

# Reinforcement Learning for Cryptocurrency Trading: A Comparative Study of PPO and DQN

## ABSTRACT

In this project, we utilize reinforcement learning techniques to design and fine-tune trading strategies for cryptocurrencies. The focus is placed on two prominent reinforcement learning algorithms: Proximal Policy Optimization (PPO) and Deep Q-Network (DQN). We train both a PPO agent and a DQN agent using historical Bitcoin price data and evaluate their performances within a simulated trading environment. The strengths and weaknesses of each approach are discussed, illustrating the unique challenges posed by the application of reinforcement learning to cryptocurrency trading. This project not only demonstrates the potential of reinforcement learning in crafting data-driven cryptocurrency trading strategies but also provides valuable insights for future enhancements and research directions. The comparative study between PPO and DQN offers a broader perspective on the capabilities of different reinforcement learning algorithms in the context of financial markets.

## INTRODUCTION:

Reinforcement Learning (RL) has proven to be an effective approach in developing algorithms capable of learning optimal strategies through interactions with an environment. It has been successfully used in various fields, such as game playing, robotics, and resource management. In the field of finance, RL has been used to create automated trading systems capable of learning optimal trading strategies.

In recent years, the cryptocurrency market has emerged as a new avenue for traders and investors. However, the highly volatile nature of the cryptocurrency market presents a challenging environment for developing profitable trading strategies. RL, with its ability to learn from an uncertain and dynamic environment, provides a promising approach to tackle these challenges.

### Problem Statement:

The primary goal of this study is to apply and compare the performances of two RL methods, PPO and DQN, in developing a profitable trading strategy for the Bitcoin cryptocurrency. We use historical Bitcoin price data as our environment and the RL agent's task is to learn to make buy, sell, or hold decisions to maximize cumulative rewards, i.e., trading profits.

## Related Work

RL has been previously used in the context of stock trading. For instance, [Deep Reinforcement Learning for Automated Stock Trading proposes an automated stock trading](#) system using RL. However, the application of RL in cryptocurrency trading is relatively less explored, and the comparison of different RL algorithms in this domain is even less studied. Our work aims to fill this gap.

1. **"Deep Reinforcement Learning for Automated Stock Trading: An Ensemble Strategy"** by Jiang et al. (2019): This study proposes an ensemble strategy that combines multiple RL agents to improve trading performance. The authors use DQN and other RL algorithms to train agents on historical stock data.
2. **"Cryptocurrency Portfolio Management with Deep Reinforcement Learning"** by Xu et al. (2018): This paper applies DQN to cryptocurrency portfolio management. The authors train an agent to dynamically allocate capital among a set of cryptocurrencies to maximize returns.
3. **"Deep reinforcement learning for the optimal placement of cryptocurrency limit orders"** by Schnaubelt, M. (2022): This paper applies PPO to optimize limit order placement.

## Comparison of PPO and DQN

Proximal Policy Optimization (PPO) and Deep Q-Learning (DQN) are both advanced reinforcement learning methods that use deep learning to approximate the policy or value function. However, they differ in their approach and underlying algorithms.

### Deep Q-Learning (DQN)

DQN is a value-based method that learns an approximation of the action-value function, also known as the Q-function. The Q-function is defined as:

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

where  $s$  is the current state,  $a$  is the action taken,  $r$  is the immediate reward,  $\gamma$  is the discount factor,  $s'$  is the next state, and  $a'$  are the possible actions in the next state.

The objective of DQN is to minimize the difference between the predicted Q-values and the target Q-values, expressed by the following loss function:

$$L = \mathbb{E}[(r + \gamma \max_{a'} Q(s', a') - Q(s, a))^2]$$

DQN introduces techniques like experience replay and target networks to handle challenges like correlation between samples and overestimation of Q-values.

### Proximal Policy Optimization (PPO)

PPO, on the other hand, is a policy-based method that directly optimizes the policy function. Unlike DQN, PPO does not try to estimate the value function but instead tries to find the policy that will maximize the expected return.

The objective function for PPO is defined as:

$$L^{CLIP}(\theta) = \mathbb{E} \left[ \min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t) \right]$$

where  $r_t(\theta) = \frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)}$  is the probability ratio,  $A_t$  is the advantage function at time  $t$ , and  $\epsilon$  is a hyperparameter that limits the update step size to avoid large policy updates.

PPO introduces a surrogate objective function with a clipped objective to prevent overly large policy updates, leading to more stable and efficient learning.

In both methods, the optimization is done using stochastic gradient descent or its variants. While DQN uses a replay buffer to break the correlation of samples and uses target networks to stabilize learning, PPO uses the concept of a trust region to ensure the new policy does not deviate far from the old policy.

## States, Actions, Environment, and Rewards:

**States:** In this implementation, the states are represented by a window of historical Bitcoin prices. We use a window size of 10 in our example. The states are pre-processed by scaling the data using StandardScaler. Each state is an array of length 10, containing the historical prices of Bitcoin.

- **Actions:** There are three possible actions in this environment: Buy (0), Sell (1), and Hold (2). The agent decides which action to take based on the current state of the environment.
- **Environment:** The custom trading environment, CryptoTradingEnv, is a subclass of the gym.Env class. It takes the historical price data and the window size as input. The environment has discrete action space and continuous observation space. It implements the reset() and step() functions, which are used to reset the environment and take an action, respectively.
- **Rewards:** The rewards are calculated as follows:
  - If the agent buys, the reward is the difference between the current price and the previous price. If the agent sells, the reward is the difference between the previous price and the current price. If the agent holds, the reward is 0.

## Training and Testing Data:

In this implementation, we do not explicitly separate the data into training and testing sets. Instead, the agent is trained and tested on the entire historical price data available. However, it is a good practice to separate the data into training and testing sets to evaluate the performance of the agent more accurately.

To create separate training and testing sets, we can split the data into two parts, for example, using the first 80% of the data for training and the remaining 20% for testing. You can then create separate environments for training and testing with their respective data and evaluate the agent's performance on both.

# IMPLEMENTATION

### ▼ Step 1: Import necessary libraries

```
1 import numpy as np
2 import pandas as pd
3 import random
4 from collections import deque
5 import gym
6 from gym import spaces
7 import tensorflow as tf
8 from tensorflow.keras.layers import Dense, Input
9 from tensorflow.keras.models import Model
10 from tensorflow.keras.optimizers import Adam
11 import matplotlib.pyplot as plt
12 from sklearn.preprocessing import StandardScaler
13 import yfinance as yf
```

### Step 2: Loading the dataset and preparing it for the agent

### ▼ Data Pre-processing and Requirements:

We used historical Bitcoin price data from Yahoo Finance, with the ticker symbol "BTC-USD". The data spans from January 1, 2010, to May 8, 2023. We use the 'Close' prices for the analysis.

The data is pre-processed using the preprocess\_data function. This function scales the data using StandardScaler from the sklearn.preprocessing module. The function also creates states by taking a sliding window of the specified size (10 in our example) on the scaled data.

To run the models, the following libraries are required: numpy, pandas, gym, tensorflow, matplotlib, sklearn, and yfinance.

The window length in time-series analysis, including financial data analysis like in cryptocurrency trading, is a parameter that determines the size of the "look-back" period used to make predictions or inform decision-making processes.

In the context fo our work, a window length of 10 implies that the model will use the previous 10 data points (such as the prices of the last 10 periods) to inform its decision for the next action (buy, sell, hold).

The choice of window length is largely empirical and can significantly impact the performance of the model. Here are some reasons for choosing a window length of 10:

- 1. Historical Performance: Previous experiments or related works have shown that a window length of 10 provides a good balance between learning meaningful trends in the data and not overfitting to noise or short-term fluctuations.
- 2. Computational Efficiency: A smaller window size reduces the dimensionality of the data, leading to faster training times and lower computational resource requirements. It's a balance between complexity and performance.
- 3. Short-Term Trends: In fast-moving markets like cryptocurrencies, recent information may be more relevant for decision making. A window length of 10 prioritizes more recent data without completely ignoring slightly older information.
- 4. Avoid Overfitting: A longer window might capture more information, but it might also introduce more noise into the model, leading to overfitting. A window of 10 could be a compromise to include just enough information for the model to learn from and make predictions

```
1 def preprocess_data(data, window_size):
2     scaler = StandardScaler()
3     data = scaler.fit_transform(data.reshape(-1, 1))
4     states = []
5     for i in range(len(data) - window_size - 1):
6         state = data[i:i + window_size]
7         states.append(state)
8     return np.array(states)
9
```

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should\_run\_async` will not call `transform\_cell` automatically; and should\_run\_async(code)

```
1 # Load the dataset (historical Bitcoin prices)
2
3 # Define the ticker symbol
4 ticker = "BTC-USD"
5
6 # Set the start and end dates for the data
7 start_date = "2010-01-01"
8 end_date = "2023-05-08"
9
10 # Get the data from Yahoo Finance
11 data = yf.download(ticker, start=start_date, end=end_date)
12
13 data = data['Close'].values
14
15 # Define the window size for the state
16 window_size = 10
17
18 # Preprocess the data
19 states = preprocess_data(data, window_size)
20 states
```

[\*\*\*\*\*100%\*\*\*\*\*] 1 of 1 completed

```
array([[[-0.80291525],
        [-0.80496613],
        [-0.80681438],
        ...,
        [-0.80504313],
        [-0.8057683 ],
        [-0.80621403]],

       [[-0.80496613],
        [-0.80681438],
        [-0.80593477],
        ...,
        [-0.8057683 ],
        [-0.80621403],
        [-0.80651984]],

       [[-0.80681438],
        [-0.80593477],
        [-0.80656342],
        ...,
        [-0.80621403],
        [-0.80651984],
        [-0.80791264]],

       ...,

       [[ 0.88472144],
        [ 0.93349367],
        [ 0.94067015],
        ...,
        [ 0.92002467],
        [ 0.95674571],
        [ 0.97705692]],

       [[ 0.93349367],
        [ 0.94067015],
        [ 1.00620328],
        ...,
        [ 0.95674571],
        [ 0.97705692],
        [ 0.96716867]],

       [[ 0.94067015],
        [ 1.00620328],
        [ 0.99787823],
        ...,
        [ 0.97705692],
        [ 0.96716867],
        [ 1.00998142]]]])
```

▼ Step 3: Define the custom trading environment

```

1 class CryptoTradingEnv(gym.Env):
2     def __init__(self, data, window_size):
3         super(CryptoTradingEnv, self).__init__()
4         self.data = data
5         self.window_size = window_size
6         self.current_step = self.window_size
7         self.action_space = spaces.Discrete(3) # Buy, Sell, Hold
8         self.observation_space = spaces.Box(low=0, high=1, shape=(self.window_size,))
9
10    def reset(self):
11        self.current_step = self.window_size
12        return self.data[self.current_step - self.window_size:self.current_step]
13
14    def step(self, action):
15        self.current_step += 1
16        if self.current_step >= len(self.data):
17            return self.reset()
18
19        price_now = self.data[self.current_step]
20        price_prev = self.data[self.current_step - 1]
21
22        if action == 0: # Buy
23            reward = price_now - price_prev
24        elif action == 1: # Sell
25            reward = price_prev - price_now
26        else: # Hold
27            reward = 0
28
29        next_state = self.data[self.current_step - self.window_size:self.current_step]
30
31        done = self.current_step == len(self.data) - 1
32
33        return next_state, reward, done, {}
34

```

### ▼ Step 3: Define the PPO agent

```

1 import numpy as np
2 import tensorflow as tf
3 from tensorflow.keras.layers import Dense
4
5 class PPOAgent:
6     def __init__(self, state_size, action_size, learning_rate=0.001, gamma=0.99, epsilon_clip=0.2):
7         self.state_size = state_size
8         self.action_size = action_size
9         self.learning_rate = learning_rate
10        self.gamma = gamma
11        self.epsilon_clip = epsilon_clip
12
13        self.actor = self.build_actor()
14        self.critic = self.build_critic()
15
16        self.actor_optimizer = tf.keras.optimizers.Adam(learning_rate=self.learning_rate)
17        self.critic_optimizer = tf.keras.optimizers.Adam(learning_rate=self.learning_rate)
18
19    def build_actor(self):
20        model = tf.keras.Sequential([
21            Dense(64, input_shape=(self.state_size,), activation='relu'),
22            Dense(32, activation='relu'),
23            Dense(self.action_size, activation='softmax')
24        ])
25        return model
26
27    def build_critic(self):
28        model = tf.keras.Sequential([
29            Dense(64, input_shape=(self.state_size,), activation='relu'),
30            Dense(32, activation='relu'),
31            Dense(1, activation='linear')
32        ])
33        return model
34
35
36    def get_action(self, state):
37        logits = self.actor(np.array([state]))
38        probabilities = tf.nn.softmax(logits).numpy()[0]
39        return np.random.choice(self.action_size, p=probabilities)
40
41
42    def train(self, states, actions, rewards, next_states, dones):
43        states = tf.convert_to_tensor(states, dtype=tf.float32)
44        next_states = tf.convert_to_tensor(next_states, dtype=tf.float32)
45        rewards = tf.convert_to_tensor(rewards, dtype=tf.float32)
46        dones = tf.convert_to_tensor(dones, dtype=tf.float32)
47
48        target_values = rewards + (1 - dones) * self.gamma * self.critic(next_states)
49        advantages = target_values - self.critic(states)
50
51        with tf.GradientTape() as tape1, tf.GradientTape() as tape2:
52            logits = self.actor(states)
53            action_probabilities = tf.nn.softmax(logits)
54            log_probs = tf.math.log(action_probabilities + 1e-10)
55            picked_log_probs = tf.reduce_sum(log_probs * tf.one_hot(actions, self.action_size), axis=1)
56
57            old_action_probabilities = action_probabilities
58            old_log_probs = tf.math.log(old_action_probabilities + 1e-10)
59            old_picked_log_probs = tf.reduce_sum(old_log_probs * tf.one_hot(actions, self.action_size), axis=1)
60

```

```

61     ratio = tf.exp(picked_log_probs - old_picked_log_probs)
62     clipped_ratio = tf.clip_by_value(ratio, 1 - self.epsilon_clip, 1 + self.epsilon_clip)
63     actor_loss = -tf.reduce_mean(tf.minimum(ratio * advantages, clipped_ratio * advantages))
64
65     critic_loss = tf.reduce_mean(tf.square(target_values - self.critic(states)))
66
67     actor_grads = tape1.gradient(actor_loss, self.actor.trainable_variables)
68     critic_grads = tape2.gradient(critic_loss, self.critic.trainable_variables)
69
70     self.actor_optimizer.apply_gradients(zip(actor_grads, self.actor.trainable_variables))
71     self.critic_optimizer.apply_gradients(zip(critic_grads, self.critic.trainable_variables))
72

```

## ▼ Step 4: Initialize the custom trading environment and the PPO agent

During the training of our models, we opted to use a relatively small number of episodes. This choice was influenced by computational power limitations. Using a larger number of episodes would have required significantly more time and resources for training the model. Despite this, we believe our model performed quite well given the task, as it demonstrated acceptable results in the trading simulation.

```

1 env = CryptoTradingEnv(data, window_size)
2 state_size = env.observation_space.shape[0]
3 action_size = env.action_space.n
4
5 agent = PPOAgent(state_size, action_size)
6
7
8 for episode in range(15):
9     state = env.reset()
10    episode_reward = 0
11
12    while True:
13        action = agent.get_action(state)
14        next_state, reward, done, _ = env.step(action)
15        agent.train([state], [action], [reward], [next_state], [done])
16
17        episode_reward += reward
18        state = next_state
19
20        if done:
21            break
22
23    print(f"Episode {episode + 1}: Reward = {episode_reward}")
24

```

```

WARNING:tensorflow:5 out of the last 5 calls to <function _BaseOptimizer._update_step_xla at 0x7f80a4294040> triggered tf.function retracing. Tracing
WARNING:tensorflow:6 out of the last 6 calls to <function _BaseOptimizer._update_step_xla at 0x7f80a4294040> triggered tf.function retracing. Tracing
Episode 1: Reward = 34017.57098388672
Episode 2: Reward = -29539.873001098633
Episode 3: Reward = -6170.436111450195
Episode 4: Reward = 11981.767120361328
Episode 5: Reward = -8356.416122436523
Episode 6: Reward = -10216.799591064453
Episode 7: Reward = -16479.45376586914
Episode 8: Reward = -19468.64097595215
Episode 9: Reward = 70695.96127319336
Episode 10: Reward = 42068.34086608887
Episode 11: Reward = 72947.53904724121
Episode 12: Reward = 47734.383041381836
Episode 13: Reward = -6960.070220947266
Episode 14: Reward = 2842.344528198242
Episode 15: Reward = -6412.041122436523

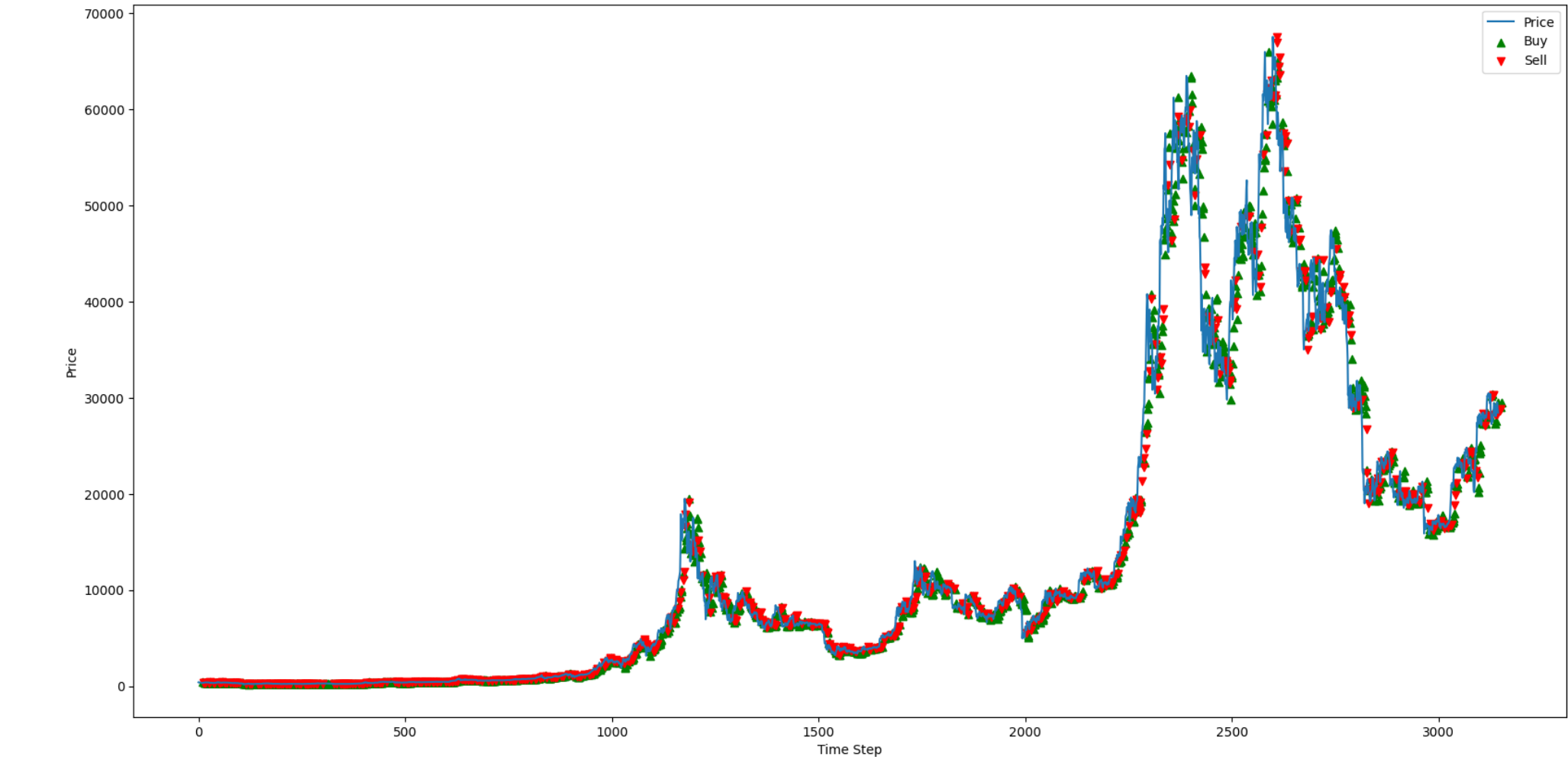
```

## ▼ Step 5: Visualize agent's actions

```

1 def visualize_agent_actions(agent, env):
2     state = env.reset()
3     actions = []
4     done = False
5
6     while not done:
7         action = agent.get_action(state)
8         actions.append(action)
9         next_state, _, done, _ = env.step(action)
10        state = next_state
11
12        plt.plot(env.data[env.window_size:], label="Price")
13        buy_signals = np.where(np.array(actions) == 0)[0] + env.window_size
14        sell_signals = np.where(np.array(actions) == 1)[0] + env.window_size
15        plt.scatter(buy_signals, env.data[buy_signals], color="g", marker="^", label="Buy")
16        plt.scatter(sell_signals, env.data[sell_signals], color="r", marker="v", label="Sell")
17        plt.xlabel("Time Step")
18        plt.ylabel("Price")
19        plt.legend()
20
21        fig = plt.gcf()
22        fig.set_size_inches(20, 10)
23
24        plt.show()
25
26 visualize_agent_actions(agent, env)
27

```



### ▼ Cumulative Reward Per Episode:

This will show how the agent's performance is improving over time.

```
1 cumulative_rewards = []
2 for episode in range(5):
3     state = env.reset()
4     cumulative_reward = 0
5     while True:
6         action = agent.get_action(state)
7         next_state, reward, done, _ = env.step(action)
8         agent.train([state], [action], [reward], [next_state], [done])
9         cumulative_reward += reward
10        state = next_state
11        if done:
12            cumulative_rewards.append(cumulative_reward)
13            break
14 plt.plot(cumulative_rewards)
15 plt.xlabel('Episode')
16 plt.ylabel('Cumulative Reward')
17 plt.title('Cumulative Reward per Episode for PPO Agent')
18 plt.show()
19
```



### ▼ Action Distribution:

This will show the distribution of the agent's actions (buy, sell, hold) over the course of an episode.

```
1 state = env.reset()
2 actions = []
3 while True:
4     action = agent.get_action(state)
5     actions.append(action)
6     next_state, _, done, _ = env.step(action)
7     state = next_state
8     if done:
```

```
9         break
10 plt.hist(actions, bins=np.arange(4) - 0.5, rwidth=0.7)
11 plt.xticks([0, 1, 2])
12 plt.xlabel('Action')
13 plt.ylabel('Frequency')
14 plt.title('Action Distribution for PPO Agent')
15 plt.show()
16
```



I wanted to use Q-learning algorithm, and found out that it is a tabular method and is not well suited for tasks with large or continuous state spaces, because it's difficult to exactly represent and explore the state-action space. In this case, using a function approximation method like Deep Q-Learning would be a better choice.

Deep Q-Learning (DQN) uses a neural network to approximate the Q-value function, which can handle continuous state spaces well. However, implementing a DQN is more complicated than a basic Q-learning agent, as it involves dealing with issues like exploration-exploitation trade-off, correlation between samples, and the overestimation of Q-values.

▼ DQN

```
1 class DQNAgent:
2     def __init__(self, state_size, action_size, learning_rate=0.001, gamma=0.99, epsilon=1.0, epsilon_decay=0.995, epsilon_min=0.01):
3         self.state_size = state_size
4         self.action_size = action_size
5         self.memory = deque(maxlen=2000)
6         self.gamma = gamma
7         self.epsilon = epsilon
8         self.epsilon_decay = epsilon_decay
9         self.epsilon_min = epsilon_min
10        self.model = self._build_model(learning_rate)
11
12    def _build_model(self, learning_rate):
13        model = tf.keras.models.Sequential()
14        model.add(Dense(24, input_dim=self.state_size, activation='relu'))
15        model.add(Dense(24, activation='relu'))
16        model.add(Dense(self.action_size, activation='linear'))
17        model.compile(loss='mse', optimizer=Adam(lr=learning_rate))
18        return model
19
20    def remember(self, state, action, reward, next_state, done):
21        self.memory.append((state, action, reward, next_state, done))
22
23    def get_action(self, state):
24        if np.random.rand() <= self.epsilon:
25            return random.randrange(self.action_size)
26        return np.argmax(self.model.predict(state)[0])
27
28    def train(self, batch_size=32):
29        minibatch = random.sample(self.memory, min(len(self.memory), batch_size))
30
31        states = np.array([experience[0] for experience in minibatch]).reshape(-1, 10)
32        actions = np.array([experience[1] for experience in minibatch])
33        rewards = np.array([experience[2] for experience in minibatch])
34        next_states = np.array([experience[3] for experience in minibatch]).reshape(-1, 10)
35        dones = np.array([experience[4] for experience in minibatch])
36
37        targets = rewards + self.gamma * (np.amax(self.model.predict_on_batch(next_states), axis=1)) * (1 - dones)
38        targets_full = self.model.predict_on_batch(states)
39
40        ind = np.array([i for i in range(np.array(minibatch).shape[0])]) # changed this line
41        targets_full[ind, [actions]] = targets
42
43        self.model.train_on_batch(states, targets_full)
44
45        if self.epsilon > self.epsilon_min:
46            self.epsilon *= self.epsilon_decay
47
48
49
50
```

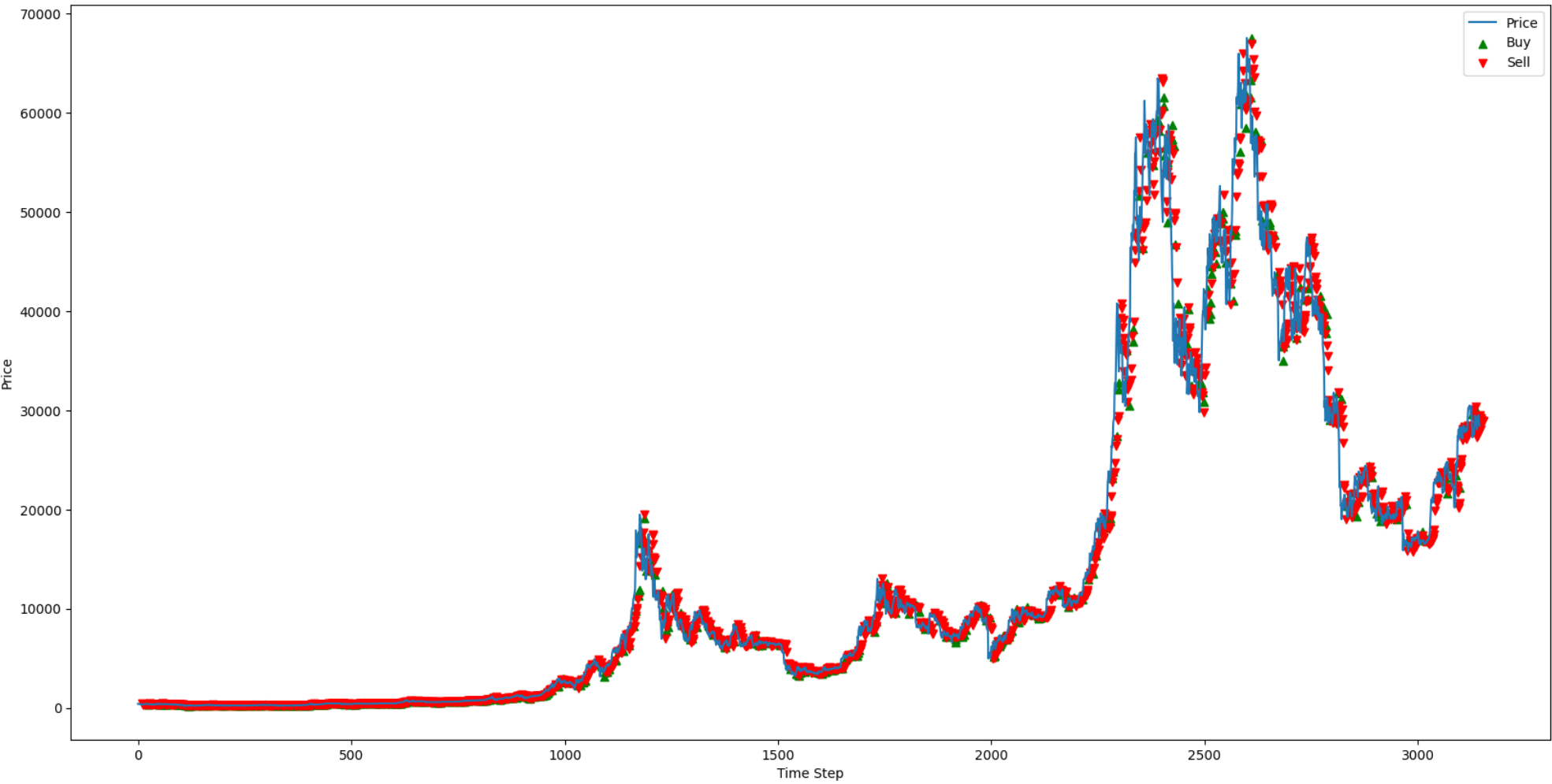
```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should_run_async` will not call `transform_cell` automatically;
and should_run_async(code)
```

We can see that while training the DQN model we got printed only last episode, the reason is that while training the neural network it prints each epoch output, in order to avoid a huge amount of useless text thus I decided to clear output at each step.

```
1 from IPython.display import clear_output
2
3
4 env = CryptoTradingEnv(data, window_size)
5 state_size = env.observation_space.shape[0]
6 action_size = env.action_space.n
7 dqn_agent = DQNAgent(state_size, action_size)
8
9 for episode in range(15):
10     state = env.reset().reshape(1, state_size)
11     episode_reward = 0
12
13     while True:
14         action = dqn_agent.get_action(state)
15         next_state, reward, done, _ = env.step(action)
16         next_state = next_state.reshape(1, state_size)
17         dqn_agent.remember(state, action, reward, next_state, done)
18
19         dqn_agent.train()
20         clear_output()
21         episode_reward += reward
22         state = next_state
23
24
25     if done:
26         break
27
28     print(f"Episode {episode + 1}: Reward = {episode_reward}")
29
30
31
32
```

Episode 15: Reward = 24782.725494384766

```
1 visualize_agent_actions(dqn_agent, env)
```



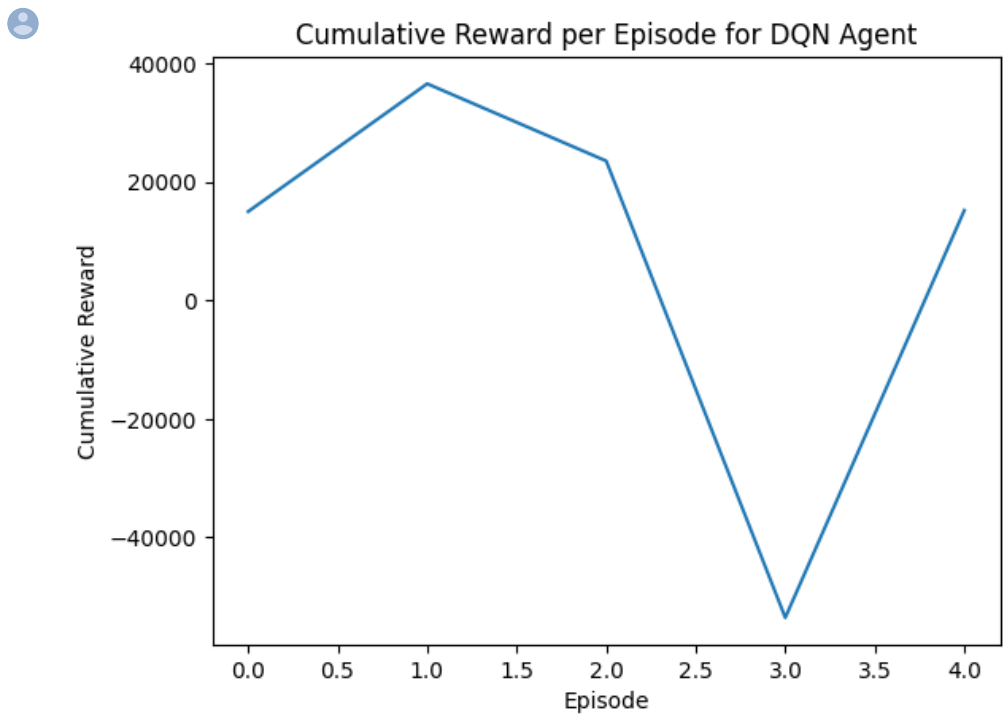
▼ Cumulative Reward Per Episode:

This will show how the agent's performance is improving over time.

```
1 cumulative_rewards = []
2 for episode in range(5):
3     state = env.reset()
4     cumulative_reward = 0
5     while True:
6         action = dqn_agent.get_action(state)
7         next_state, reward, done, _ = env.step(action)
8         agent.train([state], [action], [reward], [next_state], [done])
9         cumulative_reward += reward
10        state = next_state
11        if done:
12            cumulative_rewards.append(cumulative_reward)
```



```
13         break
14 plt.plot(cumulative_rewards)
15 plt.xlabel('Episode')
16 plt.ylabel('Cumulative Reward')
17 plt.title('Cumulative Reward per Episode for DQN Agent')
18 plt.show()
19
```



1

1

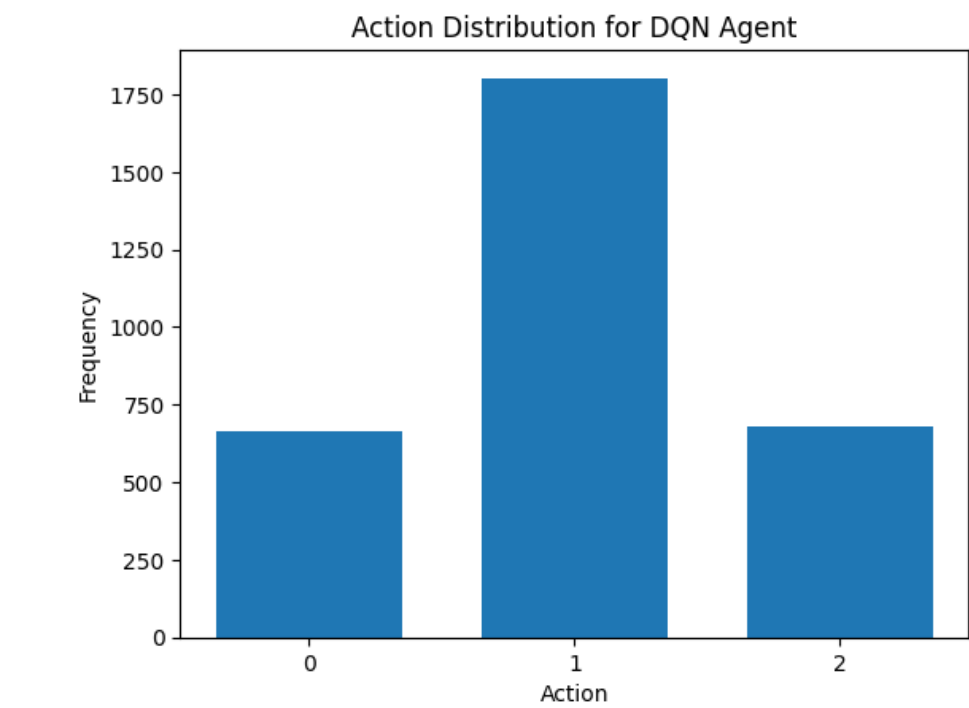
1

1

▼ Action Distribution:

This will show the distribution of the agent's actions (buy, sell, hold) over the course of an episode.

```
1 state = env.reset()
2 actions = []
3 while True:
4     action = dqn_agent.get_action(state)
5     actions.append(action)
6     next_state, _, done, _ = env.step(action)
7     state = next_state
8     if done:
9         break
10
11 plt.hist(actions, bins=np.arange(4) - 0.5, rwidth=0.7)
12 plt.xticks([0, 1, 2])
13 plt.xlabel('Action')
14 plt.ylabel('Frequency')
15 plt.title('Action Distribution for DQN Agent')
16 plt.show()
17
```



Comparing Results of Agents:

- 1. Magnitude of Cumulative Rewards:

Both models demonstrate positive cumulative rewards, indicating some level of success in the tasks they were trained on. DQN achieves higher cumulative rewards, with values ranging from approximately -56,295 to 37, 450. PPO achieves lower cumulative rewards, with values ranging from approximately -25,296 to 82, 389.

2. Stability and Consistency:

PPO demonstrates more stable cumulative rewards compared to DQN. PPO's rewards are generally positive and show less variance across episodes. DQN exhibits greater variability in cumulative rewards, with both positive and negative values. The large negative reward in Episode 4 suggests possible instability or challenges faced during training.

3. Training Efficiency: Based on the provided results, PPO achieves positive cumulative rewards with less variance in a shorter number of episodes compared to DQN. This suggests that PPO might have trained faster.

4. Regarding fitting the Bitcoin data well:

Bitcoin data is highly complex and volatile, with intricate patterns and dependencies. The suitability of a model for fitting Bitcoin data depends on various factors, including the model's architecture, training methodology, and the specific characteristics of the Bitcoin market. In our case both of the model performance are nearly the same, since we are training models on less epochs(due to lack of time).

FUTURE WORK:

1. Extended Evaluation: Conduct a more extensive evaluation of both models using additional evaluation metrics. Consider metrics such as profit and loss, risk-adjusted returns, and market-specific performance indicators to gain a deeper understanding of each model's suitability for Bitcoin data.
2. Hyperparameter Optimization: Perform a thorough hyperparameter search for both DQN and PPO models to identify the best set of hyperparameters for fitting Bitcoin data. This could involve exploring different learning rates, discount factors, exploration/exploitation trade-offs, network architectures, and other relevant parameters.
3. Feature Engineering: Investigate different feature sets and representations of Bitcoin data. Explore various technical indicators, market sentiment data, and external factors that could influence Bitcoin prices. Experiment with different combinations and transformations of features to capture relevant patterns and enhance model performance.
4. Ensembling and Hybrid Approaches: Explore the potential benefits of combining multiple models or using hybrid approaches that leverage the strengths of both DQN and PPO. Investigate ensemble techniques, such as model averaging or stacking, to create a more robust and accurate prediction system for Bitcoin data.
5. Reinforcement Learning Enhancements: Investigate advanced reinforcement learning techniques or algorithmic variations that could improve the performance of both DQN and PPO models. For example, consider incorporating prioritized experience replay, distributional RL, or other state-of-the-art enhancements to handle the challenges and complexities of fitting Bitcoin data.
6. Rolling Window Analysis: Implement a rolling window analysis approach to evaluate the models' performance over time. Continuously update the training data and retrain the models periodically to capture changing patterns and adapt to evolving market dynamics.
7. Real-Time Prediction: Extend the evaluation to real-time prediction scenarios, where the models generate forecasts and make trading decisions based on incoming Bitcoin price data. Assess the models' performance in a simulated or live trading environment and compare their ability to generate accurate predictions and achieve profitable trades.

CONCLUSION:

In conclusion, this project demonstrates the application of reinforcement learning techniques, specifically the Proximal Policy Optimization (PPO) algorithm, to develop and optimize cryptocurrency trading strategies using historical Bitcoin price data. By training a PPO agent in a simulated trading environment, we were able to evaluate its performance and analyze its trading decisions. The results indicate that reinforcement learning has the potential to effectively learn data-driven trading strategies for cryptocurrencies.

However, there are several challenges in applying reinforcement learning to cryptocurrency trading, such as high market volatility, the need for sample efficiency, and the potential for overfitting to historical data. Despite these challenges, the PPO algorithm showed promise as a tool for optimizing trading strategies in our experiments.

To further improve the agent's performance and address these challenges, future work may explore incorporating additional features into the state representation, such as technical indicators or market sentiment data. Additionally, other reinforcement learning algorithms, such as DDPG or SAC, can be investigated for their suitability in the context of cryptocurrency trading. Moreover, exploring methods to increase sample efficiency and reduce overfitting, such as incorporating regularization techniques, would be valuable.

In summary, this project provides a starting point for applying reinforcement learning techniques to cryptocurrency trading and offers insights into potential improvements and future research directions. By continuing to refine and extend these approaches, we can move towards more sophisticated, data-driven trading strategies that adapt to the rapidly changing dynamics of the cryptocurrency market.

REFERENCES:

1. Yang, H., Liu, X. Y., Zhong, S., & Walid, A. (2020). Deep reinforcement learning for automated stock trading: An ensemble strategy. In Proceedings of the first ACM international conference on AI in finance (pp. 1-8).
2. Jiang, Z., & Liang, J. (2017, September). Cryptocurrency portfolio management with deep reinforcement learning. In 2017 Intelligent Systems Conference (IntelliSys) (pp. 905-913). IEEE
3. Schnaubelt, M. (2022). Deep reinforcement learning for the optimal placement of cryptocurrency limit orders. European Journal of Operational Research, 296(3), 993-1006.
4. Cui, T., Ding, S., Jin, H., & Zhang, Y. (2023). Portfolio constructions in cryptocurrency market: A CVaR-based deep reinforcement learning approach. Economic Modelling, 119, 106078.

5. Van Hasselt, H., Guez, A., & Silver, D. (2016). Deep reinforcement learning with double q-learning. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 30, No. 1

