# Survey on Real-Time Data Processing in Finance Using Machine Learning Techniques

**Pushkar Mehendale**

Troy, MI, USA
Email: *pushkar.mehendale[at]yahoo.com*

**Abstract:** *Real-time data processing is transforming the financial industry by enabling applications that demand immediate insights and decisions. Machine learning (ML) techniques play a pivotal role in this transformation, with algorithms like Support Vector Machines (SVMs), Deep Neural Networks (DNNs), and Random Forests (RFs) widely used for tasks such as stock price prediction, credit risk assessment, fraud detection, and algorithmic trading. Frameworks like Apache Hadoop and Apache Spark facilitate efficient data handling and analysis, but challenges such as data volume and velocity, data variety, data quality, and latency need to be addressed for successful real-time data processing in finance.*

**Keywords:** Machine Learning, Real-Time Data Processing, Finance, Apache Hadoop, Apache Spark

## 1. Introduction

Real-time data processing in financial markets is crucial for maintaining a competitive edge in today's fast-paced environment [2]. The financial sector generates enormous amounts of data every second, and the ability to analyze this data in real-time can significantly enhance decision-making processes. Machine learning (ML) techniques provide a powerful means to analyze real-time data, as they can learn complex patterns and relationships in data and make accurate predictions [6].

One of the key applications of real-time data processing and ML in finance is stock market prediction. Stock prices are highly volatile, and even small changes in price can have a significant impact on investment decisions. ML algorithms can be used to analyze real-time market data, such as stock prices, trading volume, and news sentiment, to predict future price movements. This information can be invaluable for investors and traders, as it allows them to make informed decisions about when to buy, sell, or hold stocks.

Another important application of real-time data processing and ML in finance is credit card fraud detection. Credit card fraud is a major problem, and it costs financial institutions billions of dollars each year. ML algorithms can be used to analyze real-time credit card transaction data to identify fraudulent transactions. This information can be used to prevent fraudulent transactions from occurring, and it can also help to identify the perpetrators of credit card fraud.

## 2. Background

Real-time data processing involves the processing of data as it is generated or received, which is crucial for time-sensitive financial applications. Traditional batch processing methods, where data is collected over a period and processed later, are no longer adequate for handling the high speed and massive volume of financial data. Real-time data processing techniques are essential to keep up with the pace and complexity of modern financial transactions [5].

The financial industry heavily relies on accurate and timely data analysis for informed decision-making. Real-time data processing enables financial institutions to analyze data as it becomes available, allowing them to identify trends, patterns, and potential risks in real-time. This empowers decision-makers with up-to-date insights, enabling them to make quick and informed adjustments to their strategies, risk management, and investment decisions [1].

Real-time data processing also plays a vital role in fraud detection and prevention. It allows financial institutions to monitor transactions as they occur, enabling them to detect suspicious activities and fraudulent patterns in real-time. By leveraging real-time data processing, financial institutions can minimize losses, protect customer accounts, and maintain the integrity of their financial systems.

Overall, real-time data processing is an indispensable component of modern financial systems. It enables financial institutions to harness the power of data in real-time, gain valuable insights, make informed decisions, mitigate risks, and enhance the overall efficiency and accuracy of their operations.

### a) Importance of Real-Time Data Processing
Real-time data processing empowers financial institutions with the ability to make swift and informed decisions. It allows for the prompt detection of fraudulent activities and enables institutions to adapt quickly to changing market conditions. This leads to enhanced operational efficiency and effectiveness, as up-to-date information provides valuable insights for decision-making. Real-time data processing provides financial institutions with a competitive advantage by enabling them to respond swiftly to market opportunities and challenges.

### b) Traditional Data Processing Methods
Traditional batch processing methods, which gather and store data over a period before processing it in batches, are no longer adequate for the dynamic and high-velocity financial data streams of today. The latency inherent in batch processing makes it unsuitable for scenarios where immediate data analysis is essential, such as real-time risk assessment,

fraud detection, and algorithmic trading. In these instances, the time lag between data collection and processing can have significant financial implications, emphasizing the need for real-time data processing solutions that can handle large volumes of data with minimal latency.

## 3. Machine Learning in Finance

Machine learning (ML) algorithms play a crucial role in financial applications, enabling the analysis of vast datasets, pattern identification, and predictive modeling. One of the most well-known ML algorithms used in finance is linear regression, which establishes a linear relationship between input variables (such as historical stock prices) and an output variable (such as future stock prices). By fitting a linear equation to the data, linear regression can make predictions about future values based on the observed trends. Another commonly used ML algorithm in finance is decision trees, which create a hierarchical structure to classify data points based on their attributes [9]. Decision trees are often employed for fraud detection, as they can identify anomalous patterns in financial transactions that may indicate fraudulent activity [3].

In addition to linear regression and decision trees, other ML algorithms such as neural networks, support vector machines, and ensemble methods are also utilized in financial applications. Neural networks, inspired by the human brain, consist of interconnected layers of nodes that can learn complex relationships within data. Support vector machines are powerful classification algorithms that can separate data points into distinct classes by finding the optimal hyperplane that maximizes the margin between them. Ensemble methods, such as random forests and gradient boosting, combine multiple ML models to improve overall predictive performance [7]. These various ML algorithms provide a versatile toolkit for financial professionals, enabling them to tackle complex challenges such as stock price prediction, credit risk assessment, and portfolio optimization [8].

### 1) Support Vector Machines (SVM)
Support Vector Machines (SVMs) are supervised learning models highly effective in classification and regression tasks, particularly when the goal is to separate data into distinct classes. SVM works by finding the optimal hyperplane that best separates the classes in the feature space. This makes them robust against overfitting, especially in high-dimensional spaces. SVMs excel in environments involving stock price prediction or identifying fraudulent transactions, where the ability to accurately distinguish between different classes is crucial.

a) **Application in Stock Market Prediction:** SVMs are utilized in stock market prediction to classify stock movements into categories such as "up" or "down." By training on historical stock price data, SVMs can identify patterns and trends that indicate future movements, providing traders with actionable insights.

b) **Application in Fraud Detection:** In fraud detection, SVMs are used to classify transactions as either fraudulent or legitimate. The model is trained on a dataset of historical transactions, labeled as fraudulent or not, enabling it to recognize anomalies and unusual patterns indicative of fraud.

### 2) Deep Neural Networks (DNN)
Deep Neural Networks (DNNs) have revolutionized the field of financial data analysis. These models, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), possess exceptional capabilities in handling complex structures and relationships within financial data. Their ability to learn from vast historical datasets allows them to identify subtle patterns that might evade detection by simpler algorithms. DNNs find extensive applications in the financial sector, particularly in areas such as credit card fraud detection. By processing transaction data in real-time, DNNs can detect anomalies and potential frauds, offering a powerful tool for safeguarding financial institutions and consumers [5].

a) **Convolutional Neural Networks (CNNs):** CNNs are effective in processing grid-like data structures, making them suitable for analyzing time-series data like stock prices. By applying convolutional layers, CNNs can capture temporal dependencies and patterns within the data.

b) **Recurrent Neural Networks (RNNs):** RNNs are designed for sequential data, making them ideal for tasks such as predicting stock prices based on historical data. RNNs, and particularly Long Short-Term Memory (LSTM) networks, can capture long-term dependencies in the data, enhancing their predictive accuracy.

### 3) Random Forests
Random Forests is an ensemble learning technique that utilizes multiple decision trees during training to reduce overfitting and improve prediction accuracy. It combines the outputs of individual trees, providing a robust and interpretable model capable of handling large datasets with many features. In finance, Random Forests has proven valuable for tasks such as credit scoring, stock returns prediction, and other applications that require reliable and comprehensible models.

## 4. Real-Time Data Processing Frameworks
### 1) Apache Hadoop
Apache Hadoop is an open-source framework that allows for the distributed storage and processing of large datasets. Hadoop's HDFS (Hadoop Distributed File System) and MapReduce programming model are designed to scale up from single servers to thousands of machines, each offering local computation and storage. This makes Hadoop suitable for processing massive datasets generated in financial markets.

HDFS provides a scalable and reliable storage solution that can handle petabytes of data. It distributes data across multiple machines, ensuring fault tolerance and high availability.

MapReduce is a programming model that processes large datasets in parallel across a Hadoop cluster. It simplifies the data processing task by dividing it into smaller sub-tasks, which are processed independently and in parallel, significantly speeding up data processing times [10].

### 2) Evaluation of Explanations
Apache Spark is an open-source distributed computing system that provides an interface for programming entire

clusters with implicit data parallelism and fault tolerance. Spark's in-memory processing capabilities make it significantly faster than Hadoop MapReduce for certain types of data processing tasks. Spark's MLlib library provides a wide array of machine learning algorithms, which can be applied to real-time data streams for tasks such as predictive analytics and anomaly detection.

Spark's ability to perform in-memory processing drastically reduces the time taken to process data, as it eliminates the need for frequent read and write operations to disk. This makes Spark particularly suitable for real-time data analytics where speed is crucial.

MLlib is Spark's machine learning library that includes various algorithms for classification, regression, clustering, and collaborative filtering. It allows seamless integration of machine learning algorithms with Spark's data processing capabilities, enabling the development of real-time analytics applications.

## 5. Case Studies

### 1) *Stock Market Prediction*
Stock market prediction involves forecasting future stock prices using historical data. By utilizing real-time data processing frameworks like Spark and machine learning algorithms such as SVMs and DNNs, financial institutions can predict stock price movements with higher accuracy. For instance, a real-time prediction system can continuously ingest data from financial APIs, update models, and provide actionable insights for traders [4].
- *Methodology:* A real-time stock market prediction system typically involves the following steps, data ingestion, data preprocessing, model training, and real-time analysis.
- *Result:* Implementing real-time stock market prediction systems has shown improved accuracy in predicting stock movements, enabling traders to make informed decisions and optimize their trading strategies.

### 2) *Credit Card Fraud Detection*
Credit card fraud detection is another critical application of real-time data processing. Implementing a DNN-based autoencoder model allows for the continuous monitoring of transaction data streams. This model can learn the normal behavior of cardholders and detect deviations that may indicate fraud. Financial institutions can thereby minimize losses and enhance security by promptly flagging suspicious transactions for further investigation.
- *Methodology:* A real-time stock market prediction system typically involves the following steps, data collection, feature engineering, model training, real-time monitoring.
- *Result:* Real-time fraud detection systems using DNNs have demonstrated high accuracy in identifying fraudulent transactions, significantly reducing the incidence of fraud and associated financial losses.

## 6. Conclusion

Real-time data processing, empowered by machine learning, is revolutionizing the financial industry. This dynamic combination allows for faster and more precise decision-making, providing financial institutions with a competitive edge. The integration of robust frameworks like Apache Hadoop and Spark with advanced ML algorithms, such as Support Vector Machines (SVMs) and Deep Neural Networks (DNNs), offers significant advantages. These frameworks enable the efficient handling and processing of vast amounts of data in real-time, making it possible to extract valuable insights and make informed decisions promptly.

One of the key advantages of real-time data processing using machine learning in finance is the ability to identify and capitalize on market opportunities. By leveraging ML algorithms, financial institutions can analyze market data, including stock prices, economic indicators, and news sentiment, in real-time to identify potential trading opportunities. This enables them to make timely investment decisions, potentially leading to increased profitability. Additionally, real-time data processing can enhance risk management by allowing financial institutions to monitor their portfolios and identify potential risks in real-time. This proactive approach to risk management can help mitigate losses and protect investor capital.

Future research in the field of real-time data processing using machine learning should focus on improving model accuracy and scalability. As the volume and complexity of financial data continue to grow, ML models need to be developed that can handle these challenges while maintaining high levels of accuracy. Additionally, research should address the challenges associated with real-time data processing, such as latency and data consistency. By overcoming these challenges and further refining ML algorithms, financial institutions can harness the full potential of real-time data processing to make even more informed decisions and achieve even greater success.

## References

[1] Ruta, Dymitr. 2014. "Automated Trading with Machine Learning on Big Data." *In 2014 IEEE International Congress on Big Data*, 824–30. Anchorage, AK, USA: IEEE.

[2] Perera, Srinath, and Suhothayan Sriskandarajah. 2015. "Solution Patterns for Realtime Streaming Analytics." *In Proceedings of the 9th ACM International Conference on Distributed Event-Based Systems*, 334–35. Oslo, Norway: ACM.

[3] Shi, Xiang, Peng Zhang, and Samee U. Khan. 2016. "Quantitative Data Analysis in Finance." *In Handbook of Financial Data and Risk Information I*,, 441–64. Cham: Springer International Publishing.

[4] Kranthi Sai Reddy, V. 2018. "Stock Market Prediction Using Machine Learning." International Research *Journal of Engineering and Technology (IRJET)* 5 (10): 1032–35.

[5] Abakarim, Youness, Mohamed Lahby, and Abdelbaki Attioui. 2018. "An Efficient Real Time Model For Credit Card Fraud Detection Based On Deep Learning." *In Proceedings of the 2018 International Conference on Smart Information and Technology Applications*, 1–7. Rabat, Morocco: ACM.

[6] Habeeb, Riyaz Ahamed Ariyaluran, Fariza Nasaruddin, Abdullah Gani, Ibrahim Abaker Targio Hashem, Ejaz Ahmed, and Muhammad Imran. 2019. "Real-Time Big

Data Processing for Anomaly Detection: A Survey." *International Journal of Information Management* 45: 289–307.

[7] Rundo, Francesco, Francesca Trenta, Agatino Luigi di Stallo, and Sebastiano Battiato. 2019. "Machine Learning for Quantitative Finance Applications: A Survey." *Applied Sciences* 9 (24): 5574.

[8] Lv, Dongdong, Shuhan Yuan, Meizi Li, and Yang Xiang. 2019. "An Empirical Study of Machine Learning Algorithms for Stock Daily Trading Strategy." *Mathematical Problems in Engineering* 2019 (April): 7816154.

[9] Peng, Zhihao. 2019. "Stocks Analysis and Prediction Using Big Data Analytics." *In 2019 International Conference on Intelligent Transportation, Big Data & Smart City* (ICITBS), 309–12. Dalian, China: IEEE.

[10] Zhang, Peng, and Yuanyuan Gao. 2016. "QuantCloud: Big Data Infrastructure for Quantitative Finance on the Cloud." *IEEE Transactions on Big Data*.