

---

# Reinforcement Learning - Spring Semester 2023

## Assignment 3. Policy-based RL

---

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
Antonia Christodoulou (s3614794), Kyriakos Aristidou (s3510123), Shreyansh Sharma (s3772241)

### Abstract

Our work focuses on the investigation of different *policy-based* methods in reinforcement learning to solve the Catch environment. We implement **REINFORCE** and **Actor-Critic** algorithms as well as some variations of Actor-Critic which include bootstrapping or/and baseline subtraction techniques. After concluding to simple Actor-Critic to be the best agent, we make several experiments using it, changing many different configurations of the environment such as the size, speed, etc so as to examine its performance.

### 1. Introduction

Reinforcement learning (RL) is a strong machine-learning technique that has received a lot of attention over the last several years. In the previous assignment, we investigated the value-based methodology of RL, using algorithms such as *SARSA*, *Q-learning*, and *DQN*, with the objective of learning the values of actions in different states. In this assignment though, we will examine the policy-based approach. The policy-based approach, as opposed to the value-based approach, optimizes the policy directly instead of acquiring it from a representation of the value function.

The goal of this assignment is to build and compare several policy-based RL algorithms. To accomplish this, we will use the Catch environment ([Moerland](#)), a task in which the purpose is to operate a paddle to catch falling balls.

Initially, we implement various algorithms from scratch: **REINFORCE**, **ActorCritic** with bootstrapping, with *baseline subtraction*, and a combination of the two. We then tune the hyperparameters of each algorithm in order to make it a fair comparison. After which we compare their performance and select the best configuration, which is later utilized for more experimentation on the environment variations. More specifically, using the best agent, we will examine different configurations of the environment (like the ball's speed, the grid's size, etc.) to see how our agent's performance is affected.

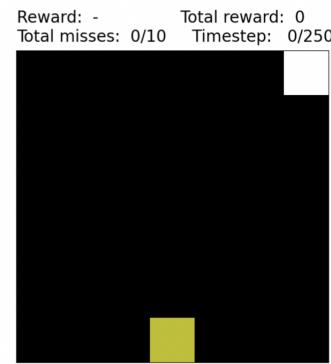


Figure 1. Illustration of the Catch environment provided by the authors of Assignment 3. The yellow block at the bottom represents the paddle and the white block above is the falling ball. The black background presents the space within which the game takes place.

### 2. Catch Environment

The Catch environment, as illustrated in *Figure 1*, is a game in which you use a paddle to catch falling balls. To get rewards, you have to move the paddle(the yellow block) and catch as many balls(the white block) as possible. The balls fall at various speeds from the top of the screen. At each timestep, you get a +1 reward if the position of the yellow block is the same as the white block's, you get -1 if the white block reaches the bottom but is found in a different position than the yellow one, and, finally, you get +0 reward in any other case. The state space includes the height and width of the grid (number of rows and columns) while the action space includes three different possible moves of the paddle: left, right, or stay idle.

The task finishes either if we reach the maximum steps set by the user, or if the paddle reaches the maximum number of misses set by the user as well.

To adjust the difficulty level, initialization configurations can be used to personalize the environment. The number of rows and columns on the game board, for example, can be changed to broaden or narrow the space that the paddle must travel. The speed rate at which the balls fall can also be adjusted to make the game quicker or slower. There are, also, different observation types, such as raw pixel images

055 or preprocessed vector representations of the game state.  
 056 Because of its ease of use and adaptability, the Catch environment  
 057 is widely utilized. It is simple to use and may be  
 058 tailored to a variety of experimental and testing scenarios.  
 059

### 060 3. Methodology

#### 061 3.1. REINFORCE

064 Reinforce is a famous algorithm in reinforcement learning  
 065 and is used to train agents to execute tasks in situations that  
 066 are unpredictable like the catch environment in this case.  
 067 The agent in Reinforce engages with the environment and  
 068 is rewarded for each action it takes. The aim of the agent  
 069 is to gradually develop a policy that maximizes the pre-  
 070 dicted reward. The agent accomplishes this by employing a  
 071 stochastic policy, which generates a probability distribution  
 072 over actions at each time step. The agent selects its actions  
 073 by sampling from this distribution.

074 The main concept of Reinforce is to update the agent's policy  
 075 using the policy gradient theorem. The policy gradient  
 076 theorem assumes that the product of the expected reward and  
 077 the gradient of the log probability of the action chosen with  
 078 respect to the policy parameters determines the expected  
 079 reward's gradient with respect to the policy parameters. In a  
 080 few words, we can adjust the policy configurations so that  
 081 the expected reward is higher by using the gradient of the log  
 082 probability of the taken action.

083 Using stochastic gradient ascent, the Reinforce algorithm  
 084 updates the network parameters based on the gradient of the  
 085 predicted return. The sum of the rewards obtained for each  
 086 action in the current episode, weighted by the probability  
 087 that the action will be taken, is used to compute the expected  
 088 return, this is what we call the Monte Carlo method.

089 First, we compute the discounted returns  $G_t$  for each time  
 090 step  $t$  using the rewards  $r_k$  received at time steps  $k \geq t$   
 091 and the discount factor  $\gamma$ . This equation computes the total  
 092 discounted return from each time step  $t$  to the end of the  
 093 episode at time step  $T - 1$ , where  $\gamma \in [0, 1)$  is the discount  
 094 factor that determines the weight of future rewards relative  
 095 to immediate rewards:  
 096

$$G_t = \sum_{k=t}^{T-1} \gamma^{k-t} r_k \quad (1)$$

097 Next, we compute the baseline  $b$ , which is the expected  
 098 value of the returns across all time steps:  
 099

$$b = \frac{1}{T} \sum_{t=0}^{T-1} G_t \quad (2)$$

Then, we compute the normalized returns  $\hat{G}_t$  by subtracting  
 the baseline  $b$  from each return  $G_t$  and dividing the result by  
 the standard deviation  $\sigma$  of the returns. This equation stan-  
 dardizes the returns, which helps to reduce the variance of  
 the gradients and improve the convergence of the algorithm.:  
 060

$$\hat{G}_t = \frac{G_t - b}{\sigma} \quad (3)$$

where  $\sigma$  is the standard deviation of the returns across all  
 time steps

After computing the normalized returns, we update the pol-  
 icy network parameters  $\theta$  by taking a gradient step in the  
 direction of the policy gradient  $\nabla_\theta \log \pi_\theta(a_t | s_t)$  weighted  
 by the normalized return  $\hat{G}_t$ . This equation updates the  
 policy network parameters using the policy gradient, which  
 is the gradient of the expected return with respect to the  
 policy parameters. The learning rate  $\alpha$  determines the step  
 size of the update.:  
 061

$$\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(a_t | s_t) \hat{G}_t \quad (4)$$

Then we compute the loss, we first calculate the policy loss  
 and entropy loss for each time step  $t$ . First, we loop through  
 the saved log probabilities and the corresponding returns for  
 each time step. For each time step, we calculate the policy  
 loss as the negative log probability of the selected action  
 multiplied by the corresponding return:

$$\mathcal{L}_{policy} = - \sum_{t=1}^T \log \pi_\theta(a_t | s_t) \cdot G_t \quad (5)$$

where  $\hat{G}_t$  is the normalized return for time step  $t$ .

Next, we calculate the entropy loss by taking the negative  
 sum of the logarithm of the policy probabilities  $\pi_\theta(a_t | s_t)$   
 for each time step  $t$ . The entropy loss encourages the policy  
 to be more exploratory by penalizing it for being too certain  
 about its actions.

$$\mathcal{L}_{entropy} = - \sum_{t=1}^T \pi_\theta(a_t | s_t) \cdot \log \pi_\theta(a_t | s_t) \quad (6)$$

We then sum up the policy loss and entropy loss across all  
 time steps to get the total loss, where  $\beta$  is the entropy weight  
 that control the strength of the entropy regularization:

$$\mathcal{L} = \sum_{t=0}^{T-1} (\mathcal{L}_{policy}(t) + \mathcal{L}_{entropy}(t) \cdot \beta)$$

Finally, we calculate the gradients of the total loss with respect to the policy parameters  $\theta$  and update the parameters using gradient descent:

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}$$

where  $\alpha$  is the learning rate.

### 3.2. Actor-Critic

As explained above, in REINFORCE we use Monte Carlo to estimate the return of a whole episode which leads to high variance. The Actor-Critic algorithm uses a hybrid architecture combining both value-based and policy-based methods aiming to reduce the variance and stabilize the learning process. So, the architecture used for this algorithm consists of two parts:

- the **Actor**: which determines the action that the agent will follow (the policy-based part).
- the **Critic**: which criticizes how good or bad this action is (the value-based part).

The general idea behind Actor-Critic is that the actor executes some actions, the critic provides feedback and the policy is updated based on this feedback. So, in this way, our agent becomes better and better at playing a game (the Catch in our case). However, not only the policy has to be updated but, also, the value function which determines how effective or not the action is. So, value and policy losses are calculated, and, then, using backpropagation the weights of the networks ( $\theta, \phi$ ) are updated.

Summing up, we have two neural networks to train:

- for the **Policy function** controlled by parameters  $\theta$
- for the **Value function** controlled by parameters  $\phi$

So, now there are two techniques to reduce variance:

1) *Bootstrapping*, and 2) *Baseline Subtraction*.

#### 3.2.1. BOOTSTRAPPING

Bootstrapping is a technique that refers to the target estimation  $Q_n$ . Unlike the Monte Carlo method which computes the full expected reward at a certain timestep  $t$ , bootstrapping sums up the rewards of depth  $n$ . As  $n$  gets higher we result in the Monte Carlo method and bootstrapping is gone. In our case, we use  $n=100$  after experimenting with different values of  $n$  such as  $n=1, 40, 100, 200$ . So, the procedure is sampling a trace, and then iterating over the number of timesteps, we compute the  $n$ -step target:

$$Q_n(s_t, a_t) = \sum_{k=0}^{n-1} r_{t+k} + V_{\theta}(s_{t+n}) \quad (7)$$

Then we update the Critic network as:

$$\theta \leftarrow \theta + \eta \sum_t [Q_n(s_t, a_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)] \quad (8)$$

and the Actor Network is updated as:

$$\phi \leftarrow \phi - \eta \sum_t \nabla_{\phi} (Q_n(s_t, a_t) - V_{\phi}(s_t))^2 \quad (9)$$

#### 3.2.2. BASELINE SUBTRACTION

The learning process can be further stabilized and the variance reduced using the **Advantage function** ( $A(s,a)$ ). The idea is that we want to push up the probability of taking actions that result in a Q-value higher than the average value of that state. So, the Advantage function basically calculates the extra reward that an action results in, compared to the average reward of that state:

$$A(s, a) = Q(s, a) - V(s) \quad (10)$$

So, using the baseline subtraction technique, the algorithm first samples a trace and then computes the target which might be calculated using n-step bootstrapping ( $Q_n$ ) as shown above or the Monte Carlo method ( $Q_{MC}$ ):

$$Q_n(s_t, a_t) = \sum_{k=0}^{n-1} r_{t+k} + V_{\theta}(s_{t+n}) \quad (11)$$

or

$$Q_{MC}(s_t, a_t) = \sum_{i=t}^{\infty} \gamma^i \cdot r_i \quad (12)$$

where  $\gamma$  is a discounted factor,  $r$  indicates the reward,  $t$  indicates the timestep, and  $V_{\theta}$  is the value of that state.

The Advantage is calculated as shown in equation (4) and then the Critic network is updated as:

$$\theta \leftarrow \theta + \eta \sum_t [A_n(s_t, a_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)] \quad (13)$$

and the Actor Network is updated as:

$$\phi \leftarrow \phi - \eta \sum_t \nabla_{\phi} (Q_n(s_t, a_t) - V_{\phi}(s_t))^2 \quad (14)$$

## 4. Experiments

Our experiments are divided into two parts. Firstly, we use the default configurations of the Catch environment and check the performance of the algorithms (REINFORCE, Actor-Critic, Actor-Critic with bootstrapping, Actor-Critic with baseline subtraction, and Actor-Critic with both) comparing them and aiming to find the best agent.

The **default parameters** of the environment are:

- 165     • **Number of rows:** 7  
 166     • **Number of columns:** 7  
 167     • **Speed of the falling ball:** 1.0  
 168     • **Observation Type:** pixel  
 169     • **Maximum number of steps in one episode:** 250  
 170     • **Maximum number of misses in one episode:** 10

175 After determining the best agent, we move on to the second  
 176 part of our experiments where we use the best algorithm in  
 177 modified environments to examine if and how the per-  
 178 formance of the agent is affected. We make changes in the size  
 179 of the environment, the speed of the falling ball, and the  
 180 observation type. All the experiments and plots presented  
 181 in the following sections are averaged over 10 repetitions.  
 182

## 183 4.1. Networks and Model Architectures

### 184 4.1.1. REINFORCE ARCHITECTURE

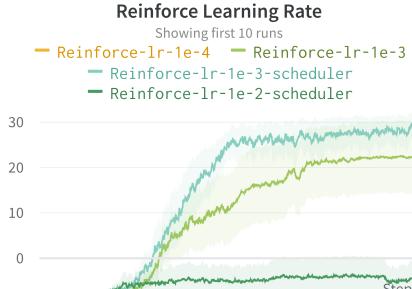
185 As mentioned before for the catch environment, we have two  
 186 options for observations: "vector" and "pixel". For "vector"  
 187 observations, we use a simple policy network that consists  
 188 of fully connected layers. The policy network takes in the  
 189 current state of the environment as an input and outputs a  
 190 probability distribution over the possible actions. We choose  
 191 this architecture because the observation space is small and  
 192 does not contain spatial information.

193 Alternatively, for "pixel" observations, we employ a con-  
 194 volutional neural network (CNN) architecture that accepts  
 195 the environment's raw pixel values as input. Since they  
 196 can learn hierarchical representations of the input by con-  
 197 volving filters over the input image, CNNs have particular  
 198 advantages for image-based observations. Two convolutional  
 199 layers with 3x3 filters and ReLU activation functions  
 200 are used in our particular architecture, followed by two fully  
 201 linked layers. The network's top layer produces a probabili-  
 202 ty distribution over all potential actions.

203 We employ ReLU activation functions in both architectures  
 204 to add non-linearity and enable the network to learn more  
 205 complex representations of the input. We employ two fully  
 206 connected layers, each with 128 units, in the policy network.  
 207 The first fully connected layer in the CNN architecture con-  
 208 sists of 448 units, followed by the first convolutional layer  
 209 with 16 filters and the second convolutional layer with 32  
 210 filters.

### 211 4.1.2. ACTOR-CRITIC ARCHITECTURE

212 As aforementioned, for Actor-Critic we need to train two  
 213 networks, the Actor and the Critic. In our experiments  
 214 though, we use one network with two heads. So, the hidden  
 215



(a)



(b)

Figure 2. Comparing various hyperparameters for tuning the Reinforce baseline. (a)Reinforce Learning Rate: alter learning rates and add scheduler (b)Reinforce Entropy Regularization: modifying the entropy weight

layers are the same for both parts but they differ in the output which for the actor has a size of three(3), giving the probabilities for each of the three possible moves: right, left, and stay idle, and for the critic has a size of one(1) which is the state's value after following those actions.

So the network is built as follows: When the observation type is set to "pixel", we use two convolutional layers and two linear layers. If we work with "vector", then we just use four linear layers. For any case, the activation function is "relu" and the optimizer is "Adam". However, a flattening is needed after the convolutional layers and before the linear layers in the case of the "pixel" type. Also, the actor's head exports the probabilities of the three actions after using the softmax function. The exact sizes of the hidden layers as well as other parameters like the learning rate will be shown and discussed in the following subsection.

220 4.1.3. HYPERPARAMETER TUNING  
 221

222 After running several experiments with different values of  
 223 the learning rate, gamma, entropy strength, depth of boot-  
 224 strapping and the use or not of a decay parameter that de-  
 225 creases the learning rate steadily, we concluded the follow-  
 226 ing values for our parameters as presented in *Table 1* and  
 227 *Table 2*. In *Figure 2* and *Figure 3*, you can see some of the ex-  
 228 periments for tuning the main parameters of the algorithms.  
 229 Specifically, you can see the tuning of the learning rate and  
 230 entropy regularization for Actor-Critic and Reinforce, and  
 231 Bootstrapping depth for Actor-Critic. For the learning rate,  
 232 we also examine if the lr\_scheduler helps increase perfor-  
 233 mance, which essentially decreases the learning rate over  
 234 time.

235 **Actor-Critic Hyperparameters**  
 236

lr	1e-2
gamma	0.85
n_steps	100
entropy_weight	0.01
lr_decay	0.99
lr_scheduler	True
num_episodes	1000

244 *Table 1.* Hyperparameters for Actor-Critic algorithm.  
 245

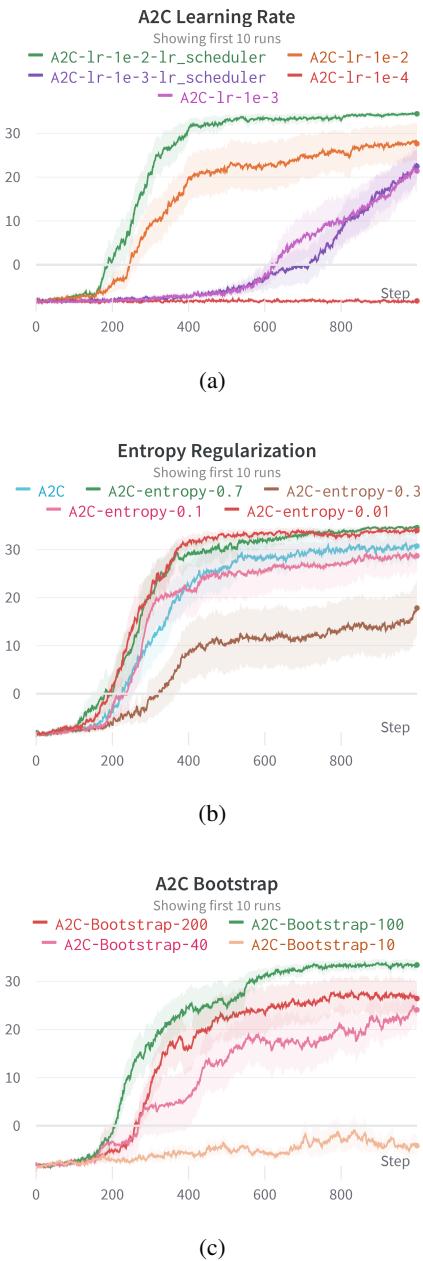
246 **REINFORCE Hyperparameters**  
 247

lr	1e-3
gamma	0.85
beta	0.01
lr_decay	0.99
lr_scheduler	True
num_episodes	1000

255 *Table 2.* Hyperparameters for REINFORCE algorithm.  
 256

257 For **Actor-Critic**, the learning rate (lr) is set to 1e-2 and  
 258 we found that decreasing it over time (scheduler) with a  
 259 standard decay parameter helps its performance. Also, the  
 260 depth of bootstrapping (n\_steps) is found to give better  
 261 results when set to 100 while lower depths were giving a low  
 262 performance. However, as shown by our experiments for  
 263 when the depth is set to 200 it starts performing relatively  
 264 bad.

265 For **REINFORCE**, the learning rate (lr) is set to 1e-3 and  
 266 we found that decreasing it over time (scheduler) with a  
 267 standard decay parameter helps the performance of the al-  
 268 gorithm same as the actor-critic. We set beta to 0.01 which  
 269 represents the entropy weight to control the strength of the  
 270 entropy regularization as mentioned in the previous section.  
 271 The entropy regularization was vital in order to find the best  
 272 agent for the Reinforce algorithm. As *Figure 2(b)* shows,  
 273 specifically the orange line, reinforce algorithm with no



274 *Figure 3.* Comparing various hyperparameters for tuning the base-  
 275 line. (a)A2C Learning Rate: alter learning rates and add scheduler  
 276 (b)Entropy Regularization: modifying the entropy weight (c)A2C  
 277 Bootstrap: Experimenting with different number of steps for boot-  
 278 strapping

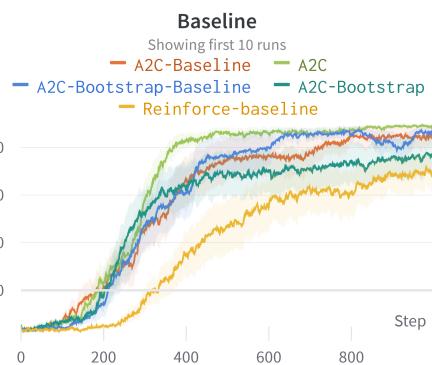
275 entropy regularization has really bad results.  
 276  
 277  
 278

## 5. Results & Discussion

### 5.1. Part I - Find the best agent

#### 5.1.1. REINFORCE vs ACTOR-CRITIC

As explained before, we run experiments using the networks for REINFORCE and Actor-Critic (A2C) and from *Figure 4*, we can come to the conclusion that Actor-Critic without bootstrapping or baseline subtraction performs better than the other algorithms. Actor-Critic with bootstrapping and baseline subtraction performs almost as well as the vanilla A2C, although it shows some instabilities and fluctuations. You can, also, observe that A2C with only bootstrapping is performing much worse than all the other variations of A2C. The reason behind this could be the depth of bootstrapping which in our case was set as 100, which although it gives better results, has a higher standard error. On the other hand, REINFORCE appears to be the worst agent as it does not even reach the maximum number of possible rewards for this environment which is 35, in 1000 epochs. It learns slower and does not converge to the optimum. It also suffers from high variance as shown in *Figure 4*. Therefore, the next section of our experiments will be performed using A2C without bootstrapping or baseline subtraction.



317 *Figure 4*. Comparisons between five different algorithms. The  
 318 main algorithms are **REINFORCE** and **Actor-Critic (A2C)** but  
 319 also some variations: **A2C-Bootstrap** which stands for Actor-Critic  
 320 using bootstrapping, **A2C-Baseline** for A2C using baseline sub-  
 321 traction and **A2C-Bootstrap-Baseline** which uses both techniques.  
 322 A2C performs the best and REINFORCE the worst.

### 5.2. Part II - Modify the environment

324 Experiments with varying environment factors can have  
 325 different maximum and minimum rewards. So in order to  
 326 understand them visually on a graph we decided all experi-  
 327 ments that have changes in the environment variables, have  
 328 their rewards normalized between a range of -1 to +1.  
 329

For our Actor-Critic architecture, we utilize a convolutional layer as the input layer, when we have the observation type as pixel. So the architecture works with square environments and does not support rectangular environment size. To overcome this issue we decided to add padding to rectangular-shaped environments, with a padding value of zero.

All experiments are performed using Wandb (Biewald, 2020) to log and make experimentation easier and Pytorch is utilized for building all the neural networks (Paszke et al., 2019).

#### 5.2.1. CHANGING THE SIZE OF THE ENVIRONMENT

The default experiments were run on an environment with rows and columns equal to 7. For this experiment, we see how our baseline performs when encountering changes in the size of the environment.

From *Figure 5* we can see as expected for the default environment, we continue to perform well and converge to the optimal reward. But when we increase the size of the environment to 10X10 we see a drop in performance and also gives an average maximum reward of 0. Similar behavior can be noticed for 25X25 environment size as we never achieve a positive reward. We believe this to be due to the current model architecture designed for default environment size and does not scale well. Further experiments can be done to increase the depth of the network as well as the hidden nodes in order to find the best parameters.

Now looking at rectangular environments, we see that 14X7 barely achieves any positive rewards and 7X14 never learns anything due to the width being bigger than the height. Whereas on the other hand for an environment size of 14X7 we see decent performance as it gets a maximum average reward of 0.5 at the end of 1000 epochs. Again we believe this to be due to our choice of architecture and additional padding added for rectangular environments.

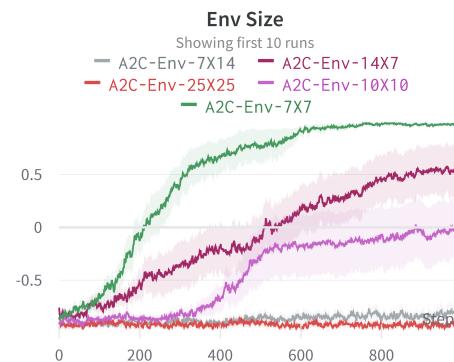


Figure 5. Changing Environment Size for A2C baseline model to see its effect on the total reward obtained.

330    5.2.2. CHANGING THE SPEED OF THE FALLING BALL  
 331   

332    The default speed of the ball for all experiments so far was  
 333    1.0. For this experiment, we alter the speed of the ball to see  
 334    its impact on the performance of A2C.

335    From *Figure 6*, we can see that for a speed of 5 A2C does  
 336    not get any rewards and this is because it cannot intercept  
 337    the ball at high speeds. For a speed value of 2.0, it performs  
 338    relatively well but does not converge to the optimal reward  
 339    in 1000 epochs. We believe this to be due to higher speeds  
 340    and a lesser probability of interception of the ball. Finally  
 341    looking at a speed of 0.5 we see a lower performance than  
 342    at speed 1. We believe this happens due to lesser rewards  
 343    obtained over the same number of maximum steps which  
 344    reduces the agent's capability to learn. If we were to increase  
 345    the number of epochs, we expect to see an increase in the  
 346    overall reward obtained, but cannot guarantee the optimal  
 347    reward.



363    *Figure 6*. Comparison between different speeds at which the ball  
 364    drops in the environment. All experiments are run on A2C for the  
 365    default environment parameters apart from the speed.  
 366   

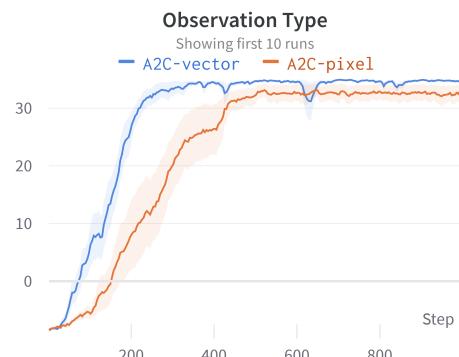
367    5.2.3. CHANGING THE OBSERVATION TYPE: PIXEL VS  
 368    VECTOR

369    As mentioned in the section about the Environment, apart  
 370    from size and speed, there are other configurations that we  
 371    can modify and examine the performance of our agent. One  
 372    of the main configurations is the observation type. We have  
 373    two options:

- 375    • **Pixel:** Pixel is set as the default observation type and it  
 376    is a binary two-channel array that "observes" and saves  
 377    the locations of the paddle in the first channel and the  
 378    ball's locations in the second one.
- 379    • **Vector:** Vector keeps the xy-coordinates of the pad-  
 380    dle's locations and the ball's location when reaching  
 381    the lowest level.

383    As shown in *Figure 7*, the agent appears to learn faster and,

384    eventually, to achieve a slightly higher reward when using  
 385    vectors as the observation type while with "pixel", it seems  
 386    to need more steps to converge. This can be caused by  
 387    the fact that Catch is a relatively easier environment and  
 388    numerical representations of the observation space are a  
 389    simple way that helps the agent to capture the relationship  
 390    between its actions and the resulting state transitions faster.



391    *Figure 7*. Comparing how A2C performs when modifying the ob-  
 392    servation type of the environment. When using "vector", the agent  
 393    learns faster and achieves higher rewards.

394    5.2.4. ENVIRONMENT VARIATIONS

395    For the final experiment, we chose to combine at least 2 vari-  
 396    ations in environment configurations and test its outcome.  
 397    From *Figure 8* we can see the different variations that we  
 398    implemented.

399    First looking at an environment with maximum misses 20  
 400    and max steps 300, it performs the best among all variations.  
 401    This was an expected result as the steps and maximum  
 402    misses increase the agent now sees more rewards and has a  
 403    longer time to learn from its mistakes.

404    Next, we ran two experiments to compare environment size  
 405    14X7 with a speed of 0.5 on both vector and pixel obser-  
 406    vation types. The vector environment reaches the optimal  
 407    reward well under 1000 epochs, whereas the pixel environ-  
 408    ment struggles to obtain positive rewards till the end. we  
 409    believe this could be due to the padding that we chose to add  
 410    to rectangular environments and due to the sparse rewards  
 411    obtained in an environment with a speed of 0.5. Addition-  
 412    ally, we can see that rectangular environment perform better  
 413    with vector observations as seen in *Figure 5*.

414    Then we ran experiments on the environment size 7X7 for  
 415    vector observations modifying the speed to 0.5 and 2.0. For  
 416    an environment with a speed of 0.5, it performs really well  
 417    and converges to the optimal reward. Whereas, when speed  
 418    is 2.0 the agent barely learns to obtain positive rewards. This  
 419    was an expected result as when we represent the environ-

385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
ment in vector observation type for speed 2.0, we have two balls that the paddle can catch one after the other. This can be easily represented in pixel observations, but in vector observation due to the fixed size of the state, the agent only sees the coordinate of the lowest ball, and cannot anticipate the second ball's position. This leads to a surprise negative rewards for the agent and it cannot learn any further.

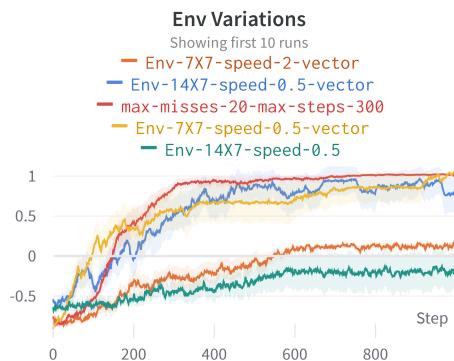


Figure 8. Comparing how A2C performs when modifying different environment variables like speed, size, max steps and max misses.

## 6. Conclusion

To sum up, we have worked with the Catch environment in order to evaluate REINFORCE and Actor-Critic algorithms as well as some variations of Actor-Critic. We compared their performance and came to a conclusion that REINFORCE indeed suffers from high variance which can be overcome by utilizing Actor-Critic instead. Bootstrapping and Baseline subtraction are two techniques used in Actor-Critic aiming to stabilize the learning process and reduce the variance but from our experiments, we understand that this does not necessarily hold in our case. Some reasons behind this could be the specific environment that we are working with which is relatively simple so the simpler version of the A2C algorithm performs well without any more complexity needed or maybe the bootstrapping depth should be lower to reduce variance. Additionally, we believe our hyperparameter tuning to have played huge role in getting high performance from the vanilla Actor-Critic algorithm.

Comparing the two techniques: bootstrapping and baseline subtraction, we can say that the learning process using bootstrapping is more stable however, using baseline convergence to the optimal reward is faster. Moreover, we have examined the effect of the entropy regularization factor which appears to play a major role in the performance as its presence increases the reward and makes the learning process faster if its strength is tuned optimally.

Finally we have run multiple experiments using A2C on

different configurations of the environment coming to the main conclusion that the agent performs poorly in very big environments due to our current architecture tuned for the default environment. Rectangular environments, pose a harder challenge as the agent learns slower. Finally, other properties like speed or observation type can also affect the agent's performance, and depending on the combinations we might get different results.

## 7. Contribution of the team members

All the team members have worked together for the best results in this Assignment. For the first part of the experiments, Kyriakos Aristidou worked on REINFORCE algorithm, and Antonia Christodoulou and Shreyansh Sharma worked together on Actor-Critic. Shreyansh Sharma worked on the "experiments.py" file while all three of us added different types of experiments on the "experiments.config.py" file. We, all, ran experiments and discussed together the results so that we could contribute to writing the final report.

## References

- Biewald, L. Experiment tracking with weights and biases, 2020. URL <https://www.wandb.com/>. Software available from wandb.com.
- Face, H. Deep reinforcement learning: Advantage actor-critic, February 2018. URL <https://huggingface.co/blog/deep-rl-a2c>.
- Moerland, T. Adapted from deepmind bsuite.
- Moerland, T. Continuous markov decision processes. *arXiv preprint arXiv:2104.05309*, 2021.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance.pdf>.
- Plaat, A. Deep reinforcement learning, a textbook, Jun 2022.
- Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. MIT press, 2nd edition edition, 2018. URL <https://web.stanford.edu/class/psych209/Readings/SuttonBartoIPRLBook2ndEd.pdf>.