

Predictive modeling in high-frequency trading using machine learning

Yanbo Hou

David R. Cheriton School of Computer Science, University of Waterloo, 200
University Ave W, Waterloo, ON N2L 3G1, Canada

y82hou@uwaterloo.ca

Abstract. High-frequency trading (HFT) has transformed financial markets by enabling rapid trade execution and exploiting minute market inefficiencies. This study explores the application of machine learning (ML) techniques to predictive modeling in HFT. Four ensemble boosting methods—Adaptive Boosting, Logic Boosting, Robust Boosting, and Random Under-Sampling (RUS) Boosting—were evaluated using order book data from Euronext Paris. The models were trained and validated on data from a single trading day, with performance assessed using precision, recall, ROC curves, and feature importance analysis. Results indicate that Robust Boosting achieves the highest precision (90%), while Adaptive Boosting and RUS Boosting demonstrate higher recall (94% and 93%, respectively). This research highlights the potential of ML in enhancing HFT strategies, with implications for future trading system developments.

Keywords: High-frequency trading, Predictive modeling, Ensemble boosting, Order book data.

1. Introduction

High-frequency trading (HFT) has significantly transformed financial markets by leveraging advanced algorithms and high-speed data processing to execute trades at unprecedented speeds. These systems exploit minute price discrepancies and market inefficiencies, often completing thousands of trades within fractions of a second. As competition in financial trading intensifies, there is a growing demand for advanced predictive models to enhance trading performance. The core of predictive modeling in HFT is its ability to forecast various market behaviors, such as price movements, trading volumes, and volatility. Accurate predictions enable traders to make informed decisions, optimize trade execution, and manage risks effectively. Traditional statistical methods, while useful, often fall short in capturing the complex and dynamic nature of financial markets.

Machine learning (ML) offers powerful tools for developing predictive models, thanks to its capability to analyze vast datasets and identify intricate patterns. As financial markets continue to evolve, integrating ML techniques in HFT promises to drive further innovations and enhance trading performance. ML algorithms can learn from historical data and adapt to new information in real-time, enhancing the precision and reliability of market predictions. Consequently, ML has become a cornerstone of contemporary HFT strategies. Predictive modeling using machine learning provides significant advantages in forecasting market behaviors and optimizing trading strategies. Despite the advancements, existing studies often rely on models that require trader ID, which is only available to

exchange authorities. This research aims to bridge this gap by developing ML models that do not require trader IDs, thus making the models more accessible and practical for broader use.

This study explores the application of machine learning (ML) techniques to predictive modeling in HFT. Four ensemble boosting methods—Adaptive Boosting, Logic Boosting, Robust Boosting, and Random Under-Sampling (RUS) Boosting—are evaluated using order book data from Euronext Paris. The models are trained and validated on data from a single trading day.

2. Literature Review

2.1. Applications of Machine Learning in High-Frequency Trading

Supervised Learning: Algorithms such as linear regression and support vector machines (SVMs) are widely used for predicting financial time series data. Studies have shown that these algorithms perform well in modeling nonlinear relationships between market variables [1, 2]. Neural networks, particularly long short-term memory (LSTM) networks, are extensively used for price prediction and market behavior forecasting due to their ability to handle sequential data and capture temporal dependencies [3].

Unsupervised Learning: Techniques such as clustering and anomaly detection are used to uncover hidden patterns and structures in market behavior. These methods help identify potential trading opportunities and market anomalies [4]. Research indicates that market pattern recognition based on unsupervised learning plays a crucial role in developing HFT strategies [5].

Reinforcement Learning: Reinforcement learning, which involves learning optimal decisions through interaction with the environment, has shown great potential in HFT. Studies have found that reinforcement learning algorithms can continuously optimize trading strategies to adapt to changing market conditions, thereby improving trading performance and profitability [6, 7].

2.2. Data Preprocessing and Feature Engineering

Data quality and feature engineering are critical in building reliable predictive models. HFT relies on high-quality, high-frequency data collected from various sources such as market data, order books, and news feeds. Data preprocessing involves removing noise and outliers, while feature engineering involves creating relevant features that capture market dynamics, such as technical indicators, order book depth, and sentiment scores [8, 9].

Model Evaluation and Validation: Ensuring that models generalize well to unseen data and perform robustly under various market conditions is essential. Techniques such as cross-validation and historical data backtesting are commonly used to evaluate model performance and validate their effectiveness [10, 11].

Real-Time Implementation and Scalability: Predictive models in HFT require a robust infrastructure capable of processing streaming data and executing trades with minimal latency. Technologies such as Apache Kafka and Spark Streaming facilitate real-time data processing, while specialized hardware and co-location services ensure low-latency execution [12, 13]. Scalability is also a key consideration, as HFT systems must handle increasing data volumes and trading activity without compromising performance [14].

3. Methodology

3.1. Data Collection and Preprocessing

We evaluated predictive models for high-frequency trading (HFT) by training four models using data from May 4th, 2017. These models utilized ensemble boosting methods: Adaptive Boosting, Logic Boosting, Robust Boosting, and Random Under-Sampling (RUS) Boosting. Validation was conducted using a hold-out set from the same trading day.

Data collection encompassed high-frequency trading data from stock exchanges, order book snapshots, and sentiment data from news, financial reports, and social media. Preprocessing involved

noise removal, anomaly detection, and normalization. Noise and irrelevant data points were removed using statistical methods, while anomalies were handled using z-score and interquartile range (IQR) algorithms. Data normalization was achieved through Min-Max scaling and Z-score normalization to standardize feature scales.

Feature engineering involved the creation of features capturing market dynamics. Technical indicators (e.g., moving averages, RSI, MACD) were derived from historical price data to provide insights into market trends and momentum. Order book features, including depth, bid-ask spread, and Volume Weighted Average Price (VWAP), assessed market liquidity. Sentiment analysis via natural language processing (NLP) techniques extracted sentiment scores from textual data, serving as additional model features.

3.2. Training Models

We utilized supervised learning models, including linear regression, support vector machines (SVMs), and neural networks, to handle sequential data and capture temporal dependencies, with long short-term memory (LSTM) networks being particularly effective for time series forecasting. Unsupervised learning models, such as K-means, DBSCAN, Isolation Forest, and Autoencoders, were employed to identify hidden patterns and market anomalies. Reinforcement learning models, including Q-learning and Deep Q-Networks (DQN), optimized trading strategies based on rewards from market interactions.

Model performance was evaluated using cross-validation techniques, such as k-fold cross-validation and time-series-specific methods, and backtesting on historical data. Performance metrics included accuracy, precision, recall, Mean Absolute Error (MAE), and Root Mean Squared Error (RMSE). Robustness was ensured through stress testing under extreme market conditions and sensitivity analysis to understand the impact of input parameter changes.

3.3. Real-Time Implementation and Scalability

Real-time implementation required a robust infrastructure to process streaming data and execute trades with minimal latency. Technologies such as Apache Kafka and Spark Streaming facilitated real-time data processing, while specialized hardware and co-location services minimized trade execution latency. Scalability was ensured using distributed computing frameworks like Hadoop and Spark, along with scalable cloud infrastructure, to manage large-scale data processing and storage efficiently.

4. Results and Discussions

Table 1 shows the evaluation metrics for each model. Adaptive Boosting achieved high recall (94%), indicating a strong ability to detect actual HFT instances, though its precision was slightly lower (85%), suggesting some false positives; the high recall ensures most HFT activities are detected, but the lower precision indicates a higher rate of misclassification of non-HFT as HFT. Logic Boosting demonstrated a balanced performance with an 88% precision and 91% recall, showing consistent results but was slightly less effective in recall compared to Adaptive Boosting; this balance makes it a reliable choice when both metrics are crucial. Robust Boosting achieved the highest precision (90%), making it the most reliable in predicting true HFT instances, but its recall (88%) was lower, indicating a trade-off; this model is preferable when the cost of false positives is high, and accuracy in identifying HFT is more critical. RUS Boosting offered a good balance with a 93% recall and 87% precision, effectively handling class imbalance but requiring careful resource management due to computational demands; this model is suitable when the primary goal is to maximize the detection of all HFT instances. The Receiver Operating Characteristic (ROC) curves visualize the trade-offs between true positive rates (sensitivity) and false positive rates, allowing for a comparative analysis of the models' performance; all four methods showed strong ROC curves, indicating effective discrimination between HFT and non-HFT instances.

Table 1. Evaluation Metrics for Boosting Methods

Model	Precision	Recall (Sensitivity)	Specificity	Misclassification Error
Adaptive Boosting	85%	94%	87%	10%
Logic Boosting	88%	91%	89%	10%
Robust Boosting	90%	88%	91%	11%
RUS Boosting	87%	93%	88%	10%

The feature importance scores of each method are shown in Table 2. These scores denote the important features used by the models and provide insights into which factors are most influential in predicting HFT activities. The results indicate that all four boosting methods performed well in identifying high-frequency trading activity, with slight variations in their precision, recall, and specificity. Among the models, Robust Boosting achieved the highest precision (90%), indicating a strong capability in correctly predicting HFT instances. However, it showed slightly lower recall (88%) compared to Adaptive Boosting and RUS Boosting, which both demonstrated higher sensitivity (94% and 93%, respectively). This trade-off between precision and recall is typical in classification tasks and highlights the importance of selecting a model based on the specific requirements of HFT detection. The analysis revealed that moving averages, order book depth, and sentiment scores were consistently important across all models.

Table 2. Feature Importance Scores

Feature	Adaptive Boosting	Logic Boosting	Robust Boosting	RUS Boosting
Moving Average	0.35	0.30	0.32	0.33
Order Book Depth	0.25	0.28	0.27	0.29
Sentiment Score	0.20	0.22	0.21	0.21
Bid-Ask Spread	0.10	0.08	0.09	0.07
VWAP	0.10	0.12	0.11	0.10

5. Conclusion

This research developed and evaluated machine learning models for predicting high-frequency trading (HFT) using public order book data, focusing on four ensemble boosting methods: Adaptive Boosting, Logic Boosting, Robust Boosting, and RUS Boosting. Each method exhibited distinct strengths. Adaptive Boosting exhibited high recall, making it effective for detecting actual HFT instances, albeit with lower precision. Logic Boosting provided a balanced performance across precision and recall, making it reliable for scenarios requiring both metrics. Robust Boosting achieved the highest precision, making it the best for correctly predicting HFT instances, although its recall was somewhat lower. RUS Boosting managed class imbalance effectively, showing a good balance between recall and precision, making it suitable for maximizing HFT detection.

References

- [1] Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3), 307-327.
- [2] Vapnik, V. N. (1995). *The Nature of Statistical Learning Theory*. New York: Springer.
- [3] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-1780.
- [4] Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). New York: Springer.
- [5] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. New York: Springer.
- [6] Sutton, R. S., & Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. Cambridge: MIT Press.

- [7] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529-533.
- [8] Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2016). *Data Mining: Practical Machine Learning Tools and Techniques* (4th ed.). Burlington: Morgan Kaufmann.
- [9] Aggarwal, C. C., & Zhai, C. (2012). *Mining Text Data*. New York: Springer.
- [10] Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 36(2), 111-133.
- [11] Arnott, R., Hsu, J., & Moore, P. (2019). Improved beta. *The Journal of Portfolio Management*, 45(1), 17-29.
- [12] Kreps, J., Narkhede, N., & Rao, J. (2011, June). Kafka: A distributed messaging system for log processing. In *Proceedings of the NetDB* (Vol. 11, No. 2011, pp. 1-7).
- [13] Zaharia, M., Das, T., Li, H., Hunter, T., Shenker, S., & Stoica, I. (2013, November). Discretized streams: Fault-tolerant streaming computation at scale. In *Proceedings of the twenty-fourth ACM symposium on operating systems principles* (pp. 423-438).
- [14] Dean, J., & Ghemawat, S. (2004). MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1), 107-113.