episode is 3 minutes.

In the gym library, Ice Hockey has an observation shape (210,160,3); it has 3 layers of RGB, each layer is shaped in height 210 and width 160.

Space: Discrete  
timesteplimit = maxEpisodeSteps  
Trials = 100  
Ram version Vs. Non-ram Version  
Deterministic? No  
NoFrameSkip?

---

<b>[0] Noop</b>	<b>[1] Fire</b>	<b>[2] Up</b>	<b>[3] Right</b>
-	(32,) ESC	(119,) W	(100,) D
<b>[4] Left</b>	<b>[5] Down</b>	<b>[6] UpRight</b>	<b>[7] UpLeft</b>
(97,) A	(115,) S	(100,119) D,W	(97,119) A,W
<b>[8] DownRight</b>	<b>[9] DownLeft</b>	<b>[10] Upfire</b>	<b>[11] RightFire</b>
(100,115) D,S	(97,115) A,S	(32,119) ESC,W	(32,100) ESC,D
<b>[12] LeftFire</b>	<b>[13] DownFire</b>	<b>[14] UpRightFire</b>	<b>[15]UpLeftFire</b>
(32,97) ESC,A	(32,115) ESC,S	(32,100,119) ESC,D,W	(32,97,119) ESC,A,W
<b>[16]DownRightFire</b>	<b>[17]DownLeftFire</b>		
(32,100,115) ESC,D,S	(32,97,115) ESC,A,S		

Table 1.1: Actions

---

# Learning from Demonstration

---

**2.1 Playing hard exploration games by watching YouTube**

**2.2 Player Experience Extraction from Gameplay Video**

**2.3 Learning Montezuma's Revenge from a Single Demonstration**

**2.4 Deep Q-learning from Demonstration**

How many chapters you have? You may have Chapter ??, Chapter ??, Chapter ??, Chapter ??, and Chapter ??.





# State Embeddings

## 3.1 Dynamic image encoder

## 3.2 Stack of difference of frame video-clip encoder

## 3.3 Grad-CAM ++

The Class Activation Mapping has been generalised into Gradient Class Activation Mapping (Grad-CAM). The interpretability of Grad-CAM in deep network can lead researchers to failures of model. For the interpretability, the Grad-CAM can discriminate the image classes by localising the region of interest (ROI) in a image.

Grad-CAM uses CNN architecture, the gradient mappings obtained from CNN can be used in the final fully connected layer to significant the regions of interest in a image.

Where to get the gradient information

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$$

$\alpha_k^c$  :

$\frac{1}{Z}$  : global average pooling

$y^c$  :

$A_{ij}^k$  : A matrix of K feature map

$\frac{\partial y^c}{\partial A_{ij}^k}$  :

$$L_{Grad\_CAM}^c = ReLU(\sum_k \alpha_k^c A^k)$$

*ReLU*: it's good to see the positive result

Grad-CAM ++ is better for video-based training.

Why "++" better for video? How differ from Grad\_CAM?

### 3.4 VCG

How many chapters you have? You may have Chapter ??, Chapter ??, Chapter ??, Chapter ??, and Chapter ??.

# RL

---

## 4.1 Definitions

Agent: An agent takes actions

Action (A): A is the set of all possible moves the agent can make

Environment: The world through which the agent moves. The environment takes the agent/ current state and action as input, and returns as output the agent/s reward and next state

State (S): A state is a concrete and immediate situation in which the agent finds itself

Reward (r): A reward is the feedback by which we measure the success or failure of an agent/ s actions

Discount factor ( $\gamma$ ): The discount factor is multiplied with future rewards as discovered by the agent in order to dampen their effect on the agent/ s choice of action. It makes future rewards worth less than immediate rewards

Policy ( $\pi$ ): The policy is the strategy the agent employs to determine the next action based on the current state. It maps states to actions

Value (V): The expected long-term return with discount, as opposed to the short-term reward r.  $V\pi(s)$  is defined as the expected long-term return of the current state under policy  $\pi$

Q-value or action-value (Q): Q-value is similar to Value, except that it takes an extra parameter, the current action a.  $Q\pi(s, a)$  refers to the long -term return of the current state/ s, taking action a under policy  $\pi$ . Q maps state-action pairs to rewards

## 4.2 Off-Policy

## 4.3 On-Policy

## 4.4 Model Based

## 4.5 Model Free

## 4.6 DQN

Use the immediate reward we receive and a value estimate of our new state to update the value estimate of original state-action pair. We only had the learned value function Q-function and the policy we followed was simply taking the action that maximised the Q-value at each step

## 4.7 A3C

Actor-critic methods combine policy gradient methods with a learned value function. we learn two different functions: the policy (or "actor"), and the value (the "critic"). The "policy" adjusts action probabilities based on the current estimated advantage of taking that action, and the value function updates that advantage based on the experience and rewards collected by following the policy

How many chapters you have? You may have Chapter ??, Chapter ??, Chapter ??, Chapter ??, and Chapter ??.

---

# IRL

---

## 5.1 Apprenticeship Learning

## 5.2 Bayesian Inverse Reinforcement Learning

## 5.3 Maximum Entropy Reinforcement Learning

## 5.4 Generative Adversarial Imitation Learning

## 5.5 Project Scope

Describe the problem your project addresses.

## 5.6 Report Outline

How many chapters you have? You may have Chapter ??, Chapter ??, Chapter ??, Chapter ??, and Chapter ??.



---

# Policy Optimisation

---

## 6.1 Dynamic image encoder

## 6.2 Stack of difference of frame video-clip encoder

Using a stack of differences of frames to capture motion and dynamics from a video clip. It was helpful to learn odd-one-out

How many chapters you have? You may have Chapter ??, Chapter ??, Chapter ??, Chapter ??, and Chapter ??.





---

# Bibliography

---