

REINFORCEMENT LEARNING-BASED CRYPTOCURRENCY PORTFOLIO MANAGEMENT USING SOFT ACTOR–CRITIC AND DEEP DETERMINISTIC POLICY GRADIENT ALGORITHMS

KAMAL PAYKAN

ABSTRACT. This study proposes a deep reinforcement learning (DRL) framework for dynamic cryptocurrency portfolio management across four major assets: Bitcoin, Ethereum, Litecoin, and Dogecoin. The trading environment is formulated as a discrete-time stochastic system in which an intelligent agent optimizes portfolio allocation through interaction with market data. Operating under partial observability, the agent utilizes daily open, high, low, and close prices together with trading volume to determine optimal asset weights for the next trading period. The framework integrates actor–critic algorithms—specifically Deep Deterministic Policy Gradient and Soft Actor–Critic—enhanced with Long Short-Term Memory networks to capture temporal dependencies within the policy and value functions. For benchmarking, the classical Markowitz mean–variance model is employed. Empirical evaluations demonstrate that the proposed DRL-based agents learn adaptive strategies that consistently outperform the Markowitz benchmark in both absolute return and risk-adjusted performance. These findings highlight the potential of reinforcement learning as a robust methodology for managing cryptocurrency portfolios in volatile and nonstationary markets.

1. INTRODUCTION

Cryptocurrency markets have emerged as a substantial component of the global financial ecosystem, attracting both retail and institutional investors. Major assets such as Bitcoin and Ethereum are traded globally, with total market capitalization frequently exceeding trillions of dollars [7]. However, cryptocurrency portfolio management remains challenging due to extreme volatility, unstable risk–return relationships, and rapidly evolving market structures. Digital assets exhibit high price fluctuations, regime shifts, and liquidity constraints, complicating both short- and long-term investment strategies [2, 4].

Classical financial models often assume normally distributed returns and stationary market conditions—assumptions that are frequently violated in cryptocurrency markets [23]. High transaction costs, thin liquidity, and structural changes further reduce the reliability of traditional models such as Markowitz’s

Key words and phrases. Portfolio management, Markov decision process, Deep reinforcement learning, Deep Deterministic Policy Gradient, Soft Actor–Critic, Long Short-Term Memory networks.

mean–variance framework. Consequently, adaptive and data-driven methods capable of learning directly from dynamic environments are required. Reinforcement learning (RL) provides such a framework by enabling agents to optimize trading strategies through interaction with the market, making RL well-suited for the complexities of digital asset portfolio management.

1.1. Motivation for Using Reinforcement Learning. Traditional portfolio optimization techniques, including the Markowitz model and CAPM, rely heavily on stable statistical assumptions and accurate estimation of expected returns and covariances. In cryptocurrency markets, where returns are non-stationary and exhibit nonlinear dependencies, such assumptions often fail [3, 28]. As a result, classical financial models struggle under sudden market regime shifts and volatile price dynamics.

Reinforcement learning, by contrast, learns optimal decision policies from experience rather than predefined models [20]. Deep reinforcement learning (DRL) integrates RL with deep neural networks, enabling effective learning in high-dimensional and continuous spaces. Algorithms such as Deep Deterministic Policy Gradient (DDPG) and Soft Actor–Critic (SAC) can learn complex trading policies directly from noisy, sequential financial data [14, 22, 24, 25]. Empirical evidence demonstrates that RL-based strategies can outperform static optimization methods under uncertainty [16, 31, 36]. These properties make RL an effective framework for cryptocurrency portfolio management.

1.2. Purpose and Scope of the Paper. This study develops and evaluates a reinforcement learning framework for cryptocurrency portfolio optimization using two DRL algorithms: DDPG and SAC. Both are actor–critic methods designed for continuous action spaces. DDPG employs deterministic policies for efficient continuous control but can exhibit instability in noisy environments [17]. SAC uses stochastic, entropy-regularized policies, promoting exploration and stability in volatile markets [5, 14].

Both algorithms are implemented within a unified trading environment to ensure a fair comparison. The classical Markowitz mean–variance model serves as a deterministic benchmark. Performance is assessed using financial metrics including the Sharpe ratio, Sortino ratio, maximum drawdown, Value-at-Risk, Conditional Value-at-Risk, and cumulative portfolio value. Results demonstrate that the DRL agents—particularly SAC—achieve superior risk-adjusted returns and portfolio stability compared with DDPG and the Markowitz baseline. This work contributes to the literature by demonstrating the effectiveness of adaptive, data-driven strategies for navigating cryptocurrency market volatility.

2. BACKGROUND

2.1. Cryptocurrency Portfolio Management. Cryptocurrency portfolio management involves selecting and rebalancing digital assets to optimize the trade-off between risk and return in a highly uncertain environment. Cryptocurrencies typically exhibit extreme price volatility, volatility clustering, and sensitivity to sentiment and regulatory developments [18]. These dynamics motivate the use

of risk-adjusted performance metrics such as the Sharpe ratio, Sortino ratio, and maximum drawdown [3, 11].

Diversification is less effective in cryptocurrency markets compared with traditional assets, as correlations tend to increase during periods of market stress [7]. Static portfolio strategies, including mean–variance optimization, risk parity, and momentum-based trading, often fail under non-stationarity and structural breaks [32]. These limitations motivate reinforcement learning, which learns from sequential market interactions and adapts to evolving market regimes.

2.2. Reinforcement Learning in Finance. Reinforcement learning is a machine learning framework for sequential decision-making, where an agent interacts with an environment and learns a policy that maximizes cumulative reward [34]. This structure aligns naturally with financial decision-making, where strategies must adapt to dynamic market conditions. Unlike static optimization frameworks, RL models temporal dependencies and delayed rewards, allowing optimization over the entire investment horizon [31].

Early RL research in finance relied on tabular methods such as Q-learning and SARSA [9], which struggle with high-dimensional data. Deep reinforcement learning alleviates this limitation by combining RL with neural networks [20]. Algorithms such as DQN, DDPG, and SAC have since been applied to portfolio optimization, algorithmic trading, and risk management, demonstrating strong adaptability under volatility and nonlinear dependencies.

2.3. Deep Deterministic Policy Gradient. The Deep Deterministic Policy Gradient (DDPG) algorithm [22] is a model-free, off-policy actor–critic method for continuous action spaces. DDPG employs an actor network that maps states to continuous portfolio weights and a critic network that evaluates the action-value function $Q(s, a)$ [33]. Training stability is improved through mechanisms such as experience replay and target networks [30]. DDPG has been shown to outperform heuristic baselines in dynamic environments [21]; however, its sensitivity to hyperparameters and limited exploration may lead to suboptimal convergence.

2.4. Soft Actor–Critic. The Soft Actor–Critic (SAC) algorithm [14] introduces entropy regularization into the objective function, encouraging exploration and improving stability. SAC learns a stochastic policy that maximizes both expected return and entropy, balancing exploration and exploitation. The architecture includes two critic networks, a stochastic actor, and a value network, which collectively mitigate overestimation bias and improve convergence [13]. SAC’s sample efficiency and robustness make it well-suited for highly volatile cryptocurrency markets [6, 21, 36].

3. PROBLEM FORMULATION AND METHODOLOGY

3.1. Portfolio Optimization Problem. Cryptocurrency portfolio management can be formulated as a sequential decision-making problem, where the investor determines daily capital allocations across assets. The portfolio consists of $n = 4$

cryptocurrencies (BTC, ETH, LTC, DOGE), with price vector

$$\mathbf{p}_t = [p_t^{(1)}, p_t^{(2)}, \dots, p_t^{(n)}]^\top.$$

The log-return of asset i is:

$$r_{t+1}^{(i)} = \ln \left(\frac{p_{t+1}^{(i)}}{p_t^{(i)}} \right).$$

The portfolio weight vector is denoted as:

$$\mathbf{w}_t = [w_t^{(1)}, \dots, w_t^{(n)}]^\top,$$

subject to:

$$\sum_{i=1}^n w_t^{(i)} = 1, \quad w_t^{(i)} \geq 0.$$

The portfolio return is:

$$R_{t+1} = \mathbf{w}_t^\top \mathbf{r}_{t+1}.$$

Accounting for transaction costs with rate c :

$$R_{t+1}^{\text{net}} = R_{t+1} - c \sum_{i=1}^n |w_t^{(i)} - w_{t-1}^{(i)}|.$$

Portfolio value evolves as:

$$v_{t+1} = v_t \exp(R_{t+1}^{\text{net}}).$$

The optimization objective is:

$$\max_{\{w_t\}_{t=1}^T} \mathbb{E} \left[\sum_{t=1}^T \gamma^{t-1} U(R_t^{\text{net}}) \right],$$

where U is a mean–variance utility function:

$$U(R_t^{\text{net}}) = \mathbb{E}[R_t^{\text{net}}] - \frac{\lambda}{2} \text{Var}(R_t^{\text{net}}).$$

However, due to the non-stationary and heavy-tailed nature of crypto returns [18, 28], static optimization methods such as Markowitz’s mean–variance model often perform poorly. Therefore, the problem is reformulated as a stochastic control task with policy $\pi(a_t | s_t)$ learned via RL methods.

3.2. Reinforcement Learning Setup.

3.2.1. *State Space.* At time t , the agent observes asset prices:

$$\mathbf{p}_t = [p_t^{(1)}, \dots, p_t^{(M)}],$$

but financial markets are partially observable. Hence, a rolling window of size $W = 50$ is used:

$$s_t = (\mathbf{p}_{t-W+1:t}, \mathbf{w}_t),$$

where $\mathbf{p}_{t-W+1:t} \in \mathbb{R}_+^{M \times W \times 5}$ represents OHLCV data.

Log-returns are computed and the remaining features are standardized. The continuous state space is:

$$\mathcal{S} \subset \mathbb{R}_+^{M \times W \times 5} \times \mathbb{R}^M.$$

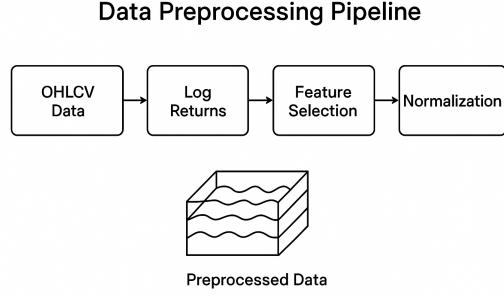


FIGURE 1. Data preprocessing workflow applied to OHLCV cryptocurrency data.

3.2.2. Action Space. The agent outputs next-day portfolio weights:

$$a_t = \mathbf{w}_{t+1} = (w_{1,t+1}, \dots, w_{M,t+1}),$$

constrained by:

$$\sum_{i=1}^M w_{i,t+1} = 1, \quad w_{i,t+1} \geq 0.$$

The policy network produces raw outputs \tilde{a}_t that are normalized via:

$$w_{i,t+1} = \frac{\exp(\tilde{a}_{i,t})}{\sum_{j=1}^M \exp(\tilde{a}_{j,t})}.$$

3.2.3. Reward Function. The simplest reward is log-return:

$$r_{t+1} = \ln \left(\frac{v_{t+1}}{v_t} \right).$$

A risk-adjusted reward based on the differential Sharpe ratio (DSR) is also used:

$$r_{t+1} = \frac{B_t(\rho_{t+1} - A_t) - \frac{1}{2}A_t(\rho_{t+1}^2 - B_t)}{(B_t - A_t^2)^{3/2}},$$

with updates:

$$A_t = A_{t-1} + \eta(\rho_t - A_{t-1}), \quad B_t = B_{t-1} + \eta(\rho_t^2 - B_{t-1}).$$

3.3. Data and Environment.

3.3.1. *Data Sources.* Experiments use four major cryptocurrencies: BTC, ETH, LTC, and DOGE. Daily OHLCV data are obtained using `yfinance`. The period spans:

Training: 2016–2022, Testing: 2023–2024.

Data preprocessing and environment interaction occur through a custom `gym`-based simulator.

3.3.2. *Time Period of Analysis.* The dataset spans January 1, 2016 to December 31, 2024, covering various market regimes [7, 18]. The split follows best practices for avoiding data leakage.

Time Period of Analysis

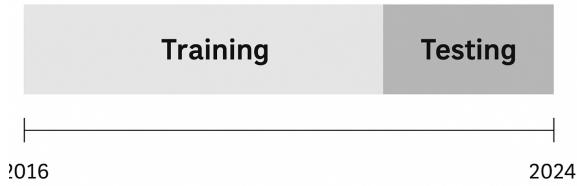


FIGURE 2. Temporal split of training and testing periods.

3.4. **Simulation Environment.** The environment simulates daily portfolio rebalancing, following the RL interaction loop. Portfolio value evolves as:

$$v_{t+1} = v_t(1 + R_{t+1}^{\text{net}}),$$

and actions are given by:

$$\mathbf{w}_{t+1} = \pi_\theta(s_t).$$

Training uses off-policy learning with replay buffers. Evaluation on the test set uses fixed parameters. Performance is assessed using Sharpe ratio, Sortino ratio, maximum drawdown, annualized volatility, and cumulative return.

This environment provides a realistic testbed for DRL-based cryptocurrency portfolio optimization.

4. METHODOLOGY: SAC vs. DDPG AND BENCHMARK INTEGRATION

4.1. **DDPG Algorithm Implementation.** The Deep Deterministic Policy Gradient (DDPG) algorithm introduced by [22] extends the deterministic policy gradient framework to high-dimensional and continuous action spaces, making it suitable for portfolio allocation problems. DDPG follows an actor–critic architecture in which two neural networks—the actor and the critic—approximate the policy and the action–value function, respectively. This design allows efficient policy optimization in stochastic and highly volatile environments such as cryptocurrency markets.

The trading environment is modeled as a Markov decision process (MDP) with state space \mathcal{S} , action space \mathcal{A} , transition dynamics $P(s_{t+1} | s_t, a_t)$, and reward function $r_t = R(s_t, a_t)$. The objective is to learn a deterministic policy

$$\mu_\theta : \mathcal{S} \rightarrow \mathcal{A},$$

parameterized by θ , that maximizes the discounted return

$$J(\theta) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t r_t \right],$$

where $\gamma \in (0, 1]$ is the discount factor.

The critic network $Q_\phi(s_t, a_t)$ estimates the action–value function:

$$Q_\phi(s_t, a_t) = \mathbb{E} [r_t + \gamma Q_{\phi'}(s_{t+1}, \mu_{\theta'}(s_{t+1}))],$$

where (ϕ', θ') denote slowly updated target-network parameters. The critic is trained by minimizing the TD error:

$$L(\phi) = \mathbb{E} [(Q_\phi(s_t, a_t) - y_t)^2], \quad y_t = r_t + \gamma Q_{\phi'}(s_{t+1}, \mu_{\theta'}(s_{t+1})).$$

The actor is updated using the deterministic policy gradient:

$$\nabla_\theta J(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[\nabla_a Q_\phi(s_t, a) \Big|_{a=\mu_\theta(s_t)} \nabla_\theta \mu_\theta(s_t) \right].$$

Algorithm 1 Deep Deterministic Policy Gradient

```

1: Initialize actor parameters  $\theta$  and critic parameters  $\phi$ 
2: Set target network parameters  $\theta' \leftarrow \theta$ ,  $\phi' \leftarrow \phi$ 
3: Initialize replay buffer  $\mathcal{D}$ 
4: for each episode do
5:   Initialize exploration noise process  $\mathcal{N}$ 
6:   Observe initial state  $s_1$ 
7:   while episode not terminated do
8:     Select action  $a_t = \mu_\theta(s_t) + \mathcal{N}_t$ 
9:     Execute  $a_t$ , observe  $r_t$  and  $s_{t+1}$ 
10:    Store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
11:    Sample minibatch of  $B$  transitions from  $\mathcal{D}$ 
12:    Compute target

$$y_t = r_t + \gamma Q_{\phi'}(s_{t+1}, \mu_{\theta'}(s_{t+1}))$$

13:    Update critic by minimizing

$$L(\phi) = \frac{1}{B} \sum_{i=1}^B (Q_\phi(s_t^i, a_t^i) - y_t^i)^2$$

14:    Update actor using

$$\nabla_\theta J(\theta) = \mathbb{E} \left[ \nabla_a Q_\phi(s, a) \Big|_{a=\mu_\theta(s)} \nabla_\theta \mu_\theta(s) \right]$$

15:    Soft-update targets:

$$\phi' \leftarrow \tau\phi + (1 - \tau)\phi', \quad \theta' \leftarrow \tau\theta + (1 - \tau)\theta'$$

16:     $s_t \leftarrow s_{t+1}$ 
17: Return: Trained actor  $\mu_\theta$  and critic  $Q_\phi$ 

```

Exploration in DDPG is achieved by adding noise to the deterministic action:

$$a_t = \mu_\theta(s_t) + \mathcal{N}_t,$$

typically modeled using an Ornstein–Uhlenbeck (OU) process [35]. This allows smooth action perturbations, appropriate for financial applications involving gradual portfolio rebalancing.

To enhance stability, DDPG incorporates target networks via soft updates:

$$\phi' \leftarrow \tau\phi + (1 - \tau)\phi', \quad \theta' \leftarrow \tau\theta + (1 - \tau)\theta',$$

and employs a replay buffer \mathcal{D} for off-policy training [30]. In this study, DDPG is applied to continuous cryptocurrency portfolio allocation where the action a_t represents the portfolio weight vector.

4.2. Soft Actor–Critic Algorithm Implementation. The Soft Actor–Critic (SAC) algorithm [14] is a maximum-entropy reinforcement learning method that improves stability and exploration by combining stochastic policies with entropy regularization. This makes SAC highly suitable for volatile and non-stationary environments such as cryptocurrency markets.

The SAC objective augments the cumulative return with a policy entropy term:

$$J(\pi) = \sum_{t=0}^T \mathbb{E} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t))], \quad (4.1)$$

where α controls the exploration-exploitation balance.

SAC employs two soft Q-networks and a value network:

$$Q_{\phi_1}(s, a), \quad Q_{\phi_2}(s, a), \quad V_{\psi}(s),$$

and a stochastic Gaussian actor

$$\pi_{\theta}(a|s) = \mathcal{N}(\mu_{\theta}(s), \sigma_{\theta}(s)).$$

Actions are sampled using the reparameterization trick:

$$a_t = \tanh (\mu_{\theta}(s_t) + \sigma_{\theta}(s_t) \odot \epsilon_t), \quad \epsilon_t \sim \mathcal{N}(0, I).$$

The soft Q-values are trained by minimizing

$$J_Q(\phi_i) = \mathbb{E} \left[(Q_{\phi_i}(s_t, a_t) - y_t)^2 \right]$$

with targets

$$y_t = r_t + \gamma \left(\min_i Q_{\phi'_i}(s_{t+1}, a_{t+1}) - \alpha \log \pi_{\theta}(a_{t+1} | s_{t+1}) \right).$$

The actor minimizes

$$J_{\pi}(\theta) = \mathbb{E} [\alpha \log \pi_{\theta}(a_t|s_t) - Q_{\phi_1}(s_t, a_t)].$$

Algorithm 2 Soft Actor-Critic Algorithm

- 1: Initialize policy parameters ϕ , critics θ_1, θ_2 , value network ψ
 - 2: Initialize target value network $\psi' \leftarrow \psi$
 - 3: Initialize replay buffer \mathcal{D}
 - 4: **for** each episode **do**
 - 5: **for** each timestep t **do**
 - 6: Sample $a_t \sim \pi_{\phi}(\cdot|s_t)$
 - 7: Execute a_t , observe r_t and s_{t+1}
 - 8: Store transition in \mathcal{D}
 - 9: Sample minibatch from \mathcal{D}
 - 10: Compute targets
 - $y_i^v = \min_j Q_{\theta_j}(s_i, a_i) - \alpha \log \pi_{\phi}(a_i|s_i)$
 - $y_i^q = r_i + \gamma V_{\psi'}(s'_i)$
 - 11: Update V_{ψ} , Q_{θ_1} , Q_{θ_2} , and π_{ϕ}
 - 12: Soft-update:
$$\psi' \leftarrow \tau \psi + (1 - \tau) \psi'$$
 - 13: **return** Policy π_{ϕ} and critics $Q_{\theta_1}, Q_{\theta_2}$
-

SAC provides strong exploration via entropy regularization and stabilizes training using double critics, reparameterization, and automatic temperature adjustment [14]. These properties yield robust performance in cryptocurrency portfolio optimization.

4.3. Benchmark Comparison Framework. To benchmark the reinforcement-learning agents, we implement the classical Markowitz mean–variance model [27]. The objective is to find portfolio weights

$$\mathbf{w} = [w_1, \dots, w_M]^T$$

that minimize variance for a target expected return. Let

$$\boldsymbol{\mu} = [\mu_1, \dots, \mu_M]^T, \quad \boldsymbol{\mathcal{E}} \text{ (covariance matrix).}$$

The optimization problem is

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \boldsymbol{\mathcal{E}} \mathbf{w} - \lambda \boldsymbol{\mu}^T \mathbf{w}$$

subject to

$$\sum_{i=1}^M w_i = 1, \quad w_i \geq 0.$$

The closed-form optimal solution is

$$\mathbf{w}^* = \frac{1}{Z} \boldsymbol{\mathcal{E}}^{-1} (\boldsymbol{\mu} - r_f \mathbf{1}),$$

with corresponding expected return and variance

$$\sigma_p^2 = \mathbf{w}^{*T} \boldsymbol{\mathcal{E}} \mathbf{w}^*, \quad \mu_p = \boldsymbol{\mu}^T \mathbf{w}^*.$$

The Markowitz model is implemented via an `MPT-Agent` subclass derived from a shared `BaseAgent` interface, alongside the `DDPG-Agent` and `SAC-Agent`. Optimization is carried out using `SciPy`. Despite its simplifying assumptions, the model serves as a strong deterministic baseline [8, 26, 29].

4.4. Network Architectures.

4.4.1. Input Network (Feature Extraction Module). Figure 3 illustrates the feature extraction module. Each asset has a multivariate time series of standardized log-differences $\boldsymbol{\rho}_{i,t-W+2:t}$ constructed from OHLCV data [12, 17]. A fully connected layer followed by three LSTM layers [1] produces a latent representation v_t .

A cash asset with constant price produces a zero-vector embedding, stabilizing training.

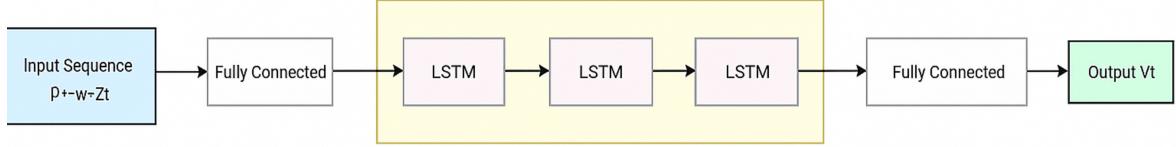


FIGURE 3. Structure of the feature extraction network. Each $p_{t-W+2:t}$ sequence is processed independently by an identical network producing latent features v_t .

4.4.2. Actor Network. The actor receives the concatenated latent features and the current portfolio weights.

For DDPG, the actor outputs deterministic portfolio weights using a softmax layer:

$$\sum_{i=1}^M w_{t+1}^{(i)} = 1.$$

For SAC, the actor outputs mean and standard deviation of a Gaussian:

$$a_t \sim \mathcal{N}(\mu_\theta(s_t), \sigma_\theta(s_t)),$$

then normalizes:

$$\hat{a}_t = \frac{a_t}{\sum_i a_t^{(i)}}.$$

4.4.3. Critic Network. The critic receives (s_t, a_t) and outputs $Q(s_t, a_t)$. DDPG uses one critic; SAC uses two to reduce overestimation [14]. All critic networks use dense layers with Leaky ReLU activations.

4.4.4. Hyperparameter Configuration. Both algorithms share comparable architectures and training settings for fair comparison. SAC adapts its temperature parameter α , while DDPG uses OU exploration noise.

Tables 1 and 2 summarize the hyperparameters used throughout training and evaluation.

4.4.5. Forecasting Network. In the reinforcement learning framework employed in this study, the primary features influencing the reward signal are the closing prices of the assets. Accordingly, the associated forecasting task is designed to predict the vector of logarithmic returns of these closing prices at the next time step, denoted by $\rho_{t+1}^{\text{close}}$, based on the historical observation window $\rho_{i,t-W+2:t}$.

To facilitate this task, the deterministic actor network originally developed for portfolio management was adapted for forecasting purposes. In this adaptation, the components responsible for processing current portfolio weights were removed, along with the final *softmax* transformation layer. This modification was essential, as the prediction of logarithmic returns does not require the normalization constraint imposed on portfolio weights, which must sum to one.

TABLE 1. Hyperparameters for the Soft Actor–Critic Algorithm

Hyperparameter	Value
Actor Network	Adam, LR 3×10^{-4} , 64 hidden units
Critic Network	Adam, LR 3×10^{-4} , 64 hidden units
Temperature α	Adam, LR 10^{-3} , initial value 1
Discount Factor γ	0.99
Soft Update Rate τ	10^{-3}
Replay Memory Size	10^5
Batch Size	64
Episodes	1000
Steps per Episode	10000

TABLE 2. Hyperparameters for the Deep Deterministic Policy Gradient Algorithm

Hyperparameter	Value
Actor Network	Adam, LR 3×10^{-4} , 64 hidden units
Critic Network	Adam, LR 3×10^{-4} , 64 hidden units
Exploration Noise	OU process, $\mu = 0$, $\sigma = 0.3$, $\theta = 0.20$
Discount Factor γ	0.99
Soft Update Rate τ	10^{-3}
Replay Memory Size	10^5
Batch Size	64
Episodes	1000
Steps per Episode	10000

The forecasting model was trained using the same historical data as the reinforcement learning environment. Training followed a supervised learning paradigm: after preprocessing, the input at each time step t was defined as

$$X_t = \boldsymbol{\rho}_{t-W+2:t},$$

and the corresponding target output as

$$y_t = \boldsymbol{\rho}_{t+1}^{\text{close}}.$$

Let \hat{y}_t denote the model’s prediction at time t . Since the task involves multivariate time-series forecasting, both y_t and \hat{y}_t are multidimensional. To ensure balanced accuracy across all assets, the loss function was formulated as the mean Root

Mean Square Error (RMSE):

$$\mathcal{L} = \frac{1}{M} \sum_{i=1}^M \sqrt{\frac{1}{T} \sum_{t=1}^T (y_t^{(i)} - \hat{y}_t^{(i)})^2},$$

where M denotes the number of assets and T the total number of time steps.

Optimization was performed using stochastic gradient descent with the Adam optimizer [19]. Beyond its predictive role, the forecasting model also serves as a diagnostic tool within the reinforcement learning framework: if it can reliably forecast future returns using the same input features as the RL agent, this indicates that the data contain meaningful predictive signals and that the shared feature extraction module is effective. Thus, the forecasting network functions not as a decision-maker but as a supporting model validating the representational capacity of the architecture.

The hyperparameter configurations used for the forecasting model are summarized in Table 3.

TABLE 3. Hyperparameters for the Forecasting Network

Hyperparameter	Value
Forecasting Network	Optimizer: Adam Learning Rate: 3×10^{-4} Hidden Size: 64
Batch Size	128
Number of Episodes	1000

4.5. Comparison Criteria. To evaluate the performance and practicality of the Soft Actor–Critic (SAC) and Deep Deterministic Policy Gradient (DDPG) algorithms in cryptocurrency portfolio management, a comprehensive set of comparison criteria is considered. These include convergence speed, final portfolio value, risk-adjusted performance metrics, and robustness under market volatility.

Convergence speed measures how rapidly an algorithm reaches a stable and near-optimal policy. In financial reinforcement learning, faster convergence implies efficient learning from limited data. The entropy-regularized objective of SAC generally yields smoother and more stable convergence compared to deterministic methods like DDPG [14].

Final portfolio value quantifies cumulative profitability. However, profitability alone does not capture risk exposure; therefore, several risk-adjusted metrics are employed.

The *Sharpe ratio* is defined as

$$\text{Sharpe} = \frac{\mathbb{E}[R_p - R_f]}{\sigma_p},$$

where R_p denotes the portfolio return, R_f is the risk-free rate, and σ_p is the standard deviation of returns.

To focus specifically on downside fluctuations, the *Sortino ratio* is also used:

$$\text{Sortino} = \frac{\mathbb{E}[R_p - R_f]}{\sigma_{\text{down}}},$$

where σ_{down} measures only the negative deviations below a minimum acceptable return.

Volatility is assessed using the *standard deviation of returns*, while tail-risk is analyzed through the Value at Risk (VaR) and Conditional Value at Risk (CVaR):

$$\text{VaR}_\alpha = \inf\{x \in \mathbb{R} : \mathbb{P}(L > x) \leq 1 - \alpha\}, \quad \text{CVaR}_\alpha = \mathbb{E}[L | L > \text{VaR}_\alpha],$$

where L denotes portfolio loss.

Maximum Drawdown (MDD) measures the worst peak-to-trough loss. The *Calmar ratio* evaluates annualized return relative to MDD:

$$\text{Calmar} = \frac{R_{\text{annual}}}{\text{MDD}}.$$

Robustness reflects an algorithm's ability to maintain performance under shifting market regimes. SAC, with its stochastic exploration mechanism, generally exhibits superior robustness compared to DDPG, whose deterministic policy may overfit or become unstable under noisy market conditions [5, 14, 36].

5. EMPIRICAL EVALUATION AND DISCUSSION

This section presents the experimental design, evaluation methodology, and empirical analysis of the Soft Actor–Critic (SAC) and Deep Deterministic Policy Gradient (DDPG) algorithms for cryptocurrency portfolio management. The goal is to assess the learning dynamics, robustness, and financial viability of the proposed reinforcement learning methods under realistic market conditions.

5.1. Experimental Setup. All models were trained and evaluated using historical cryptocurrency data. Hyperparameters for SAC and DDPG (Tables 1 and 2) were kept consistent across experiments to ensure a fair comparison. Cross-validation procedures were applied where appropriate to reduce overfitting and improve generalization.

A forecasting network was also trained using the same historical data as the reinforcement learning agents. The input at time t consisted of historical return sequences

$$X_t = \rho_{t-W+2:t},$$

and the target output was the next-day log return

$$y_t = \rho_{t+1}^{\text{close}}.$$

Although the forecasting module does not participate directly in portfolio allocation, it serves as an auxiliary tool for validating feature extraction and temporal modeling. Figure 4 displays its training and testing phases, with prediction accuracy reaching approximately 70%. This moderate accuracy reflects the strong volatility and nonstationarity inherent in cryptocurrency markets.

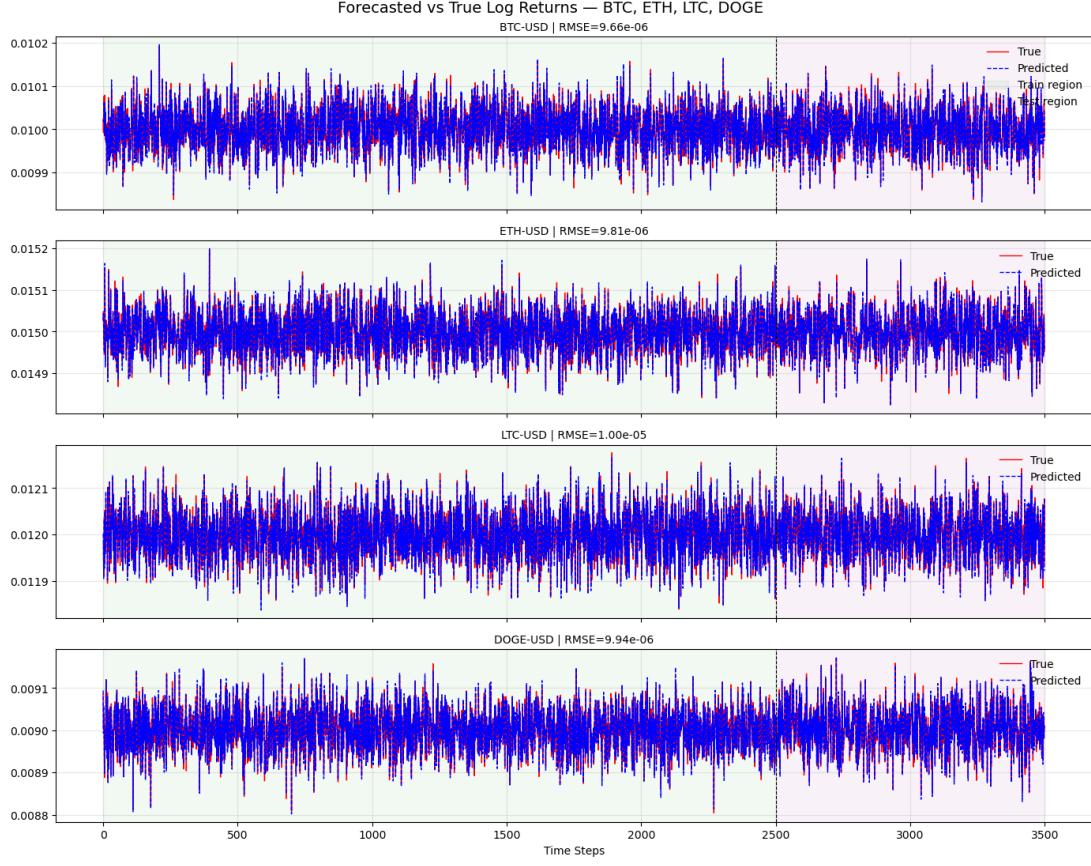


FIGURE 4. Training and testing performance of the forecasting network. The green and purple regions denote the training and test periods, respectively. Red lines show true values; blue lines show predicted values.

5.2. Asset-Level Risk and Return Analysis. To characterize the underlying assets considered by the RL agents, Figure 5 shows normalized price trajectories and 20-day rolling Sharpe ratios for BTC-USD, ETH-USD, LTC-USD, and DOGE-USD. Normalizing prices enables a direct comparison of growth across assets with different nominal price scales, while the rolling Sharpe ratio captures fluctuations in risk-adjusted performance.

Bitcoin and Ethereum exhibit comparatively stable return dynamics and smoother Sharpe ratio trajectories, indicating stronger efficiency and lower volatility. In contrast, Litecoin and Dogecoin show higher fluctuations and less consistent performance, reflecting greater exposure to market instability. Within a reinforcement learning framework, prioritizing assets with stable Sharpe behavior can improve portfolio robustness, whereas highly volatile assets should be managed cautiously to mitigate drawdowns.

5.3. Evaluation Metrics and Comparative Results. Portfolio performance was assessed using standard financial and risk-based metrics, including Sharpe

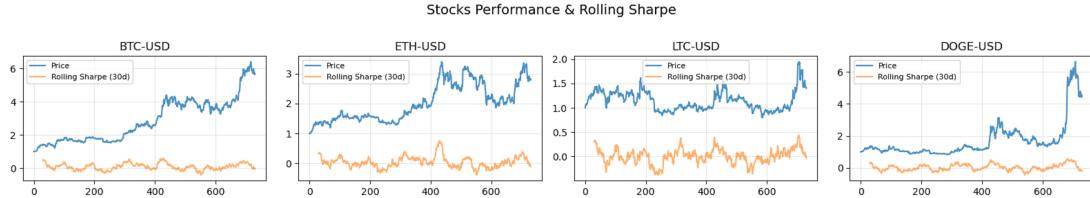


FIGURE 5. Normalized prices (blue) and 20-day rolling Sharpe ratios (red) for BTC-USD, ETH-USD, LTC-USD, and DOGE-USD.

and Sortino ratios, Maximum Drawdown, Value-at-Risk (VaR), Conditional Value-at-Risk (CVaR), and cumulative portfolio value. This comprehensive evaluation captures both profitability and downside risk.

Figure 6 compares the SAC and DDPG agents with a Markowitz Mean–Variance Portfolio Theory (MPT) baseline. The MPT model was implemented as a deterministic benchmark using SciPy to compute optimal asset weights at each step under classical mean–variance assumptions.

Both reinforcement learning agents outperform the MPT benchmark, with SAC consistently delivering superior risk-adjusted performance. SAC’s entropy-regularized stochastic policy enables continual exploration and adaptation to changing volatility regimes, resulting in smoother convergence and more stable returns. DDPG performs well during stable periods but is more sensitive to hyperparameters and market noise, occasionally leading to unstable behavior.

These results highlight the advantage of RL-based strategies in highly volatile and nonstationary cryptocurrency markets.

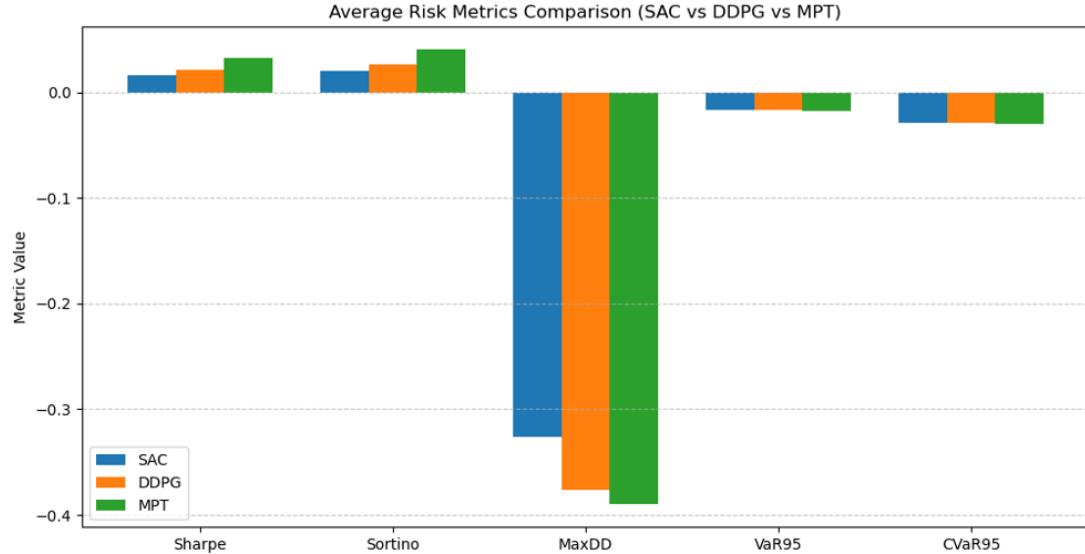


FIGURE 6. Average portfolio performance metrics for SAC, DDPG, and the Markowitz MPT benchmark.

Figure 7 illustrates episodic learning curves for Sharpe ratio, Sortino ratio, Maximum Drawdown, VaR, and CVaR. The DDPG agent improves gradually

but exhibits higher tail risk, as shown by lower VaR and CVaR values. SAC, on the other hand, achieves smoother convergence with consistent downside risk control, demonstrating the stabilizing benefits of entropy regularization.

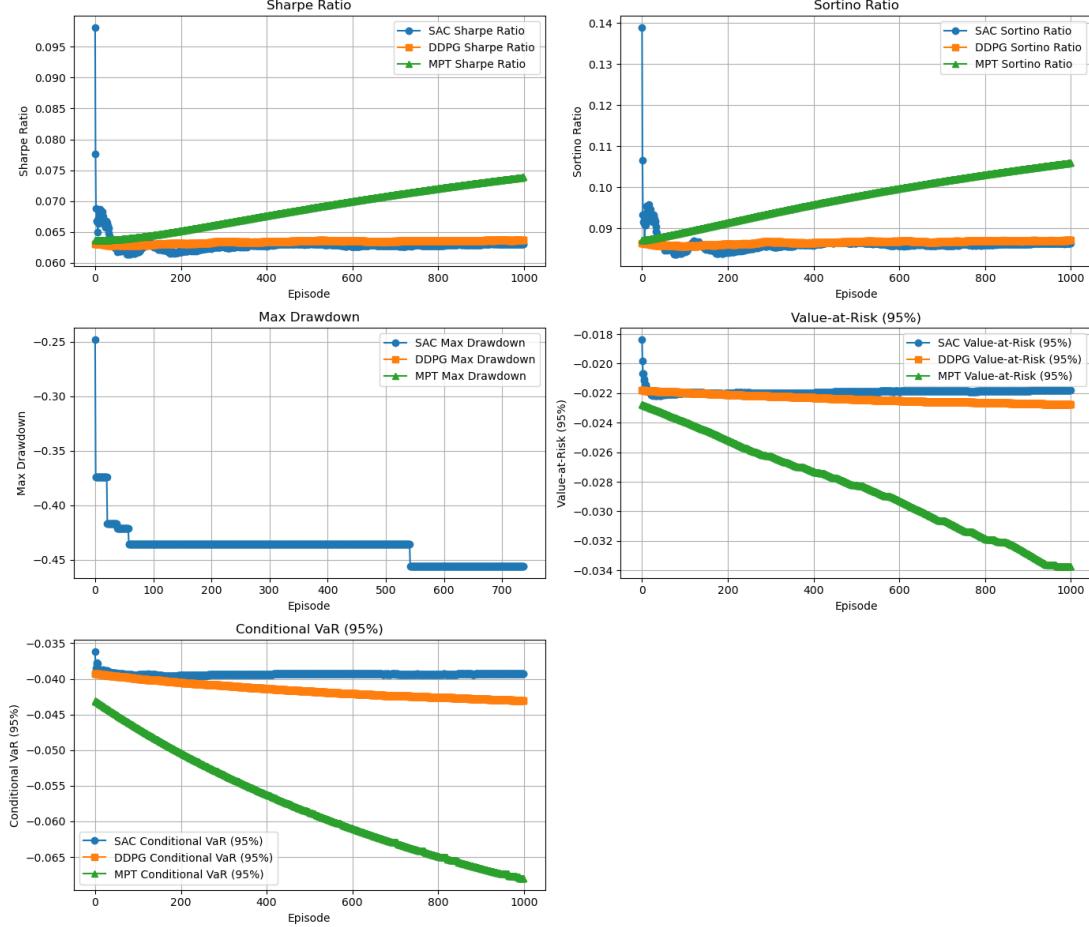


FIGURE 7. Episodic comparison of Sharpe ratio, Sortino ratio, Maximum Drawdown, VaR (95%), and CVaR (95%) for SAC and DDPG algorithms.

Figure 8 shows normalized portfolio values for SAC, DDPG, and the MPT strategy over the test period, benchmarked against Bitcoin. SAC demonstrates the strongest cumulative return, nearly tripling the initial investment. DDPG exhibits higher volatility and less consistent performance. The MPT strategy produces stable but slower growth. The results emphasize the adaptability of SAC and the effectiveness of reinforcement learning in dynamic market conditions.

Table 4 summarizes the test-period results. SAC achieves the highest final portfolio value (2.7627), the largest mean log return, and the best risk-adjusted metrics (Sharpe = 0.0673, Sortino = 0.1093). It also exhibits lower maximum drawdown and reduced tail risk compared to DDPG and MPT.

DDPG shows negative mean returns and higher drawdowns, indicating instability in volatile environments. MPT produces conservative but steady results.

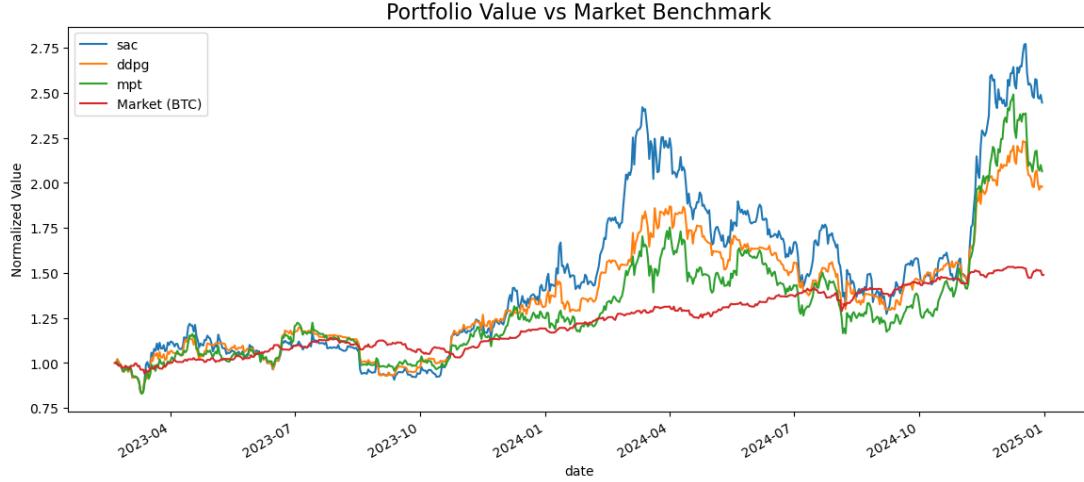


FIGURE 8. Normalized portfolio values for SAC, DDPG, and MPT across BTC, ETH, LTC, and DOGE, compared with Bitcoin during the test period.

SAC's diversified asset allocation further supports its robustness in nonstationary markets.

TABLE 4. Performance comparison of DDPG, SAC, and MPT on the test dataset.

Metric	DDPG	SAC	MPT
Final Portfolio Value	1.969214	2.762704	2.035221
Mean Log Return	-0.001463	0.001492	0.001043
Standard Deviation	0.022354	0.027567	0.023023
Sharpe Ratio	0.013126	0.067340	0.026794
Sortino Ratio	0.018515	0.109281	0.037859
Maximum Drawdown	-0.681602	-0.409029	-0.719820
VaR (95%)	-0.040353	-0.041248	-0.039025
CVaR (95%)	-0.060043	-0.057676	-0.057710
Avg. Weight (BTC)	0.515419	0.100008	0.202570
Avg. Weight (ETH)	0.010000	0.393539	0.233725
Avg. Weight (LTC)	0.048458	0.606461	0.224882
Avg. Weight (DOGE)	0.436123	0.101000	0.189328

5.4. Conclusion and Future Work. This study proposed a deep reinforcement learning framework for dynamic cryptocurrency portfolio management using SAC and DDPG, enhanced with LSTM networks for temporal feature extraction. A

Markowitz MPT baseline was implemented for comparison. Historical data from 2016–2024 were used to train and test the models under realistic conditions.

Empirical evaluations showed that both RL-based agents outperform the classical MPT benchmark, with SAC demonstrating superior risk-adjusted performance, stability, and adaptability in volatile environments. DDPG achieved competitive returns but displayed greater sensitivity to market noise and hyperparameter initialization.

Challenges included overfitting in low-liquidity assets, data quality issues, and computational demands—particularly for SAC due to entropy-regularized updates. Despite these challenges, the results highlight the potential of reinforcement learning for robust and adaptive portfolio management in nonstationary markets.

Future research may incorporate multimodal data sources (e.g., sentiment, blockchain metrics), explore advanced architectures such as transformers, evaluate scalability to larger asset universes, and improve model interpretability using explainable AI techniques.

DISCLOSURE STATEMENT

The author declares no conflicts of interest.

FUNDING

This research received no external funding.

REFERENCES

- [1] Bao, W., Yue, J., and Rao, Y. (2017). A deep learning framework for financial time series using stacked autoencoders and LSTMs. *Neurocomputing*, 356, 284–296.
- [2] Baur, D. G., Hong, K., and Lee, A. D. (2018). Bitcoin: Medium of exchange or speculative assets? *Journal of International Financial Markets, Institutions and Money*, 54, 177–189.
- [3] Brière, M., Oosterlinck, K., and Szafarz, A. (2015). Virtual currency, tangible return: Portfolio diversification with Bitcoin. *Journal of Asset Management*, 16(6), 365–373.
- [4] Cheah, E. T., and Fry, J. (2015). Speculative bubbles in Bitcoin markets? An empirical investigation. *Economics Letters*, 130, 32–36.
- [5] Chen, S., Zhang, R., and Chen, H. (2021). Soft Actor–Critic algorithm for portfolio management. *Expert Systems with Applications*, 168, 114384.
- [6] Christodoulou, P. (2019). Soft Actor–Critic for discrete action settings. *arXiv:1910.07207*.
- [7] Corbet, S., Lucey, B., Urquhart, A., and Yarovaya, L. (2018). Cryptocurrencies as a financial asset. *International Review of Financial Analysis*, 62, 182–199.
- [8] De Miguel, V., Garlappi, L., and Uppal, R. (2009). Optimal versus naive diversification. *Review of Financial Studies*, 22(5), 1915–1953.

- [9] Dempster, M. A. H., and Leemans, V. (2006). An automated FX trading system using reinforcement learning. *Expert Systems with Applications*, 30(3), 543–552.
- [10] Deng, Y., Bao, F., Kong, Y., Ren, Z., and Dai, Q. (2017). Deep direct reinforcement learning for trading. *IEEE Trans. Neural Networks*, 28(3), 653–668.
- [11] Dyrberg, A. H. (2016). Bitcoin, gold and the dollar. *Finance Research Letters*, 16, 85–92.
- [12] Fischer, T., and Krauss, C. (2018). LSTM models for financial prediction. *European Journal of Operational Research*, 270(2), 654–669.
- [13] Fujimoto, S., van Hoof, H., and Meger, D. (2018). Addressing function approximation error in actor–critic methods. ICML, 1587–1596.
- [14] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. (2018). Soft Actor–Critic. ICML, 1861–1870.
- [15] Henderson, P., Islam, R., Bachman, P., et al. (2018). Deep reinforcement learning that matters. AAAI, 32(1), 3207–3214.
- [16] Huang, G., Zhou, X., and Song, Q. (2020). DRL for portfolio management. *arXiv:2012.13773*.
- [17] Jiang, Z., Xu, D., and Liang, J. (2017). DRL framework for portfolio management. *arXiv:1706.10059*.
- [18] Katsiampa, P. (2017). Volatility estimation for Bitcoin. *Economics Letters*, 158, 3–6.
- [19] Kingma, D. P., and Ba, J. (2015). Adam optimizer. ICLR.
- [20] Li, Y., Jiang, B., Li, X., and Qin, M. (2019). DRL for portfolio management. *Neurocomputing*, 358, 256–270.
- [21] Liang, Z., Chen, X., Zhu, Y., et al. (2018). Adversarial DRL in portfolio management. *arXiv:1808.09940*.
- [22] Lillicrap, T. P., Hunt, J. J., et al. (2016). DDPG algorithm. ICLR.
- [23] Liu, Y., and Tsyvinski, A. (2021). Risks and returns of cryptocurrency. *Review of Financial Studies*, 34(6), 2689–2727.
- [24] Liu, F., Li, Y., Li, B., et al. (2021). Bitcoin trading via DRL. *Applied Soft Computing*, 113, 107952.
- [25] Lucarelli, G., and Borrotti, M. (2020). Deep Q-learning portfolio management. *Neural Computing and Applications*, 32, 17229–17244.
- [26] Maringer, D. (2008). *Financial Engineering*. Springer.
- [27] Markowitz, H. (1952). Portfolio selection. *Journal of Finance*, 7(1), 77–91.
- [28] Mensi, W., Sensoy, A., and Kang, S. H. (2019). Crypto–financial asset dependence. *Finance Research Letters*, 29, 68–78.
- [29] Michaud, R. (1989). Is optimized optimal? *Financial Analysts Journal*, 45(1), 31–42.
- [30] Mnih, V., Kavukcuoglu, K., Silver, D., et al. (2015). Human-level control through DRL. *Nature*, 518(7540), 529–533.
- [31] Moody, J., and Saffell, M. (2001). Learning to trade via reinforcement. *IEEE Trans. Neural Networks*, 12(4), 875–889.
- [32] Shen, D., Urquhart, A., and Wang, P. (2020). Bitcoin safe-haven? *Finance Research Letters*, 33, 101317.

- [33] Silver, D., Lever, G., Heess, N., et al. (2014). Deterministic policy gradients. ICML, 387–395.
- [34] Sutton, R. S., and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*, 2nd ed. MIT Press.
- [35] Uhlenbeck, G. E., and Ornstein, L. S. (1930). Brownian motion theory. *Physical Review*, 36(5), 823–841.
- [36] Zhang, Z., Zohren, S., and Roberts, S. (2020). DRL for trading. *Journal of Financial Data Science*, 2(2), 25–40.

DEPARTMENT OF MATHEMATICS, TAFRESH UNIVERSITY, TAFRESH, IRAN

Email address: k.paykan@gmail.com