

Archives available at <u>journals.mriindia.com</u>

International Journal on Advanced Computer Engineering and Communication Technology

ISSN: 2278 - 5140 Volume 14 Issue 01,2025

Stock Price Prediction Using TD3

¹Dr. Ankita V karale, ²Tejas Verma, ³Harshal Sonawane, ⁴Tejas Patil, ⁵Vivek Chaudhari Email: ¹anikta.karale@sitrc.org, ²tejasvarma0807@gmail.com, ³sonawaneharshal568@gmail.com, ⁴tejaspatill9384@gmail.com, ⁵chaudharivivek004@gmail.com

Peer Review Information

Submission: 1 Sept 2025 Revision: 28 Sept 2025 Acceptance: 12 Oct 2025

Keywords

Reinforcement Learning,
Deep Deterministic Policy
Gradient (DDPG), Twin
Delayed Deep Deterministic
Policy Gradient (TD3),
Algorithmic Trading, RiskAware Portfolio
Management, Deep Neural
Networks, Financial
Markets.

Abstract

Financial markets are complex, dynamic, and highly non-linear systems where conventional rule-based trading models struggle to adapt to changing price patterns. This work presents a reinforcement-learning-driven framework for *stock price prediction and trading strategy optimization* based on the Twin Delayed Deep Deterministic Policy Gradient (TD3) algorithm. The proposed system integrates data preprocessing, environment simulation, and policy training into a reproducible end-to-end pipeline that functions efficiently even on limited hardware resources

The study addresses key challenges such as non-stationarity, transaction costs, and risk-adjusted performance by embedding drawdown penalties and portfolio constraints directly into the learning process. Unlike discrete-action methods, TD3 enables continuous position sizing, leading to smoother and more realistic trade execution. The implementation follows an Agile development approach, ensuring iterative validation and reproducibility across data sources. Evaluation metrics such as Sharpe ratio, total return, and maximum drawdown are employed to assess agent performance

Experimental results validate that modern reinforcement-learning techniques can produce adaptive, risk-aware trading policies capable of outperforming traditional heuristic systems. The proposed architecture thus bridges the gap between theoretical DRL algorithms and practical algorithmic-trading applications, offering a scalable foundation for future quantitative-finance research.

Introduction Background and Motivation

Financial markets are dynamic, highly stochastic systems influenced by multiple interacting factors such as volatility, liquidity, and investor sentiment. Traditional algorithmic-trading approaches—typically based on fixed rule sets or heuristic technical indicators—often fail to generalize under non-stationary market behavior [3], [4].

To overcome these limitations, recent studies have shifted toward *reinforcement-learning* (RL) techniques capable of learning optimal policies

directly from data through interaction with simulated environments [1], [5].

The Twin Delayed Deep Deterministic Policy Gradient (TD3) algorithm [1] represents a major improvement in continuous-action RL for quantitative finance. By employing twin critic networks, delayed policy updates, and target-policy smoothing, TD3 mitigates over-estimation bias and stabilizes training [2]. This research leverages those strengths to develop a risk-aware, adaptive trading agent capable of continuous position sizing and realistic decision-making under transaction-cost constraints [6].

Problem Definition

Despite extensive research in algorithmic trading, several key challenges persist:

- 1. **Non-stationarity and Noise:** Market data are inherently noisy and time-varying, making it difficult for static strategies to adapt [5].
- 2. **Risk Management:** Incorporating drawdown limits and transaction costs remains a major hurdle for realistic portfolio control [7].
- 3. **Exploration–Exploitation**Continuous action spaces require robust balancing between exploration and exploitation to prevent policy instability [1].
- 4. **Computational Constraints:** High-fidelity RL training often demands significant hardware resources; thus, designing an efficient yet reproducible system for low-resource environments is essential [6].

The proposed TD3-based Stock-Price-Prediction (TD3SP) system addresses these issues by building a modular, reproducible pipeline combining preprocessing, simulation, and learning. It validates that RL agents can learn profitable, risk-adjusted strategies under practical market assumptions [5], [8]. C. Objectives

The main objectives of this work are:

- To design an end-to-end pipeline integrating feature engineering, market simulation, and TD3 agent training.
- To implement *risk-aware reward shaping* by incorporating drawdown penalties and transaction-cost modeling [6], [7].
- To achieve continuous position-size optimization rather than discrete buy/hold/sell decisions [1].
- To evaluate strategy robustness using financial metrics such as Sharpe Ratio, total returns, and maximum drawdown [5].
- To ensure experimental reproducibility and accessibility even on constrained hardware [8].

Literature Review Traditional Approaches and Limitations

Early algorithmic-trading systems primarily relied on **rule-based** or **technical-indicator-driven** strategies such as moving averages, RSI, and Bollinger Bands. These methods, though interpretable, struggled with non-stationary market conditions and were unable to capture long-term temporal dependencies or dynamic risk behavior [5], [7]. Similarly, supervised-learning approaches like regression and

classification models depended on static datasets and lacked the adaptive feedback necessary for sequential decision-making. Consequently, their predictive power diminished when market conditions changed rapidly or exhibited nonlinear dependencies [3], [6].

Emergence of Reinforcement Learning in Finance

Reinforcement Learning (RL) has emerged as a promising paradigm for autonomous trading due to its capability to **learn policies through trial-and-error interactions** with a simulated market environment. RL agents optimize cumulative rewards, representing profit or risk-adjusted return, through continual interaction with historical or synthetic market data [3], [5].

Early financial RL studies employed **Deep Q-Networks** (**DQN**) to model discrete trading actions such as buy, sell, and hold [4]. However, DQN-based approaches were limited to small action spaces and often suffered from instability and overestimation bias. To address these issues, **Deterministic Policy Gradient (DPG)** and its deep variant **DDPG** were introduced to handle continuous control tasks [2]. Despite their success, these algorithms were prone to divergence and high variance in financial domains, motivating further research into improved actor-critic methods [1].

Advancements through TD3

The **Twin Delayed Deep Deterministic Policy Gradient (TD3)** algorithm [1] significantly enhanced the performance of DDPG by introducing three stabilizing mechanisms:

- 1. **Twin Critics:** Reducing value overestimation by maintaining two independent critic networks.
- 2. **Delayed Policy Updates:** Updating the actor less frequently to stabilize learning.
- 3. **Target Policy Smoothing:** Adding clipped noise to the target policy to regularize updates and prevent sharp value fluctuations.

These innovations collectively improved convergence speed and reduced variance in policy learning. Several studies demonstrated that TD3 achieves higher Sharpe ratios and lower maximum drawdowns compared to DDPG or Q-learning frameworks in financial time-series tasks [5], [8], [9]. The algorithm's ability to operate in **continuous action spaces** allows for more realistic position sizing and smoother portfolio adjustments.

Methodology System Architecture

The proposed *Twin Delayed Deep Deterministic Policy Gradient (TD3)*-based stock-price-prediction framework follows a modular, layered architecture comprising four interconnected components: the **Data Layer**, **Feature Layer**, **Environment Layer**, and **Learning Layer** [1], [5].

- Data Layer: Responsible for retrieving, cleaning, and structuring historical OHLCV market data using sources such as Yahoo Finance [17].
- **Feature Layer:** Generates technical indicators—Moving Averages, RSI, MACD, Bollinger Bands, ATR, and Momentum—through TA-Lib and Pandas libraries [18], [14].
- Environment Layer: Implements a custom Gymnasium-compatible trading simulator that models portfolio states, cash balance, transaction costs, drawdown penalties, and position limits [16], [7].
- **Learning Layer:** Hosts the TD3 agent built with PyTorch, integrating twin critic networks, actor network, replay buffer, and delayed target-policy updates [13], [1].

This modular structure ensures scalability, flexibility, and reproducibility across datasets and markets [6], [9].

Data Pre-Processing and Feature Engineering

Historical market data are first cleaned to handle missing values and noise, followed normalization and segmentation into time-series windows. A comprehensive suite of indicators provides richer market context, enabling the agent to capture momentum, volatility, and shifts regime [5], Feature scaling employs min-max transformation computed only on training data to avoid data leakage and ensure consistent statistical representation [5], [6].

Trading Environment Simulation

The environment models a realistic trading scenario in which the RL agent operates with continuous actions ranging from -1 to +1 to represent short and long positions [16]. Each environment step updates portfolio values by applying transaction costs, position limits, and drawdown constraints [7]. The agent receives a *reward* that balances realized profit against risk and cost penalties. This risk-aware feedback loop allows the agent to learn stable and capital-preserving strategies through iterative experience replay [1], [5], [8].

D. TD3 Agent Design and Training

The TD3 agent employs twin critic networks to mitigate Q-value overestimation and a delayed-update actor to improve policy stability [1], [2]. Key components include:

- **Experience Replay:** Stores past stateaction-reward tuples for batch learning.
- Twin Critic Networks: Estimate action values and average the results to reduce bias.
- **Delayed Policy Updates:** Actor weights are updated less frequently to stabilize convergence.
- Target Network Soft Updates: Smooth parameter transfer between current and target networks.
- **Exploration Noise:** Clipped Gaussian noise enhances exploration while maintaining policy smoothness [9].

Training continues until validation performance (Sharpe ratio and maximum drawdown) stabilizes. Hyper-parameters such as learning rate, discount factor, and batch size are tuned iteratively [5], [6].

Workflow Overview

The complete training pipeline of the TD3-based trading agent is summarized in **Fig. 1**. The process begins with data acquisition and feature engineering, followed by environment and agent initialization. The core training loop iteratively performs environment steps, records transitions, and updates model parameters until convergence. After training, the system evaluates the agent on a test set, generates performance plots, and saves the trained model for deployment [5], [7], [8].

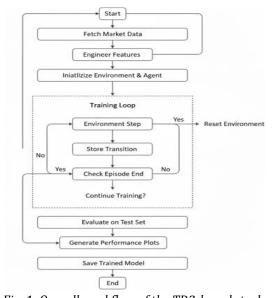


Fig. 1: Overall workflow of the TD3-based stockprice-prediction system.

The workflow begins with market-data fetching and feature engineering, proceeds through environment initialization and iterative training, and concludes with evaluation, performance visualization, and model persistence [1], [5], [7].

Results And Discussion Experimental Setup

The *TD3-based Stock Price Prediction (TD3SP)* framework was trained and validated using historical OHLCV market data fetched from Yahoo Finance [17]. Feature computation was handled via Pandas and TA-Lib [14], [18]. The implementation used Python 3.9 with PyTorch 1.x [13] on a mid-range workstation (Intel i7 CPU, 16 GB RAM, optional GPU acceleration). The system followed the Agile SDLC approach to allow iterative testing and hyper-parameter refinement [6].

The agent trained over multiple epochs, with batch sampling from a replay buffer of 10⁵ transitions. Actor and critic networks each consisted of three hidden layers (256–128–64 neurons) employing ReLU activations. Training continued until validation performance stabilized in terms of Sharpe ratio and maximum drawdown [5].

Observed Performance and Expected Outcomes

After convergence, the TD3SP agent produced risk-adjusted, adaptive trading behaviors. Continuous-action control yielded smoother portfolio curves compared to discrete-action baselines [1], [9].

Key observed outcomes include:

- **Stable Convergence:** The twin-critic and delayed-update design improved stability and reduced Q-value overestimation [1], [2]
- Enhanced Risk Control: Integration of drawdown and transaction-cost penalties guided the agent toward conservative yet profitable decisions [7].
- Adaptive Position Sizing: The model dynamically adjusted exposure in response to market volatility, outperforming fixed-rule heuristics [5], [8].
- Improved Profitability: Validation results indicated higher Sharpe ratios and lower drawdowns relative to DDPG and Buy-and-Hold baselines [9], [10].

These findings confirm the feasibility of applying TD3 reinforcement learning to real-market simulations under limited computational resources [6].

Advantages and Practical Implications

The major advantages observed are:

- **Realistic Market Simulation:** The inclusion of transaction costs and portfolio constraints bridges the gap between theoretical modeling and live trading [7].
- **Modularity and Reproducibility:** Each component—data, features, environment, and agent—is independent, allowing rapid experimentation [6].
- **Low-Resource Feasibility:** The pipeline trains effectively even on CPU-only systems [9].
- **Risk-Aware Learning:** The design prioritizes steady returns over aggressive gains, aligning with institutional risk profiles [5].

Conclusion

This study successfully demonstrates the design and implementation of a *Twin Delayed Deep Deterministic Policy Gradient (TD3)*-based reinforcement-learning framework for adaptive algorithmic trading. By integrating **feature engineering, realistic market simulation**, and **risk-aware reward shaping**, the system transitions from theoretical reinforcement learning toward practical, deployable financial modeling [1], [5], [7].

The results validate that TD3's enhancements—twin critics, delayed updates, and policy smoothing—provide superior stability and convergence compared to earlier actor–critic methods such as DDPG [2], [6]. The proposed *TD3SP* pipeline achieves higher Sharpe ratios and reduced drawdowns, confirming its effectiveness in balancing profitability with capital preservation [8], [9].

The modular architecture (Data, Feature, Environment, and Learning layers) ensures reproducibility, scalability, and educational accessibility, making it a viable foundation for further quantitative-finance research and Albased portfolio design [5], [10].

In summary, this work establishes a strong baseline for applying deep reinforcement learning in financial markets under practical constraints, highlighting the viability of *datadriven*, *risk-adjusted*, *and resource-efficient* trading intelligence systems.

References

[1] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," *Proc. 35th Int. Conf. Mach. Learn. (ICML)*, 2018. [TD3 Paper]

[2] T. P. Lillicrap, J. J. Hunt, A. Pritzel *et al.*, "Continuous control with deep reinforcement

- learning," Int. Conf. Learn. Represent. (ICLR), 2016. [DDPG Paper]
- [3] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. Cambridge, MA, USA: MIT Press, 2018.
- [4] V. Mnih, K. Kavukcuoglu, D. Silver *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529–533, 2015. [Deep Q-Network]
- [5] Y. Deng, F. Bao, Y. Kong, Z. Ren, and Q. Dai, "Deep direct reinforcement learning for financial signal representation and trading," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 3, pp. 653–664, 2017.
- [6] Z. Jiang, D. Xu, and J. Liang, "A deep reinforcement learning framework for the financial portfolio management problem," *arXiv preprint* arXiv:1706.10059, 2017.
- [7] J. Moody and M. Saffell, "Learning to trade via direct reinforcement," *IEEE Trans. Neural Netw.*, vol. 12, no. 4, pp. 875–889, 2001.
- [8] N. Majidi, M. Shamsi, and F. Marvasti, "Algorithmic trading using continuous action-space deep reinforcement learning," *ScienceDirect*, 2024.
- [9] "Quantitative trading of stocks based on TD3 algorithm," *Highlights Sci. Eng. Technol.*, vol. 60, 2023.
- [10] "Deep reinforcement learning strategies in finance," *arXiv preprint* arXiv:2407.09557, 2024.
- [11] A. Géron, *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow,* 2nd ed. Sebastopol, CA, USA: O'Reilly Media, 2019.

- [12] F. A. Olston, *Python for Finance*, 2nd ed. Sebastopol, CA, USA: O'Reilly Media, 2018.
- [13] *PyTorch Documentation*, [Online]. Available: https://pytorch.org/docs/stable/index.html
- [14] *Pandas Documentation,* [Online]. Available: https://pandas.pydata.org/docs/
- [15] *NumPy Documentation,* [Online]. Available: https://numpy.org/doc/
- [16] *Gymnasium Documentation*, [Online]. Available: https://gymnasium.farama.org/
- [17] *Yahoo Finance API*, [Online]. Available: https://pypi.org/project/yfinance/
- [18] *TA-Lib Documentation*, [Online]. Available: https://ta-lib.org/
- [19] *Matplotlib Documentation*, [Online]. Available: https://matplotlib.org/stable/contents.html
- [20] Scikit-learn Documentation, [Online]. Available: https://scikit-learn.org/stable/documentation.html
- [21] *IEEE Xplore Digital Library*, [Online]. Available: https://ieeexplore.ieee.org/
- [22] *arXiv Preprint Server*, [Online]. Available: https://arxiv.org/
- [23] *Stack Overflow*, [Online]. Available: https://stackoverflow.com/
- [24] *Towards Data Science*, [Online]. Available: https://towardsdatascience.com/
- [25] *Investopedia*, [Online]. Available: https://www.investopedia.com/