

# DRL for Tennis

## 1. Introduction

## 2. Environment

## 3. Dependency

## 4. Implementation

### 4.1. Model Description

#### 4.1.1. The Agent

##### 4.1.1.1. *The Actors*

##### 4.1.1.2. *The Critics*

##### 4.1.1.3. *The Replay Buffers*

The replay buffers have been generated using a simple `deque`. The replay buffer is a class that has been defined within the file `src/utils/memory.py` by a class called `ReplayBuffer`. This exposes several methods that allow the replay buffer to be populated (`append`, and `appendMany`), sampled (`sample`) and also has methods that allow the `ReplayBuffer` to be saved (`save`) and reloaded (`load`) to and from disk at any point. This is especially useful because it is possible that we have already had a great set of memories that we can use. It is typically a total waste to throw away the entire buffer at one go. Hence, these buffers may be saved, so that they may be used in the future if needed.

There are several characteristics of the specific part of the game that makes the design of the memory buffer efficient.

1. The game is episodic.
2. Points are awarded every time the racquet of an agent is able to hit the ball across the net.
3. Each hit is fairly independent of the other hits. Hence the probability that a particular agent is going to be able to hit the ball across the net depends mostly upon the set of actions that the agent takes a few time-points (let's say 5 time points) before actually hitting the ball, at which point the racquet actually lobs the ball over the fence. Hence, it would be most ideal if the agent learns around this segment, rather than during all segments that the agent does not actually make a dent at the learning experiment.

- Note that it is quite easily possible to calculate the actual cumulative reward for this problem for a fairly large number of episodes.

Hence, in this specific case, the memory buffer contains a tuple of the following form:

```
(state, action, reward, next_state, done, cumRewards, numHits)
```

The cumulative rewards (represented by `cumRewards`) are calculated with a  $\gamma$  of 1. However, the choice of the value of  $\gamma$  that should be used will be discussed in the next section. `numHits` represents the number of hits that a particular agent is having for a particular hit. This is easier to visualize with the following graphs

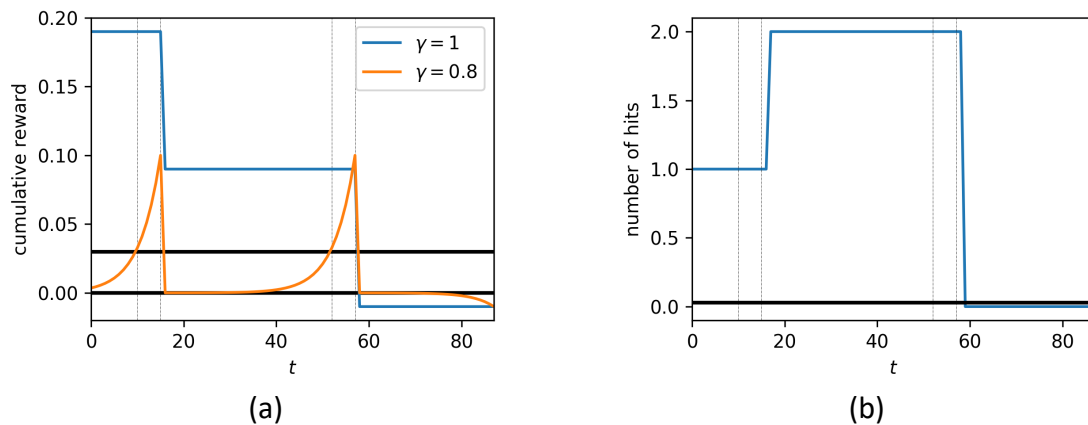


Fig. <>. Different characteristics of an agent is shown for an entire episode. In this episode, the agent starts, and successfully lobs the ball over the net (approximately at a time point of 10 units), the other agent hits the ball back, and this agent is able to hit it one more time approximately around 60 time units. The cumulative rewards for the first agent is shown in (a) for two different values of  $\gamma$ , while the hit number is represented by the figure in (b).

**Generating memories:** Refreshing a memory buffer is a needed very often. Hence a function has been presented that will allow a set of two list of tuples to be generated. The function called `memories` is present in the file `src/utils/generateMemories.py`. This function and the ideas that went into designing this function will be briefly discussed in the following paragraphs.

Since we are able to calculate the cumulative reward for every point in an episode, this is exactly what has been done. An episode is played using a given `policy`<sup>1</sup> that will last at most  $N_{st}$  steps (or when a game stops).

Now, note that there is very little to be learnt from regions where the agent is unable to get a score. Hence, it would be great if we were able to learn from the region close to where the agent actually gets a point. For finding this reason, the cumulative rewards with a cumulative rewards is generated with a  $\gamma$  of 0.8 as shown in the Figure above. This shown that the cumulative rewards decreases exponentially, as we get farther away from the location of a successful hit. Here, we have found regions within the episode wherein this value of the

---

<sup>1</sup> Generation of policies for both generating memories and training is discussed in [Section <Update this>](#).

```
rewards      1 :  
[ 0.    0.    0.    0.    0.    0.    0.    0.    0.    0.    0.  
   0.    0.    0.    0.    0.1  0.    0.    0.    0.    0.    0.  
   0.    0.    0.    0.    0.    0.    0.    0.    0.    0.    0.  
   0.    0.    0.    0.    0.    0.    0.    0.    0.    0.    0.  
   0.    0.    0.    0.    0.    0.    0.    0.    0.    0.    0.1  
   0.    0.    0.    0.    0.    0.    0.    0.    0.    0.    0.  
   0.    0.    0.    0.    0.    0.    0.    0.    0.    0.    0.  
   0.    0.    0.    0.    -0.01]
```

