

# Practical Directions for the Integration of Financial Big Data and Quantitative Risk Control

Xun Hong

Faculty of Business, Macau University of Science and Technology, Macao, 999078, China

## ABSTRACT

With the all-round development of digital technology in recent years, so too have people's risk-bearing habits changed. Systematically explore the practical directions of applying large-scale financial big data and high-end quantitative risk control frameworks in this paper. Utilise a large number of organised and unorganised alternative data to improve the accuracy of prediction and operation. Introduce new kinds of high-performance, frequently updated quantitative algorithms and move away from the old credit-scoring system. The first few practical applications are shown below: real-time fraud detection, dynamic credit assessment and comprehensive market risk monitoring. This paper will also discuss the issues that have arisen due to such a combination, such as data privacy laws, algorithmic transparency, and a large-scale computational infrastructure. In short, now that large-scale big data analysis and high-end quantitative methods have been introduced, all-encompassing risk management systems for large banks can be constructed to face the turbulence of the economy, credit losses due to defaults can be reduced, and continuous profit generation in the era of deep globalisation can be achieved.

## KEYWORDS

Financial big data; Quantitative risk control; Machine learning; Alternative data; Credit scoring; Algorithmic integration

## 1. INTRODUCTION

At present, the financial system has undergone changes due to a new-age technological wave and the rapid development of digital transaction data and computing power. Due to the old bureaucratic system recently, many financial institutions have had to increase the proportion of subjective and manual credit assessment for enterprises and individuals. Although these traditional ways have provided a certain level of operating stability during a time of relatively weak economic fluctuation, they are far too inadequate for dealing with the extremely high complexity, high speed and large scale of the world's financial system today. Given that human analysts cannot promptly process the large volume of diverse data, a high-precision, rule-based risk control system has begun to be developed.

Quantitative risk control is an all-or-nothing improvement in financial safety that uses high-level mathematical models, deep statistical analysis and complex algorithmic systems to determine the amount of money that could be lost in the event of a risk, and takes corresponding strong measures to reduce this amount. In other words, even if a quantitative model is good, it will not be practical if the foundation data for its construction are poor. In recent years, with the widespread deployment of digital payment gateways, global e-commerce platforms and ubiquitous mobile banking applications, an exceptionally large and completely new reservoir of financial data has been generated and is now classified as big data. This immense amount of new information from the internet includes not only the old structured financial data that has always been available, but also huge quantities of

unstructured alternative data, such as real-time location information, complex behavioural biometrics, and rich social network connections.

The very top of today's financial engineering is now a very high-level application of large-scale financial big data in combination with advanced quantitative risk control models. Systematically feed a large number of very diverse datasets into dynamic, continuously learning algorithms to actively discover extremely subtle and deep-seated patterns of imminent default or highly complex fraudulent behaviour that would be entirely hidden from traditional analysis methods. The two are working together to shift the entire mode of risk management from the old, passive-reactive and historically focused type to a new, proactive-predictive operating system. Due to the good integration of its branches, the bank will have an advantage and be able to extend loans to undeveloped areas safely; at the same time, funds will be used more efficiently and the risk of a widespread financial crisis will be reduced.

The first purpose of this all-encompassing research paper is to conduct an in-depth study of very practical and easily implemented paths for successfully carrying out the integration of financial big data and advanced quantitative risk control. Section 2 will first present the theoretical foundation of Big Data architecture and then introduce the development of high-end data mining from traditional statistical learning. Section 3: Development of quantitative credit scoring models and their application in algorithms. Section 4 introduces the specific application paths for the above ideas, such as real-time streaming analysis and alternative data underwriting. Section 5 concludes with an all-encompassing statement on the future trend of digital financial risk management. Through this systematical investigation, the purpose of this study is to offer financial managers, quantitative analysts and regulatory authorities a very practical and scientific plan for constructing high-resilience digital financial infrastructure.

## **2. THEORETICAL FOUNDATIONS OF BIG DATA IN FINANCE**

To know how to apply this rich wealth of financial data effectively in recent years, we first need to learn about its internal structure. At the bottom of this technological system is an application of advanced statistical learning and very complex data mining. In the past, the limits of traditional relational databases had restricted the scope of financial risk analysis and were unable to handle large amounts of unstructured or irregularly structured numerical data. However, current financial data is inherently very disorganized and highly diverse in form. At present, extremely sophisticated data mining tools can be employed in practice to extract valuable and actionable predictive intelligence from the vast amount of messy data stored in large-scale data lakes with great success [4].

The foundation of linear regression has been further extended through the development of other models of machine learning and now has high accuracy. The foundation of statistical learning clearly indicates that as the absolute quantity and basic dimensionality of the training data increase exponentially, highly flexible and strongly non-linear algorithms are absolutely required to successfully capture the true underlying variance without suffering from severe mathematical overfitting [5]. Given the specific situation of financial risk, the behaviour of borrowers often changes in a large-scale and unpredictable manner due to abrupt economic downturns, and thus traditional models frequently break down. Therefore, many quantitative researchers have been using large-scale ensemble methods to stabilise their extremely complex risk forecasts. For example, a strategic deployment of high-end random forest algorithms has successfully combined the outputs of thousands of independent decision trees to substantially reduce the overall model variance and improve the absolute robust predictive accuracy of the top-level risk management framework [3].

In addition, the large-scale integration of big data has led to an all-new approach for the initial data exploration. First build a full set of exploratory data analysis (EDA), then build a high-quality predictive model. The first is a general exploration step to uncover all kinds of serious data problems,

hidden structures and large areas of missing data that might make the quantitative model invalid. Deeply embrace large-scale exploratory techniques that have been specifically designed for incredibly vast auditing and financial datasets to ensure the fundamental mathematical integrity of the raw input layer [9]. A first-class discipline in preparing the raw material for model construction is required, and otherwise all-weather high-end quantitative risk models cannot be built reliably. Without the above-mentioned large-scale foundation discipline, the whole complex superstructure of big data risk control will inevitably be subject to catastrophic algorithmic failure.

### **3. EVOLUTION OF QUANTITATIVE RISK MODELS**

With the long history of quantitative risk models, so too has the development of consumer credit scoring gone through a period of expansion. At the end of the 20th century, the first use of simple statistical division in the retail credit industry completely changed how it was run. The first round of credit-scoring systems was well-structured and allowed for the quick, full-scale, and systematic evaluation of many millions of loan applications by banks; at the same time, this reduced human error and a huge amount of manual work in credit approval [1]. Basic discriminant analysis and standard logistic regression were employed to determine which historical variables of a financial institution's future financial delinquency are most closely associated with it; among these, previous loan defaults and a basic income ratio were identified [2].

Although the first statistical studies on credit and behaviour scoring provided an early foundation for predicting large-scale financial risks, their basic structure was also very limited due to a small number of organised financial indicators [2]. With the deep digitalisation of the world's economy, a large number of people around the world, especially young adults and new immigrants, have never had credit histories before. These 'thin-file' consumers were unfairly penalized by the old scoring system and thus lost many very profitable loans that should have been approved. To actively address the severe structural deficiencies mentioned above, the leading quantitative scholars have completely pioneered the aggressive application of highly complex machine learning algorithms that are suitable for processing extremely large and highly irregular datasets [6].

Systematically collect a large number of all kinds of unconventional alternative data, such as very detailed digital transaction records, high-volume mobile application usage behaviour, and highly complex behavioural biometrics, and build exceptionally accurate and robust credit risk models for previously unrateable groups of consumers using machine learning [6]. Therefore, a large number of analyses have been carried out, and they are relatively comprehensive. Contemporary financial institutions need to use all kinds of high-end measurement tools and various sophisticated enterprise software to continually observe, conduct extensive backtesting, and strictly verify the mathematical validity of their very large-scale and complicated credit portfolios. A very strict and highly continuous validation cycle has been introduced to ensure that the heavily data-driven machine learning algorithms do not deviate from the fundamental economic reality; thus, the core stability and overall substantial profits of the institutional lending system will be rigorously maintained [8].

### **4. PRACTICAL DIRECTIONS AND IMPLEMENTATION STRATEGIES**

Large-scale modern financial institutions are now using integrated big data and quantitative risk control systems in their daily operations, so high-precision and highly strategic operation is required. One of the indispensable, very practical applications of this large-scale integration is a high-performance, all-weather, real-time streaming risk engine. In the heavily traditional legacy system, the risk assessment was generally a very slow, large-scale batch process that had to be completed at the end of the business day. However, given that now we live in an era of instantaneous digital payments and extremely fast algorithmic high-frequency trading, a delay in risk assessment is virtually the same as no risk assessment at all.

The demand from current financial technology frameworks is to carry out high-volume, high-frequency mathematical calculations of continuous data streams in real time. Due to the profound research problems at present in artificial intelligence for modern financial technology risk management, there is an urgent and essential need to build highly stable, extremely high-speed algorithmic structures that can instantly detect large-scale, deeply intricate fraud networks and serious liquidity risks in real time. Directly embed extremely lightweight and high-performance machine learning models within the core transactional payment infrastructure of modern financial institutions to identify and prevent substantial capital outflows at an earlier stage, thus avoiding significant financial losses [10].

In addition, the extended incorporation of various other big data sources is also starting to show some results in the microeconomics of the consumer loan market. A number of economic studies have shown that by extending and optimising high-precision credit scoring models over time, the total cost of capital has been reduced substantially, credit has been extended to more segments of the population, and the large lending institutions have achieved stable financial development [7]. A very positive, large-scale economic feedback loop will be driven forward more rapidly by continuous updates to the precise limits of acceptable financial risk through big data analysis. Finally, the institutions will have to face a difficult time navigating the large-scale and extremely complex regulatory system for algorithmic underwriting. The world's financial authorities have begun to require the full disclosure of algorithms and stringent adherence to the principle of non-discrimination. Therefore, the most successful and practical implementation is one that combines very high predictive power with high transparency in explanation to achieve both exceptionally large commercial profits and strict adherence to all laws simultaneously [7, 10]. In addition, the large-scale deployment of these extremely high-integration big data risk frameworks fully supports the deep and extremely complex process of macroeconomic stress testing.

In practice, the application of this implementation should not be confined to the technical deployment of predictive models. An overall plan for data governance in the entire risk and control system should cover the collection and processing of data, selection of relevant features, model performance tracking and periodic reevaluation. Financial big data often contains noise, missing values and potential sampling bias; therefore, standard procedures need to be established to ensure the reliability of the input data before training the model [9]. Continuously test the deployed model based on changes in the market; that is to say, behaviour patterns and fraud characteristics among borrowers are likely to shift rapidly during economic downturns. Credit risk analysis has also found that to prevent model drift and keep the portfolio stable, backtesting and performance monitoring need to be carried out [8]. At the same time, explainable models should also be applied in the risk control system. A high-performance machine learning model is used to increase the accuracy of the prediction; however, detailed reasons for the credit decision should still be provided to regulators, customers, and other divisions of the company. Therefore, instead of substituting traditional risk managers with algorithms, a new model for the years can be built by combining machine learning-based early warning and the final decision-making authority of human analysts in cases of complex or serious risks [10].

## **5. CONCLUSION**

At present, highly aggressive, all-encompassing systems of large-scale financial big data and advanced quantitative risk control mechanisms have become indispensable requirements for the operation of modern financial institutions. As the above comprehensive, exceptionally deep-seated systematic analysis has shown, relying solely on old, extremely rigid legacy underwriting models is fundamentally out of date in a deeply hyper-connected and incredibly massive digitalised global economy. Due to a very deep transition to an all-automated, high-machine-learning-driven risk architecture, it is now feasible to handle the sheer complexity and scale of modern capital markets safely with extremely large-scale mathematical computations and rapid processing by deep learning.

High-end data mining and very complex statistical learning algorithms are employed to process a large number of vast oceans of unstructured alternative data and identify extremely subtle but highly profitable predictive signals that traditional models have been entirely unable to detect. Credit scoring has developed substantially over time, moving from relatively simple, high-constraint linear regression models to very stable and complex non-linear ensemble frameworks today, and as a result, banks are now able to safely extend necessary credit to many unbanked people around the world. Furthermore, the excellent case of a fully real-time, high-streaming risk engine offers robust defense against the large-scale and continually evolving risks from complex digital fraud and severe systemic liquidity crises.

In the future, as the growth rate of global computing infrastructure will be higher and higher, the core function of large-scale data collection, analysis and intelligent response by banks will be determined by the quality and efficiency of this infrastructure. However, given that this extensive technological capability must be continuously and profoundly balanced with the imperative for very strict model verification, deep algorithmic openness, and full compliance with all-around global data privacy laws. Financial institutions that can perfectly master this extremely delicate and complex technological balance will not only reduce their absolute exposure to catastrophic financial defaults at the fundamental level but will also lead the vast global financial industry into an extraordinarily secure, profoundly prosperous, and highly technologically advanced digital future sustainably.

## REFERENCES

- [1] Hand, D. J., & Henley, W. E. (1997). Statistical classification methods in consumer credit scoring: A review. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 160(3), 523–541. <https://doi.org/10.1111/1467-985X.00078>
- [2] Thomas, L. C. (2000). A survey of credit and behavioural scoring: Forecasting financial risk of lending to consumers. *International Journal of Forecasting*, 16(2), 149–172. [https://doi.org/10.1016/S0169-2070\(99\)00045-4](https://doi.org/10.1016/S0169-2070(99)00045-4)
- [3] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
- [4] Witten, I. H., & Frank, E. (2005). *Data mining: Practical machine learning tools and techniques* (2nd ed.). Morgan Kaufmann.
- [5] Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
- [6] Khandani, A. E., Kim, A. J., & Lo, A. W. (2010). Consumer credit-risk models via machine-learning algorithms. *Journal of Banking & Finance*, 34(11), 2767–2787. <https://doi.org/10.1016/j.jbankfin.2010.06.009>
- [7] Einav, L., Jenkins, M., & Levin, J. (2013). The impact of credit scoring on consumer lending. *The RAND Journal of Economics*, 44(2), 249–274. <https://doi.org/10.1111/1756-2171.12020>
- [8] Baesens, B., Roesch, D., & Scheule, H. (2016). *Credit risk analytics: Measurement techniques, applications, and examples in SAS*. John Wiley & Sons.
- [9] Lessmann, S., Baesens, B., Seow, H. V., & Thomas, L. C. (2015). Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. *European Journal of Operational Research*, 247(1), 124–136. <https://doi.org/10.1016/j.ejor.2015.05.030>
- [10] Giudici, P. (2018). Fintech risk management: A research challenge for artificial intelligence in finance. *Frontiers in Artificial Intelligence*, 1, 1. <https://doi.org/10.3389/frai.2018.00001>