

Reinforcement learning project

Hariharan Janardhanan

Section 1

States:

Cart Position (x): The position of the cart along the horizontal axis.

Cart Velocity (\dot{x}): The velocity of the cart along the horizontal axis.

Pole Angle (θ): The angle of the pole with respect to the vertical axis.

Pole Angular Velocity ($\dot{\theta}$): The angular velocity of the pole.

Actions:

0: Apply a force to move the cart to the left.

1: Apply a force to move the cart to the right.

Rewards:

+1: The agent receives a reward of 1 for each time step where the pole remains upright and the cart stays within defined thresholds.

0: The episode terminates, either because the cart goes beyond the defined position thresholds or the pole's angle exceeds the defined limits, or the maximum number of time steps (500) is reached.

Markov Decision Process (MDP):

State Space (S): The set of all possible states as described above.

Action Space (A): The set of all possible actions $\{0, 1\}$.

Transition Function (T): Defines the probability of transitioning from one state to another given an action. In this problem, the transition function is deterministic because the next state is completely determined by the current state and action.

Transitions between states are determined by the laws of physics as applied in the simulation. The equations of motion take into account gravitational forces, the input force from the action, and the interactions between the cart's and pole's movements. Depending on the action chosen, the environment computes the new state using semi-implicit Euler integration for numerical stability.

Reward Function (R): Defines the immediate reward received after transitioning from one state to another due to an action. As mentioned earlier, the reward is +1 for each time step where the pole remains upright and 0 otherwise.

Discount Factor (γ): The discount factor, if used, would discount future rewards in favor of immediate rewards. It is not explicitly defined in this problem.

The agent interacts with the environment in discrete time steps, choosing actions based on the current state, receiving rewards, and transitioning to the next state according to the dynamics of the system. The goal of the agent is to learn a policy that maximizes the cumulative reward over time, i.e., keeping the pole balanced for as long as possible within the constraints of the environment.

Section 2

TODO: change input and output sizes depending on the environment

`input_size = 4` #Input size is 4 because of that is the state info

`nb_actions = 2` #This is because there are two action, left or right

`action = np.random.choice(a=[0, 1], p=probabilities.detach().numpy())`

`total_reward = 14.0`

`total_reward = 24.0`

`total_reward = 15.0`

Objective:

The aim is to modify the 'ActorModelV0' class within the 'actor_v0.py' file to suit the specific requirements of the Cartpole environment, which involves configuring the input size and the number of actions correctly.

Implementation:

The Cartpole environment is characterized by a state vector consisting of four continuous variables: cart position, cart velocity, pole angle, and pole angular velocity. The action space is binary, consisting of two possible actions: moving the cart left or right.

The `ActorModelV0` class was initially set up with placeholder values for `input_size` and `nb_actions`. To align the model with the Cartpole environment, these values were updated to reflect the actual dimensions of the state and action spaces. Specifically, the following changes were made in the `__init__` method of the `ActorModelV0` class:

1. Input Size: Updated from 0 to 4 to match the four dimensions of the Cartpole state vector.
2. Number of Actions: Updated from 0 to 2 to represent the two possible actions in the Cartpole environment.

Updated Code Snippet:

```
```python
class ActorModelV0(nn.Module):
 """Deep neural network designed for policy approximation in Cartpole environment."""

 def __init__(self):
 super(ActorModelV0, self).__init__()
 input_size = 4 # Cart position, cart velocity, pole angle, pole angular velocity
 nb_actions = 2 # Actions: move left, move right

 self.fc0 = nn.Linear(input_size, 128)
 self.fc1 = nn.Linear(128, 128)
 self.policy = nn.Linear(128, nb_actions)
 ...
```

Implementation of Action Selection in 1\_test\_model.py

Objective:

To implement a method for selecting actions based on the probabilities output by the `ActorModelV0` when testing the policy on the Cartpole environment.

Implementation:

The policy network outputs the probabilities of taking each possible action given a state. The action selection process involves choosing an action based on these probabilities to ensure that actions with higher probabilities are more likely to be selected.

The action selection was implemented using a probabilistic approach, where the softmax function is applied to the network's output to obtain a valid probability distribution over actions.

An action is then randomly selected according to this distribution. This method allows the exploration of actions proportional to their estimated utility.

Updated Code Snippet in `1_test_model.py`:

```
```python
probabilities = torch.nn.functional.softmax(policy(state), dim=-1).squeeze() # Normalize logits
to probabilities
action = np.random.choice(a=[0, 1], p=probabilities.detach().numpy()) # Randomly select an
action based on probabilities
```
```

The updates to the actor model ensure that it accurately reflects the characteristics of the Cartpole environment, enhancing the relevance and effectiveness of the learned policy. The probabilistic action selection method allows for both exploitation of known good actions and exploration of less frequently chosen actions, which is crucial in dynamic environments like Cartpole.

## Performance Analysis

The scores provided suggest that on average, the pole falls after about 14 to 24 steps. This level of performance is relatively low, indicating that the policy may struggle with effectively balancing the pole over longer durations.

## Possible Reasons for This Performance

### 1. Insufficient Training:

- The neural network (actor) may not have been trained for enough iterations or with enough data to adequately learn the task.
- The model might not have encountered a sufficient variety of states during training, limiting its ability to generalize to new situations during testing.

### 2. Exploration vs. Exploitation:

- If the training process did not properly balance exploration and exploitation, the policy might be overfitting to particular scenarios or not exploring the action space thoroughly enough. This can result in a policy that performs well in known states but poorly in unseen or slightly different states.

## Steps to Improve Performance

- Longer or More Effective Training: Ensure that the model trains over a larger number of episodes or until performance stabilizes at a higher reward threshold.
- Hyperparameter Tuning: Experiment with learning rates, batch sizes, the number of neurons, layers, and other architecture details.
- Enhanced Exploration Techniques: Implement or enhance exploration mechanisms like epsilon-greedy strategies during training to ensure a wider range of state-action pairs are experienced.

## Section 3

### Implementation

#### 1. Discounted Reward Computation:

The implementation effectively handles the inherent variance of Monte Carlo methods by computing the sum of discounted rewards within the `reinforce.py` script:

```
```python
discounted_rewards = [0.0 for _ in saved_rewards]
discounted_reward = 0 # Discounted reward
for t in reversed(range(len(saved_rewards))):
    discounted_reward = saved_rewards[t] + DISCOUNT_FACTOR * discounted_reward
    discounted_rewards[t] = discounted_reward

discounted_rewards = torch.tensor(discounted_rewards)
mean, std = discounted_rewards.mean(), discounted_rewards.std()
discounted_rewards = (discounted_rewards - mean) / (std + 1e-7)
```
```

This method ensures that rewards expected in the future are appropriately weighted less than immediate rewards, adhering to the principle that uncertainties in future outcomes should be progressively discounted.

#### 2. REINFORCE Loss Calculation:

The REINFORCE loss function connects the logarithm of the policy's action probabilities to the computed discounted rewards, forming a direct feedback mechanism that promotes actions leading to higher rewards:

```

python
actor_loss = []
for p, g in zip(saved_probabilities, discounted_rewards):
 time_step_actor_loss = -torch.log(p) * g
 actor_loss.append(time_step_actor_loss)

actor_loss = torch.cat(actor_loss).sum()

```

This approach encourages the policy to increase the likelihood of actions that yield higher returns, optimizing the agent's behavior to maximize rewards.

### 3. Hyperparameter Optimization:

Adjusting the learning rate ( $\alpha$ ) and the discount factor ( $\beta$ ) is crucial for the training process. The selected values ( $\alpha = 0.001$ ), ( $\beta = 0.9999$ ) strike a balance between ensuring efficient learning and maintaining stability in updates, allowing the model to learn consistently without significant oscillations in policy behavior.

### 4. Training Termination Metric:

Using a moving average of recent rewards, this approach monitors the performance over a series of episodes, ensuring that the stopping criterion is based on consistent performance rather than a single peak reward:

```

python
moving_average_reward = sum(recent_rewards[-N:]) / N
if episode_count > MIN_EPISODES:
 if (episode_count % CHECK_INTERVAL == 0 and
 (moving_average_reward - best_mean_reward < IMPROVEMENT_THRESHOLD)):
 print(f"Training stopped after {episode_count} episodes due to lack of improvement.")
 break
if moving_average_reward >= TARGET_REWARD:
 print(f"Training stopped after reaching target reward of {TARGET_REWARD}.")
 break

```

This methodology helps ensure that the agent's learning is robust and stable across multiple episodes, rather than being influenced by fluctuations or anomalies in single episodes. It's an indicator of the agent's overall capability and learning stability.

## Testing Approaches and Evidence

### Testing Procedure:

The `'3_test_model.py'` script evaluates the effectiveness of the trained policy by testing it in varied initial conditions within the same Cartpole environment.

### Results and Analysis:

During testing, the policy consistently achieves the target cumulative reward of 500, demonstrating its robustness and efficacy. This performance signifies not only the success of training and parameter optimization but also the model's ability to generalize across different states and dynamics of the environment.

### Evidence of Effectiveness:

The consistent high rewards observed during the testing phase serve as concrete evidence of the successful application of the REINFORCE algorithm. This empirical data supports the theoretical underpinnings of the algorithm's effectiveness in real-world scenarios.

## Concluding Thoughts

The training phase showed a marked improvement in rewards, culminating consistently at the maximum reward of 500 over multiple iterations. This upward trajectory in performance reflects the robustness and adaptability of the trained policy under the REINFORCE algorithm. The final testing confirmed the model's reliability and effectiveness, achieving the optimal score: Total reward = 500.0. These results highlight the algorithm's success in real-world application, proving its ability to train an agent that performs optimally in varied conditions.

## Output

`refinorace.py`

```
iteration 825 - last reward: 376.00
iteration 830 - last reward: 353.00
iteration 835 - last reward: 341.00
iteration 840 - last reward: 360.00
iteration 845 - last reward: 354.00
iteration 850 - last reward: 424.00
iteration 855 - last reward: 500.00
```

iteration 860 - last reward: 500.00  
iteration 865 - last reward: 500.00  
iteration 870 - last reward: 500.00  
iteration 875 - last reward: 500.00  
iteration 880 - last reward: 500.00  
iteration 885 - last reward: 500.00  
iteration 890 - last reward: 500.00  
iteration 895 - last reward: 500.00  
iteration 900 - last reward: 500.00  
iteration 905 - last reward: 500.00  
iteration 910 - last reward: 500.00  
iteration 915 - last reward: 500.00  
iteration 920 - last reward: 500.00  
iteration 925 - last reward: 500.00  
iteration 930 - last reward: 500.00  
iteration 935 - last reward: 500.00  
iteration 940 - last reward: 500.00  
iteration 945 - last reward: 500.00  
iteration 950 - last reward: 500.00

Training stopped after reaching target reward of 500.

hariharanjanardhanan@HariHarans-MacBook-Pro          mlpr          %          /usr/bin/python3  
/Users/hariharanjanardhanan/Desktop/mlpr/3\_test\_model.py

ActorModelV1(  
  (fc0): Linear(in\_features=2, out\_features=128, bias=True)  
  (fc1): Linear(in\_features=128, out\_features=128, bias=True)  
  (policy): Linear(in\_features=128, out\_features=2, bias=True)  
)

Total reward = 500.0



```

iteration 855 - last reward: 500.00
iteration 860 - last reward: 500.00
iteration 865 - last reward: 500.00
iteration 870 - last reward: 500.00
iteration 875 - last reward: 500.00
iteration 880 - last reward: 500.00
iteration 885 - last reward: 500.00
iteration 890 - last reward: 500.00
iteration 895 - last reward: 500.00
iteration 900 - last reward: 500.00
iteration 905 - last reward: 500.00
iteration 910 - last reward: 500.00
iteration 915 - last reward: 500.00
iteration 920 - last reward: 500.00
iteration 925 - last reward: 500.00
iteration 930 - last reward: 500.00
iteration 935 - last reward: 500.00
iteration 940 - last reward: 500.00
iteration 945 - last reward: 500.00
iteration 950 - last reward: 500.00
Training stopped after reaching target reward of 500.
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3 /Users/hariharanjanardhanan/Desktop/mlpr/3_test_model.py
ActorModelV1(
 (fc0): Linear(in_features=2, out_features=128, bias=True)
 (fc1): Linear(in_features=128, out_features=128, bias=True)
 (policy): Linear(in_features=128, out_features=2, bias=True)
)
Total reward = 500.0

```

## Section 4

### Implementation

#### Modifications in CartpoleEnvV1

In CartpoleEnvV1, accelerations are omitted from state data, reducing the observation size to 2. This requires a corresponding update in ActorModelV1:

```

Modify the state processing to work without accelerations, making the environment
effectively Markovian by using the rate of change for position and angle.
This change helps capture the dynamics necessary for the REINFORCE algorithm to learn
effectively.

```

```

if self.last_state is not None:

```

```

 state_diff = self.state - self.last_state

```

```

 processed_state = processed_state + (state_diff[[0, 2]] / 0.2)

```

```

self.last_state = self.state

```

```

return processed_state

```

#### Actor Model Adjustment in ActorModelV1

The ActorModelV1 is adapted to handle a state input size of 2, reflecting the reduced state information:

```

Initialize model with adjusted input and output sizes based on the new environment's state size

```

```
input_size = 2 # Updated to match the reduced state size of CartpoleEnvV1
nb_actions = 2 # Matches the number of possible actions in the environment
```

### Training and Testing the Model

With these adjustments, scripts 1 and 2 train and test the model to evaluate convergence:

#### Initial Model Convergence Before Modifications

Does it converge?

- No, the model does not converge initially.

Why?

- The primary reason for the initial lack of convergence is the removal of acceleration components from the state representation in `CartpoleEnvV1`. This results in a state that lacks critical dynamic information necessary for predicting future movements effectively. The agent struggles to anticipate and balance the cartpole because it lacks insights into how fast the pole is moving or how its acceleration affects stability. This insufficient state representation prevents the model from developing a robust policy that can navigate the complexities of balancing the cartpole.

#### Model Convergence After Adjusting the State

To address the initial convergence issue, the state representation in `CartpoleEnvV1` is adjusted to include the rate of change of position and angle, compensating for the absence of direct acceleration data. This change aims to restore the Markov property, where the state includes sufficient information for predicting future states based on the current state alone.

Does it converge?

- Yes, the model converges after these adjustments.

Why?

- By including the rate of change (computed as the difference between the current and previous states normalized by a factor) in the state representation, the environment effectively captures the necessary dynamics that were initially missing. This enhancement allows the agent to infer velocities and trends in movement, which are crucial for making informed decisions. The REINFORCE algorithm leverages this enriched state information to learn a policy that can successfully navigate the cartpole to achieve and maintain balance. The continuous updates and richer state feedback loop lead to improved learning and eventually convergence, as evidenced by consistent high rewards during training and testing.

## Results from Script Runs

Initial Test with 1\_test\_model.py: The model achieves a total reward of 19.0, indicating struggles due to incomplete state information.

Training with 2\_reinforce.py: Post-modification, the training shows consistent achievement of 500.0 rewards, demonstrating successful convergence with the adjusted state dynamics.

Final Test with 3\_test\_model.py: The trained model reaches a total reward of 500.0, confirming its robust performance post-adjustments.

## Conclusion

Initially, the model struggles to converge due to the lack of acceleration data in the state, which impairs its ability to predict effective actions. By modifying the state in CartpoleEnvV1 to include the rate of change, the environment regains the Markov property, enabling the REINFORCE algorithm to effectively train the model. This leads to convergence, as the modified states provide sufficient information for decision-making.

## Output before

```
1_test_model.py
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3
/Users/hariharanjanardhanan/Desktop/mlpr/1_test_model.py
ActorModelV1(
 (fc0): Linear(in_features=2, out_features=128, bias=True)
 (fc1): Linear(in_features=128, out_features=128, bias=True)
 (policy): Linear(in_features=128, out_features=2, bias=True)
)
total_reward = 19.0
```

```
refinorice.py
iteration 990 - last reward: 38.00
iteration 995 - last reward: 38.00
iteration 1000 - last reward: 32.00
iteration 1005 - last reward: 41.00
iteration 1010 - last reward: 36.00
iteration 1015 - last reward: 39.00
iteration 1020 - last reward: 13.00
iteration 1025 - last reward: 37.00
iteration 1030 - last reward: 69.00
```

iteration 1035 - last reward: 58.00

iteration 1040 - last reward: 67.00

iteration 1045 - last reward: 37.00

iteration 1050 - last reward: 23.00

Training stopped after 1050 episodes due to lack of improvement.

hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3

/Users/hariharanjanardhanan/Desktop/mlpr/3\_test\_model.py

ActorModelV1(

(fc0): Linear(in\_features=2, out\_features=128, bias=True)

(fc1): Linear(in\_features=128, out\_features=128, bias=True)

(policy): Linear(in\_features=128, out\_features=2, bias=True)

)

Total reward = 43.0

```
iteration 990 - last reward: 38.00
iteration 995 - last reward: 38.00
iteration 1000 - last reward: 32.00
iteration 1005 - last reward: 41.00
iteration 1010 - last reward: 36.00
iteration 1015 - last reward: 39.00
iteration 1020 - last reward: 13.00
iteration 1025 - last reward: 37.00
iteration 1030 - last reward: 69.00
iteration 1035 - last reward: 58.00
iteration 1040 - last reward: 67.00
iteration 1045 - last reward: 37.00
iteration 1050 - last reward: 23.00
Training stopped after 1050 episodes due to lack of improvement.
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3 /Users/hariharanjanardhanan/Desktop/mlpr/3_test_model.py
ActorModelV1(
 (fc0): Linear(in_features=2, out_features=128, bias=True)
 (fc1): Linear(in_features=128, out_features=128, bias=True)
 (policy): Linear(in_features=128, out_features=2, bias=True)
)
Total reward = 43.0
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3 /Users/hariharanjanardhanan/Desktop/mlpr/1_test_model.py
ActorModelV1(
 (fc0): Linear(in_features=2, out_features=128, bias=True)
 (fc1): Linear(in_features=128, out_features=128, bias=True)
 (policy): Linear(in_features=128, out_features=2, bias=True)
)
total_reward = 19.0
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr %
```

Output after

1\_test\_model.py

hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3

/Users/hariharanjanardhanan/Desktop/mlpr/1\_test\_model.py

ActorModelV1(

(fc0): Linear(in\_features=2, out\_features=128, bias=True)

(fc1): Linear(in\_features=128, out\_features=128, bias=True)

(policy): Linear(in\_features=128, out\_features=2, bias=True)

)

total\_reward = 23.0

./2\_reinforce.py

iteration 940 - last reward: 500.00

iteration 945 - last reward: 279.00

iteration 950 - last reward: 163.00

iteration 955 - last reward: 150.00

iteration 960 - last reward: 250.00

iteration 965 - last reward: 500.00

iteration 970 - last reward: 500.00

iteration 975 - last reward: 500.00

iteration 980 - last reward: 500.00

iteration 985 - last reward: 500.00

iteration 990 - last reward: 500.00

iteration 995 - last reward: 500.00

iteration 1000 - last reward: 500.00

iteration 1005 - last reward: 500.00

iteration 1010 - last reward: 500.00

iteration 1015 - last reward: 500.00

iteration 1020 - last reward: 500.00

iteration 1025 - last reward: 500.00

iteration 1030 - last reward: 500.00

iteration 1035 - last reward: 500.00

iteration 1040 - last reward: 500.00

iteration 1045 - last reward: 500.00

iteration 1050 - last reward: 500.00

Training stopped after 1050 episodes due to lack of improvement.

3\_test\_model.py

hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3

/Users/hariharanjanardhanan/Desktop/mlpr/3\_test\_model.py

ActorModelV1(

(fc0): Linear(in\_features=2, out\_features=128, bias=True)

(fc1): Linear(in\_features=128, out\_features=128, bias=True)

(policy): Linear(in\_features=128, out\_features=2, bias=True)

)

Total reward = 500.0

```

iteration 925 - last reward: 500.00
iteration 930 - last reward: 500.00
iteration 935 - last reward: 500.00
iteration 940 - last reward: 500.00
iteration 945 - last reward: 279.00
iteration 950 - last reward: 163.00
iteration 955 - last reward: 150.00
iteration 960 - last reward: 250.00
iteration 965 - last reward: 500.00
iteration 970 - last reward: 500.00
iteration 975 - last reward: 500.00
iteration 980 - last reward: 500.00
iteration 985 - last reward: 500.00
iteration 990 - last reward: 500.00
iteration 995 - last reward: 500.00
iteration 1000 - last reward: 500.00
iteration 1005 - last reward: 500.00
iteration 1010 - last reward: 500.00
iteration 1015 - last reward: 500.00
iteration 1020 - last reward: 500.00
iteration 1025 - last reward: 500.00
iteration 1030 - last reward: 500.00
iteration 1035 - last reward: 500.00
iteration 1040 - last reward: 500.00
iteration 1045 - last reward: 500.00
iteration 1050 - last reward: 500.00
Training stopped after 1050 episodes due to lack of improvement.
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr % /usr/bin/python3 /Users/hariharanjanardhanan/Desktop/mlpr/3_test_model.py
ActorModelV1(
 (fc0): Linear(in_features=2, out_features=128, bias=True)
 (fc1): Linear(in_features=128, out_features=128, bias=True)
 (policy): Linear(in_features=128, out_features=2, bias=True)
)
Total reward = 500.0
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr %

```

## Section 5

### Implementing and Adapting the REINFORCE Algorithm for the GridWorld Environment

#### Overview

The REINFORCE algorithm is a Monte Carlo policy gradient method used in reinforcement learning to optimize policy directly. This outlines the implementation of the REINFORCE algorithm specifically designed for a GridWorld environment, where an agent navigates a grid to reach a goal. Additionally, it discusses how this implementation can be adapted to other types of environments or reinforcement learning tasks.

#### Implementation Details

##### GridWorld Environment

##### Description:

The GridWorld environment is a simple and commonly used testbed in reinforcement learning. The agent's objective is to navigate through a grid from a start position to a goal position, minimizing the number of moves and avoiding possible penalties.

##### Key Components:

State: Defined as the agent's current position in the grid, typically represented as coordinates (x, y).

Actions: Four possible movements—up, down, left, right.

Rewards: Typically, a small penalty (e.g., -1) for each move to encourage shortest path solutions and a larger positive reward (e.g., +100) upon reaching the goal.

Actor Model

Structure:

A simple neural network that inputs the state of the environment (often one-hot encoded to represent grid positions) and outputs a probability distribution over possible actions.

Input Layer: Matches the number of possible states (e.g., 25 for a 5x5 grid).

Output Layer: Corresponds to the number of possible actions (4 in this scenario).

Functionality:

The model predicts which action to take at each state by providing probabilities for each action, which guides the agent's decision-making process.

REINFORCE Training Script

Algorithm Workflow:

Initialize the Environment and Model: Set up the GridWorld environment and the Actor Model with appropriate parameters.

Episode Iteration: For a set number of episodes, the agent interacts with the environment:

State Initialization: Begin from a start position.

Action Selection: Use the model to determine the best action based on the current state.

Environment Interaction: Execute the action in the environment, observe the new state and reward.

Reward Accumulation and Log Probability Recording: Store rewards and log probabilities of actions for later gradient calculation.

Policy Update:

At the end of each episode, calculate the return (total discounted rewards) for each state-action pair.

Update the policy by maximizing expected returns using gradient ascent on the logged probabilities.

Loss Calculation:

The loss is calculated as the negative of the sum of the product of log probabilities and discounted rewards, promoting actions that lead to higher rewards.

## Output Analysis from an Example Execution

From the recorded episode outputs:

**Actions and States:** The agent navigates through the grid, with actions and subsequent state changes logged, such as moving from (2, 3) to (3, 3) and eventually to (4, 4), the goal position.

**Rewards:** The agent incurs a -1 penalty for each move but gains a significant reward of 100 upon reaching the goal.

**Policy Behavior:** Action probabilities show the agent's decision-making process, highlighting the model's learning and preference over different actions in given states.

```
New State: (4, 3), Reward: -1, Done: False
Action Probabilities: tensor([[0.0985, 0.5421, 0.0705, 0.2889]], grad_fn=<SoftmaxBackward0>)
Action taken: 2
New State: (3, 3), Reward: -1, Done: False
Action Probabilities: tensor([[0.0843, 0.2706, 0.0887, 0.5565]], grad_fn=<SoftmaxBackward0>)
Action taken: 3
New State: (4, 3), Reward: -1, Done: False
Action Probabilities: tensor([[0.0985, 0.5421, 0.0705, 0.2889]], grad_fn=<SoftmaxBackward0>)
Action taken: 2
New State: (3, 3), Reward: -1, Done: False
Action Probabilities: tensor([[0.0843, 0.2706, 0.0887, 0.5565]], grad_fn=<SoftmaxBackward0>)
Action taken: 1
New State: (3, 4), Reward: -1, Done: False
Action Probabilities: tensor([[0.0984, 0.3493, 0.0624, 0.4898]], grad_fn=<SoftmaxBackward0>)
Action taken: 3
New State: (4, 4), Reward: 100, Done: True
Episode 1000: Total Reward = -8.881784197001252e-16, Steps = 12, Loss = -2.7902092933654785
hariharanjanardhanan@HariHarans-MacBook-Pro mlpr %
```

## Thoughts on REINFORCE in GridWorld

REINFORCE's strength lies in its simplicity and direct approach to optimizing policies. However, its efficiency can be compromised by high variance in the estimated gradients, often requiring careful tuning of the learning rate and reward signals.

## Approaches for Enhancing REINFORCE

**Reward Shaping:** Adjusting the reward structure, such as scaling or adding intermediate rewards, can significantly stabilize training.

**Baseline Implementation:** Introducing a baseline subtraction from the rewards can reduce gradient variance, thus stabilizing learning.

**Exploration Techniques:** Beyond epsilon-greedy, methods like entropy regularization can encourage exploration by penalizing low entropy in action distributions.

The adaptation of the REINFORCE algorithm to the GridWorld environment demonstrates the flexibility of reinforcement learning techniques in addressing varied types of problems. Through this implementation, foundational concepts of state representation, policy-based learning, and reward structuring are employed to effectively train an agent in a simple yet illustrative task.



This exercise not only highlights the practical application of the REINFORCE algorithm but also serves as a baseline for further exploration into more complex environments and learning algorithms.

2  
...