

# Upper Confidence Bound (UCB) Algorithm for CartPole-v1

---

## Project Overview

This project demonstrates the application of the **Upper Confidence Bound (UCB)** algorithm, an exploration strategy, to the classic reinforcement learning environment **CartPole-v1** from OpenAI's Gym. The objective is to train an agent to balance a pole on a moving cart by applying forces to the cart, preventing the pole from falling over.

---

## Why Upper Confidence Bound (UCB)?

The UCB algorithm is based on the principle of "optimism in the face of uncertainty." It selects actions that maximize the upper confidence bound of the expected reward, effectively balancing exploration and exploitation. This approach is particularly useful in scenarios like the multi-armed bandit problem, where an agent must choose between multiple options with uncertain rewards. ([GeeksforGeeks](#))

---

## Prerequisites

### Required Libraries

- `numpy` : For numerical computations.
- `gym` : For the CartPole-v1 environment.

### Installation

Install the necessary libraries using pip:

```
pip install numpy gym
```

---

## Code Description

The implementation consists of the following key components:

### 1. Importing Libraries

```
import numpy as np
import gym
```

### 2. Environment Setup

Initialize the CartPole-v1 environment:

```
env = gym.make('CartPole-v1')
```

### 3. Hyperparameters

Define the hyperparameters for the UCB algorithm:

```
alpha = 0.1 # Learning rate
gamma = 0.99 # Discount factor
c = 2 # Exploration factor in UCB
```

## 4. Q-Table and Action Count Initialization

Initialize the Q-table and a counter for action selections:

```
Q = np.zeros((env.observation_space.shape[0], env.action_space.n))
N = np.zeros((env.observation_space.shape[0], env.action_space.n)) # Action count
```

## 5. UCB Policy Function

Define the policy for selecting actions based on the UCB algorithm:

```
def ucb_policy(state):
    total_actions = np.sum(N[state])
    if total_actions == 0: # No actions taken yet, explore
        return np.random.choice(env.action_space.n)
    else:
        ucb_values = Q[state] + c * np.sqrt(np.log(total_actions + 1) / (N[state] + 1))
        return np.argmax(ucb_values)
```

## 6. Training Loop

Implement the training loop to update the Q-values and action counts based on the agent's interactions with the environment:

```
def ucb_q_learning(env, n_episodes=1000):
    for episode in range(n_episodes):
        state = env.reset()
        done = False
        while not done:
            action = ucb_policy(state) # Select action based on UCB policy
            next_state, reward, done, _ = env.step(action)

            # Update Q-value and action count
            N[state, action] += 1
            Q[state, action] += alpha * (reward + gamma * np.max(Q[next_state]) - Q[state, action])

            state = next_state
```

## 7. Running the Algorithm

Execute the training function:

```
ucb_q_learning(env)
```

---

## Expected Outcomes

- **Policy Development:** The agent learns a policy to balance the pole by selecting actions that maximize the upper confidence bound of the expected reward.

- **Efficient Exploration:** The UCB algorithm effectively balances exploration and exploitation, leading to efficient learning in the CartPole-v1 environment.
- 

## Use Cases

- **Reinforcement Learning:** Applying exploration strategies in environments where balancing exploration and exploitation is crucial.
  - **Multi-Armed Bandit Problems:** Situations where decisions need to be made under uncertainty with the goal of maximizing cumulative rewards.
- 

## Future Enhancements

- **State Discretization:** Implement state discretization techniques to handle the continuous state space of CartPole-v1.
  - **Reward Shaping:** Modify the reward structure to provide more informative feedback to the agent.
  - **Comparison with Other Algorithms:** Compare the performance of the UCB algorithm with other exploration strategies like  $\epsilon$ -greedy or Thompson Sampling.
- 

## References

- [Upper Confidence Bound Algorithm in Reinforcement Learning](#)
  - [The Upper Confidence Bound Algorithm – Bandit Algorithms](#)
-