Volume 2, Issue 2, 2025

Print ISSN:3007-6951 Online ISSN: 3007-696X

# Journal of Trends in Financial and Economics



**Copyright© Upubscience Publisher** 

# Journal of Trends in Financial and Economics

Volume 2, Issue 2, 2025



**Published by Upubscience Publisher** 

#### **Copyright**<sup>©</sup> The Authors

Upubscience Publisher adheres to the principles of Creative Commons, meaning that we do not claim copyright of the work we publish. We only ask people using one of our publications to respect the integrity of the work and to refer to the original location, title and author(s).

Copyright on any article is retained by the author(s) under the Creative Commons Attribution license, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Authors grant us a license to publish the article and identify us as the original publisher.

Authors also grant any third party the right to use, distribute and reproduce the article in any medium, provided the original work is properly cited.

Journal of Trends in Financial and Economics Print ISSN: 3007-6951 Online ISSN: 3007-696 Email: info@upubscience.com Website: http://www.upubscience.com/

# **Table of Content**

DEEP LEARNING-BASED CREDIT RISK MODELING: ADDRESSING DATA IMBALANCE AND INVARIANCE LiangYu Chen, Hao Lin*	1-8
ENHANCING GLOBAL FINANCE: A BLOCKCHAIN-BASED SOLUTION FOR EFFICIENT AND COST-EFFECTIVE CROSS-BORDER PAYMENTS Gang Min	9-12
BANK CUSTOMER DEPOSIT PRODUCT PURCHASE ANALYSIS AND PREDICTION YuanJia Guo	13-22
FEEDBACK TRADING, DELAYED ARBITRAGE, AND ASSET PRICE BUBBLES: THEORETICAL AND EMPIRICAL STUDIES BASED ON DIFFERENT CYCLES ShiQi Wang*, Xin Pan, Cong Chen	23-26
LITERATURE REVIEW OF SELF-IMPROVEMENT PRODUCTS DongPing Zheng*, Pei Feng, JinFu Huang	27-31
CHARACTERISTICS, ISSUES, AND COUNTERMEASURES OF CHINA'S DIGITAL ECONOMY DEVELOPMENT ZhenJie Guo	32-35
<b>STOCK PRICE RESEARCH BASED ON ARIMA-GARCH-LSTM HYBRID MODEL</b> ChaoYan Wei*, LanLan Li, PangLeYi Chen, MeiHui Huang, HuiLin Wei, KunYao Yao, XuYang Wang, Xin Ya, ChaoHai Wei	36-43
THE APPLICATION OF DEEP REINFORCEMENT LEARNING IN ASSET ALLOCATION: A THEORETICAL FRAMEWORK AND EMPIRICAL ANALYSIS ZiLin Zhou	44-50
"TARO MEETS NEW FOOD GENERATION" — MARKET INVESTIGATION OF DEEP-PROCESSED TARO PRODUCTS IN LIPU RongJin Li	51-72
MACHINE LEARNING APPROACHES FOR ACCURATE DEMAND FORECASTING IN SUPPLY CHAIN MANAGEMENT	73-78

Liu Zhen, Yang Lin\*

# DEEP LEARNING-BASED CREDIT RISK MODELING: ADDRESSING DATA IMBALANCE AND INVARIANCE

LiangYu Chen, Hao Lin\*

School of Management, Sun Yat-sen University, Guangzhou 510275, Guangdong, China. Corresponding Author: Hao Lin, Email: 627709133@qq.com

**Abstract:** Credit risk modeling plays a crucial role in financial decision-making, helping lenders assess the likelihood of default and optimize lending strategies. Traditional credit risk assessment models, including logistic regression and decision trees, often struggle with data imbalance and invariance issues, leading to biased risk predictions and reduced generalization. The rapid advancement of deep learning (DL) techniques has introduced more sophisticated models capable of learning complex credit risk patterns. However, most DL-based credit scoring models still suffer from class imbalance in default prediction and fail to maintain fairness and stability across different demographic and economic conditions.

This study proposes a DL-based credit risk modeling framework designed to address data imbalance through advanced resampling techniques and generative modeling, while also incorporating adversarial learning to improve model invariance across diverse borrower segments. The proposed framework utilizes autoencoders, generative adversarial networks (GANs), and cost-sensitive learning techniques to enhance risk assessment accuracy while reducing bias. Additionally, domain adaptation techniques are introduced to ensure that the model remains robust across different financial environments.

Experiments on real-world credit datasets demonstrate that the proposed framework significantly improves credit risk prediction accuracy, enhances model fairness, and reduces sensitivity to class imbalance compared to traditional credit scoring approaches. The findings highlight the importance of integrating data-centric augmentation techniques with fairness-aware deep learning to improve the reliability of credit risk modeling in modern financial applications.

**Keywords:** Credit risk modeling; Deep learning; Data imbalance; Invariance; Fairness; Generative models; Adversarial learning

#### **1 INTRODUCTION**

Credit risk modeling is a critical component of financial decision-making, enabling lenders and financial institutions to assess the likelihood of borrower default [1]. An accurate credit risk assessment system ensures that loans are granted to creditworthy individuals while minimizing the risk of financial losses. Traditional credit risk models, such as logistic regression (LR) and decision trees (DTs), have been widely used for decades to predict borrower default probabilities. However, these models rely on manually engineered features and linear relationships, making them limited in capturing complex credit risk patterns present in real-world financial data [2].

One of the major challenges in credit risk modeling is class imbalance, where the proportion of borrowers who default on their loans is significantly lower than those who repay their debts [3]. This imbalance often leads to biased credit scoring models, where machine learning algorithms favor the majority class (non-defaulters) while misclassifying minority class samples (defaulters). Traditional machine learning (ML) models trained on imbalanced datasets tend to exhibit poor recall for high-risk borrowers, leading to inaccurate risk estimation and suboptimal lending decisions [4]. Addressing class imbalance is crucial to improving the performance and fairness of credit risk models.

Another key issue in credit risk modeling is invariance, referring to a model's ability to maintain consistent predictive performance across different demographic and economic conditions [5]. Many existing credit scoring models exhibit distributional bias, where certain borrower groups receive systematically different credit risk scores due to underlying imbalances in the dataset. This raises ethical and regulatory concerns, as unfair credit assessments can lead to discriminatory lending practices and potential legal consequences. Improving invariance in credit risk models is essential to ensure fairness and regulatory compliance in financial decision-making [6].

Recent advancements in deep learning (DL) have introduced more sophisticated methods for credit risk assessment, enabling models to learn non-linear and hierarchical representations from large-scale financial datasets [7]. Techniques such as artificial neural networks (ANNs) and long short-term memory (LSTM) networks have been employed to improve credit risk prediction by leveraging complex borrower-lender interaction patterns [8]. However, despite their advantages, DL-based credit scoring models still face significant challenges related to class imbalance and distributional invariance.

This study proposes a DL-based credit risk modeling framework designed to address these limitations by integrating data augmentation, generative adversarial networks (GANs), and adversarial training techniques. The framework enhances credit risk assessment by improving recall for high-risk borrowers and ensuring model stability across different borrower distributions. By incorporating autoencoders (AEs) for feature learning, cost-sensitive learning techniques for imbalanced classification, and domain adaptation strategies for fairness enhancement, the proposed approach improves both the accuracy and fairness of credit risk predictions.

The framework is evaluated using real-world credit risk datasets, demonstrating its effectiveness in mitigating data imbalance, reducing discriminatory biases, and improving model generalization. The findings highlight the necessity of combining data-centric resampling techniques with fairness-aware deep learning approaches to develop robust and equitable credit risk assessment systems for modern financial applications.

#### **2 LITERATURE REVIEW**

Credit risk modeling has been extensively studied in financial research, with traditional methods focusing on statistical and ML-based techniques to predict borrower default probabilities [9]. While LR and DTs have been widely used due to their interpretability and regulatory acceptance, these models often struggle with complex borrower behaviors and non-linear credit risk patterns [10]. The advent of DL has provided new opportunities for financial institutions to enhance risk assessment accuracy by leveraging large-scale datasets and learning intricate borrower-lender interactions. However, DL models also present challenges, particularly in addressing class imbalance and ensuring invariance in credit risk classification [11].

Class imbalance remains a fundamental issue in credit risk modeling [12]. Many credit datasets contain a significantly lower proportion of defaulters compared to non-defaulters, leading to skewed model predictions. Traditional ML models trained on such imbalanced datasets tend to favor the majority class, resulting in high overall accuracy but poor recall for high-risk borrowers. To address this issue, various resampling techniques, including oversampling methods such as Synthetic Minority Over-sampling Technique (SMOTE) and undersampling methods, have been introduced to balance class distributions. Cost-sensitive learning has also been explored as an alternative, assigning higher misclassification penalties to the minority class to improve model sensitivity to defaulters [13-16]. While these approaches mitigate imbalance to some extent, they do not fully capture the complex borrower relationships that contribute to default risk [17].

DL-based models, particularly ANNs and LSTMs, have demonstrated superior predictive performance in credit risk modeling by capturing hierarchical borrower representations and sequential financial behaviors [18]. However, their reliance on large, imbalanced datasets can lead to biased risk predictions [19]. One solution to this challenge is the use of GANs to generate synthetic borrower profiles, expanding the representation of high-risk borrowers in training data. Additionally, autoencoders have been utilized for feature extraction, improving model generalization and reducing overfitting in imbalanced credit datasets. These data augmentation techniques enhance the robustness of DL models, enabling them to learn better from minority class samples and improving recall for high-risk borrowers [20].

Another key challenge in credit risk modeling is ensuring invariance across different borrower demographics and economic conditions [21]. Many traditional credit scoring models exhibit distributional bias, where borrowers from specific demographic or socioeconomic backgrounds receive systematically lower or higher credit scores due to underlying dataset imbalances. Addressing this issue requires adversarial learning strategies that enforce fairness constraints during model training [22-26]. Domain adaptation techniques, such as adversarial domain alignment, have been introduced to ensure that risk predictions remain stable across different borrower groups [8]. These techniques help mitigate the effects of biased credit assessments, ensuring that DL-based credit risk models are both accurate and fair [27-29].

Despite the advancements in DL-based credit risk modeling [30,31], there remain open questions regarding explainability and regulatory compliance. Financial institutions are required to provide transparent justifications for credit decisions, which can be challenging when using complex neural networks. Research into explainable AI methods, including feature attribution techniques and interpretable DL architectures, aims to bridge this gap by improving model interpretability while maintaining high predictive performance [32].

This study builds on these advancements by integrating GAN-based data augmentation, adversarial learning for fairness enhancement, and autoencoder-driven feature learning into a unified DL-based credit risk modeling framework. The proposed approach aims to improve credit risk assessment by addressing class imbalance and ensuring distributional invariance, enhancing both model performance and fairness in credit lending decisions. The next section outlines the methodology used to implement and evaluate the proposed framework.

#### **3 METHODOLOGY**

#### 3.1 Data Preprocessing and Feature Engineering

Credit risk modeling relies on high-dimensional financial datasets containing diverse borrower attributes, transaction histories, and economic indicators. Ensuring data quality and consistency is crucial before training DL-based models. The first step in preprocessing involves handling missing values, which are common in credit datasets due to incomplete borrower information. Missing values are addressed using imputation techniques, including mean imputation for numerical features and mode imputation for categorical attributes. More complex techniques such as k-nearest neighbors imputation and autoencoder-based imputations are also employed to enhance data consistency.

Outlier detection is another important preprocessing step, as extreme values in financial data can bias model predictions. Borrowers with unrealistic credit scores, income levels, or transaction frequencies are identified using statistical anomaly detection methods, including interquartile range and Mahalanobis distance. Data normalization is applied to ensure that all numerical variables are scaled appropriately, preventing models from being influenced by large-magnitude financial values. Feature engineering plays a crucial role in improving the predictive power of credit risk models. Derived financial attributes, such as debt-to-income ratio, revolving credit utilization, and historical loan repayment behavior, are extracted to provide more informative borrower representations. Temporal features are also introduced by analyzing borrower behavior over different time windows, allowing the model to capture financial trends and repayment consistency. Dimensionality reduction techniques, including principal component analysis and autoencoder-based feature selection, are applied to eliminate redundant information while preserving critical credit risk indicators.

#### 3.2 Handling Class Imbalance with Generative and Cost-Sensitive Approaches

Class imbalance is one of the most significant challenges in credit risk modeling, where the number of default cases is substantially lower than non-default cases. Training models on imbalanced datasets leads to biased predictions, where the model learns to favor the majority class, resulting in high precision but low recall for high-risk borrowers. To address this issue, multiple techniques are integrated into the proposed framework to improve model sensitivity to defaulters while maintaining overall classification accuracy.

Resampling techniques, including SMOTE-based oversampling and random undersampling, are employed to rebalance class distributions. While these methods improve recall for the minority class, they can introduce noise and redundancy in training data. To mitigate this, GAN-based data augmentation is implemented, generating synthetic borrower profiles that mimic the statistical properties of real defaulters. The GAN framework consists of a generator that learns to create realistic borrower data and a discriminator that distinguishes between real and synthetic profiles, resulting in more diverse and representative training samples.

Cost-sensitive learning is incorporated into the DL framework to assign higher misclassification penalties for false negatives, ensuring that defaulters are correctly identified. The model optimizes a weighted loss function that prioritizes minimizing the impact of incorrectly classified high-risk borrowers. Adaptive threshold tuning is also applied, where the decision boundary for classifying defaulters is adjusted dynamically based on dataset imbalance ratios. These techniques collectively enhance the model's ability to recognize high-risk borrowers, leading to more reliable credit risk assessments.

#### 3.3 Adversarial Training for Invariance and Fairness Enhancement

Ensuring that the credit risk model remains fair and unbiased across different borrower demographics is essential for regulatory compliance and ethical lending practices. Many traditional credit scoring models exhibit disparities in loan approval rates due to dataset biases, where certain demographic groups receive systematically different risk assessments. To mitigate this issue, adversarial training is integrated into the DL framework to enhance fairness and improve model invariance across different borrower segments.

The adversarial learning framework consists of two competing models: the primary credit risk classifier and an adversary trained to detect disparities in credit risk predictions. During training, the adversary attempts to identify which demographic group a borrower belongs to based on the classifier's predictions. If the adversary successfully distinguishes borrower groups, the classifier is penalized, forcing it to learn risk assessment criteria that are independent of demographic attributes. This process ensures that the credit risk model learns unbiased decision-making rules, reducing the influence of sensitive attributes such as age, gender, or ethnicity.

Domain adaptation techniques are also introduced to ensure that the model maintains stability across different economic conditions and borrower distributions. The model is trained on credit datasets from multiple financial institutions and economic cycles, improving its ability to generalize across varied lending environments. Transfer learning is employed to fine-tune the model on new borrower datasets, ensuring that it remains adaptable to evolving financial conditions without requiring full retraining.

#### 3.4 Model Training, Optimization, and Evaluation Metrics

The proposed DL-based credit risk framework is implemented using a deep feedforward neural network architecture combined with LSTMs for capturing sequential borrower behaviors. The model is trained using a hybrid loss function that balances classification accuracy, fairness constraints, and cost-sensitive penalties for high-risk borrowers. The optimization process is conducted using the Adam optimizer with dynamic learning rate adjustments, ensuring that the model converges efficiently without overfitting.

Hyperparameter tuning is performed using Bayesian optimization and grid search techniques to identify optimal model configurations, including the number of hidden layers, activation functions, dropout rates, and regularization parameters. To further improve model generalization, ensemble learning techniques such as stacked neural networks and boosting-based deep learning are explored. These techniques enhance prediction robustness by combining multiple neural network architectures to minimize variance and bias.

The evaluation of the proposed framework is conducted using multiple performance metrics to ensure a comprehensive assessment of its effectiveness. Precision, recall, and F1-score are used to measure classification accuracy, while AUC-ROC curves assess the model's ability to distinguish between defaulters and non-defaulters. The fairness of credit risk assessments is evaluated using disparate impact ratios and equality of opportunity metrics, ensuring that loan approval decisions do not disproportionately disadvantage any demographic group.

Scalability and computational efficiency are also key considerations, as credit risk models must be capable of processing large borrower datasets in real time. The model's inference speed, memory consumption, and scalability across different hardware environments are analyzed to determine its suitability for deployment in large-scale financial systems. The results confirm that the proposed framework achieves superior credit risk assessment accuracy, fairness, and computational efficiency compared to traditional credit scoring models.

By integrating advanced data augmentation, fairness-aware adversarial learning, and cost-sensitive classification techniques, the proposed DL-based credit risk framework effectively addresses data imbalance and invariance challenges, providing financial institutions with a more accurate, ethical, and scalable approach to credit risk modeling. The next section presents experimental results and discusses the impact of combining these techniques on credit risk prediction performance.

#### **4 RESULTS AND DISCUSSION**

#### 4.1 Credit Risk Prediction Accuracy and Model Performance

The effectiveness of the proposed DL-based credit risk modeling framework was evaluated using real-world credit datasets, comparing its performance with traditional ML models such as LR, DTs, and gradient boosting classifiers. The evaluation metrics included precision, recall, F1-score, and AUC-ROC, providing a comprehensive assessment of the model's ability to classify borrowers accurately. The results demonstrated that integrating deep learning techniques, particularly generative models and adversarial training, significantly enhanced credit risk assessment performance.

The model achieved superior recall for high-risk borrowers compared to conventional ML-based credit scoring models, which tend to misclassify defaulters due to class imbalance. By incorporating GAN-based data augmentation and cost-sensitive learning, the model successfully improved recall for defaulters by 22% while maintaining high precision, ensuring that non-defaulters were not incorrectly classified as high-risk borrowers. The improvement in recall is particularly beneficial for financial institutions, as it reduces the likelihood of misclassifying potential defaulters, thereby minimizing loan default risks.

The overall AUC-ROC score of the proposed model was consistently higher across different credit datasets, indicating better discriminatory power between default and non-default cases. Compared to traditional ML models, which often fail to capture non-linear borrower patterns, the proposed DL-based framework demonstrated higher adaptability to complex credit risk scenarios, allowing it to generalize effectively across diverse borrower groups.

Figure 1 presents a comparative analysis of credit risk prediction accuracy across different models, illustrating the performance advantages of the proposed DL-based framework.



Figure 1 Credit Risk Prediction Accuracy Comparison

#### 4.2 Impact of Generative Modeling and Resampling Techniques on Class Imbalance

Addressing class imbalance is crucial for improving the effectiveness of credit risk models, as imbalanced datasets often lead to biased predictions that favor the majority class. Traditional resampling techniques, such as SMOTE, have been widely used to balance class distributions, but they often introduce synthetic noise, leading to model overfitting. The proposed framework integrates GAN-based data augmentation and cost-sensitive learning to enhance model robustness against class imbalance.

The impact of generative modeling was assessed by training the credit risk classifier on datasets enhanced with synthetic borrower profiles generated by GANs. The results showed that the use of synthetic defaulter data significantly improved model generalization, particularly in cases where default rates were extremely low. Unlike oversampling

methods that simply duplicate minority class instances, GAN-based augmentation introduced diverse yet realistic borrower profiles, allowing the model to learn more representative credit risk features.

In addition to generative modeling, cost-sensitive learning played a crucial role in optimizing the classification process. By assigning higher misclassification penalties for false negatives, the framework effectively reduced the number of misclassified defaulters without substantially increasing false positives. This approach ensures that the model prioritizes risk management without over-penalizing creditworthy borrowers.

Figure 2 illustrates the impact of GAN-based data augmentation and cost-sensitive learning on class imbalance reduction, highlighting the effectiveness of these techniques in improving borrower risk classification.



#### 4.3 Fairness and Invariance Across Borrower Demographics

Ensuring fairness and invariance in credit risk modeling is critical for regulatory compliance and ethical lending practices. Many traditional credit scoring models exhibit biases related to demographic attributes, where specific borrower groups receive systematically different risk scores due to historical imbalances in training data. The proposed framework incorporates adversarial training to mitigate these biases, ensuring that the credit risk model remains invariant across different borrower segments.

To evaluate fairness, disparate impact ratios and equality of opportunity metrics were computed for different borrower groups, analyzing whether the model disproportionately assigned higher risk scores to specific demographics. The results confirmed that the adversarial learning framework significantly reduced bias in risk assessments, ensuring that credit scores were determined primarily by financial indicators rather than demographic characteristics.

The adversarial training component effectively prevented the model from learning biased correlations by enforcing fairness constraints during the training process. Additionally, domain adaptation techniques were employed to fine-tune the model on credit datasets from different economic environments, ensuring that its predictions remained stable across varying financial conditions. The findings demonstrate that integrating fairness-aware training strategies not only improves the ethical considerations of credit risk modeling but also enhances model robustness.

Figure 3 presents an analysis of fairness-enhancing adversarial training, showcasing its impact on reducing credit risk classification bias across borrower demographics.



#### 4.4 Computational Efficiency and Scalability of the Credit Risk Model

For credit risk models to be deployed in large-scale financial applications, they must be capable of processing vast amounts of borrower data in real time while maintaining high classification accuracy. The computational performance of the proposed DL-based framework was analyzed by measuring inference speed, memory consumption, and scalability across increasing dataset sizes.

The results showed that the framework maintained low latency inference, processing thousands of borrower applications per second without significant computational overhead. Compared to traditional ML models, which often require manual feature engineering and retraining, the proposed DL framework leveraged parallel processing and GPU acceleration to achieve higher processing efficiency. The use of autoencoder-based feature extraction further reduced dimensionality, optimizing computation while preserving critical risk assessment features.

Scalability tests were conducted by evaluating the framework's performance on datasets containing varying numbers of borrower records, ranging from 100,000 to over 10 million transactions. The model's performance remained stable, confirming its ability to handle large-scale credit risk assessment tasks. The integration of transfer learning further improved scalability, allowing the model to be fine-tuned on new financial datasets without requiring full retraining.

The findings confirm that the proposed framework achieves high computational efficiency and scalability, making it suitable for deployment in financial institutions requiring real-time credit risk assessments for large borrower pools. The results also indicate that the model's ability to adapt to new credit data while maintaining efficiency makes it a viable solution for dynamic financial environments.

Figure 4 presents an evaluation of computational performance and scalability, highlighting the framework's efficiency in large-scale credit risk modeling.



#### Figure 4 Computational Efficiency and Salability

#### **5 CONCLUSION**

Credit risk modeling is a fundamental component of financial risk management, enabling lenders to assess borrower creditworthiness and optimize lending strategies. While traditional ML models such as LR and DTs have been widely used in credit scoring, they exhibit limitations in handling class imbalance and demographic invariance, leading to biased predictions and suboptimal risk assessments. The emergence of DL has introduced more powerful modeling techniques capable of capturing non-linear borrower-lender interactions, yet challenges remain in ensuring fairness, mitigating imbalance, and maintaining model stability across different financial environments.

This study proposed a DL-based credit risk modeling framework designed to address these challenges by integrating GAN-based data augmentation, adversarial fairness learning, and cost-sensitive classification techniques. The results demonstrated that incorporating generative models significantly improved the classification recall for high-risk borrowers, reducing bias caused by class imbalance. Unlike traditional resampling methods, which introduce synthetic noise, GANs provided realistic borrower profiles, allowing the model to learn more representative credit risk features. The introduction of adversarial learning further enhanced model invariance, ensuring that credit risk assessments remained consistent across different borrower demographics.

The experimental evaluation confirmed that the proposed GNN-GAN-adversarial framework outperforms conventional credit risk models in terms of accuracy, fairness, and computational efficiency. The model achieved higher recall for defaulters, improved credit risk classification precision, and reduced false positives. The integration of adversarial learning also reduced disparities in credit scoring across borrower segments, mitigating potential regulatory concerns and ensuring compliance with ethical lending standards.

Scalability and computational efficiency were key factors in determining the practicality of the proposed approach for large-scale credit risk modeling. The model was tested on credit datasets containing millions of borrower records, demonstrating that the proposed DL-based framework maintained stable inference speed and memory efficiency, making it suitable for deployment in financial institutions with high transaction volumes. The results confirmed that the combination of autoencoder-based feature selection, GPU acceleration, and transfer learning techniques allowed the model to process borrower data in real-time while maintaining predictive accuracy.

Despite its advantages, the proposed framework presents certain limitations that warrant future research. One key challenge is the computational cost associated with training GANs and adversarial networks on large-scale financial datasets. While the current implementation ensures scalability, further optimizations such as model compression techniques and distributed training strategies should be explored to enhance efficiency. Another limitation is the explainability of deep learning-based credit risk assessments, as financial institutions require transparency in risk prediction for regulatory compliance. Future work should focus on interpretable AI methods, integrating techniques such as SHAP values and attention-based explanations to improve trust in credit risk predictions.

Additionally, future research should investigate the integration of multi-source financial data, including alternative credit scoring indicators such as transactional behaviors, spending patterns, and social network interactions, to enhance borrower profiling accuracy. Extending the framework to cross-border credit risk modeling will also improve its applicability in global financial markets.

This study highlights the importance of combining DL with fairness-aware learning and data augmentation strategies to develop robust, equitable, and scalable credit risk models. By incorporating generative models, adversarial fairness techniques, and cost-sensitive classification, the proposed framework provides a comprehensive solution to credit risk assessment, ensuring that financial institutions can make fair, data-driven lending decisions while minimizing credit losses and regulatory risks.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCES

- [1] Mashrur A, Luo W, Zaidi N A, et al. Machine learning for financial risk management: a survey. IEEE Access, 2020, 8: 203203-203223.
- [2] Sudjianto A, Zhang A. Model validation practice in banking: A structured approach for predictive models. 2024. DOI: 10.48550/arXiv.2410.13877.
- [3] Addo P M, Guegan D, Hassani B. Credit risk analysis using machine and deep learning models. Risks, 2018, 6(2): 38.
- [4] Bhatore S, Mohan L, Reddy Y R. Machine learning techniques for credit risk evaluation: a systematic literature review. Journal of Banking and Financial Technology, 2020, 4(1): 111-138.
- [5] Hamori S, Kawai M, Kume T, et al. Ensemble learning or deep learning? Application to default risk analysis. Journal of Risk and Financial Management, 2018, 11(1): 12.
- [6] Han X, Yang Y, Chen J, et al. Symmetry-Aware Credit Risk Modeling: A Deep Learning Framework Exploiting Financial Data Balance and Invariance. Symmetry, 2025, 17(3): 341.
- [7] Bussmann N, Giudici P, Marinelli D, et al. Explainable machine learning in credit risk management. Computational Economics, 2021, 57(1): 203–216.

- [8] Mhlanga D. Financial inclusion in emerging economies: The application of machine learning and artificial intelligence in credit risk assessment. International Journal of Financial Studies, 2021, 9(3): 39.
- [9] Bazarbash M. Fintech in financial inclusion: machine learning applications in assessing credit risk. International Monetary Fund, 2019.
- [10] Leo M, Sharma S, Maddulety K. Machine learning in banking risk management: A literature review. Risks, 2019, 7(1): 29.
- [11] Moscato V, Picariello A, Sperlí G. A benchmark of machine learning approaches for credit score prediction. Expert Systems with Applications, 2021, 165: 113986.
- [12] Munkhdalai L, Munkhdalai T, Namsrai O E, et al. An empirical comparison of machine-learning methods on bank client credit assessments. Sustainability, 2019, 11(3): 699.
- [13] Berhane T, Melese T, Seid A M. Performance Evaluation of Hybrid Machine Learning Algorithms for Online Lending Credit Risk Prediction. Applied Artificial Intelligence, 2024, 38(1): 2358661.
- [14] Kaisar S, Sifat S T. Explainable Machine Learning Models for Credit Risk Analysis: A Survey. Data Analytics for Management, Banking and Finance: Theories and Application. Cham: Springer Nature Switzerland, 2023: 51-72.
- [15] Talaei Khoei T, Ould Slimane H, Kaabouch N. Deep learning: systematic review, models, challenges, and research directions. Neural Computing and Applications, 2023, 35(31): 23103-23124.
- [16] Kim E, Lee J, Shin H, et al. Champion-challenger analysis for credit card fraud detection: Hybrid ensemble and deep learning. Expert Systems with Applications, 2019, 128: 214-224.
- [17] Guo H, Ma Z, Chen X, et al. Generating artistic portraits from face photos with feature disentanglement and reconstruction. Electronics, 2024, 13(5): 955.
- [18] Lee Z, Wu Y C, Wang X. Automated Machine Learning in Waste Classification: A Revolutionary Approach to Efficiency and Accuracy. Proceedings of the 2023 12th International Conference on Computing and Pattern Recognition, 2023, 299-303.
- [19] Wang X, Wu Y C, Ma Z. Blockchain in the courtroom: exploring its evidentiary significance and procedural implications in US judicial processes. Frontiers in Blockchain, 2024, 7: 1306058.
- [20] Seera M, Lim C P, Kumar A, et al. An intelligent payment card fraud detection system. Annals of Operations Research, 2024, 334(1): 445-467.
- [21] Wang X, Wu Y C, Zhou M, et al. Beyond surveillance: privacy, ethics, and regulations in face recognition technology. Frontiers in Big Data, 2024, 7: 1337465.
- [22] Liu Y, Wu Y C, Fu H, et al. Digital intervention in improving the outcomes of mental health among LGBTQ+ youth: a systematic review. Frontiers in Psychology, 2023, 14: 1242928.
- [23] Ejiofor O E. A comprehensive framework for strengthening USA financial cybersecurity: integrating machine learning and AI in fraud detection systems. European Journal of Computer Science and Information Technology, 2023, 11(6): 62-83.
- [24] Carneiro N, Figueira G, Costa M. A data mining based system for credit-card fraud detection in e-tail. Decision Support Systems, 2017, 95: 91-101.
- [25] Al-Hashedi K G, Magalingam P. Financial fraud detection applying data mining techniques: A comprehensive review from 2009 to 2019. Computer Science Review, 2021, 40: 100402.
- [26] Liang Y, Wang X, Wu Y C, et al. A study on blockchain sandwich attack strategies based on mechanism design game theory. Electronics, 2023, 12(21): 4417.
- [27] Li X, Wang X, Chen X, et al. Unlabeled data selection for active learning in image classification. Scientific Reports, 2024, 14(1): 424.
- [28] Kalluri K. Optimizing Financial Services Implementing Pega's Decisioning Capabilities for Fraud Detection. International Journal of Innovative Research in Engineering & Multidisciplinary Physical Sciences, 2022, 10(1): 1-9.
- [29] Chen S, Liu Y, Zhang Q, et al. Multi-Distance Spatial-Temporal Graph Neural Network for Anomaly Detection in Blockchain Transactions. Advanced Intelligent Systems, 2025, 2400898.
- [30] Sailusha R, Gnaneswar V, Ramesh R, et al. Credit card fraud detection using machine learning. 2020 4th International Conference on Intelligent Computing and Control Systems (ICICCS), 2020: 1264-1270. IEEE.
- [31] Cui Y, Han X, Chen J, et al. FraudGNN-RL: A Graph Neural Network With Reinforcement Learning for Adaptive Financial Fraud Detection. IEEE Open Journal of the Computer Society, 2025.
- [32] Thennakoon A, Bhagyani C, Premadasa S, et al. Real-time credit card fraud detection using machine learning. 2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence), 2019: 488-493. IEEE.

## ENHANCING GLOBAL FINANCE: A BLOCKCHAIN-BASED SOLUTION FOR EFFICIENT AND COST-EFFECTIVE CROSS-BORDER PAYMENTS

Gang Min

*Beijing Ruirong Technology Co., Ltd., Beijing 100080, China. Corresponding Email: mingang8732@outlook.com* 

**Abstract:** This study explores the transformative potential of blockchain technology in revolutionizing cross-border payment systems. Traditional methods are hindered by inefficiencies such as high transaction fees, prolonged processing times, and opaque operations, which impede seamless global financial interactions. Blockchain, with its decentralized and immutable ledger, offers a secure and transparent alternative that can significantly streamline payment processes. This paper examines how blockchain can facilitate real-time settlements, eliminate intermediaries, and enhance data integrity, thereby reducing costs and improving efficiency. Further, it addresses the practical applications and regulatory challenges associated with integrating blockchain into existing payment infrastructures. Ultimately, this research aims to provide actionable insights for developing a more efficient, transparent, and cost-effective cross-border payment ecosystem.

Keywords: Blockchain technology; Cross-border payments; Financial efficiency; Cost reduction

#### **1 INTRODUCTION**

In an era of increasing global economic interdependence, the demand for efficient cross-border financial transactions has never been higher. However, traditional cross-border payment systems are plagued by inefficiencies that hinder their ability to meet modern demands. These systems are characterized by complex processes involving multiple intermediaries, which lead to extended processing times, high costs, and a lack of transparency. These challenges not only burden businesses with increased operational costs but also limit the accessibility of cross-border transactions for small and medium-sized enterprises (SMEs) and individual users. In response, blockchain technology has emerged as a promising solution, offering a decentralized, transparent, and immutable ledger that can revolutionize the way cross-border payments are processed. This study aims to investigate how blockchain can optimize payment efficiency and reduce costs, providing a theoretical foundation for the development of a more efficient and cost-effective cross-border payment system.

With the deepening development of global economic integration, the demand for cross-border trade and financial transactions is increasing. However, the existing cross-border payment system is facing many restrictions and has become a "bottleneck" factor in international financial transactions. Traditional cross-border payment processes are cumbersome and involve multiple intermediaries, resulting in long processing time, high cost and lack of transparency. This situation not only increases the operating costs of enterprises, but also limits the convenience of cross-border transactions between small and medium-sized enterprises and individuals. In this context, how to optimize the efficiency of the cross-border payment system and reduce the transaction cost has become one of the research priorities in the global fintech field[1].

Blockchain technology, as a distributed ledger technology that has emerged in recent years, has brought new possibilities for cross-border payments. Its decentralization, transparency and imtamable features can significantly reduce the mediation of the payment system and improve the processing efficiency and security. Therefore, cross-border payment systems based on blockchain technology have gradually become the focus of attention in both academia and industry. This study aims to explore the potential of blockchain technology in cross-border payment systems, focusing on how to improve payment efficiency and optimize costs through this technology, so as to provide theoretical support for the construction of a more efficient and low-cost cross-border payment system[2-3].

#### 2 CURRENT SITUATION AND CHALLENGES OF CROSS-BORDER PAYMENTS

#### 2.1 Overview of the Traditional Cross-Border Payment System

Traditional cross-border payments rely heavily on established financial intermediaries and networks such as SWIFT, RTGS, and ACH. While these systems have long been the backbone of international finance, they are not without their drawbacks. SWIFT, for instance, excels at transmitting transaction information globally but relies on correspondent banking relationships for actual fund settlements, leading to delays and high costs. Similarly, RTGS systems, designed for high-priority payments, struggle with cross-border transactions due to their localized nature and the need for international bank participation. These complexities result in slow clearing speeds and increased transaction costs, highlighting the urgent need for innovation in cross-border payment systems.

Traditional cross-border payment systems often rely on banks and other financial intermediaries to complete the transnational transfer of funds. The main systems include SWIFT (Global Association for Banking, Finance and Telecommunications), RTGS (real-time full settlement system) and ACH (automatic clearing system), etc. SWIFT system is a widely used cross-border payment information transmission network in the global banking industry. It helps banks to transmit cross-border transaction information through its standardized information format and global network support. However, SWIFT does not directly liquidate the funds, and still needs to complete the actual capital flow through the account transfer between the corresponding bank and the agent bank. This process is usually more complex, involving multi-party participation, and a long information transmission process, leading to the actual settlement speed of funds[4].

The RTGS system is mainly used for large, high-priority payments, which is characterized by full settlement within the same day and is suitable for a single large cross-border transaction. However, RTGS systems mostly operate in a single country or region, while cross-border RTGS payments rely on financial institutions from other countries to join in. Therefore, the efficiency of RTGS in cross-border payment is limited. At the same time, traditional cross-border payment systems usually rely on multiple account structures, with banks having agent accounts in different countries, which makes the clearing speed in the cross-border payment process slow and increases the transaction costs[5].

#### 2.2 Major Challenges of Cross-Border Payments

The inefficiencies in traditional cross-border payment systems manifest in several critical areas:

Delays: The involvement of multiple intermediaries in cross-border transactions often results in significant delays. Even with advanced networks like SWIFT, settlements can take days or even longer when multiple countries and banks are involved. Time zone differences and varying operational hours further exacerbate these delays, making real-time payments a distant goal.

High Costs: Each intermediary in the payment chain charges fees, leading to cumulative costs that can be prohibitively high, especially for SMEs and individual users. Additionally, exchange rate fluctuations and bank conversion fees add to the financial burden, limiting the economic feasibility of cross-border transactions for many.

Lack of Transparency: Tracking the flow of funds and understanding fee allocations in traditional systems is challenging due to the involvement of multiple layers of intermediaries. This lack of transparency not only affects user experience but also increases the risk of fraud and reduces trust in the payment process.

Traditional cross-border payment system faces multiple challenges in its practical operation, mainly reflected in the following aspects:

Delay: Cross-border payment processes usually require multiple intermediaries, resulting in a long time to transfer funds. Even with SWIFT, funds can be settled for days, especially when transfers between multiple countries and banks, often extending further. In addition, differences in operating times and time zones of different banks may also lead to payment delays, making it difficult to meet the demand of real-time payment.

High fees: Cross-border payment involves multi-party participation, and each intermediary agency will charge a certain handling fee, and the total cost is relatively high after superposition. In the traditional system, intermediary banks and agent banks will charge intermediary fees in the transaction process, and the exchange rate difference and the conversion fee of the bank will also increase the transaction cost, and bring additional economic burden to both parties. Especially for small and medium-sized enterprises or individual users, this high cost is an important obstacle to their participation in cross-border transactions.

Lack of transparency: In traditional cross-border payment processes, capital flow, real-time state and fee allocation are often difficult to track because transfers involve multi-layer intermediaries and multinational banks. This lack of transparency not only affects the user experience, but also increases the risk of money being withheld or falling under fraud, and reduces the security and trust of transactions[6-7].

#### **3** OVERVIEW OF BLOCKCHAIN TECHNOLOGY

Blockchain technology represents a paradigm shift in the way data is stored and shared. By leveraging cryptography and consensus algorithms, blockchain creates a decentralized and immutable ledger that records transactions in a transparent and tamper-proof manner. Each block in the chain contains a unique cryptographic hash linking it to the previous block, ensuring data integrity and security. This decentralized structure eliminates the need for intermediaries, allowing for direct peer-to-peer transactions. The transparency and immutability of blockchain not only enhance trust but also provide a robust foundation for secure and efficient financial transactions.

Blockchain is a distributed ledger technology that enables the decentralized storage and sharing of data through cryptography and consensus algorithms. The core idea is to package transaction data into blocks and connect them in chronological order to form an immutable and traceable chain. Each block contains a collection of data records, which includes the transaction records and the hash values of the previous block. All the data in the blockchain are verified and transmitted through the encryption algorithm, and are jointly maintained by each node through the consensus mechanism, which ensures the integrity and security of the data[8].

Blockchain is decentralized, transparent, and tamper-proof. This decentralized structure reduces the reliance on third-party intermediaries, allowing all participants to reach an agreement without full trust. The imtamability of the blockchain stems from its chain structure: each block contains the hash value of the previous block, and tampering with

any block will lead to a mismatch in the hash value on the chain, thus revealing traces of changes. Therefore, blockchain provides a technical guarantee for the realization of the high security and reliability of data[9].

#### 4 CROSS-BORDER PAYMENT SYSTEM DESIGN BASED ON BLOCKCHAIN

A blockchain-based cross-border payment system aims to address the inefficiencies of traditional methods by offering a streamlined, low-cost, and transparent solution. The system architecture typically comprises several layers: user interface, smart contracts, consensus mechanisms, and distributed data storage. By leveraging real-time data synchronization and eliminating intermediaries, the system significantly reduces transaction times and costs. The payment process, from initiation to settlement, is automated through smart contracts, ensuring transparency and reducing the risk of errors or fraud.

Blockchain-based cross-border payment systems aim to achieve low-cost, real-time payment solutions. The system architecture typically includes a user layer, a smart contract layer, a consensus layer, and a data storage layer. The user layer is responsible for receiving payment instructions, the smart contract layer sets payment conditions and execution logic, the consensus layer ensures the validity of the transaction, and the data storage layer saves the transaction data through the distributed ledger. Real-time data synchronization between the nodes eliminates the intermediary dependence, simplifies the payment process, and improves the anti-aggression and fault tolerance of the system, and ensures the data security[10].

The payment process generally includes payment initiation, transaction verification, settlement execution, and result confirmation. After the user initiates the payment, the system will trigger the smart contract, check the account balance and authentication, confirm it through the consensus mechanism, and broadcast to the blockchain network to complete the settlement. The decentralization of the payment process significantly shortens transaction time, and users can track payment status in real time, improving transparency.

#### 4.1 Consensus Mechanism Selection and Optimization

Consensus mechanism affects the efficiency and security of cross-border payment systems. Byzantine fault tolerance (BFT) and proof of equity (PoS) are commonly used mechanisms, among which BFT is suitable for small-scale and high-trust networks, with fast confirmation speed, while PoS is suitable for large payment systems, with low and stable energy consumption. In order to improve the transaction efficiency, the system can adopt a multi-level consensus mechanism, such as PoS in a large range of nodes and BFT in a small range of nodes. At the same time, combined with the penalty score mechanism to deal with malicious nodes, to improve the security and stability of the system.

#### 4.2 Smart Contract Design and Application

Smart contract is an important tool for automated transaction, which can be used for the setting and execution of cross-border payment conditions, such as exchange rate conversion, account balance verification, etc., which can be automatically executed after the conditions are met to reduce intermediary intervention and manual operation. In order to ensure the accuracy and safety of the contract, strict tests need to be conducted, and a double confirmation mechanism is introduced when necessary to reduce the risk of contract execution.

#### 4.3 System Security and Compliance

Blockchain payment systems should pay attention to data security and compliance. Prevent data tampering through multiple signatures and encryption algorithms, and comply with anti-money laundering (AML) and KYC (understand customer) requirements to ensure the authenticity of user identity and transaction compliance. The combination of on-chain and off-chain data helps to achieve transparent management and facilitate regulatory review.

#### **5** CONCLUSIONS

Through decentralized and transparent operations, blockchain can significantly enhance payment efficiency and reduce costs. The proposed blockchain-based payment system, supported by advanced consensus mechanisms and smart contracts, offers a secure and efficient alternative to existing methods. However, challenges related to scalability and regulatory compliance must be addressed to fully realize the benefits of blockchain in cross-border payments. Future research should focus on overcoming these hurdles to pave the way for a more efficient, transparent, and cost-effective global payment ecosystem.

This paper studies the application of blockchain technology in cross-border payment systems, aiming to solve the problems of high cost, low efficiency and insufficient transparency in traditional cross-border payment systems. Blockchain offers new solutions for cross-border payments through its decentralization, immutability and high transparency. This paper designs a blockchain-based cross-border payment system, including system architecture, payment process, and key technology implementation. Through multi-level consensus mechanisms, smart contracts, and multiple security guarantees, the system automates and reduces low cost of payment processes, and improves payment efficiency and security. At the same time, this paper discusses the compliance and application challenges of the

blockchain payment system. Despite scalability and legal compliance issues, blockchain has broad prospects in the cross-border payment field, bringing efficient and secure solutions to the global payment system.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCE

- [1] Wang Xiaolan. Research on the new model of cross-border payment based on blockchain technology. China market, 2022(18).
- [2] Deng Zhongsheng, Wang Honghai, Deng Tingwen, et al. Research on the application of blockchain technology in the field of "Belt and Road" cross-border payment. Land Bridge View, 2022(02).
- [3] Wang Yu. Research on cross-border payment and settlement innovation based on blockchain technology. Hainan Finance, 2021(10).
- [4] Zheng Bugao. Research on cross-border payment issues based on blockchain technology. Modern Finance Guide, 2021(02).
- [5] Wang Xiaoguang, Hu Jingjing. Construction and evaluation of the cross-border payment system for bulk commodities based on blockchain. Supply Chain Management, 2020(12).
- [6] Li Yimeng. Research on the application of cross-border payment in the "One Belt And One Road" region. Technology and Industry, 2019(08).
- [7] Ba Shusong, Zhu Yuanqian, Wang Ke. Blockchain promotes the path of financial change. China Finance, 2019(08).
- [8] Li Bo, Zheng Bo, Guo Ziyang, et al. Development and prospect of the application of blockchain technology in the financial direction. Journal of Applied Science, 2019(02).
- [9] Fei Wang. Research on blockchain technology and new ideas for promoting the development of cross-border e-commerce in China. Theoretical Monthly, 2019(03).
- [10] Hu Qilei. Research on the application of blockchain technology in the financial field —— is based on the "One Belt and One Road" perspective. Friends of Accounting, 2019(05).

## BANK CUSTOMER DEPOSIT PRODUCT PURCHASE ANALYSIS AND PREDICTION

#### YuanJia Guo

Management of National Economics, Renmin University of China, Beijing 100872, China. Corresponding Email: christine228@126.com

**Abstract:** This study examines the differences in characteristics between customers who purchase deposit products and those who do not, using relevant information from a bank's customer dataset. The key features analyzed include customer ID, age, occupation, marital status, credit card default history, mortgage status, contact method, last contact month and duration, the three-month interbank lending rate, previous marketing campaign results, the number of contacts before the current campaign, the number of days since the last contact, employment variation rate, consumer price index, consumer confidence index, number of employees, and whether the customer purchased a deposit product. Python is employed for descriptive analysis and classification analysis, and decision tree, logistic regression, and random forest models are used to predict whether a customer will purchase a deposit product. The analysis results reveal key factors influencing customer decisions, providing insights for banks to conduct targeted marketing within limited time frames, increase the likelihood of customer purchases, and ultimately improve overall deposit performance. **Keywords:** Bank customers; Deposit products; Descriptive analysis; Classification analysis; Decision tree, Logistic regression; Random forest; Model prediction

#### **INTRODUCTION**

Deposit business is a key source of funding for commercial banks, playing a decisive role in their development. It is also an essential guarantee for profitability and risk management. Stable and low-cost deposits not only enable banks to generate more profit through lending but also help mitigate liquidity risks and reduce incentives to pursue high-risk assets, thereby ensuring stable operations. Regularly promoting deposit products to customers is an important task for banks. However, with continuous innovation in financial models and the increasing number of commercial banks, attracting customer deposits has become increasingly challenging. This study aims to analyze bank customer information to identify the factors influencing the purchase of fixed deposit products and predict whether customers will choose to purchase them. The goal is to enable banks to conduct more targeted marketing campaigns for potential customers.

#### **1 RESEARCH BACKGROUND AND OBJECTIVES**

#### 1.1 Research Background

Deposit business is a crucial source of funds for commercial banks, significantly impacting their long-term development. It serves as a fundamental safeguard for both profitability and risk management. Stable and low-cost deposits not only support banks in generating higher profits through lending but also help mitigate liquidity risks and reduce the incentive to pursue high-risk assets, thereby ensuring sound financial operations.

However, with the continuous innovation of financial models and the rapid growth in the number of commercial banks, the competition for attracting customer deposits has intensified. In this evolving landscape, banks must gain deeper insights into customer needs and behaviors to design effective marketing strategies that enhance the willingness to purchase deposit products. Prior research emphasizes that customer behavior is influenced by various psychological and social factors, including observational learning and promotional incentives [1].

Moreover, the application of machine learning algorithms in financial services has shown great potential in understanding consumer patterns and making data-driven marketing decisions. As reviewed by Mahesh [2], machine learning offers robust tools for classification and prediction tasks, making it highly applicable in customer segmentation and targeting. Supervised learning techniques, in particular, such as decision trees and logistic regression, have been widely adopted in predictive modeling for consumer behavior analysis [3]. These methods can support banks in identifying key variables that influence deposit product adoption and in building adaptive marketing models that respond to real-time behavioral data.

#### **1.2 Research Objectives**

The objective of this study is to analyze the relevant information of bank customers regarding the purchase of deposit products, identify key factors influencing customer decisions, and predict whether a customer will purchase a deposit product using classification models such as decision trees, logistic regression, and random forests. Specifically, this study aims to:

• Conduct a descriptive analysis of various characteristics of bank customers to understand their basic profiles and behavioral patterns.

• Apply classification analysis methods to identify the main factors influencing customers' decisions to purchase deposit products.

• Build and evaluate decision tree, logistic regression, and random forest models to predict customer purchasing behavior.

• Provide a basis for banks to develop targeted marketing strategies based on the analysis and prediction results, enabling them to effectively increase customer purchase rates within a limited timeframe and ultimately improve overall deposit performance.

#### **2 DATA COLLECTION AND PREPROCESSING**

#### 2.1 Data Collection

The data used in this study is sourced from Tianchi, containing information on bank customers, including basic demographic details, communication records, and customer identifiers. The dataset consists of 22,500 entries and includes 22 columns. The field names and their corresponding meanings are listed in the table1 below:

Field Name	Meaning		
id	Customer identifier		
age	Age		
job	Occupation		
marital	Marital status		
education	Education level		
default	Credit card default status		
housing	Mortgage status		
loan	Other loan status		
contact	Contact method		
month	Last contact month		
day_of_week	Last contact day of the week		
duration	Last contact duration		
campaign	Number of contacts before the current campaign		
Table 2The Field	Names and Their Corresponding Meanings (Descriptive and Categorical Analysis)		
Field Name	Meaning		
pdays	Number of days since the last contact		
previous	Number of previous marketing contacts		
poutcome	Outcome of the previous marketing campaign		
emp_var_rate	Employment variation rate		
cons_price_index	Consumer price index		
cons_conf_index	Consumer confidence index		
lending_rate3m	Three-month interbank lending rate		
nr_employed	Number of employees		
subscribe	Whether the customer purchased the deposit product		

Table 1 The Field Names and Their Corresponding Meanings (Basic Demographic Information)

This table 2 will be used for **descriptive and classification analysis** to identify key factors influencing bank customers' decisions to purchase deposit products. Additionally, it will be leveraged to build predictive models to determine the likelihood of a customer purchasing a deposit product.

#### 2.1.1 Missing data handling

Based on the data description, all fields contain 22,500 entries, meaning there are no missing values. However, some fields contain "unknown" values. The proportion of "unknown" values in each column is as follows:

Job category (job): 1.21% Marital status (marital): 1.42%

Education background (education): 4.42%

Credit card default status (default): 21.60%

Mortgage status (housing): 3.94%

Other last status (1003112): 2.050

Other loan status (loan): 3.95%

Based on the proportion of "unknown" values, the following processing measures were applied:

For fields where the proportion of "unknown" values is less than 5% (job, marital, education, housing, and loan), records containing "unknown" values were deleted.

For fields where the proportion of "unknown" values exceeds 20% (default), the missing values were filled with the mode (most frequent value) of the column.

#### 2.1.2 Duplicate data handling

After analyzing the dataset, no duplicate records were found. However, a feature analysis revealed redundancy among the following fields:

- "Number of days since last contact" (pdays)
- "Month of last contact" (month)

#### • "Day of the week of last contact" (day\_of\_week)

Since **pdays** already captures the information about the last contact, the **month** and **day\_of\_week** fields were **removed** to eliminate redundant features.

#### 2.1.3 Outlier handling

A descriptive analysis of numerical fields revealed the presence of extreme outliers in certain fields, which could negatively impact the model's performance. To ensure data accuracy and model effectiveness, these outliers were handled appropriately.

The analysis found extreme values in:

Number of contacts made during the current campaign (campaign);

Number of contacts made before the current campaign (previous);

To address these outliers, the following filtering criteria were applied:

Retained records where campaign < 32;

Retained records where previous < 23;

#### **3 DESCRIPTIVE DATA ANALYSIS**

#### 3.1 Overall Sample Data Analysis

Based on the chart "Proportion of Subscription Status", the following conclusions can be drawn:

• The proportion of customers who did not subscribe to the deposit product is significantly higher than those who subscribed.

• Specifically, customers who did not subscribe ("no") make up the vast majority of the dataset, accounting for over 90%, while those who subscribed ("yes") account for less than 10%.

• This indicates a significant data imbalance, where the number of non-subscribing customers far exceeds that of subscribing customers.

Such an imbalance may cause the model to be biased toward predicting non-subscriptions in subsequent model training. To address this, certain data balancing techniques such as oversampling or undersampling should be applied.

The low subscription rate suggests that the current marketing strategy may not be effective in attracting customers to subscribe to deposit products. Banks may need to adjust and optimize their marketing strategies to increase customer subscription rates. Some possible improvements include:

• Implementing more precise customer segmentation

- Offering more attractive deposit products
- Strengthening customer relationship management

To gain deeper insights into customer subscription behavior, further analysis can be conducted using other variables such as age, marital status, education level, and occupation. Identifying the key factors influencing customer subscriptions will help develop more targeted marketing strategies, ultimately improving the subscription rate (Figure 1).



Figure 1 Proportion of Subscription Status

#### 3.2 Univariate Descriptive Analysis

#### 3.2.1 Relationship between bank customers' age and subscription to fixed deposit products

Based on the Figure 2 "Age Distribution by Subscription Status", we can analyze the purchasing behavior of customers across different age groups.

• The majority of customers fall within the 30 to 50 age range, which represents the largest customer base, regardless of whether they subscribed to the deposit product or not.

• Across all age groups, the number of non-subscribers (blue section) is significantly higher than that of subscribers (orange section), which aligns with the overall dataset's high proportion of non-subscribing customers.

• In particular, within the 30 to 40 age group, the number of subscribing customers is relatively higher, possibly because individuals in this age range are more focused on financial planning and deposit products.

• However, within the 40 to 50 age group, the number of non-subscribers is noticeably higher than in other age groups, suggesting that customers in this age range may have alternative financial investments or show less interest in deposit products.

• For customers aged 30 to 40, banks can enhance promotional efforts, offering more personalized and attractive deposit products to encourage subscriptions.

• For customers aged 40 to 50, further investigation is needed to understand why they are not subscribing. Based on these insights, banks can optimize product design and marketing strategies to attract more customers in this segment. By analyzing the relationship between age and subscription behavior, we can identify variations in deposit product subscriptions across different age groups. These insights help banks and financial institutions better understand customer needs, refine their marketing strategies, and ultimately increase subscription rates.

Particularly for young and middle-aged customers, banks should develop more targeted marketing plans to boost their willingness to subscribe to deposit products.



#### 3.2.2 Relationship between marital status and subscription to fixed deposit products

Based on the Figure 3 "Marital Status by Subscription Status", we can observe differences in purchasing behavior among customers with different marital statuses.

Married customers make up the largest proportion of the dataset, whether they subscribed to the deposit product or not. Single customers represent the second-largest group, but the number of non-subscribers among them is significantly higher than that of subscribers.

Divorced customers have the smallest representation, with non-subscribers also significantly outnumbering subscribers. This indicates that customer behavior regarding deposit product subscriptions varies significantly by marital status. Married customers are more likely not to subscribe, while single and divorced customers have relatively lower subscription rates.



Marital Status by Subscription Status

Figure 3 Marital Status by Subscription Status

Using this information, banks and financial institutions can develop more targeted marketing strategies for customers with different marital statuses to increase deposit product subscriptions.

#### 3.2.3 Relationship between education level and subscription to fixed deposit products

Based on the Figure 4 "Education Background by Subscription Status", we can analyze differences in subscription behavior across various educational backgrounds.

Customers with a university degree (university.degree) form the largest group in the dataset, both among subscribers and non-subscribers. However, the number of non-subscribers significantly outweighs that of subscribers.

Customers with a high school diploma (high.school) are the second-largest group, with non-subscribers significantly outnumbering subscribers.

Customers with professional courses (professional.course) and basic education (basic.9y, basic.4y, basic.6y) are fewer in number, but they also follow the trend where non-subscribers outnumber subscribers.

Illiterate (illiterate) customers have the lowest representation, with non-subscribers still significantly outnumbering subscribers.

Overall, across all education levels, the proportion of non-subscribers is significantly higher than that of subscribers. This suggests that education background may have some influence on customer subscription behavior.

Banks and financial institutions can tailor marketing strategies based on customers' education levels:

For higher-educated customers, offering more sophisticated financial products and educational services may attract them to subscribe.

For lower-educated customers, enhancing product explanations and promotions can help them better understand the benefits of deposit products, making them more likely to subscribe.



#### **3.3 Correlation Analysis**

Based on the Figure 5, correlation matrix heatmap and previous analyses, the relationships between different variables and customer subscription to deposit products (subscribe) can be summarized as follows:

Call duration (duration) and previous contact count (previous) have a positive correlation, indicating that more frequent contacts are associated with longer call durations. Additionally, longer call durations tend to be linked to higher subscription rates.

Number of marketing campaigns (campaign) shows low correlation with other variables, suggesting that simply increasing the number of marketing campaigns may not significantly improve subscription rates. Instead, marketing strategies should be optimized for better results.

Days since last contact (pdays) and previous contact count (previous) have a negative correlation, meaning that longer gaps between customer contacts are associated with shorter call durations and fewer previous contacts. This indicates that a well-planned contact frequency and interval can help improve subscription rates.

Employment variation rate (emp\_var\_rate) and three-month interbank lending rate (lending\_rate3m) have a positive correlation, reflecting their linkage within the economic environment. These factors may indirectly influence customers' financial decisions and subscription behaviors.

Call duration, contact frequency, and the economic environment have a significant impact on customer subscription behavior.

Banks and financial institutions should take these factors into account when designing marketing strategies. Optimizing customer communication by extending effective call durations and scheduling contacts appropriately can significantly.



Figure 5 Correlation Matrix Heatmap

#### **4 THREE MODELS AND ANALYSIS**

#### 4.1 Application of the Decision Tree Model in Predicting Customer Subscription to Deposit Products

In this analysis, we use the Decision Tree model to predict whether bank customers will subscribe to deposit products. This choice is supported by prior research showing the Decision Tree's interpretability and adaptability to categorical data [4].

The dataset includes various customer information. In the preprocessing phase, we removed unnecessary columns (such as IDs) and applied One-Hot Encoding to transform categorical variables into a numerical format, suitable for model training. The dataset was then split into a training set (70%) and a test set (30%), a widely adopted strategy that ensures sufficient training data while keeping an independent subset for performance evaluation and preventing overfitting.

The Decision Tree model was built using the DecisionTreeClassifier, which learns the relationship between features and the target variable (i.e., whether a customer subscribed). During training, the model iteratively splits the dataset based on feature thresholds to construct a hierarchical decision structure. The test set, unseen by the model during training, was used to assess its generalization capability.

Model evaluation involved computing standard metrics: Accuracy, Confusion Matrix, Precision, Recall, and F1-Score. The Decision Tree model achieved an accuracy of 84.09%, demonstrating promising overall performance. However, it struggled to correctly identify customers who subscribed, showing relatively low precision and recall in that category. This is consistent with existing literature on customer decision-making, which points out the complexity and psychological nuances of purchase intentions, especially under the influence of external factors like promotional stimuli [1].

The observed imbalance in prediction performance highlights the importance of understanding not just model mechanics but also customer behavioral patterns. As noted by Qiu et al. [5], interactions between variables — such as socio-economic status and previous marketing contact — may significantly affect predictive accuracy and should be carefully addressed in model design and variable selection(table 3).

Table 3	Decision	Tree Model	Analysis
			~

	precision	recall	f1-score	support
0	0.91	0.91	0.91	5860
1	0.40	0.41	0.41	890
accuracy			0.84	6750
macro avg	0.66	0.66	0.66	6750
weighted avg	0.84	0.84	0.84	6750

This suggests that further improvements, such as data balancing techniques or more complex models, may be needed to enhance prediction accuracy. By conducting this analysis, we can identify key factors influencing customer subscription behavior, providing data-driven support for banks to develop more targeted marketing strategies. Additionally, by improving the model (e.g., through data balancing techniques or more advanced models), we can further enhance prediction accuracy.

#### 4.2 Application of the Logistic Regression Model in Predicting Customer Subscription to Deposit Products

This analysis uses a logistic regression model to predict whether bank customers will subscribe to deposit products. While both logistic regression and decision tree models are commonly used for binary classification, they differ in modeling principles and processing approaches.

During data preprocessing, we removed the ID column and applied One-Hot Encoding to convert categorical variables into numerical form, allowing the model to process these features effectively. One-Hot Encoding transforms categorical data into binary vectors, ensuring compatibility with both logistic regression and decision tree models.

The dataset was then split into 70% for training and 30% for testing. The training set enables the model to learn featuretarget relationships and optimize parameters, while the test set evaluates performance on unseen data. This approach ensures a balanced trade-off between training adequacy and model validation.

The logistic regression model performs a linear combination of the input features and maps the result through the sigmoid function (logistic function) to a value between 0 and 1, thereby outputting the probability of the event occurring(table 4).

	precision	recall	f1-score	support
0	0.88	0.98	0.93	5860
1	0.47	0.14	0.22	890
accuracy			0.87	6750
macro avg	0.68	0.56	0.57	6750
weighted avg	0.83	0.87	0.83	6750

Table 4 Logistic Regression Model Analysis

Its core mathematical formula is:

 $P(y=1 | x)=1+e^{-(\beta 0+\beta 1x1+\beta 2x2+\dots +\beta nxn)}$ (1)

where  $\beta 0 = 0$  is the intercept, and  $\beta 1, \beta 2, ..., \beta n = 1$ ,  $\beta 1,$ 

parameter estimates.

The logistic regression model is widely used for predictions due to its simplicity, efficiency, and ability to output probabilities. It is easy to implement and interpret, making it ideal for establishing baseline models. Additionally, it offers fast computation, even with large datasets, and allows flexible decision threshold adjustments based on business needs.

Unlike decision trees, logistic regression is a linear model suited for linearly separable data, while decision trees handle complex nonlinear relationships. Although logistic regression may not capture intricate patterns as well as decision trees, its efficiency and simplicity make it a popular choice in many applications.

In this analysis, the logistic regression model achieved 86.58% accuracy in predicting customer subscription to deposit products. While it performed well in identifying non-subscribers, it had limitations in detecting subscribers, likely due to data imbalance. Future improvements could involve data balancing techniques or more advanced models to enhance prediction accuracy and support targeted marketing strategies.

#### 4.3 Application of the Random Forest Model in Predicting Customer Subscription to Deposit Products

The random forest model enhances prediction accuracy and stability by constructing an ensemble of decision trees, each trained on different data subsets with random feature selection at each split. This ensemble learning strategy reduces variance and mitigates overfitting, offering a more robust solution than a single decision tree [6]. It is particularly well-suited for high-dimensional datasets and excels at evaluating feature importance, making it a popular choice in banking applications.

Compared to logistic regression, which primarily captures linear relationships, random forests effectively model complex, non-linear patterns. While logistic regression remains interpretable and computationally efficient, its performance is often limited in scenarios involving intricate customer behaviors. Prior studies [5] have highlighted the limitations of linear models when dealing with interaction effects between variables, underscoring the advantages of tree-based methods like random forests in such contexts.

In our analysis, the random forest model outperformed both logistic regression and decision trees in predicting whether customers would subscribe to deposit products. It achieved an accuracy of 88.18%, correctly classifying 5,730 non-subscribers and 222 subscribers, while misclassifying 668 subscribers. These results demonstrate a strong performance overall but also indicate room for improvement, particularly in terms of recall for the subscribed category.

Despite its superior performance, the random forest model still struggles with accurately identifying subscribing customers — a challenge noted in similar studies on customer behavior prediction [7,8]. Factors such as class imbalance and subtle behavioral traits might contribute to this. Future enhancements such as SMOTE (Synthetic Minority Oversampling Technique), feature engineering, and parameter optimization could further enhance predictive power and generalizability.

In the broader context of banking risk management, machine learning models like random forests are increasingly recognized for their value in early risk detection and customer segmentation [9]. As banks seek to optimize their marketing strategies and retain valuable clients, the integration of advanced AI models into CRM (Customer Relationship Management) systems will be indispensable(table 5).

Tuble 5 Random Forest Woder Finarysis				
	precision	recall	fl-score	support
0	0.90	0.98	0.93	5860
1	0.63	0.25	0.36	890
accuracy			0.88	6750
macro avg	0.76	0.61	0.65	6750
weighted avg	0.86	0.88	0.86	6750

 Table 5 Random Forest Model Analysis

Accuracy of the Random Forest model: 0.882.

#### 4.4 Comparison of the Three Models

In this analysis, the decision tree model offers strong interpretability with its intuitive structure but tends to overfit, limiting its generalization [4]. It achieved 84.09% accuracy, performing well for non-subscribers but struggling with subscribing customers, highlighting its limitations in handling complex and imbalanced data [5].

The logistic regression model, known for its simplicity and computational efficiency, reached 86.58% accuracy. However, as a linear model, it does not effectively handle non-linear relationships or class imbalance, resulting in low recall for subscribed customers [5,10]. This aligns with previous research which identifies logistic regression's weaknesses in dynamic banking environments [11].

The random forest model, leveraging ensemble learning, achieved the highest accuracy at 88.18%, significantly improving stability and reducing overfitting [6,12]. Its ability to process high-dimensional data and evaluate variable importance makes it the most robust choice for predicting customer subscription behavior. However, the model still underperforms in identifying subscribing customers, suggesting further improvement is needed through data balancing and hyperparameter tuning [7,13].

As highlighted by Guerra and Castelli [14], banking supervision increasingly relies on machine learning tools not just

for prediction but for customer-centric decision-making. These models are also widely used in insolvency prediction [9] and credit risk assessment [11], demonstrating their transferability and scalability across domains.

#### **5** CONCLUSION

This study examined customer behavior towards bank deposit product subscriptions through descriptive statistics and predictive modeling. A key insight is the existence of significant data imbalance — with over 90% of customers not subscribing, the models naturally tend to favor the majority class. Such imbalance can distort evaluation metrics and compromise model reliability. Techniques like SMOTE or under-sampling are therefore recommended for future optimization [10].

Our findings show that demographic variables, such as age, marital status, and education, significantly influence subscription behavior. Customers aged 30 - 40 display higher subscription rates, suggesting they are at a life stage conducive to long-term savings. In contrast, those aged 40 - 50 show reduced interest, which may be due to increased financial obligations. While married customers represent the largest demographic group, their non-subscription rate is the highest, indicating the need for tailored strategies. In contrast, singles and divorced individuals are more likely to subscribe. Additionally, higher education does not directly translate into increased subscriptions, indicating a potential demand for more sophisticated or diversified financial products.

Moreover, call duration and contact frequency play critical roles in influencing customer decisions. Longer calls are positively correlated with higher conversion rates, and frequent engagement—especially during promotional periods— can significantly enhance customer responsiveness [1]. These behavioral insights are crucial for banks to optimize customer relationship management (CRM) strategies.

#### **6 DISCUSSION**

To increase customer subscription rates, banks should adopt targeted, data-driven strategies. For example, customers aged 30–40 can be approached with customized, long-term financial plans, while for more educated groups, investment-linked deposit products or educational webinars could improve engagement. Additionally, further market research is needed on married and mid-aged customers to understand their hesitations and adjust offers accordingly [8].

Given the impact of data imbalance, banks should integrate balancing techniques such as oversampling minority classes or using ensemble methods like Boosting during model training. Enhancing client communication by increasing effective call duration and optimizing frequency, particularly in high-conversion windows, may also improve results.

Moreover, the integration of advanced models like Gradient Boosted Trees (GBT), XGBoost, or even deep learning neural networks is recommended to capture more nuanced patterns and latent variables influencing decisions [10, 13]. Monitoring model performance over time with continuous learning systems can ensure adaptability to shifting customer behaviors and market dynamics, thereby supporting long-term marketing and risk management objectives [9,14].

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCES

- [1] Hao Liancai, Zou Peng, Li Yijun. The impact of sales promotion based on observational learning on customer purchase intention. Journal of System Management, 2012, 21(6): 795-801.
- [2] Mahesh B. Machine learning algorithms-a review. International Journal of Science and Research (IJSR), 2020, 9(1): 381-386.
- [3] Singh A, Thakur N, Sharma A. A review of supervised machine learning algorithms//2016 3rd international conference on computing for sustainable global development (INDIACom). Ieee, 2016: 1310-1315.
- [4] Yang Xuebing, Zhang Jun. Decision tree algorithm and its core technology. Computer Technology and Development, 2007, 17(1): 43-45.
- [5] Qiu Hong, Yu Dexin, Wang Xiaorong, et al. Analysis and evaluation of interaction effects in the logistic regression model. Chinese Journal of Epidemiology, 2008, 29(9): 934-937.
- [6] Lü Hongyan, Feng Qian. A review of research on the random forest algorithm. Journal of Hebei Academy of Sciences, 2019, 36(3): 37-41.
- [7] Leo M, Sharma S, Maddulety K. Machine learning in banking risk management: A literature review. Risks, 2019, 7(1): 29.
- [8] Lu Haiyan. Research on customer satisfaction in commercial banks. China Market, 2008 (44): 31-33.
- [9] Petropoulos A, Siakoulis V, Stavroulakis E, et al. Predicting bank insolvencies using machine learning techniques. International Journal of Forecasting, 2020, 36(3): 1092-1113.
- [10] Hu L, Chen J, Vaughan J, et al. Supervised machine learning techniques: An overview with applications to banking. International Statistical Review, 2021, 89(3): 573-604.
- [11] Munkhdalai L, Munkhdalai T, Namsrai O E, et al. An empirical comparison of machine-learning methods on bank client credit assessments. Sustainability, 2019, 11(3): 699.

- [12] Carbo-Valverde S, Cuadros-Solas P, Rodríguez-Fernández F. A machine learning approach to the digitalization of bank customers: Evidence from random and causal forests. Plos one, 2020, 15(10): e0240362.
- [13] Donepudi P K. Machine learning and artificial intelligence in banking. Engineering International, 2017, 5(2): 83-86.
- [14] Guerra P, Castelli M. Machine learning applied to banking supervision a literature review. Risks, 2021, 9(7): 136.

# FEEDBACK TRADING, DELAYED ARBITRAGE, AND ASSET PRICE BUBBLES: THEORETICAL AND EMPIRICAL STUDIES BASED ON DIFFERENT CYCLES

#### ShiQi Wang<sup>\*</sup>, Xin Pan, Cong Chen

Jiangsu University School of Finance and Economics, Zhenjiang 212013, Jiangsu, China. Corresponding Author: ShiQi Wang, Email: wangshiqi1014@163.com

**Abstract:** The 20th CPC National Congress emphasized that financial security is the cornerstone of national security, and preventing asset price bubbles is the core task of financial market risk prevention and control. China's stock market has developed rapidly, but it is characterized by frequent price fluctuations, and the deviation of prices caused by irrational investor behavior and arbitrage restrictions is prominent. Existing research mostly focuses on the theoretical level and lacks empirical analysis targeting the Chinese market. Moreover, the inverse relationship between investor sentiment and stock returns contradicts the reality of the weak presence of rational traders in China's market. This paper, from the perspective of behavioral finance, combines feedback trading and delayed arbitrage theory to construct an investor sentiment index based on principal component analysis and explores its dynamic relationship with the formation of stock price bubbles in China. Through empirical tests across multiple cycles, it reveals the mechanism by which arbitrage restrictions affect the persistence of bubbles and explains the "reverse induction puzzle," providing theoretical support and policy insights for improving China's stock market risk prevention and control system. **Keywords:** Investor sentiment; Feedback trading; Delayed arbitrage; Stock price bubbles; Arbitrage restrictions

#### **1 INTRODUCTION**

#### 1.1 Research Background

The 20th CPC National Congress elevated financial security to a national strategic level, emphasizing that "preventing and resolving financial risks is the fundamental task of financial work." Data shows that the total market value of China's A-shares has exceeded 80 trillion yuan, accounting for more than 12% of the global capital market. However, its price formation mechanism still has significant particularities: compared with mature Western markets, China's stock market, with a shorter development history and distinct characteristics of institutional transformation, exhibits more frequent asset price fluctuations and bubble accumulation phenomena.

Through a systematic review of the literature, two core issues were identified: first, the disconnection between theoretical mechanisms and empirical tests. Although Berger and Turtle conducted empirical tests on the investor sentiment index using U.S. sample data, no scholars have conducted empirical analyses using Chinese sample data [1]. Second, the theoretical paradox of investor sentiment indicators and the need for localization reconstruction. China has fewer investment institutions and weaker forces of rational traders, making it more prone to stock price bubbles. Therefore, this inverse indicator cannot match the current situation of China's stock market. This also necessitates the expansion through an investor sentiment index, studying the dynamic relationship between investor sentiment and stock returns to solve the reverse induction puzzle and better explain why China is more prone to stock price bubbles.

#### **1.2 Research Significance**

This paper has three practical values: First, it systematically analyzes the internal mechanisms of stock market fluctuations. By revealing the dynamic characteristics of the interaction between investor sentiment and market cycles, it provides a theoretical basis for understanding the special operating laws of China's capital market. Second, it innovates the investor decision-making analysis framework. It constructs a multi-dimensional decision-making model based on the identification of behavioral biases, providing methodological support for individual investors to avoid cognitive traps and optimize asset allocation. Third, it strengthens the effectiveness of financial risk prevention and control. By identifying early signals of emotion-driven price deviations, it helps regulatory authorities establish a forward-looking early warning system and improve the response mechanism for market imbalances caused by irrational trading, effectively maintaining the stable operation of the capital market.

#### 2 DOMESTIC AND INTERNATIONAL RESEARCH STATUS

#### 2.1 Definition and Measurement of Investor Sentiment

Li Junpan et al. pointed out that investor sentiment originates from investors' minds and psychological states and is a subjective influencing factor in investors' decision-making and actions [2]. The definition of investor sentiment has not

(1)

(2)

yet been unified, and existing research mainly defines it from the perspective of investors' beliefs and preferences deviating from traditional rational theory. In terms of measurement, although principal component analysis is widely used, it has certain limitations, such as requiring investor sentiment components to account for the largest proportion in common factors. Some studies have attempted to construct an emotion index using partial least squares to make up for the shortcomings of principal component analysis.

#### 2.2 Feedback Trading

Chen Jian and Zeng Shiqiang studied the impact of optimistic and pessimistic emotions on the behavior of noise traders and rational traders from both theoretical and empirical perspectives [3]. They further analyzed how these emotions drive market feedback trading behavior and explored the reaction of returns to investor sentiment shocks. Chen Zijun argued that since arbitrageurs' understanding of mispricing is sequential [4], they generally do not act immediately when they discover mispricing but choose the right time to enter and arbitrage. This also indicates that the impact of behavior on prices has a short-term resistance to arbitrage.

#### 2.3 Delayed Arbitrage

He Chengying et al. pointed out in their theoretical and empirical research that limited arbitrage is one of the main reasons for investor sentiment anomalies in the market [5]. In stock portfolios with more severe arbitrage restrictions, the phenomenon of negative correlation between investor sentiment and stock returns is more evident. At the same time, investor sentiment and stock price anomalies are more pronounced under extreme market emotions and extreme optimism.

#### 2.4 Asset Price Bubbles

Numerous studies have defined and interpreted asset price bubbles from different perspectives and used various methods to identify and test the existence, cycle, frequency, and degree of bubbles.

#### **2.5 Research Review**

Existing research has achieved certain results in the impact of investor sentiment on the stock market but still has shortcomings. There is a lack of theoretical and empirical research combining feedback trading and delayed arbitrage, and the explanation for the reverse induction puzzle is not deep enough, without fully considering the characteristics of China's stock market. This study will address these deficiencies and conduct in-depth analysis of the related issues.

#### **2.6 Theoretical Research**

We assume that there is a stock D in the market and divide the market into positive feedback traders and fundamental traders. We construct a model using a sentiment index to explain the specific situation of positive feedback traders and fundamental traders considering the impact of favorable information. Through the mutual game of the two types of market traders and based on market clearing to solve the equilibrium price of the stock in different periods, we discuss the stock price mechanism from a theoretical level. For the basic price of the stock, we divide this model into four trading periods. Period 0: There is no trading in the market. Period 1: Positive feedback traders do not react to the current price, and the demand is 0. Fundamental traders buy low and sell high, and their demand function is:

$$D1 = \alpha 1 - \beta 1P$$

When the stock price becomes basic public information, the demand of positive feedback traders in Period 2 is:

$$D_2^e = \alpha(\varphi - P_2)$$

The demand of fundamental traders in Period 2 is negatively correlated with the price of the security relative to its fundamental value, and the demand is:

$$D_2^f = \beta(P_1 - P_0) \tag{3}$$

The demand of fundamental traders in Period 2 is negatively correlated with the price of the security relative to its fundamental value, and the demand is:

$$D_2^e = \alpha(\phi - P_2) \tag{4}$$

After a favorable signal for stock D appears in the stock market, this project studies the game process of positive feedback traders and fundamental traders and combines the supply and demand theory to obtain the fundamental value and price sequence at different times. Through numerical simulation, the changes in variables are obtained, and the price time series change chart is drawn.

#### **3** EMPIRICAL RESEARCH

#### 3.1 Constructing Daily, Weekly, and Monthly Investor Sentiment Indexes

First, several original investor sentiment indicators are determined. Principal component analysis is used to construct investor sentiment, which can effectively remove noise factors in the original data and retain the collinear components of these data.

#### 3.1.1 Selection of investor sentiment proxy indicators

Turnover Rate (ATR): Reflects market liquidity and can be used to distinguish between optimistic and pessimistic emotions.

Closed-End Fund Discount Rate (CEFD): The change in the discount rate of closed-end funds reflects the impact of investor sentiment changes.

Initial Public Offering Return (IPOR): The return on the first day of listing can well reflect the degree of investor enthusiasm and is a positive indicator of sentiment.

Consumer Confidence Index (CCI): It reflects consumers' satisfaction with the current economic situation and their expectations for future economic development. This indicator is monthly data and reflects consumers' confidence in the market this month. When studying daily and weekly data, the data will be adjusted. CCI is the only subjective indicator. To construct the daily, weekly, and monthly composite sentiment indexes, we first take the first-order difference of the sentiment proxy indicators and then perform principal component extraction on the processed variables. When conducting principal component analysis, we extract the first five principal components based on the method of eigenvalues greater than 1. These five principal components are then weighted averaged according to their respective eigenvalues to obtain the preliminary composite investor sentiment index (INSI1). The calculation method is as follows:

$$INSI_1 = P * \frac{VEL}{EVEL}$$
(5)

#### 3.1.2 Constructing the daily investor sentiment change index

To address the timeliness issue of the classic BW monthly frequency index, this paper constructs a daily sentiment change index using principal component analysis. A 1% winsorization is applied to this index to exclude the influence of outliers. Subsequently, the continuous positive parts are accumulated and summed up. It should be noted that the original data are all positive values. However, after standardization by mean and standard deviation, negative values may appear, which represent low sentiment scenarios.

#### 3.1.3 Constructing the weekly investor sentiment change index

For the weekly index, after taking the first-order difference of the sentiment indicators, we apply the Partial Least Squares (PLS) method to effectively extract the sentiment S (sentiment change index). This process filters out irrelevant components and captures the consistent relationship between investor sentiment and expected stock returns across periods.

#### 3.1.4 Constructing the monthly investor sentiment change index

When constructing the monthly index, after obtaining the preliminary composite investor sentiment indicator using the aforementioned methods, we need to take the first-order difference again. In the second construction of the sentiment composite index, only the first and second principal components are extracted, and then the coefficients of their linear combination are calculated.

#### 3.2 Establishing the Relevant Linear Regression Model

In this empirical study, we need to test two hypotheses related to the accumulation of sentiment changes. We create a new indicator, Sumt-1,pos , to represent the accumulation of high sentiment. For these values, if the sentiment change in a period is positive, then the value of Sumt-1,pos is accumulated with the sentiment of that period; otherwise, it is reset to zero. The purpose of this treatment is to ensure that this indicator measures sentiment that has been continuously accumulating. In this way, we can derive the investor sentiment accumulation index at different frequencies and levels. As the bubble persists and sentiment increases, we predict that the relationship between positive sentiment and subsequent regression will weaken. The growth rate of the bubble will slow down when fundamental traders engage in arbitrage, and the bubble will burst when it inflates to a certain extent. We test this through the construction of a linear regression equation:

$$R_{i:t} - R_f = \alpha_i + \theta_i SUM_{t-1,pos} + \varphi_i SUM_{t-1,pos}^2 + e_{i,t}$$
(6)

t represents the period. Depending on the sentiment index of different cycles, the length of each period will vary. Rj,t denotes the excess return of asset j in period t, Sumt-1, pos is the accumulation of high sentiment, and Sumt-1,pos2 is the square of this accumulation. According to Baker et al. [1], we do not include risk factors that may weaken the

empirical results.

We predict that the  $\theta$  coefficient will be positive because, in the short term, a surge in sentiment will lead to an increase in future stock returns. We also predict that the  $\phi$  coefficient will be negative because, as the bubble expands, the selling pressure from rational arbitrageurs increases, thereby suppressing the growth rate of the bubble. Ultimately, we can verify that as the accumulation of high sentiment has a diminishing marginal effect on future excess returns, the growth rate of excess returns gradually slows down. When the accumulation of high sentiment reaches a sufficiently high level, a price correction will occur.

# **3.3** The Accumulation of Investor Sentiment and the Behavior of Positive Feedback Traders and Fundamental Traders

The diminishing marginal effect of high sentiment accumulation on future excess returns leads to a gradual slowdown in the growth rate of excess returns. When the accumulation of high sentiment reaches a sufficiently high level, a price correction occurs. Through the fitted quadratic relationship, we obtain a graph that intuitively shows the process of mutual game between positive feedback traders and fundamental traders in the market[6-7].

#### 4 RESULTS

#### 4.1 Validation of Sentiment Index Effectiveness

The daily sentiment volatility (standard deviation of 0.38) is significantly higher than the monthly sentiment volatility (standard deviation of 0.21). Moreover, the correlation coefficient between daily sentiment and turnover rate reaches 0.67 (p<0.01), indicating that short-term market sentiment has a significant driving effect[8].

#### 4.2 Bubble Formation and Bursting Thresholds

The regression results show coefficients of 0.15 (p<0.05) and 0.12 (p<0.01), confirming the nonlinear impact of sentiment accumulation on returns. When the sentiment accumulation SUM<sub>t,pos</sub>>4.7, the probability of price correction increases to 68% (95% Confidence Interval: 63%-73%), and the bubble cycle is shortened to 12-15 trading days.

#### **5** DISSCUSION AND CONCLUSION

By measuring investor sentiment across different cycles (daily, weekly, and monthly), this paper captures the rapid changes in investor sentiment, which can better prevent asset price bubbles caused by investor sentiment. Analyzing the impact of investor sentiment fluctuations on stock market volatility can assess the potential impact of sentiment volatility on economic growth and market stability. It also helps us to further clarify the patterns of stock market fluctuations in China, providing guidance for policy-making and risk management to some extent. This enables financial regulators to promptly detect the formation of "irrational bubbles" in the market, thereby preventing them and maintaining normal stock price fluctuations. It also aids in making correct macroeconomic decisions in a timely manner to stabilize the stock market. Analyzing the impact of investor sentiment fluctuations on market trading behavior can assess the degree of impact of sentiment volatility on market efficiency. It also provides a new perspective for investors and financial institutions to make decisions, enhancing investors' risk management capabilities and contributing to the stability of the financial market.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCES

- [1] Berger D, Turtle J H. Sentiment bubbles. Journal of Financial Markets, 2015: 2359-74.
- [2] Li Junpan, Wang Zhenjia. The impact of investor sentiment on investor behavior. Liaoning Economic, 2024(02): 73–77. DOI: 10.14041/j.cnki.1003-4617.2024.02.009.
- [3] Chen Jian, Zeng Shiqiang. How do optimism and pessimism drive market feedback trading behavior? ——From the perspective of the joint effect of noise traders and rational traders. Investment Research, 2023, 42(06): 127–141.
- [4] Chen Zijun. Synchronous cooperation risk and delayed arbitrage. China Foreign Investment, 2012(10): 212–213.
- [5] He Chengying, Chen Rui, Xue Bing, et al. Investor sentiment, limited arbitrage, and stock price anomalies. Economic Research Journal, 2021, 56(01): 58–73.
- [6] Hu Changsheng, Chi Yangchun. Research on Investor Sentiment and Abnormal Fluctuations in Asset Prices. Wuhan University Press, 2014.
- [7] Hu Changsheng, Xia Guoxiao. The cumulative effect of investor sentiment and stock price volatility. Statistics and Decision, 2023, 39(05): 152-157. DOI: 10.13546/j.cnki.tjyjc.2023.05.028.
- [8] Xu H-C, Zhou WX. A weekly sentiment index and the cross-section of stock returns. Finance Research Letters, 2018, 27: 135–139.

# LITERATURE REVIEW OF SELF-IMPROVEMENT PRODUCTS

DongPing Zheng\*, Pei Feng, JinFu Huang

School of Business Administration, Guangxi University, Nanning 530004, Guangxi, China. Corresponding Author: DongPing Zheng, Email: zdpgxu@163.com

Abstract: The rise of self-improvement products is closely related to lifelong learning, skill iteration, and mental health needs in modern society. This paper integrates the theories of evolutionary psychology and sociology to reveal that they drive individual development through the dual path of compensatory psychological compensation and developmental multidimensional optimization. Products cover physical, cognitive, economic and spiritual categories, and consumption behavior is influenced by emotional perception (e.g., guilt triggers immediate consumption, and awe reinforces long-term commitment), cognitive motivation (goal setting and social comparison), and social environment (resource competition and group norms). Although technological empowerment has eased class anxiety, it has led to digital exclusion and ethical risks. In the future, it is necessary to strengthen interdisciplinary integration and long-term evaluation, and to balance technological rationality and humanistic care at the practical level in order to promote sustainable development.

Keywords: Self-improvement; Self-improvement products; Compensatory motivation; Emotional perception; Social structure

#### **1 INTRODUCTION**

The rise of Self-Improvement Products (SIPs) in the 21st century is closely related to structural changes in modern society. Increased global competition and the transformation of the knowledge economy have created a need for lifelong learning, skill iteration, and mental health[1]. In postmodern values, individuals have shifted from "collective attachment" to "self-empowerment", and the pursuit of individualized growth has become the mainstream ideology[2]. Self-improvement products have become an important growth pole of the new consumer economy, but their rapid expansion has also triggered controversy in the academic community: on the one hand, research has confirmed that such products promote individual class mobility and social equity by lowering the threshold of knowledge acquisition and providing psychological adjustment tools; on the other hand, over-marketing may exacerbate anxiety (e.g., the "knowledge payment trap"), and technology-enabled products have become the most popular form of consumer products in the world. On the other hand, excessive marketing may exacerbate anxiety (e.g., the "knowledge payment trap"), and the ethical risks of data behind technological empowerment (e.g., self-objectification triggered by face recognition) are gradually emerging. Currently, the fields of psychology, sociology and economics have conducted research from the perspectives of consumption motivation, social context and technological intervention, etc. However, the lack of interdisciplinary dialogues, the absence of long-term impact assessment, and the comparative weakness of cultural differences still constrain the deepening of theories. The purpose of this paper is to systematically review the research results at home and abroad, to reveal the dual influence mechanism of self-improvement products on individual behavior and social structure, and to provide theoretical references for optimizing product design, avoiding technological alienation, and promoting sustainable development.

#### **2** SELF-IMPROVEMENT

Self-improvement mechanisms, as an important area of research in evolutionary psychology, reveal the psychological regulatory systems that humans have developed in the process of adapting to their natural and social environments. Evolutionary psychologists point out that this adaptive mechanism prompts individuals to respond to environmental challenges through continuous self-optimization, thus ensuring the continuation of survival advantages[3]. Existing theoretical systems mainly present a dichotomous divide: compensatory self-enhancement (CSE) refers to the psychological compensatory mechanism by which individuals compensate for negative self-evaluations through positive feedback, while basic self-improvement (BSE) emphasizes the process of self-improvement driven by purely developmental motivations[4]. Despite the differences in their paths of action, both are rooted in the fundamental human motivation to pursue self-worth enhancement, which essentially consists of three core dimensions: enhancing personal value perceptions, maintaining a positive self-schema, and averting the threat of negative evaluations.

At the level of theoretical constructs, the Self-Concept Enhancement Strategies Model (SCENT) provides an integrative framework for understanding this mechanism[5]. The model distinguishes two paths of action: direct self-enhancement achieves enhancement by reinforcing positive self-perceptions, while strategic self-improvement encompasses the secondary dimensions of self-improvement, self-validation, and self-assessment. Given the structural complexity of broad self-improvement theory, this article focuses on the narrower dimension of self-improvement, which is the process by which an individual systematically strives to achieve multidimensional development in key areas such as skill refinement, image management, health optimization, and wealth accumulation[5].

Research in developmental psychology suggests that self-improvement motivation shows significant activation during adolescence to early adulthood (18-25 years old). Individuals at this age are faced with critical developmental tasks such as independent development, career exploration, and social relationship building, and their future time perspective is characterized by significant stretching, a cognitive trait that forms a two-way reinforcement mechanism with self-enhancement motivation[6]. Notably, the level of achievement motivation has a significant moderating effect on self-improvement efficacy-individuals with high achievement motivation are more likely to develop goal-directed, sustained improvement behaviors. Evidence from localized studies suggests that when individuals activate self-improvement motivation, they systematically optimize core dimensions such as personality traits, health status, and subjective well-being, thereby building a virtuous psychological ecosystem[7].

#### **3 SELF-IMPROVEMENT PRODUCTS**

The concept of self-improvement products is based on the idea of self-improvement and specifically refers to goods and services that can help individuals advance in self-related dimensions such as physical functions (like enhanced physical fitness), psychological capital (such as cognitive improvement), economic capabilities (such as financial intelligence development), or spiritual realms (like spiritual growth)[8]. According to their functional orientation, self-improvement products can be classified into four major types: the first type is physical function enhancement, including fitness equipment and health monitoring devices; the second type is cognitive ability development, such as professional books and educational applications; the third type is economic capital accumulation, for instance, financial intelligence training courses; and the fourth type is spiritual realm improvement, represented by meditation aids[9]. This classification system not only clearly reflects the functional positioning of the products but also profoundly reveals their intrinsic connection with human self-perfection motives. According to evolutionary psychology theory, consuming such products is essentially a modern manifestation of an adaptive behavior - individuals enhance their survival competitiveness through resource investment, a mechanism that can be traced back to the psychological ecological balance theory [5], where product consumption becomes an important strategy for maintaining the self-esteem system. In the development of teenagers, the core functions of self-improvement products focus on the orientation of ability breakthrough, especially emphasizing the assistance in academic competition and career preparation through skill enhancement (such as intelligent learning devices) and physical optimization (such as sports monitoring equipment)[8]. Its mechanism of action is manifested in a three-stage path: knowledge internalization (skill acquisition), task efficiency optimization (performance improvement), and resource accumulation to break through bottlenecks (development empowerment), thereby building a closed-loop improvement model of "input - transformation - output"[10].

#### 4 INFLUENCES OF SELF-IMPROVEMENT PRODUCTS

#### 4.1 Emotion Perception

The driving effect of emotions on the consumption of self-improvement products shows significant compensatory and reinforcing characteristics. Negative emotions (such as guilt and shame) drive consumption by stimulating compensatory needs, while positive emotions (such as awe) enhance the persistence of behavior by strengthening self-control resources.

#### 4.1.1 Compensatory mechanism of negative emotions

Guilt, as a core emotion for cross-domain behavioral regulation, significantly drives consumers to repair their self-concept through the consumption of self-improvement products. For instance, Allard and White found that guilt resulting from failed healthy eating prompts individuals to purchase self-improvement books to alleviate cognitive dissonance[8]. Similarly, research on Chinese high school students indicates that guilt significantly enhances the preference for self-improvement products such as English courses by triggering self-improvement motivation, while shame, due to its self-deprecating tendency, does not have the same effect[11]. Further research shows that rejection-type social exclusion (such as workplace ostracism) drives compensatory consumption by lowering self-esteem, and independent self-construal individuals are more sensitive to this[12]. This suggests that the effect of negative emotions is moderated by the structure of an individual's self-concept.

#### 4.1.2 Reinforcement of self-control by positive emotions

Positive emotions promote the long-term use of self-improvement products by enhancing self-efficacy. Zhao Jianbin's experimental study shows that awe (such as watching nature documentaries) activates the "small self" cognition, significantly enhancing an individual's self-control resources and thereby increasing the preference for self-improvement products like fitness courses. This effect is more pronounced among incremental theory individuals (who believe abilities are malleable)[13]. Additionally, technological means can enhance the effect of positive emotions by increasing self-focus. For example, facial recognition technology, by providing real-time feedback on facial features, activates self-focus and thereby increases consumers' willingness to pay for self-improvement products such as beauty devices[14].

#### 4.1.3 Cultural heterogeneity

Gregg and Sedikides cross-cultural study found that East Asian consumers are more likely to view the consumption of self-improvement products as fulfilling social obligations (such as "improving abilities for the family"), while Western

consumers emphasize personal growth[15]. This difference leads to the differentiation of the driving effects of the same emotions in different cultures. For example, guilt is more likely to trigger the consumption of self-improvement products for social obligations in collectivist cultures.

#### 4.2 Cognitive Motivation

The cognitive motivation system of consumers regulates the adoption and continuous use of self-improvement products through goal-oriented mechanisms and social comparison strategies, and its effect is complexly influenced by ability beliefs and political ideologies.

#### 4.2.1 Goal setting and feedback mechanism

The difficulty of goals and the design of feedback mechanisms are key factors affecting the utility of self-improvement products. Diefenbach found that smart devices (such as water-saving timers) enhance users' self-control motivation by providing immediate behavioral feedback, but their long-term effects depend on the degree of goal internalization[16]. Xiao et al. further pointed out that time cues (such as "limited-time offers") promote the consumption of self-improvement products by activating a locomotion orientation, but overly difficult goals may trigger self-efficacy doubts, leading to behavioral reversals[17]. Therefore, dynamic adaptation algorithms (such as adjusting course difficulty based on user ability) have become technical solutions for optimizing goal setting.

#### 4.2.2 Bidirectional effects of social comparison

The application of social comparison strategies in the marketing of self-improvement products has directional differences. Ma Pu and Li Ji demonstrated through experiments in the Chinese market that positive comparisons (such as "surpass 80% of users") significantly increase the willingness to use fitness apps by stimulating competitive awareness, while negative comparisons (such as "avoid falling behind") may trigger defensive avoidance. However, low self-esteem groups are prone to falling into a negative comparison cycle[18]. Johnson and Angelo found that low self-esteem consumers tend to choose low-quality self-improvement products due to self-verification needs and need to break this cycle through identity reconstruction interventions (such as reinforcing the "quality consumer" label)[19].

#### 4.3 Social Environment

The social environment shapes the collective logic of the consumption of self-improvement products through the pressure of resource competition and group norms. While technological empowerment alleviates class anxiety, it also gives rise to new ethical risks.

#### 4.3.1 Social congestion and status threats

The competition for resources in the process of urbanization has intensified the demand for self-improvement. Social congestion is perceived as a "territorial invasion threat", prompting individuals to obtain class mobility capital through the consumption of self-improvement products, and this effect is more significant in areas with low employment rates and high perception of social equity. Due to the dual high characteristics of resource competition intensity and perception of social equity, first-tier cities in China have become the core scenarios for the consumption of self-improvement products, and the population density is significantly positively correlated with self-improvement motivation[9]. When consumers are faced with status threats, they will perceive loss of control and then purchase self-improvement products to compensate for the sense of loss of control[20].

#### 4.3.2 Group pressure and identity

Group pressure drives compensatory consumption through the identity mechanism. The sense of financial limitation prompts consumers to choose career development and self-improvement products such as MBA courses by triggering self-identity threats (such as "devaluation of professional identity")[21]. Self-identity threats (such as midlife career crises) promote the consumption of self-improvement products by activating problem-focused coping strategies, and the sense of social fairness positively moderates this path[22].

#### 4.3.3 Ethical dilemmas of technology empowerment

Although technologies such as AR and big data have enhanced the personalization level of self-improvement products, excessive technicalization may exacerbate social inequality. AR technology enhances user engagement through immersive experiences (such as virtual language environments), but low-income groups face new digital exclusion due to device access restrictions[23]. Furthermore, the immediate feedback from smart devices may give rise to "clock-in self-improvement", leading users to fall into a false sense of achievement[24], and weakening the sustainability of intrinsic motivation.

#### 5 CONCLUSION

This study systematizes the theoretical connotation, mechanism of action and influencing factors of self-improvement products, revealing their dual impact in individual development and social structure. At the theoretical level, the self-improvement mechanism is rooted in the framework of evolutionary psychology's explanation of adaptive behavior, and constructs a theoretical spectrum of psychological compensation and self-improvement through the dichotomy of compensatory and developmental motivation. The study confirms that emotional perception, cognitive motivation,

social environment and technological empowerment constitute a multidimensional driving system: negative emotions trigger immediate consumption through compensatory demand, while positive emotions promote long-term commitment through self-control reinforcement; the two-way regulation mechanism of goal-setting and social comparison is influenced by the interaction between beliefs about competence and ideology; the pressure of competition for resources and group norms shape the collective logic of consumption, and technological empowerment both alleviates class anxiety and generates digital exclusion. Technological empowerment both relieves class anxiety and creates ethical risks such as digital exclusion.

In practice, the design of self-improvement products needs to balance instrumental rationality and humanistic concern. Technical means such as dynamic adaptation algorithms and immersive experiences can optimize user participation, but the tendency of self-objectification triggered by technological dominance should be avoided. Marketing strategies need to emphasize cultural heterogeneity, such as strengthening the social obligation narrative in collectivist cultures and highlighting the value of personal growth in individualist cultures. Policymakers need to pay attention to the inequality behind technological empowerment, and lower the barriers to participation for disadvantaged groups through the popularization of digital infrastructure and the regulation of data ethics.

Existing research still has three limitations: first, interdisciplinary integration is insufficient, and the research on behavioral mechanisms from the perspective of psychology and the social criticism of technology ethics have not yet formed an effective dialogue; second, long-term impact assessment is missing, and the existing conclusions are mostly based on cross-sectional data, which makes it difficult to reveal the dynamic evolution law of consumption behaviors; and third, the research on cultural differences is weak, and the localized theoretical construction in non-Western contexts still needs to be deepened. Future research can focus on the following directions: (1) constructing an integrated model of emotion-cognition-environment to analyze the mechanism of multifactorial synergy; (2) conducting tracking studies to assess the long-term effects of technological interventions, especially the potential depletion effect on intrinsic motivation; and (3) strengthening cross-cultural comparative analyses to explore the interaction patterns between traditional values and consumerism in the process of modernization. These explorations will provide theoretical support for the construction of a responsible self-improvement ecosystem and promote the transformation of consumption practices from instrumental rationality to sustainable development paradigm.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### FUNDING

Supported by 2023 National Guangxi University College Students' Innovation and Entrepreneurship Training Program (Project No. 202310593104).

#### REFERENCES

- [1] A Giddens. Modernity and Self-Identity. In: Social Theory Re-Wired, 3rd ed. Routledge, 2023.
- [2] Z Bauman. Social Issues of Law and Order. The British Journal of Criminology, 2000, 40(2): 205–221. DOI:10.1093/bjc/40.2.205.
- [3] J Veldhuis, F Te Poel, R Pepping, et al. "Skinny Is Prettier and Normal: I Want to Be Normal"—Perceived Body Image of Non-Western Ethnic Minority Children in the Netherlands. Body Image, 2017, 20: 74–86. DOI:10.1016/j.bodyim.2016.11.006.
- [4] J S Shrauger, M Schohn. Self-Confidence in College Students: Conceptualization, Measurement, and Behavioral Implications. Assessment, 1995, 2(3): 255–278. DOI:10.1177/1073191195002003006.
- [5] C Sedikides, M J Strube. Self-Evaluation: To Thine Own Self Be Good, To Thine Own Self Be Sure, To Thine Own Self Be True, and To Thine Own Self Be Better. In: M P Zanna (ed.). Advances in Experimental Social Psychology, Vol. 29. Academic Press, 1997: 209–269. DOI:10.1016/S0065-2601(08)60018-0.
- [6] S C Levinson. Language and Space. Annual Review of Anthropology, 1996, 25: 353–382. DOI:10.1146/annurev.anthro.25.1.353.
- [7] X-C Liu, S Biao, D-Hui Dong. Pros and Cons of Self-Enhancement: Theory, Empirical Research, and Application. Advances in Psychological Science, 2011, 19(6): 883.
- [8] T Allard, K White. Cross-Domain Effects of Guilt on Desire for Self-Improvement Products. Journal of Consumer Research, 2015, 42(3): 401–419. DOI:10.1093/jcr/ucv024.
- [9] Y Ding, J Zhong. The Effect of Social Crowding on Individual Preference for Self-Improvement Products. Acta Psychologica Sinica, 2020, 52(2): 216–228. DOI:10.3724/SP.J.1041.2019.00216.
- [10] X Sun, X Li. Self-Improvement or Self-Enhancement? A Construal Level Theory Perspective. Psychological Science, 2012, 35(2): 264–269.
- [11] Y An. The Influence of Guilt Emotions on High School Students' Preference for Self-Improvement Products: The Mediating Role of Self-Enhancement Motivation. Sichuan Normal University, 2021.
- [12] L Tong. The Impact of Types of Social Exclusion on Preference for Self-Improvement Products. Zhongnan University of Economics and Law, 2021.
- [13] J Zhao. The Effect of Awe on Consumers' Preference for Self-Improvement Products. Journal of Marketing Science, 2019(4): 38–51.
- [14] J Li, X Li, Y Huang, et al. The Impact of Facial Recognition on Consumers' Preference for Self-Improvement Products. Nankai Business Review, 2023, 26(5): 181–189.
- [15] A Gregg, C Sedikides. Self-Enhancement. In: F Maggino (ed.). Encyclopedia of Quality of Life and Well-Being Research. Cham: Springer International Publishing, 2023: 6245–6248.
- [16] S Diefenbach. The Potential and Challenges of Digital Well-Being Interventions: Positive Technology Research and Design in Light of the Bitter-Sweet Ambivalence of Change. Frontiers in Psychology, 2018, 9. DOI:10.3389/fpsyg.2018.00331.
- [17] H Xiao, G Li, Y Chen, et al. The Time Effect on Desire for Self-Improvement Products. Current Psychology, 2023, 42(26): 23003–23017. DOI:10.1007/s12144-022-03405-3.
- [18] P Ma. The Influence of Comparative Direction and Form of Data Feedback in Self-Improvement Applications on Consumers' Intent to Use. Central University of Finance and Economics, 2022.
- [19] C S Johnson, S D'Angelo. Low Self-Esteem Leads to Low-Quality Purchases. (2021-08-10). https://business.cornell.edu/hub/2021/08/10/low-self-esteem-leads-to-low-quality-purchases/.
- [20] Y Wang, X Wang, H Chen (Allan), et al. Effect of Status Threat on Preference for Cross-Domain Self-Improvement Products: The Moderation of Trade-Off Beliefs. Journal of Business Research, 2024, 172: 114400. DOI:10.1016/j.jbusres.2023.114400.
- [21] J Yang, X Li, X Zhou. The Impact of Financial Constraints on Consumers' Choice of Self-Improvement Products: The Moderating Role of Self-Value and Financial Relatedness. China Soft Science, 2019(1): 136–145.
- [22] R Gong. The Effect of Threat to Self-Identity on Preference for Self-Improvement Products. Shanghai University of Finance and Economics, 2025. DOI:10.27296/d.cnki.gshcu.2023.001272.
- [23] W Zhou, X Lin, J Lei, et al. MFFENet: Multiscale Feature Fusion and Enhancement Network for RGB–Thermal Urban Road Scene Parsing. IEEE Transactions on Multimedia, 2022, 24: 2526–2538. DOI:10.1109/TMM.2021.3086618.
- [24] F Kehr, M Hassenzahl, M Laschke, et al. A Transformational Product to Improve Self-Control Strength: The Chocolate Machine. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. New York, NY, USA: Association for Computing Machinery, 2012: 689–694. DOI:10.1145/2207676.2207774.

### CHARACTERISTICS, ISSUES, AND COUNTERMEASURES OF CHINA'S DIGITAL ECONOMY DEVELOPMENT

ZhenJie Guo

Al-Farabi Business School, Al-Farabi KazNU, Almaty, Kazakhstan. Corresponding Email: 1144697673@qq.com

**Abstract:** Under the wave of the new technological revolution, the digital economy is reshaping the allocation of production factors, making it a pivotal arena in global competition. While China's digital economy is advancing rapidly, it faces challenges such as platform monopolies, data security risks, and the digital divide. Digital productivity not only enhances wealth creation efficiency but also profoundly transforms production relations, necessitating the establishment of a new governance framework. Accordingly, this study proposes targeted governance strategies from the perspective of government intelligence to support high-quality economic development.

Keywords: Digital economy development; Characteristics; Government

### **1 INTRODUCTION**

The rapid advancement of digital technologies has ushered in a new era of economic transformation, redefining traditional production models and reshaping global competitive dynamics. As a key driver of this transformation, the digital economy is not only revolutionizing the allocation of resources but also introducing unprecedented opportunities and challenges. In China, the digital economy has emerged as a vital force for high-quality growth, yet its expansion is accompanied by structural issues such as market concentration, data governance dilemmas, and uneven digital access[1]. These challenges underscore the need to re-examine existing governance frameworks and explore adaptive strategies that align with the evolving nature of digital productivity. This study delves into the characteristics, obstacles, and policy responses surrounding China's digital economic development, aiming to provide actionable insights for sustainable and inclusive progress in the digital age.

Nevertheless, China's digital economy development still faces multiple challenges, including lagging regulatory frameworks, data security vulnerabilities, and uneven regional development. In response to these issues, General Secretary Xi Jinping emphasized in the report to the 20th CPC National Congress the need to strengthen government guidance and improve institutional design, aiming to cultivate internationally competitive digital industry clusters that can provide new momentum for high-quality economic development. Currently, China is adopting a dual-drive approach combining technological innovation and institutional innovation to propel its digital economy toward a higher-quality development phase.

### 2 KEY CHARACTERISTICS OF DIGITAL ECONOMY DEVELOPMENT

### 2.1 Technology-Driven Foundations

Digital productivity fundamentally transforms economic paradigms by digitizing human understanding of objective laws through "digital-real integration," shifting resource allocation from localized to global optimization while generating innovative economic organizations and business models. Research indicates that achieving 80% adoption of intelligent computing centers during China's 14th Five-Year Plan period could drive nearly 300% growth in core AI industries. This growth stems from algorithms that digitally express physical laws, enabling software platforms to optimize data circulation. The synergistic combination of data, computing power and algorithms is revolutionizing both operational tools and decision-making processes, creating essential infrastructure for comprehensive socioeconomic digital transformation through real-time resource allocation, intelligent automation of complex decisions, and continuous optimization of industrial ecosystems via closed-loop data flows[2].

### 2.2 Innovation-Driven Development Models

The digital revolution is fundamentally transforming traditional industrial upgrading paradigms by overcoming the limitations of physical prototyping - long development cycles, high costs, and significant risks associated with conventional trial-and-error methods. This shift is enabled by disruptive digital twin and simulation technologies that create dynamic virtual replicas of physical systems, establishing a new innovation validation framework with three distinctive advantages: precise digital mirroring of real-world entities through similarity principles, rapid design iteration in near-zero marginal cost virtual environments, and closed-loop systems integrating description, diagnosis, prediction, and decision-making capabilities. This digital transformation has demonstrated remarkable efficiency gains across aerospace, automotive, and biopharmaceutical sectors, where virtual validation now accelerates R&D processes while dramatically reducing resource expenditures compared to traditional physical testing approaches.

### 2.3 Enhanced Economies of Scale

The digital economy creates increasing returns through network effects - as more users join, the value grows exponentially while costs approach zero. This stems from digital networks' unique properties: each new user adds connections that boost overall value (Metcalfe's Law), data accumulates to improve services, and ecosystem synergies multiply benefits[3]. These dynamics enable digital platforms to scale faster and achieve higher valuations than traditional businesses, as seen with major tech firms that leverage network effects for rapid growth. Essentially, digital networks rewrite traditional economic rules about scaling and value creation.

### 2.4 The Digital Economy Exhibits Distinct Long-Tail Characteristics

The digital economy substantially reduces enterprises' marginal production costs and diversification expenses, enabling cost-effective servicing of fragmented demand in long-tail markets. This economic paradigm not only expands market boundaries but also creates new opportunities for SMEs to pursue differentiated development strategies. By precisely capturing niche market demands, small and medium enterprises can achieve breakthrough innovations in specialized domains. Concurrently, the release of consumers' personalized needs continues to drive product market differentiation, fostering increasingly diversified industrial ecosystems.

### 2.5 Comprehensive Information Aggregation

The operation of market economies has long been constrained by fragmented transaction information due to technological limitations. The digital economy, leveraging its networked connectivity, real-time interaction capabilities, and advanced data processing, has successfully overcome geographical and physical barriers, making commercial transactions significantly more efficient and accessible. This innovative economic model establishes round-the-clock data collection systems that ensure comprehensive transaction records while employing sophisticated analytics to maximize data value extraction. By creating intelligent trading ecosystems, digital platforms have achieved end-to-end digital transformation - converting raw transaction data into actionable business intelligence through seamless technological integration.

### **3 KEY CHALLENGES IN DIGITAL ECONOMY DEVELOPMENT**

### 3.1 Growing Concentration and Monopolistic Trends

The digital economy has fostered a competition landscape where technological capabilities and data assets become decisive competitive advantages, exhibiting strong increasing returns to scale that reinforce market dominance for leading firms. The blurring industry boundaries further amplify this effect, as digital-native companies easily cross sectors through innovation, reshaping traditional competition patterns while raising entry barriers through technological complexity and data accessibility challenges[4]. This "winner-takes-more" dynamic creates new antitrust dilemmas, requiring regulators to carefully balance innovation incentives with fair competition preservation in an environment where scale advantages and market power increasingly concentrate in few players.

### 3.2 Data Security and Privacy Protection Challenges

We are witnessing an unprecedented era of data explosion, where information generated by various entities grows exponentially in both volume and complexity. These data assets, carrying immense commercial value and privacy implications, have elevated security concerns to strategic importance. Currently, data security faces three fundamental challenges: technologically, emerging architectures complicate data flows, rendering traditional protection systems inadequate; legally, severe information asymmetry between data subjects and processors undermines meaningful consent; economically, extreme data concentration fosters monopolistic practices and new forms of market failure. Particularly concerning is how data monopolists leverage their dominance through algorithmic black boxes to implement discriminatory pricing a practice that harms consumers and distorts market competition. Developing governance frameworks that balance innovation with security has become pivotal for sustaining healthy digital economic growth[5].

### **3.3 Digital Divide and Inequality Challenges**

The global digital transformation has revealed striking disparities in technological adoption across regions and demographics. These inequalities stem from both objective conditions like geographical constraints and economic development levels, as well as social factors including cultural traditions and educational attainment, collectively creating a multifaceted digital divide. This phenomenon manifests through uneven infrastructure coverage, varying technical competencies, and unequal access to digital services - fundamentally altering resource allocation patterns. While digitally-advantaged groups enter a virtuous cycle of opportunity accumulation, vulnerable populations risk exclusion from digital dividends. In China's western regions particularly, residents in remote mountainous areas face compounded disadvantages due to inadequate broadband infrastructure and digital illiteracy, limiting their access to

e-government and online education services, thereby exacerbating regional development gaps[6].

### **4 STRATEGIES FOR ADVANCING THE DIGITAL ECONOMY**

### 4.1 Enhancing Technological Innovation

Cultivating digital productivity requires advancing both fundamental research and applied innovation through sustained technological breakthroughs. This drives deeper development and industrial application of digital technologies, enabling new business models and economic growth—particularly through integrating digital and physical economies. Key to this effort is focusing on cutting-edge digital research, strengthening national scientific capabilities, optimizing innovation platforms, and fostering cross-sector collaboration to develop top research talent—ensuring technology translates into real economic transformation[7].

### 4.2 Advancing Digital Governance

Building a digital government is pivotal to modernizing governance capabilities. This requires innovatively applying digital technologies in public services to enable intelligent restructuring of governance processes and comprehensive improvement in service efficiency. Systematically advancing "digital governance" initiatives within a legal framework will deeply integrate data governance principles into social management systems.

### 4.3 Promoting Data Resource Sharing

To unlock the full potential of data as a key production factor, we must systematically break down data silos and establish standardized sharing mechanisms through unified exchange platforms and scientific protocols, enabling secure and orderly public data openness while prioritizing critical public services like healthcare and education. By leveraging intelligent technologies to enhance service accessibility and convenience, we can create a virtuous cycle where cross-domain data integration directly translates into tangible social benefits - transforming data flows into improved digital services that generate higher-quality data, ultimately ensuring citizens reap the real dividends of digital transformation through more equitable and user-friendly public services.

### 4.4 Strengthening Data Security Protection

To comprehensively strengthen the data governance system, a coordinated three-dimensional approach must be prioritized: Legally, accelerating the improvement of data-related legislation to clarify ownership rights and delineate responsibilities among all entities, while establishing fundamental systems for data classification, grading, and cross-border flows; Regulatorily, constructing a professional data security supervision framework by establishing dedicated agencies with technical enforcement capabilities and implementing graduated penalty mechanisms encompassing fines, operational suspensions, and criminal liabilities; Corporately, strictly enforcing platform accountability by mandating comprehensive lifecycle security management systems covering critical aspects like encrypted storage, granular access control, and geographically-distributed disaster recovery, supplemented by third-party audits and real-time compliance reporting[8]. These three dimensions form an interconnected, closed-loop system - legal frameworks provide the basis for supervision, regulatory enforcement drives corporate compliance, and enterprise-level implementation in turn validates the completeness of legislation, ultimately creating a new trinity governance paradigm integrating legal, regulatory, and technological elements.

### **5 CONCLUSION**

The digital economy, with its inherent attributes of data-centric operations, collaborative openness, and equitable accessibility, serves as a catalyst for both disseminating and optimally deploying informational, intellectual, and technological assets across economic systems.conventional economic development models. This requires establishing integrated digital industrial clusters that incorporate technological innovation, digital governance, data sharing mechanisms and security safeguards into a unified framework.

### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

### REFERENCES

- [1] Gao Taishan, Ma Yuan. Issues, Opportunities and Recommendations for China's Digital Economy Development. China Economic Report, 2020(02).
- [2] Wang Yiming. New Trends and Characteristics of Digital Economy Development. New Urbanization, 2023(07).
- [3] Zhen Qian, Zhang Yong. Thoughts on Promoting the Development of Digital Economy. Gansu Science and Technology, 2021(07).

- [4] Yu Yewen, Chen Gengxuan. Issues and Policy Recommendations for China's Digital Economy Development. Southwest Finance, 2021(07).
- [5] BUKHT R, HEEKS R. Defining, Conceptualising and Measuring the Digital Economy. International Organisations Research Journal, 2018, 13(2): 143-172.
- [6] SATO S, HAWKINS J. Electronic Finance: An Overview of the Issues. Electronic Finance: A New Perspective and Challenges, 2001, 7: 1-12.
- [7] Yu Yewen, Chen Gengxuan. Challenges and Policy Recommendations for the Development of China's Digital Economy. Southwest Finance, 2021(07).
- [8] GOOLSBEE A. In a World Without Borders: The Im-pact of Taxes on Internet Commercel. The Quarterly Journal of Economics, 2000, 115(2): 561-576.

# STOCK PRICE RESEARCH BASED ON ARIMA-GARCH-LSTM HYBRID MODEL

ChaoYan Wei<sup>\*</sup>, LanLan Li, PangLeYi Chen, MeiHui Huang, HuiLin Wei, KunYao Yao, XuYang Wang, Xin Ya, ChaoHai Wei

Institute of Information Technology, Guangxi Police College, Nanning 530028, Guangxi, China. Corresponding Author: ChaoYan Wei, Email: wei\_chaoyan@163.com

**Abstract:** As financial markets become increasingly complex, the demand for stock price forecasting is growing. To capture both linear trends and volatility in sequences as well as nonlinear dependencies, this paper proposes an ARIMA-GARCH-LSTM hybrid model. First, ARIMA is used to extract linear factors, followed by GARCH to express residual volatility conditions, and finally LSTM to capture deep nonlinear features. Based on the closing prices of the Shanghai Composite Index over 1,027 trading days from 2021 to 2025, RMSE, MAE, and MAPE were used for moving forecasts and multi-indicator estimates. The experiments show that the hybrid model outperforms individual ARIMA, GARCH, or LSTM models in all metrics, confirming its accuracy and robustness. Additionally, the hybrid model demonstrates strong adaptability during periods of high volatility.

Keywords: Hybrid model; Stock price forecast; ARIMA model; GARCH family model; LSTM model

### **1 INTRODUCTION**

The stock market is a crucial component of the modern financial market system, playing a vital role in resource allocation and reflecting economic activities. However, stock prices are influenced by various factors such as macroeconomic indicators, industrial policies, and market sentiment. Therefore, stock prices exhibit not only clear linear patterns but also volatility[1]. Moreover, there are complex nonlinear dependencies, which pose significant challenges to the accuracy of stock price predictions. Thus, effectively modeling and analyzing multiple features has become an important task both scientifically and practically.

The Autoregressive Integrated Moving Average (ARIMA) model, due to its powerful linear modeling capabilities and ease of use, is widely applied in short-term predictions of stock prices and profits. Yu Ting analyzed the adaptability of stock price series modeling based on white noise characteristics testing[2]. The study shows that some data are not suitable for direct construction of ARIMA models; using GARCH models or exponential smoothing methods can more accurately capture data features, enhancing modeling and prediction performance, thus providing a scientifically effective modeling process[3].

The Generalized Autoregressive Conditional Heteroskedasticity (GARCH) model is widely used to describe the volatility characteristics of financial time series, aiming for more specific descriptions of financial time series data. Liu Qi utilized the ARIMA-GARCH model to preprocess stock prices and then combined it with the CNN-BiLSTM-AT model and XGBoost algorithm to establish a hybrid prediction model[4], validating its effectiveness in achieving stock price prediction and return target goals; Xu Shuya et al. also found that joint modeling based on ARIMAARCH improves prediction stability and accuracy[5]. However, traditional GARCH models can only represent second-order matrix dynamics, with limited ability to describe higher-order nonlinear structures.

With the rise of deep learning technology, recursive neural networks (RNNs) and their derived short-term and long-term memory networks (LSTMs), due to their empirical ability to simulate long-term correlations and complex nonlinear relationships, are often used for predicting financial time series. Wang Xiaorui proposed a model that combines ARFIMA-GARCH and LSTM for disease prediction[6], which can improve feature extraction and predictive performance. The combination of GARCH and LSTM provides a model with volatility and nonlinear dynamics for disease prediction. Jiang Min et al. described the PCA-GARCH-LSTM framework, using PCA to reduce computational complexity, thereby enhancing prediction stability and accuracy[7]. Song Zhifan proposed the decomposition of Ceemdan and GARCH-LSTM methods[8], which utilize capturing multidimensional nonlinear features to significantly enhance inflation prediction performance.

This paper proposes a hybrid ARIMA-GARCH-LSTM model to address the aforementioned issues. In this method, an ARIMA model is first used to remove the linear components of time series data. Next, an appropriate GARCH model is employed to represent the temporal changes in balance[9]. Finally, the preprocessed multi-source features are input into LSTM to extract deep nonlinear relationships and predict closing prices. Empirical comparisons using historical closing values and combining sliding window prediction with multi-indicator estimation reveal that the hybrid model outperforms ARIMA, GARCH, and LSTM models in RMSE, MAE, and MAPE, demonstrating excellent predictive performance and robustness.

### **2 RESEARCH METHODS**

In order to fully describe the linear trend, volatility clustering effect and nonlinear dependence characteristics in the

stock price sequence, this paper constructs a hybrid prediction framework based on ARIMA, GARCH family and LSTM. The specific research methods are divided into the following five parts:

### 2.1 Autoregressive Integrated Moving Average Model

The ARIMA (p, d, q) model eliminates the non-stationarity of the series through difference operation and fits the linear trend by using the autoregressive and moving average components.

### 2.1.1 Autoregressive terms

The autoregressive term represents the relationship  $\bar{p}$  between the current observation of a time series and the previous lagged observations, and can be used to capture the dependence and trend of a time series. A model with the following structure is called a p-autoregressive model, or AR (p) for short.

$$\begin{cases} x_{t} = \phi_{0} + \phi_{1}x_{t-1} + \phi_{2}x_{t-2} + \dots + \phi_{p}x_{t-p} + \omega_{t} \\ \phi_{p} \neq 0 \\ E(x_{t}) = 0, \text{Var}(\omega_{t}) = \delta_{\varepsilon}^{2}, E(\omega_{t}\omega_{s}) = 0, s \neq t \\ E(x_{s}\omega_{t}) = 0, \forall s < t \end{cases}$$
(1)

Among them, is the observation value  $x_t \phi_0, \phi_1, \phi_2, \dots \phi_p$  of the current sample of the time series, represents the autoregressive coefficient, and represents the error term.

### 2.1.2 Differential item (I)

The primary purpose of differentiation is to ensure the stability of time series data. When the data in a time series shows trends and seasonality, these can be eliminated by differencing one or more time series data points, thereby stabilizing the time series data. The notation i(d) represents this, with the formula:

$$I(d) = \Delta x_t = x_t - x_{t-d}$$
<sup>(2)</sup>

The first difference represents the time seriesd, and the order of the difference.

### 2.1.3 Moving average term

This section introduces the random residual used to capture time series data. In ARIMA, it compares the current observation with past observations of the error term, helping the model correct for random fluctuations that the autoregressive part cannot capture. It is assumed that the error term is independently and identically distributed, with no autocorrelation. The moving average term is denoted as MA(q), abbreviated as MA(q):

$$\begin{cases} x_{t} = \mu + \omega_{t} + \theta_{1}\omega_{t-1} + \cdots + \omega_{t-q} \\ \theta_{q} \neq 0 \\ E(\omega_{t}) = 0, Var(\omega_{t}) = \delta_{\omega}^{2}, E(\omega_{t}\omega_{s}) = 0, s \neq t \end{cases}$$
(3)

Among them, represents the error  $\bar{\omega}_r \mu_1, \mu_2, \cdots \mu_q$  term and represents the moving average coefficient. The ARIMA model is obtained by combining the above three terms, that is:

$$ARIMA(p, d, q) = AR(p) + I(d) + MA(q)$$
(4)

The basic steps are as follows: ① Stationarity test and differencing. Perform an ADF test on the original closing price series. If unstable, apply D-order differencing until the series stabilizes. ② Model identification: Combine the autocorrelation function (ACF) and partial autocorrelation function (PACF) of the differenced series to preliminarily determine the AR order p and MA order q parameters. Parameter estimation: Use maximum likelihood estimation (MLE) or the YUIE-WAIkER equation to fit the model parameters. Model diagnostic tests: Conduct the LJUNG-BOX test on the residual series to ensure no serial correlation; ⑤ Rolling prediction: Predict the next closing price using a sliding window on the training set, obtain the residual series, and prepare it for subsequent volatility modeling and deep learning input.

### 2.2 Generalized Autoregressive Conditional Heteroscedasticity model

The GARCH family model is used to describe the time-varying volatility of financial time series residuals. The typical form is GARCH (p, q):

$$\begin{aligned} & \epsilon_{t} = \beta_{s}^{1/2} v_{s}, \quad v_{s} \sim i.i. d. N(0, 1), \\ & , \beta_{s} = \alpha_{0} + \sum_{i=1}^{q} \alpha_{i} \epsilon_{s-i}^{2} + \sum_{j=1}^{p} \beta_{j} \beta_{s-j} \end{aligned}$$
(5)  
In wich  $\alpha_{0} > 0, \ \alpha_{i} \ge 0, \ \beta_{j} \ge 0, \ \sum_{i=1}^{q} \alpha_{i} + \sum_{j=1}^{p} \beta_{j} < 1. \end{aligned}$ 

In this paper, GARCH(1,1) and its variants are selected for comparison in the model construction to  $\{\beta_s\}$  extract the conditional volatility of the residual sequence and use it as the auxiliary input feature in the deep learning stage.

### 2.3 Long Short-Term Memory Network

LSTM overcomes the problem of long-term dependence  $f_t i_t O_t C_t$  that is difficult to capture in traditional RNN by introducing forgetting gate, input gate, output gate and cell state transmission. The core calculation process of forgetting gate, input gate and output gate in long-term short-term memory model is shown in Figure 1: (1) forget gate:

$$f_t = \sigma(W_f[x_t, h_{t-1}] + b_f),$$
 (6)

② Input gates and candidate memories:

$$i_{t} = \sigma(W_{i}[x_{t}, h_{t-1}] + b_{i}),$$
  

$$\tilde{C}_{t} = \tanh(W_{C}[x_{t}, h_{t-1}] + b_{C})$$
(7)

③ Status update:

$$C_{t} = f_{t} \odot C_{t-1} + i_{t} \odot C_{t}, \tag{8}$$

④ Output gate and hidden state:

$$o_{t} = \sigma(W_{o}[x_{t}, h_{t-1}] + b_{o}),$$

$$h_{t} = o_{t} \odot \tanh(C_{t}).$$

$$\epsilon_{t} = h_{t}^{1/2} v_{t}, v_{t} \sim i.i. d. N(0,1),$$

$$h_{t} = \alpha_{0} + \sum_{i=1}^{q} \alpha_{i} \epsilon_{t-i}^{2} + \sum_{i=1}^{p} \beta_{i} h_{t-i},$$
(10)



Figure 1 LSTM Gated Information Flow Diagram

LSTM is good at learning complex nonlinear dynamic relationships from multi-dimensional inputs, and is a key tool for capturing deep patterns of stock prices in this paper.

### 2.4 ARIMA-GARCH-LSTM Hybrid Model

### 2.4.1 Steps of the combined model

(1) Data preprocessing: Fill in missing values, normalize, and differencing the original closing prices to obtain a stationary series; (2) ARIMA modeling: Fit an ARIMA model on the training set  $\hat{y}_t^{ARIMA} \epsilon_t \{\epsilon_t\} \{h_t\}$  to extract linear predictions for the next time step and corresponding residuals. (3) Volatility extraction: Fit GARCH family models on the residual series to calculate conditional variance (4) LSTM training: Use the ARIMA predicted values within the time window, residual volatility, and the normalized original prices as multi-dimensional features to construct the training set, which is then input into the LSTM network for deep learning; (5) Rolling prediction and ensemble: Generate final prediction values Hybrid on the test set using the same process, and compare them with the outputs of each base model. The specific steps of ARIMA-GARCH-LSTM hybrid model are shown in Figure 2.



Volume 2, Issue 2, Pp 36-43, 2025

### 2.4.2 Advantages of hybrid models

(1) The multi-source information fusion organically combines the linear trend, conditional fluctuation and nonlinear characteristics, making up for the deficiency of a single model;

(2) After ARIMA eliminates the trend, GARCH supplements the volatility, and LSTM captures the deep dynamics, which significantly reduces the prediction error;

(3) Over-fitting control uses GARCH volatility and ARIMA residual as auxiliary features to reduce the noise fitting of LSTM and improve its generalization ability;

(4) Robustness is enhanced, and the hybrid framework performs well in different market environments, especially during periods of high volatility, where reliable forecasts can be produced.

### 2.5 Model Evaluation Index

To evaluate the prediction performance of each model, the following indicators are adopted in this paper:

$$RMSE = \sqrt{\frac{1}{N} \sum_{t=1}^{N} (y_t - \hat{y}_t)^2},$$

$$MAE = \frac{1}{N} \sum_{t=1}^{N} |y_t - \hat{y}_t|,$$

$$MAPE = \frac{100\%}{N} \sum_{t=1}^{N} \left| \frac{y_t - \hat{y}_t}{y_t} \right|.$$
(11)

Among them, and are the  $y_t \hat{y}_t$ Nreal value and predicted value respectively, and is the sample number. The effectiveness and superiority of the hybrid model are verified through the comprehensive comparison of the above indicators.

### **3 EXPERIMENT AND RESULT ANALYSIS**

Based on the ARIMA-GARCH-LSTM hybrid model, the historical closing price data of Shanghai Stock Index were analyzed to verify the prediction performance of the model.

#### 3.1 Data Selection and Description

The experimental data are from Yahoo Finance, covering 1027 trading days of Shanghai Composite Index from January 4, 2021 to April 25,2025. The data fields include: date, opening price, highest price, lowest price, closing price, adjusted opening price, and trading volume. This paper takes the closing price as the main object of analysis. In the experimental design:

(1) The closing price data of Shanghai 1027 index is divided into training set and test set in chronological order. The first 900 trading days are the training set, and the last 127 trading days are the test set.

(2) In the GARCH modeling stage, the logarithmic return series is calculated according to the closing price of the training set, and the residual is modeled.

(3) All input features (ARIMA forecast values, GARCH conditional volatility, and normalized closing prices) were standardized from minimum to maximum.

(4) In GARCH modeling, the logarithmic return is calculated using the closing price of the training set:

$$\mathbf{r}_{t} = \ln \frac{\mathbf{P}_{t}}{\mathbf{P}_{t-1}},\tag{12}$$

And fit the volatility model with its residual sequence;

Before LSTM training, Min-Max normalization is used for all features, such as normalized closing price, ARIMA prediction value and GARCH volatility:

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}.$$
 (13)

### 3.2 Model Environment and Parameter Setting

Development environment: Python 3.8.10;

Integrated development environment: Visual Studio Code;

Key libraries: statsmodels (ARIMA), arch GARCH family models, TensorFlow, pandas, numpy, scikit-learn, matplotlib; ARIMA model: The optimal order (5,1,0) is automatically selected by AIC criterion, and MLE estimation is used on the training set;

GARCH model: based on GARCH(1,1), compared with EGARCH and APARCH, and finally took the residual volatility of GARCH(1,1) as the input;

Number of iterations: 50 epochs;

Loss function: mean square error (MSE);

Batch size: 32;

3D feature matrix: normalized price, ARIMA residual, GARCH volatility;

Data loading and preprocessing: read Excel data, process missing values, calculate logarithmic return rate, normalize

### closing price;

Training/test division: the first 900 trading days for training, and the last 127 trading days for testing, training set and test set.

### 3.3 Analysis of Empirical Results

### 3.3.1 Prediction analysis of Shanghai Stock Index data set



As shown in Figure 3, the data set partitioning diagram intuitively illustrates the time series segmentation method used in the study. After normalizing the test set data of the Shanghai Composite Index, the prediction results of ARIMA, GARCH(1,1), single LSTM, and hybrid models were obtained according to the aforementioned process, and RMSE, MAE, and MAPE were calculated.



Figure 4 Log Returns and GARCH Volatility

Based on the closing price of the Shanghai Composite Index on that day, a GARCH(1,1) model was constructed using its actual logarithmic return as the dependent variable for volatility prediction. As shown in Figure 4, it can be seen that when the predicted volatility is high, the empirical logarithmic return is also high. The high consistency between these two trends indicates that the GARCH(1,1) model can capture information about yield volatility in the Shanghai Composite Index.

In order to verify the normalization effect, the fluctuations before and after the normalization of the closing price are compared. As shown in Figure 5, the gray line represents the violent fluctuations of the closing price, and the blue line represents the stable trend after normalization. The dimensional difference is eliminated by normalization, making it easier to carry out subsequent processing.



In order to comprehensively demonstrate the prediction ability of the hybrid model on the closing price of the Shanghai Composite Index, we conducted a visualization of the comparison of the hybrid model's prediction effect as shown in Figure 6. In the figure, the gray line represents the actual normalized closing price, while the blue dashed line and the green solid line respectively show the prediction results of ARIMA and LSTM.



Figure 6 Mixed Model Prediction Effect

As shown in Figure 6, the ARIMA model can only closely follow trends, but there is a deviation in price fluctuations during the mid-to-late stages of 2023; the LSTM model can track volatility locally, but it will match at the end of 2024 in extreme regions. From the comparison results, a single model has limited capability in predicting complex financial time series. Traditional statistical models are slow to respond to changes, while deep learning models are affected by noise and shocks. Therefore, extracting features from ARIMA, GARCH, and LSTM at multiple scales and modeling volatility can enhance their predictive strength and accuracy.

In order to clearly reflect the prediction performance of ARIMA-GARCH-LSTM, the data are visualized as shown in Figure 7.



Table 1 Comparison of Prediction Performance of Different Models of Shanghai Composite Index

model	RMSE	MAE	MAPE (%)
ARIMA	9.47	7.49	0.24
GARCH(1,1)	8.15	6.50	0.21
LSTM	5.78	4.63	0.15
hybrid model	4.09	3.26	0.11

As shown in Table 1, the prediction results of the hybrid model are 29.2% lower than that of the single LSTM model in RMSE, 29.6% lower than that of the single LSTM model in MAE, and 26.7% lower than that of the single LSTM model in MAPE. In order to clearly reflect the prediction effect of the hybrid model, the line chart comparing the actual closing price and the predicted value during the test period from 2021 to 2025 is shown in Figure 7.

### **4 CONCLUSION**

In response to the limitations of traditional single models that cannot simultaneously consider linear trends, volatility clustering, and nonlinear dynamics, an ARIMA-GARCH-LSTM hybrid model is proposed. First, the ARIMA algorithm is used to extract the linear component from the stock price series. Then, the linear part is removed. Next, GARCH(1,1) is employed to characterize the time-varying volatility of resistance. Finally, LSTM is trained with the expected value from ARIMA, the conditional variance from GARCH, and normalized prices as multi-dimensional inputs to learn a deep nonlinear function. Based on empirical data from the Shanghai Composite Index covering 1027 trading days from 2021 to 2025, the RMSE, MAE, and MAPE indicators show that the proposed model significantly outperforms single models, demonstrating the rationality and stability of the model[9-11].

Compared with traditional methods, the ARIMA-GARCH-LSTM hybrid model family has the following advantages: (1)Integrate multi-source information to fully capture price dynamics;

(2)The prediction accuracy is significantly improved, and the error index decreases by two digits;

(3)Reduce the risk of overadaptation and improve the generalization ability;

(4)It has remained stable in different market environments, especially during periods of high volatility.

However, there are also shortcomings: first, the overall framework is overly complex, with extremely high demands on computer hardware resources and parameter optimization; second, the GARCH and LSTM structures along with hyperparameter tuning need further examination; third, the ability to adapt to extreme market shocks requires further investigation. Future research could focus on introducing asymmetric fat-tail distributions, combining attention mechanisms with multi-scale degradation techniques, and incorporating algorithms such as XGBoost and Transformer to build more efficient integrated frameworks.

In a word, ARIMA-GARCH-LSTM hybrid model can provide an effective way in stock prediction. It contains linear, fluctuation and non-linear characteristics in the hybrid model, which has great application possibilities and further improvement possibilities[12].

### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

### REFERENCES

- [1] Wang J, Zhang T, Lu T, et al. A Hybrid Forecast Model of EEMD-CNN-ILSTM for Crude Oil Futures Price. Electronics, 2023, 12(11): 2521.
- [2] Ting Y H, Long Y B, Lin D, et al. Application of a Hybrid Model Based on ICEEMDAN, Bayesian Hyperparameter Optimization GRU and the ARIMA in Nonferrous Metal Price Prediction. Cybernetics and Systems, 2023, 54(1): 27-59.
- [3] Sweeti S, Surendiran B, Dhanalakshmi R, et al. Forecasting COVID-19 Pandemic Using Prophet, ARIMA, and Hybrid Stacked LSTM-GRU Models in India. Computational and mathematical methods in medicine,2022, 1556025. DOI: https://doi.org/10.1155/2022/1556025.
- [4] Liu Qi. Stock Forecasting Based on GARCH-BiLSTM Network. Changchun University of Science and Technology, 2024. DOI: 10.26977/d.cnki.gccgc.2024.000613.
- [5] Xu Shuya, Liang Xiaoying. Research on Stock Price Prediction Based on ARIMA-GARCH Model. Journal of Henan Institute of Education (Natural Science Edition), 2019, 28(04): 20-24.
- [6] Wang Xiaorui. Application of ARFIMA-GARCH-LSTM Model in the Prediction of Hand, Foot and Mouth Disease in Shanxi Province. Shanxi Medical University, 2023. DOI: 10.27288/d.cnki.gsxyu.2023.000726.
- [7] Jiang Min, Zhang Chuyi, Sun Deshan. Stock price prediction based on PCA-GARCH-LSTM model. Software Guide, 2025, 24(01): 43-48.
- [8] Song Zhifang. Research on the Prediction of China's Inflation Rate Based on CEEMDAN-GARCH-LSTM Model. North Minzu University, 2023. DOI: 10.27754/d.cnki.gbfmz.2023.000181.
- [9] Jiang Mincong. Research on Option Trading Strategies Based on LSTM Model to Predict Volatility. University of Electronic Science and Technology of China, 2024. DOI: 10.27005/d.cnki.gdzku.2024.003958.
- [10] Hu Yamei. Stock Price Prediction Based on GARCH Family Models and LSTM Models [C]// Proceedings of the Comprehensive Innovation and Development Academic Forum, Chongqing Dingyun Culture Communication Co., Ltd. School of Statistics and Mathematics, Guangdong University of Finance and Economics; 2023, 106-109. DOI: 10.26914/c.cnkihy.2023.023020.
- [11] Yongchao J, Renfang W, Xiaodie Z, et al. Prediction of COVID-19 Data Using an ARIMA-LSTM Hybrid Forecast Model. Mathematics, 2022, 10(21): 4001.
- [12] Ma C, Wu J, Hu H, et al. Predicting Stock Prices Using Hybrid LSTM and ARIMA Model. IAENG International Journal of Applied Mathematics, 2024, 54(3): 424-432.

### THE APPLICATION OF DEEP REINFORCEMENT LEARNING IN ASSET ALLOCATION: A THEORETICAL FRAMEWORK AND EMPIRICAL ANALYSIS

ZiLin Zhou

La Salle College Preparatory High School, Pasadena 91107, United States. Corresponding Email: zilinzhou08@gmail.com

Abstract: Asset allocation is a fundamental challenge in investment management, traditionally addressed through models such as mean-variance optimization. However, dynamic market environments and multi-period investment horizons limit the effectiveness of static methods. In recent years, the emergence of deep reinforcement learning (DRL) has provided a powerful tool for addressing complex sequential decision-making problems in finance. This paper conducts a comprehensive academic analysis of the application of DRL in asset allocation. First, we introduce the asset allocation problem and its challenges, then review the basic concepts of DRL and its relevance to financial decision-making. Next, we propose a theoretical framework for transforming the asset allocation problem into a Markov decision process and describe in detail how DRL agents learn optimal investment strategies under various assumptions and structures within this framework. Subsequently, through a review of foreign academic literature, this paper examines existing findings on the application of DRL in asset allocation from a qualitative perspective, including the superior performance of DRL strategies relative to traditional methods in certain scenarios and cautionary results where DRL remains competitive even under simple benchmarks. We discuss the current limitations of DRL methods, high transaction costs, and potential directions for improvement and future research priorities. The study concludes that while DRL holds great potential for enhancing asset allocation theory and practice, several key practical challenges must be addressed before its full potential can be realized.

Keywords: Asset allocation; Deep Reinforcement Learning (DRL); Markov decision process

### **1 INTRODUCTION**

Asset allocation is the process of distributing capital across different asset classes and is a core component of portfolio management. Its objective is to achieve the optimal balance between expected returns and risk. The modern portfolio theory proposed by Markowitz (1991) first provided a rigorous quantitative framework for this problem. Markowitz's mean-variance model aims to maximize expected returns given a certain level of risk, and this model remains the foundation of finance to this day[1]. However, Markowitz's method is essentially a single-period model, assuming that asset returns and investor preferences remain unchanged throughout the decision-making period. In practice, investors face multi-period decision-making and evolving markets[1]. Merton (1975) extended portfolio theory to continuous-time dynamic environments, introducing the concept of dynamic asset allocation[2]. Merton's framework emphasizes that optimal asset weights need to adjust over time as conditions change; his research demonstrated that solutions to multi-period asset allocation problems may differ significantly from those of single-period problems[2].

The dynamic characteristics of real-world markets—including changes in correlations, transitions between market states, and the presence of transaction costs—add complexity to asset allocation. Traditional analytical solutions for multi-period portfolio optimization often require strong assumptions or simplifications to maintain solvability. This has sparked interest in numerical computation and algorithmic methods to address more realistic complex scenarios. Reinforcement learning (RL), a branch of machine learning, enables agents to learn sequential decisions through trial-and-error feedback, offering a promising framework for addressing such problems. Unlike traditional optimization, RL can, in principle, learn strategies adaptively through interaction with the environment without requiring an analytical solution to the problem, thereby potentially addressing the intractable aspects of investor optimization problems.

This paper examines the application of DRL in asset allocation from both theoretical and empirical perspectives. We first outline the core concepts of DRL and its role in financial decision-making. Then, we construct a theoretical framework for the application of DRL in asset allocation, discussing the definition of states, actions, and rewards, as well as model assumptions. Next, we review the qualitative empirical findings of existing academic research on the application of DRL in asset allocation, summarizing the research approaches and main findings. Finally, we discuss the limitations of current methods and propose directions for future improvements and research. Through a comprehensive literature review, we aim to clarify how and to what extent DRL can advance the practice and theory of asset allocation.

## **2** OVERVIEW OF DEEP REINFORCEMENT LEARNING AND ITS APPLICATION IN FINANCIAL DECISION-MAKING

Reinforcement learning is a machine learning paradigm in which an agent learns optimal policies by interacting with an environment and receiving cumulative rewards[3]. Formally, such problems are often modeled as Markov decision

processes (MDPs), consisting of states, actions, state transition probabilities, and reward functions. Unlike supervised learning, which relies on labeled data for learning, RL learns from the outcomes of actions: the agent receives reward signals based on the actions taken at each step and adjusts its strategy accordingly to gradually improve the long-term cumulative reward. In a financial context, the RL agent can be viewed as a trading or allocation strategy, where the state represents market information, the action corresponds to portfolio adjustments, and the reward is a measure of investment performance. This framework is naturally suited to sequential decision-making problems, making it well-suited to multi-period asset allocation, where investment decisions at different time points interact with each other. The application of DRL in financial decision-making has been explored in multiple domains, including portfolio allocation, trading (market timing), option pricing, and order execution. Specifically in asset portfolio management, numerous studies have shown that DRL agents outperform some heuristic or static strategies due to their ability to time decisions or exploit complex dependencies. DRL algorithms applied to learning trading strategies include deep Q-networks, policy gradient methods, and actor-critic architectures. These algorithms must address unique challenges in financial settings, including the low signal-to-noise ratio in financial data and risk management requirements. However, the flexibility of RL frameworks—such as the ability to incorporate risk considerations directly into the reward function or as constraints—offers significant opportunities for innovation in asset allocation methods.

### 3 THEORETICAL FRAMEWORK: APPLICATION OF DRL IN ASSET ALLOCATION

Applying deep reinforcement learning to asset allocation requires transforming the problem into a RL form. This involves defining the environment, state space, action space, reward function, and learning mechanism:

State (S): The state should contain the information needed to make the optimal allocation decision. In an asset allocation environment, the state at time t can include recent asset prices or returns, macroeconomic indicators, volatility estimates, and any other relevant market characteristics. For example, the state can be a vector containing the returns of each asset over the past N days, plus some macroeconomic variables or market sentiment indicators. Some frameworks distinguish between fully observed and partially observed states; in finance, we typically face partial observation, but we construct the state vector based on the available information. Recent research often uses deep networks such as CNNs or recurrent neural networks (RNNs) to process historical price data and extract useful features as the state representation for DRL agents.

Action (A): An action refers to the portfolio decision made by the agent. In simplified cases, actions can be discrete choices. However, in real asset allocation, actions are more naturally represented as a continuous vector—the investment weights of each asset. Many DRL methods treat actions as continuous variables. This typically requires the use of policy gradients or actor-critic algorithms. The action space may be constrained: for example, weights may be required to be non-negative and sum to 1, or leverage levels may be restricted. These constraints can be satisfied in implementation by transforming the network outputs.

Reward (R): The definition of reward is a key design element that influences the learned strategy. In portfolio management, a natural reward is the change in portfolio value over each period. For example, after executing an action at time t, the portfolio return from t to t+1 can serve as the reward at time t. Thus, the cumulative reward over a complete episode corresponds to total return or compound growth. Some studies adopt risk-adjusted rewards: for example, penalizing volatility or maximum drawdown in the reward function to implicitly consider risk. It is worth noting that Benhamou et al. (2024) provide theoretical insights showing that if an RL agent is myopic and its reward is determined solely by the first and second moments of returns, the converged strategy of the agent is the Markowitz mean-variance portfolio[4]. In other words, when the reward is appropriately chosen, the classical Markowitz portfolio can be viewed as a special case of RL one-step optimization. By extending the reward to multiple periods and incorporating more information, agents can theoretically find strategies that outperform the static Markowitz solution[1].

Environmental dynamics: In RL, the state transition rules of the environment do not require prior knowledge by the agent but are critical for the simulation process during training. In asset allocation, the environment is essentially the market. For model-free DRL, we do not need an explicit market model; instead, we can sample state transitions using historical or simulated data. A common approach is to use rolling windows of historical time series as training episodes. Each episode may correspond to a fixed time interval, during which the agent "virtually trades." State transitions are determined by the actual market trends: the asset price in the next state is obtained from actual observations. In simulation-based methods, a generative model can also be used to generate synthetic asset paths for training to increase the number of training samples. However, if the simulation generator is not accurate enough, the intelligent agent may learn to exploit biases in the simulated environment rather than effective patterns in the real market, introducing model risk.

Assumptions: When modeling asset allocation as RL, a key assumption is the Markov property—that is, future state transitions depend only on the current state and action, and are independent of more distant history. Financial markets are not strictly Markovian, but we select appropriate state representations to satisfy the Markov condition as much as possible. Another assumption is that the statistical properties of the environment are stationary during training: the agent typically assumes that the data distribution experienced during training and execution remains consistent. However, in reality, markets evolve, meaning that an RL strategy trained during one period may perform poorly when market conditions change, unless the strategy can adapt on its own. Some theoretical frameworks partially address non-stationarity by explicitly incorporating time or state indicators into the state space, such as incorporating labels

representing the macroeconomic environment.

Learning algorithms: Under the above assumptions, various algorithms can be used to train DRL agents. Since the action space for asset allocation is mostly continuous, policy gradient methods are commonly used. The agent's policy is represented by a neural network that maps states to actions. In each training step, the agent takes an action in a given state, observes the reward and the next state, and then updates the network parameters to improve the expected future cumulative reward. The objective function is typically to maximize the expected cumulative return. Some implementations use a discount factor to reduce the weight of future rewards, although in investment, a more natural objective is often the undiscounted total return or terminal wealth.

The above theoretical framework essentially transforms asset allocation into a "game" scenario: the agent "plays" against the market, with the goal of maximizing its own wealth. If training is successful, the resulting strategy can exhibit complex behavioral patterns, such as timing and dynamic rebalancing. These behaviors are not pre-programmed through human rules but are autonomously discovered by the agent through learning if they can improve rewards. For example, even without explicitly instructing the agent on momentum or hedging strategies, it may independently learn to increase allocations to strong assets during upward trends or hold more risk-averse assets during periods of high correlation based on reward feedback.

On this basic framework, researchers have proposed several improvements. One is to incorporate transaction costs into the environment and rewards. Transaction costs can be modeled by imposing penalties on portfolio rebalancing behavior. This is crucial because strategies that appear profitable when costs are ignored may become unfeasible once costs are factored in. Another improvement is allowing the setting of risk aversion levels: for example, using logarithmic utility or mean-variance utility functions as rewards. Adjusting the reward function yields a range of strategies with different risk preferences, analogous to points on the efficient frontier. Jiang, Olmo, and Atwi (2025) explicitly incorporated risk aversion parameters into the DRL framework and demonstrated how strategy aggressiveness varies with different parameter settings[5].

In summary, the theoretical framework for modeling asset allocation as a DRL problem is as follows: an agent observes market states and decides portfolio weights to maximize cumulative returns. Within this framework, traditional strategies can be viewed as special cases. The flexibility of this framework allows for the incorporation of numerous real-world factors, making DRL a powerful theoretical tool for searching for improved allocation strategies in complex environments.

### 4 EMPIRICAL ANALYSIS: A REVIEW OF RESEARCH ON THE APPLICATION OF DRL IN ASSET ALLOCATION

Utilizing alternative data and cross-sectional information: Some studies have incorporated data beyond prices into the DRL framework. Aboussalah, Xu, and Lee (2021) explored the value of cross-sectional learning methods. In their paper published in Quantitative Finance, the agents learn not only from the time series of a single portfolio but also from the overall market cross-sectional data of numerous assets, thereby enabling the strategy to be more generalizable across a wider range of market conditions and assets[6]. They found that this cross-sectional training improved performance, indicating that information extracted from a broader market can benefit allocation agents. Compared to agents trained solely on the historical data of a single asset, agents trained using cross-sectional data demonstrated superior risk-adjusted performance in terms of returns.

Chen and Ge (2021) proposed a "learning-based strategy" for portfolio selection[7]. At its core, they introduced an algorithm in the International Journal of Economics and Finance that can adapt to changing market conditions. Their research emphasizes that a data-driven strategy is more resilient than static models under changing market conditions. Although the paper details are somewhat vague here, the key conclusion is that machine learning/DRL-driven strategies can continue to learn from new data, adjust themselves when traditional models fail, and remain effective[7].

Incorporating market sentiment and macroeconomic context: Financial markets are not only influenced by historical prices but also driven by investor sentiment and macroeconomic news. Recognizing this, Wei et al. (2021) incorporated asymmetric investor sentiment as part of the state in their RL portfolio model[8]. By using sentiment indicators, their agents could anticipate market changes that could not be predicted solely based on price history. Their research in the Journal of Expert Systems and Applications showed that DRL agents using sentiment data achieved better performance, especially during periods of market stress or euphoria: the agents were able to reduce positions before a decline and increase positions during a recovery. This result highlights the flexibility of DRL in integrating diverse data types; in principle, RL agents can learn to interpret sentiment indicators and price trends comprehensively to make better allocation decisions.

Macroeconomic variables or regime indicators are also important contextual information. Some studies provide agents with contextual data beyond asset prices, such as interest rates, volatility indices, or macroeconomic cycle labels. Benhamou et al. found that providing this additional information improves the performance of DRL models relative to traditional mean-variance strategies. In fact, it has been demonstrated that agents can utilize macroeconomic environment data to adjust their strategies. However, they also note that this comes with increased complexity—the more information the agent considers and the more forward-looking it is, the more challenging the training process becomes, though the potential benefits are greater.

Performance comparison with traditional methods: A common theme in many empirical studies is that DRL-driven asset allocation strategies often outperform traditional strategies in backtesting. For example, Jiang, Olmo, and Atwi (2025)

designed a DRL agent combining CNN and WaveNet components to handle high-dimensional portfolios and multi-period problems[5]. They tested the strategy under various market conditions, risk aversion levels, and with real transaction costs factored in. The results showed that DRL strategies outperformed multiple benchmark methods in terms of both returns and adaptability. This suggests that, under conditions of ample training data and carefully designed models, DRL can identify trading patterns and portfolio adjustments overlooked by static models, thereby achieving a better risk-return tradeoff.

Challenges and Differentiating Results: While many reports are positive, some important studies have highlighted potential shortcomings of DRL. Kruthof and Müller (2025) provide a cautionary example. These authors conducted a rigorous evaluation of state-of-the-art DRL algorithms (Soft Actor-Critic, SAC) in the Financial Research Letters, using a sample spanning seven stock markets and a total of 300 years. Interestingly, they found that in a no-transaction-cost scenario, DRL agents did exhibit some timing ability—outperforming market benchmarks during certain periods—but overall, they did not systematically outperform a simple 1/N equal-weight strategy. More notably, when introducing a mild transaction cost of 0.1%, the DRL strategy, due to its high turnover rate, resulted in negative net returns and performed worse than the 1/N strategy. The latter, which rarely adjusts its portfolio, is largely unaffected by transaction costs and thus significantly outperforms the DRL strategy when costs are considered. They also examined return distributions and tail risk metrics and found no consistent advantage for the DRL agents. This study highlights that a complex AI strategy does not necessarily outperform simple rules, especially when real-world frictions are incorporated. Its conclusions emphasize that without careful design, DRL strategies may lose their advantage due to over-trading or overfitting to noise, and require cost awareness and rigorous validation to enhance their practicality.

In summary, the empirical literature to date paints an encouraging but cautious picture. Many studies have documented performance improvements of DRL asset allocation strategies relative to traditional methods, which are attributed to DRL's ability to leverage complex data and continuously optimize strategies. However, the best results often come from carefully designed models that incorporate domain knowledge. At the same time, reliable evaluation is crucial—simple DRL models, if mishandled, may suffer from overfitting or over-trading issues, as revealed by some rigorous tests. Therefore, it is necessary to gain a deeper understanding of these limitations, which we will further discuss and propose future improvement directions in the next section.

### 5 CONCLUSION AND DISCUSSION: LIMITATIONS, POTENTIAL, AND FUTURE DIRECTIONS

The application of deep reinforcement learning to asset allocation remains an emerging field with significant potential but also clear limitations. Based on the theoretical framework and empirical findings discussed above, we summarize several key issues and propose possible directions for future research.

### 5.1 Limitations and Challenges

Market non-stationarity: Financial markets are non-stationary—their statistical properties change over time. DRL algorithms typically assume that the statistical properties of the environment remain constant during training. Strategies learned under one market regime may perform poorly when conditions change. This introduces the risk of overfitting to historical data. Many published DRL strategies perform well in-sample or during specific backtesting periods but may not generalize to new data or different market environments. Hambly et al. (2023) review that fat-tailed return distributions and regime switching pose fundamental challenges for RL in financial settings. To mitigate this issue, future methods can incorporate regime detection techniques or adopt meta-learning to enable agents to adjust themselves when detecting new regimes. Continuous learning frameworks or periodically retraining agents when new data arrives may be necessary to maintain the effectiveness of strategies.

Sample efficiency and data scarcity: Unlike games, financial data is constrained by historical length—markets have only a limited number of past years available for learning, and underlying processes are complex and variable. DRL algorithms typically require a large number of training samples to converge to a good strategy, which is problematic in finance because each episode is unique and not independently and identically distributed. Researchers have attempted to expand data through bootstrapping or simulation, but simulation data must be used cautiously; if the simulation model is overly simplified, the agent may learn strategies that rely more on the characteristics of the simulation environment rather than effective signals from the real market. This requires the development of RL algorithms with higher sample efficiency or algorithms that can incorporate prior knowledge. Model-based RL is a potential direction: the agent first learns a model of market dynamics and then uses this model for planning to find optimal strategies with fewer real samples. However, learning an accurate market model itself is also quite challenging.

Exploration and exploitation in real trading: In typical RL, exploration is necessary to discover the optimal strategy. However, in real trading, exploration means intentionally taking suboptimal actions to learn environmental characteristics, which is unacceptable to investors. Therefore, in practice, DRL often relies on offline training followed by online execution of learned strategies. This means that the agent may not have encountered scenarios that perfectly match those it will face, and it cannot "trial and error" extensively in real trading without incurring financial losses. Future research could explore safe exploration methods or construct highly realistic market simulation environments to allow agents to practice and refine strategies without directly exposing real capital to risk. For example, generative adversarial networks could be used to simulate realistic market scenarios for agent training, better preparing them for real-world trading environments. High Turnover Rate and Transaction Costs: As emphasized by Kruthof and Müller (2025), DRL strategies may suffer from over-trading issues. An unconstrained RL agent may pursue every perceived short-term opportunity, frequently rebalancing its portfolio[9]. This can lead to extremely high turnover rates, which, when accounting for transaction costs, may erode or even offset all gross returns. High turnover rates may also introduce tax burdens and other frictions that may not be adequately modeled in research. For practical applications, it is essential to incorporate transaction cost constraints or penalties into DRL strategies during training. Future research could integrate more refined cost models into RL environments. Additionally, regularization or adding costs to actions in the reward function could encourage agents to avoid frequent portfolio rebalancing unless the expected return is significant.

Risk Management and Tail Risk: Many DRL implementations optimize for average returns, but in investing, downside risk is a primary concern. Agents may unintentionally learn strategies that yield high average returns but expose the portfolio to rare, massive losses. Traditional portfolio management often sets risk limits or uses risk measures such as VaR. Current DRL frameworks may not address these concerns unless explicitly incorporated through reward functions or constraints. This requires the development of risk-sensitive DRL that incorporates risk elements into the agent's objectives. Some approaches may include optimizing CVaR in the objective or using multi-objective RL to optimize both return and risk metrics simultaneously. It is also necessary to ensure that agents adhere to risk limits even under extreme market conditions. This is critical for the robustness of strategy implementation in practice.

Interpretability and trust: Deep learning models are typically "black boxes." In the financial sector, a completely unexplainable strategy is difficult to trust, especially when fund managers or investment committees need to understand and trust the decision-making process, and regulators require explanations of trading logic. Unlike the Markowitz model, which explicitly shows the trade-offs between risk and return, the rationale behind the investment decisions output by DRL strategies is difficult to explain in human language. This opacity may hinder its application. Therefore, the application of explainable artificial intelligence (XAI) in finance has become an important direction. Cong et al. (2022) incorporated explainability into DRL models in their AlphaPortfolio study, using techniques to make the decision-making basis of agents more transparent[10]. Future research can focus on extracting rules or important factors from trained DRL agents—for example, using sensitivity analysis to determine which input features have the greatest impact on agent decisions. In summary, improving the explainability of DRL decisions while maintaining its flexibility is crucial for its adoption in real-world investment.

### **5.2 Potential and Future Directions**

Despite the challenges, DRL holds significant potential for asset allocation. With increasing computational power and access to more data, DRL methods are likely to continue improving. We anticipate the following directions for enhancing DRL in asset allocation:

Benchmarking and Reproducibility: The academic community will benefit from establishing standard benchmark data and environments for portfolio management tasks. This will enable fair comparisons between different algorithms and drive progress. Consistent evaluation protocols need to be developed, including testing under multiple market scenarios and considering transaction costs. As emphasized by Kruthof and Müller (2025), rigorous and robust validation protocols are essential[9]. Through community-agreed benchmarks, researchers can better identify which improvements truly lead to performance gains.

Contextualization and meta-reinforcement learning: Making DRL agents more "context-aware" can improve their robustness. For example, contextual reinforcement learning involves providing agents with information about the current market state. Agents can learn different sub-strategies for different contexts. Another example is meta-learning algorithms, which enable agents to learn how to quickly adjust their learning process when encountering new environments—essentially "learning how to learn." This is analogous to an investment strategy that knows when its conventional methods may fail and adjusts its behavior accordingly. Through meta-learning, an agent might adapt to entirely new market characteristics after observing only a small amount of new data, which would be highly valuable for handling sudden shifts in market structure.

Multi-agent reinforcement learning: Markets are composed of numerous interacting participants. Single-agent RL frameworks treat the market as part of the environment, but an interesting extension involves including multiple RL agents in the model that compete against each other or interact with models of other agents. Multi-agent RL can be used to simulate a market where some participants are also RL-driven, observing how they co-evolve. Although complex, this may reveal some equilibrium behaviors or demonstrate how RL agents can exploit or coordinate with other strategies. It can also be used to characterize the game-theoretic aspects of markets. Additionally, multi-agent frameworks can simulate investor-environment interactions, such as agents interacting with market makers or with several typical trading strategies, which may help in understanding strategy robustness.

Integration with traditional methods: Rather than viewing DRL as an alternative to traditional portfolio optimization, it is more productive to consider their integration. For example, a hierarchical approach could be adopted: RL agents determine high-level allocations, while traditional methods handle more granular-level allocations. Conversely, traditional models could guide the exploration of RL agents. This hybrid approach may combine the strengths of both—the theoretical robustness of classical models and the adaptability of RL.

Improved training algorithms: From an algorithmic perspective, methodological innovations in the field of reinforcement learning can also be applied to financial scenarios. For example, more stable training methods, more effective exploration strategies, or distributed RL are all worth exploring. For example, distributed RL can be used to

characterize the entire distribution of portfolio returns, thereby adapting to risk management objectives. Another example is safety-constrained reinforcement learning, a hot topic in recent years that aims to ensure that RL agents adhere to safety constraints during training. When mapped to asset allocation, this means ensuring that risks do not exceed certain thresholds or losses do not exceed certain thresholds during training, which helps to obtain strategies that are optimized under constraints and is beneficial for financial applications.

Actual deployment and feedback: Finally, an important development step is to deploy DRL strategies in real trading and obtain feedback. Real-world performance can inform research: Analyzing cases of strategy success and failure can help improve models. The growing interest in artificial intelligence within the financial industry may spur more collaboration between academia and industry, providing data, computing power, and expertise to advance the application of DRL in asset allocation. With real-market validation and data, researchers can calibrate models to better align with practical needs.

### 5.3 Conclusions

In summary, applying deep reinforcement learning to asset allocation provides a rich and flexible framework for this field, addressing the uncertainty and dynamics that traditional methods struggle with from a sequential decision-making perspective. Modeling asset allocation as an MDP enables us to unify classical methods with modern machine learning techniques, revealing that traditional solutions are merely special cases within the broader RL paradigm. Empirical studies provide encouraging evidence: carefully designed DRL agents can adapt to complex market patterns and sometimes outperform static or short-sighted strategies. These agents are capable of processing large information sets—including price history, asset interrelationships, and even sentiment data—and continuously adjusting their decisions as the environment changes, thereby endowing asset allocation decisions with context sensitivity and continuous optimization characteristics.

However, the current state of research also cautions us to remain cautious. Some of the successes reported in the literature may reflect learning from historical data and may not hold true in new data or different market environments. Issues such as high turnover rates, sensitivity to hyperparameters, and the need for large amounts of training data mean that a DRL strategy may require careful calibration and risk control before being deployed in real-money applications. Some studies have found that once real-world frictions are considered, a complex DRL algorithm does not outperform a simple equal-weight strategy—a thought-provoking reminder that in finance, strategy complexity does not automatically guarantee success; the true test lies in its robust generalization performance and ability to navigate market cycles.

Ongoing developments in this field hold promise for addressing many of the aforementioned challenges. As researchers integrate more financial domain knowledge into DRL frameworks and design algorithms specifically tailored to financial scenarios, we can anticipate more reliable and interpretable outcomes. The intersection of finance and reinforcement learning is highly promising: in the future, DRL-driven asset allocation systems may be able to adapt in real-time to market changes, strictly adhere to predefined risk profiles, and provide human managers with insights previously difficult to obtain.

Overall, the application of deep reinforcement learning in asset allocation represents a prime example of interdisciplinary innovation at the intersection of finance and artificial intelligence—blending financial theory, economic principles, and cutting-edge AI technologies. Despite ongoing challenges, the existing theoretical foundations and empirical evidence suggest that, with further refinement, DRL has the potential to emerge as a powerful tool for portfolio management, driving more adaptive and effective investment strategies in the coming years.

### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

### REFERENCES

- [1] Markowitz H M. Foundations of portfolio theory. The journal of finance, 1991, 46(2): 469-477.
- [2] Merton R C. Optimum consumption and portfolio rules in a continuous-time model. In Stochastic optimization models in finance. Academic Press, 1975: 621-661.
- [3] Sutton R S, Barto A G. Reinforcement learning: An introduction. Cambridge: MIT press, 1998, 1(1): 9-11.
- [4] Benhamou E, Guez B, Ohana J J. Deep Reinforcement Learning: Extending Traditional Financial Portfolio Methods. Available at SSRN, 2024.
- [5] Jiang Y, Olmo J, Atwi M. High-dimensional multi-period portfolio allocation using deep reinforcement learning. International Review of Economics & Finance, 20025, 98: 103996.
- [6] Aboussalah A M, Xu Z, Lee C G. What is the value of the cross-sectional approach to deep reinforcement learning? Quantitative Finance, 2021, 22(6): 1091–1111.
- [7] Chen S, Ge L. A learning-based strategy for portfolio selection. International Review of Economics & Finance, 2021, 71: 936-942.
- [8] Wei J, Yang Y X, Jiang M, et al. Dynamic multi-period sparse portfolio selection model with asymmetric investors' sentiments. Expert Systems with Applications, 2021, 177: 114945.

- [9] Kruthof G, Müller S. Can deep reinforcement learning beat 1N. Finance Research Letters, 2025: 106866.
- [10] Cong L W, Tang K, Wang J, et al. AlphaPortfolio: Direct construction through deep reinforcement learning and interpretable AI. Available at SSRN, 2021: 3554486.

### "TARO MEETS NEW FOOD GENERATION" —— MARKET INVESTIGATION OF DEEP-PROCESSED TARO PRODUCTS IN LIPU

RongJin Li

Department of Statistics, Guangxi Normal University, Guilin 541006, China. Corresponding Email: 1602647498@qq.com

**Abstract:** The old proverb "A steamed taro makes all the neighbors fragrant" proves the unique charm of Lipu taro. In the cloud mountain area of Lipu County, Guangxi, the century-old planting technology is realizing the modernization transformation of traditional technology through technological innovation. Nowadays, the deep processing of Lipu taro has formed a complete industrial chain. In recent years, the public awareness of the deep-processed products of Lipu taro has been very high, but the homogeneity of the products is serious, and the innovation is insufficient. There are still some people who have never bought or rarely bought them, and there is still great development potential in the consumer market. Therefore, this paper analyzes the factors that affect consumers to buy such products, and provides suggestions and future development marketing strategies for Lipu taro deep-processing products merchants.

In this study, by using SPSS statistical analysis, this paper investigates and analyzes the basic information of consumers of Lipu taro deep-processed products, as well as the factors influencing consumers' satisfaction and recognition of Lipu taro, users' emotion and stickiness. The questionnaire survey was carried out with the help of Tencent's questionnaire platform, and the quality of the questionnaire was controlled by reliability and validity analysis in the pre-survey stage, and the Cronbach coefficient was 0.969, which showed that the questionnaire data had good reliability and validity. Finally, 550 questionnaires were collected, 93 invalid questionnaires were eliminated, and 457 valid questionnaires were obtained, with an effective rate of about 83.09%.

Structural equation model was used to explore the influencing factors of consumption willingness of taro in Lipu. By using the grey relational analysis, it is concluded that there is a good correlation between the selected variables. Exploratory factor analysis is used to put forward path hypothesis for three different perceptions. According to the fitting results of the structural equation model, the fitness index has reached the ideal standard, indicating that the overall model has a good path fitting degree, in which two path assumptions, product quality and purchase perception, are established, and new strategies are formulated to influence customers' purchase intention from these two aspects. The K-means cluster analysis is used to classify the consumer groups, and the consumer groups are divided into developmental consumers, key consumers and potential consumers, which shows that different groups have different consumption levels and attitudes towards the deep-processed products of Lipu taro. According to their consumption habits and psychological characteristics, personalized marketing strategies are formulated respectively.

From the market point of view, this paper gives the reference strategy of product marketing model from three aspects: product sales promotion, target consumer groups and product elements promotion, and proposes to deepen market research to accurately locate demand, strengthen product innovation to create differential advantages, pay attention to brand building to enhance influence, strengthen marketing promotion to expand share, strengthen industrial chain cooperation and coordinated development, and inherit and carry forward the social responsibility and sustainable development of Lipu taro deep-processed products.

Keywords: Lipu taro; Deep processing products; Semantic network analysis; Structural equation model; K-means cluster analysis

### **1 INTRODUCTION**

### 1.1 Research Background and Significance

Lipu taro is a symbolic specialty of Lipu City, Guangxi Zhuang Autonomous[1], enjoying the reputation of "a treasure in the taro", plays a decisive role in the local agricultural economy. By 2023, the planting area of taro in Lipu City has reached 50,000 mu, the annual output has exceeded 100,000 tons, and the annual output value has exceeded 2 billion yuan. This industry has not only promoted the economic income growth of tens of thousands of farmers, but also become a key force to promote the rural revitalization strategy[2].

In Guangxi and South China, Lipu taro, as a necessary delicacy on the banquet table, carries profound cultural implications. However, a large number of counterfeits advertised as "Lipu taro" but not authentic origin have emerged in the market, which makes it difficult for consumers to identify the authenticity only by appearance or taste. What's more, some merchants confuse the genuine with the ordinary taro at a low price, which seriously infringes on the trust and rights of consumers. Lipu taro was once famous as a tribute in the Qing Dynasty, and it was well known to the public through the wide spread of the TV series Prime Minister Liu Luoguo[3]. Nevertheless, its rich historical and cultural connotation has not been fully inherited and carried forward. Among young consumers, Lipu taro is regarded as

a kind of "online celebrity food", but its value as a traditional symbol of local culture has not been given due attention and recognition.

With the upgrading of consumption and the development of agricultural industrialization, Lipu taro has extended from traditional fresh food to deep processing, and gradually formed a food industry chain with taro as the core raw material, which has become an important pillar industry for the revitalization of local villages. Enterprises such as "Guangxi Lipu Yuwang" and "Guilin Mingdian" have sprung up in Lipu, and integrated planting and processing through the model of "company+base+farmer"[4]. Some enterprises have introduced automatic production lines to improve processing efficiency, but the overall industrial concentration is still low, and small and medium-sized workshop-style processing plants account for a relatively high proportion.

Policy support and technical cooperation have contributed to the deep processing of taro products in Lipu. Guangxi has incorporated Lipu taro into the planning of characteristic agricultural industrial clusters, and local governments have provided subsidies to encourage research and development of deep processing technology, and built agricultural products processing parks (such as Lipu Food Industrial Park)[5]. Some enterprises cooperated with Guangxi Academy of Agricultural Sciences and South China Agricultural University to develop taro preservation, freeze-drying technology and extraction of functional components (such as polyphenols and polysaccharides), but the technical conversion rate needs to be improved. In order to support the development of the taro industry in Lipu, in 2023, the Lipu municipal government issued the policy of benefiting the people with the green demonstration planting award of Lipu in 2023, and farmers of Lipu taro can enjoy the 600 yuan award per mu. This policy stimulated farmers' enthusiasm for planting and expanded the planting area of Lipu taro.

### **1.2 Research Significance**

The significance of studying the deep-processed products of Lipu taro is far beyond the category of single industry, and it is a multi-dimensional integration of economic value, rural revitalization, cultural inheritance and green sustainable development. Deep processing can break through the limitation of short shelf life and high transportation cost, and transform Lipu taro into high value-added products, such as ready-to-eat food and functional raw materials, and improve the unit output value. It can also optimize the industrial chain structure, promote the standardization of upstream planting, drive downstream supporting industries, form an integrated model, and enhance risk resistance. As the "hometown of taro", the development of deep processing in Lipu City can attract investment and create jobs, such as factory workers and e-commerce operators, which will inject impetus into the high-quality development of county economy and promote local economic growth. Its taro industry accounts for over 30% of the agricultural output value.

Deep processing helps rural revitalization: guarantee the purchase price and increase farmers' income with the model of "company+cooperative+farmer"; Develop taro theme projects in combination with cultural tourism to promote the integration of rural three industries; Attract talents to return home to start businesses, improve rural infrastructure and narrow the gap between urban and rural areas.

Deep processing promotes sustainable development: recycling by-products, such as skin dregs and defective products, extracting polysaccharides for health care products, producing organic fertilizers or biofuels; Standardize planting, popularize ecological technology and protect local ecological diversity.

### **1.3 Literature Review**

Under the background of globalization, challenges such as population growth, resource constraints and climate change force the transformation and upgrading of agricultural products processing and food industry, and promote the iterative development of agricultural products processing and food innovation in the direction of green, intelligence and high added value by taking technological innovation as the path to drive industrial chain coordination, standardize production system construction and accurate adaptation of consumer demand[6]. Lei Yuliang, Xiao Lin (2024) through the spatial Dobbin model analysis shows that the agglomeration of agricultural products processing industry has a significant positive role in promoting agricultural modernization, and there is regional heterogeneity, and at the same time, it promotes the level of agricultural modernization in neighboring areas through the positive spatial spillover effect[7]. Guan Yonghua and Zhou Qiuhong (2023) think that the agricultural products processing industry in Guangxi has some problems, such as short industrial chain, low added value of technology and insufficient regional coordination. By strengthening the drive of scientific and technological innovation, deepening the integration of the three industries and improving the policy support system, the level of agricultural modernization in Guangxi can be effectively improved, and the rural revitalization and the high-quality development of the real economy can be helped[8].

With the development of Lipu taro industry, it is the premise of product promotion to explore the current public attitudes and views on Lipu taro and its deep-processed products. Bai Jian and Hong Xiaojuan (2022) through the emotional classification and theme analysis of online public opinion barrage[9]. And said that this method can show the emotional tendency and focus of attention of netizens in multiple dimensions; In addition, Ding Liuhua et al. (2023), based on the mining of online comments, used text analysis and social network analysis to conduct qualitative and quantitative research, thus constructing semantic networks and dividing them into different dimensions for further exploration[10].

How to combine Lipu taro with modern food and spread it to a larger market, and explore the factors affecting consumers of Lipu taro deep-processed products has become a top priority. Shi Xiaochen et al. (2023) set various

variables and established SEM structural equation model to verify the influence of various variables on consumers' purchase intention, and then explored the future development model of products and formulated corresponding strategies[11]; Cui Hongcheng and Chen Qingguo (2024) explored users' willingness to continue using products, tested the influencing variables by using the structural equation model, and optimized products and services according to the test results, and put forward relevant suggestions to improve users' satisfaction and willingness to continue using products[12].

Some scholars also have related research on product development and market extension. Qingliang Meng et al. (2014) explored the relationship between customer satisfaction and express service performance by integrating the improved Kano model and IPA analysis method, built the process model of express service quality detection, and gave the decision-making scheme of express service improvement[13]; Dujili et al. (2023) took the park as the research object, determined the priority of park landscape improvement order according to Kano-IPA analysis method through questionnaire data collection and information feedback, and put forward improvement suggestions respectively[14].

In the current era of big data, data mining technology can build a portrait of consumer groups, thus achieving accurate marketing of products. Cai Shaolin et al. (2022) based on K-Means clustering analysis, classified customers and mined information, and constructed a portrait of agricultural products consumers, thus accurately defining target customers, core products and marketing models[15]; Shi Lemeng et al. (2022), in order to explore the influencing factors of young consumers' purchasing behavior of sugar-reduced products, conducted preliminary market research based on questionnaire survey and literature research, obtained consumer portraits through K-Means cluster analysis, and then analyzed the factors and paths influencing consumers' purchasing tendency of sugar-reduced products through structural equation model, thus giving relevant suggestions on product sales methods and industry development trends[16].

### 1.4 Research Contents and Ideas

In order to understand the market situation of Lipu taro deep-processed products and put forward the future development direction, this paper is divided into six chapters, aiming at discussing the consumer research and market promotion of Lipu taro deep-processed products.

Firstly, the background and significance of Lipu taro are introduced. In order to provide direction and factual basis for the follow-up questionnaire design, this paper uses text mining to conduct a preliminary study on consumer demand. Through Python, Weibo's content with Lipu taro as the key word is crawled, and the word cloud image analysis and emotion analysis are carried out. The LDA theme is constructed to analyze the semantic network of high-frequency words, observe the PMI values between words and judge the core points that consumers pay attention to.

Secondly, based on the above conclusions, we designed the contents of the questionnaire, completed the distribution and collection of the questionnaire, and conducted a pre-investigation. Clean the collected data and analyze its reliability and validity to prepare for further statistical analysis. Determine the sample size, do a good job in quality control of questionnaire survey and strictly implement it. In field research, we use stratified sampling. Firstly, according to the grey relational analysis, we observed the correlation between the selected variables in the heat map, classified the characteristics through exploratory factor analysis, and put forward the path hypothesis for the structural equation model, so as to test the consumption intention. Finally, according to K-means clustering, the potential customers are classified and their characteristics are mined.

Finally, based on the research results and related marketing theories, we put forward the future development direction of Lipu taro deep-processing products, and at the same time provide relevant suggestions for merchants in product innovation, marketing strategy, channel expansion and so on.

### **1.5 Characteristics and Innovation**

### 1.5.1 Method innovation

The innovation of research methods is reflected in the application of big data and machine learning technology (mining online and offline behavior data of consumers) and visual analysis tools (such as Tableau and Power BI). These methods not only fully capture the characteristics of consumer behavior, but also reveal the key influencing factors and dynamic changes of the consumption willingness of Lipu taro deep-processed products through the structural equation (SEM) model. This multi-dimensional innovative research method provides scientific and accurate data support for research, and at the same time provides a new methodological paradigm for the study of agricultural product consumption behavior, which helps market strategy formulation and product innovation.

### 1.5.2 Perspective innovation

From the user's perspective, combined with consumer behavior characteristics, this paper focuses on analyzing consumers' personal experience (such as taste and health attributes) and preference trends (such as diversified products, cultural identity and environmental awareness) of Lipu taro processed products, and points out the shortcomings in the current market (such as homogenization, low brand awareness and insufficient promotion).

The potential users of non-Lipu taro deep-processed products were explored. In the process of data analysis, this paper not only analyzes the consumers' willingness to buy the deep-processed products of Lipu taro, but also makes a cluster study on the non-consumer groups who have not bought the deep-processed products of Lipu taro at present, so as to explore their future purchase willingness of the deep-processed products of Lipu taro and provide a series of improvement suggestions for expanding the product consumption market.

### 1.5.3 The concept of sustainability

Exploring the successful transformation of Lipu taro into deep processing industry not only optimizes the agricultural industrial structure, but also provides a new impetus for sustainable development for rural revitalization. Dynamic market observation and real-time data, through field research, consumer interviews and sales data analysis, capture the changes of market dynamics and consumption trends, and ensure the timeliness and practicality of research results.

### 2 RESEARCH SCHEME DESIGN

### 2.1 The Purpose of Research

In order to fully implement the spirit of the 20th National Congress of the Communist Party of China and thoroughly implement the spirit of General Secretary President Xi's important speech when he visited Guangxi, we should firmly implement the rural revitalization strategy, make overall plans to promote the rural revitalization strategy, and take the transformation and upgrading of agriculture, the prosperity of rural industries and the increase of farmers' income as the goal. As a national geographical indication agricultural product, Lipu taro is not only the most distinctive crop in Lipu city, but also an important driving force to promote the revitalization of local villages and increase farmers' income. Its unique quality and geographical advantages have injected vitality into the development of agricultural economy in Lipu City and become a key industry for farmers to get rich. The development of Lipu taro industry can not only enhance the market competitiveness of local agricultural products, but also provide strong support for the implementation of rural revitalization strategy and help farmers achieve sustained income increase.

Therefore, in order to realize industrial transformation and upgrading, we will speed up the adjustment of industrial structure and further strengthen the development of characteristic industries in Lipu taro around the principle of "market-led, government-driven, diversified investment and characteristic industries". By collecting consumers' ratings on the factors influencing the cognition and consumption habits of Lipu taro deep-processed products, we can formulate more targeted product development and marketing strategies and continuously improve the market competitiveness of Lipu taro price products. The results of the questionnaire survey are also helpful to provide enterprises with the direction of optimizing product quality, improving sales channels and enhancing brand value.

### 2.2 Research Object and Scope

The subjects of the questionnaire are people with different consumption, people of different age groups and people in different regions. Through the investigation of different consumer groups, we can understand the acceptance, consumption habits and purchasing power of Lipu taro deep-processed products. To investigate the preferences of different age groups for taste, nutritional value and price of Lipu taro, and to investigate the influence of regional culture, eating habits and other factors on the deep-processed products of Lipu taro.

This survey will cover the major consumption areas of Lipu taro in China, including 21 major first-tier cities including Beijing, Shanghai, Guangzhou and Shenzhen, 30 second-tier cities including Nanjing, Chengdu, Wuhan and Changsha, and 43 third-tier cities including Guilin, Dali and Yichang. In addition, Guilin, Guangxi, the specific production place of Lipu taro, will be the key research area of this survey.

### 2.3 Research Arrangements

The main objective of this survey is to deeply explore the present situation, advantages, market demand and potential in rural revitalization of the deep-processed taro industry in Lipu, Guangxi, especially to make substantial progress on the key issues of upgrading the brand of deep-processed taro products in Lipu, expanding market share and promoting industrial development. Through this investigation, we hope to provide feasible strategic suggestions for the future development of Lipu taro deep processing industry, so as to better promote local economic growth and increase farmers' income, and provide support for the implementation of rural revitalization strategy.

In recent years, with the continuous promotion of the country's rural revitalization strategy, the rural economy is undergoing profound changes. Agricultural products industry, especially characteristic agricultural products, has gradually become an important force to promote rural revitalization. As a local characteristic agricultural product, Guangxi Lipu taro has unique variety advantages and market potential, which has become an important part of industrial revitalization. Therefore, through this survey, we hope to deeply understand the production, processing, sales and opportunities and challenges of Lipu taro deep-processed products, and provide data support and theoretical basis for the improvement and branding of its industrial chain.

In the initial stage of the investigation, we first understand the overall development trend of agricultural products industry under the background of rural revitalization through the hot evaluation and analysis of the network platform. In particular, consumers' interest and demand for local specialty agricultural products are also increasing year by year. This kind of information helps us to determine the focus of research on the deep-processed products of Lipu taro, thus exploring the market demand, consumer awareness and consumption habits of the deep-processed products of Lipu taro, and providing effective suggestions for the future development of Lipu taro industry.

In the initial stage of investigation, we collected and analyzed all kinds of documents related to the taro industry in Lipu, which provided us with an overview of the development status of the deep processing industry of Lipu taro, including the distribution of producing areas, main production enterprises, processing technology and market sales. By analyzing the existing data of Lipu taro deep processing industry, we made clear the research direction and specific objectives, and ensured the smooth development of the follow-up research work. At the same time, we also actively paid attention to the policy documents related to rural revitalization and the state's support policies for the development of agricultural industry. These policies provide a good environment for the sustainable development of agricultural products industry, and also provide an opportunity for the upgrading of Lipu taro deep processing industry.

In the middle stage of investigation, we went deep into the production base of Lipu taro in Guilin, Guangxi, and visited the main planting areas of Lipu taro and related agricultural cooperatives. Through face-to-face communication with local farmers and enterprises, we learned about the production of Lipu taro, including planting scale, yield, quality control, harvesting and processing. In addition, we also conducted a detailed investigation on the upstream and downstream links in the industrial chain of Lipu taro, and learned the whole process from planting, processing to sales. In order to further understand the demand and consumption trend of consumers, we designed a questionnaire survey on the consumption market of Lipu taro and its deep-processed products. By distributing questionnaires in different regions and different consumer groups, we have obtained a lot of data about consumers' preferences and purchase frequency of Lipu taro deep-processed products. These data provide valuable reference for our subsequent market analysis and product positioning.

In the final stage of the investigation, we made a detailed analysis of the collected data. First of all, through the statistical analysis of the questionnaire data, we understand the consumer's awareness, purchase perception, community influence, purchase intention and the changing trend of market demand for Lipu taro deep-processed products. Secondly, using data analysis tools such as SPSS and Python, we made a multi-dimensional analysis on the production cost, market price fluctuation and consumer behavior of Lipu taro, and revealed the main challenges and development opportunities faced by the deep processing industry of Lipu taro at present.

Based on the preliminary understanding of the market of Lipu taro deep-processed products, we have made plans for the follow-up questionnaire design, distribution and data analysis.

### **3** INVESTIGATION AND IMPLEMENTATION

### **3.1 Preparatory Preparation**

### 3.1.1 The market scale of modern Lipu taro deep-processed products is good

According to statistics, by 2023, the planting area of taro in Lipu City, Guangxi Province exceeded 200,000 mu, and the annual output exceeded 300,000 tons, accounting for more than 15% of the total national taro production. A "planting-processing-selling" integrated industrial chain has been formed in the local area, and the number of deep processing enterprises has increased to more than 200, with an annual processing capacity of 100,000 tons, driving the output value of the whole industrial chain to exceed 5 billion yuan. In order to tap the consumption trend and capture the emerging demand, through Baidu index The platform analyzes the behavior data of netizens.



Figure 1 Trend Chart of Search Volume for "Lipu Taro" Image source: Baidu Index (https://index.baidu.com/v2/index.html#/)

In Baidu index, search, filter and analyze with "Lipu taro" as the key word to obtain its search index trend chart, such as Figure 1. As shown. As can be seen from the figure, the search index of "Lipu Taro" was relatively high in December 2024. By January-March 2025, although the search index was not so high, there was a relatively stable search index every day, which indicated that people continued to pay attention to the Lipu Taro market steadily.

Similarly, through the search and analysis of "Lipu taro" as the key word, the demand map is obtained. The correlation between netizens' needs and keywords is divided into three layers, and the intensity gradually weakens from the inside out, in which the circle size represents the search index size, green represents the index decline, and red represents the index increase. It can be seen that the words with the strongest correlation with "Lipu taro" are Xiangsha taro, Yuhuazhai street, accounting service company, etc., which shows that netizens are extremely accidental in searching for

this keyword, and we need to further understand its development status.

In the wave of upgrading the deep processing industry of agricultural products in the past 20 years, the taro industry in Lipu has undergone a transformation from traditional rough processing to high value-added deep processing, and gradually built a modern industrial chain system covering raw material supply, product innovation and brand marketing, guided by the evolution of consumers' demand for health, convenience and multiple scenarios. Nowadays, consumers' attention to "low sugar and low calorie", "instant convenience" and "functional attributes" continues to heat up, which promotes the extension of Lipu taro deep-processed products from single raw materials to snack, prepared food, healthy meal replacement and other sub-fields, which not only activates the market potential of traditional taro products, but also analogizes innovative products such as quick-frozen taro paste and taro chips into new tea, baking and healthy food tracks. Under this background, relying on the rural revitalization policy to support characteristic agricultural products, Lipu taro industry is exploring a new path of "geographical indications+industrial innovation" to drive the nationalization of regional brands through the integration of deep processing technology empowerment and consumption scenarios, providing a demonstration sample for the transformation and upgrading of Guangxi characteristic agriculture.

### 3.1.2 Lagging of the development status of Lipu taro deep-processed products

In this research, consumers who are willing to buy deep-processed products of Lipu taro are taken as the research object. Because of limited funds and short time, we choose to distribute questionnaires directly online, including online media such as Questionnaires, QQ, WeChat, etc., so as to obtain questionnaire data quickly. In recent years, although the deep processing industry of Lipu taro has made some progress under the support of large-scale planting and policies, its production infrastructure construction is still lagging behind, which is difficult to meet the needs of high-quality development of the whole industrial chain. The concrete manifestations are: insufficient iteration of processing technology, weak standardization and quality control[17]. And lack of market plan and brand premium. In addition, small and medium-sized processing enterprises are generally faced with the problem of shortage of funds. Nearly 80% of enterprises have fallen into the "low-end OEM" mode because they are unable to invest in the renovation of aseptic workshops and the upgrading of intelligent production lines, and their new product research and development and brand marketing capabilities are seriously insufficient. According to the industry survey data in 2023, the output value of deep processing of taro in Lipu only accounts for 9.2% of the total scale of taro products in China, which is seriously out of balance with its position of raw material output accounting for 15% in China. In the Baidu index, search for the term "Lipu taro" to get the regional and age distribution map.



**Figure 2** Regional distribution of search index for "Lipu Taro" Image Source: Baidu Index (https://index.baidu.com/v2/index.html#/)

By Figure 2, it can be seen that the areas with high search index of "Lipu Taro" are mainly Guangdong, Guangxi, Shandong, Jiangsu, Zhejiang, Fujian and other places, mostly in the eastern coastal areas, which shows that the propaganda of Lipu Taro in Guangxi is not in place and its popularity is not high.



**Figure 3** Age distribution of search index of "Lipu taro" Image source: Baidu Index (https://index.baidu.com/v2/index.html#/)

By Figure 3, it can be seen that most of the search groups of "Lipu Taro" are 30-39 years old, followed by young people aged 20-29. Because when TGI is equal to 100, this kind of users' attention to a problem is at an average level, and when TGI is greater than 100, it is higher than the overall level, so it can be seen that people aged 20-39 pay more attention to Lipu taro than the average level, which can be used as our research object.

### **3.2 Questionnaire Design**

In order to make the questionnaire achieve ideal results and convey detailed and accurate data information, in the statistical investigation report, the questionnaire design should follow the following principles: First, the functional principles, including the principles of consistency, completeness, accuracy and feasibility, ensure that the questionnaire questions are clearly expressed, the questions are consistent, avoid ambiguity and make the respondents easy to understand. Second, the principle of neutrality, question design should remain neutral, avoid guiding respondents' answers, and ensure the objectivity of data. Third, the principle of orderliness, the order of arranging questions should be organized, from general to specific or from simple to complex, in order to improve the participation of respondents. Fourth, comprehensive coverage, the questionnaire should comprehensively cover the main aspects within the scope of research to ensure comprehensive information. Fifth, the principle of adaptability, according to the purpose of the study and the characteristics of the object, design questions to ensure the adaptability and effectiveness of the questionnaire, and use the simplest inquiry method under the condition of ensuring the same information. Sixth, the principle of pre-test, pre-test before the formal implementation, test the rationality and effectiveness of the questionnaire through small-scale investigation, and adjust the questions in time.

Combined with the specific investigation, this paper explores consumers' willingness to consume the deep-processed products of Lipu taro to determine the questionnaire structure.

### **3.3 Pre-Investigation**

After the completion of the questionnaire, in order to ensure the effectiveness and quality of the questionnaire, the investigation team conducted a pre-survey on the questionnaire before it was officially distributed. A total of 155 questionnaires were collected through one-on-one pre-survey, of which 150 were valid.

### 3.3.1 Reliability analysis

Reliability is reliability, which refers to the consistency of the results when the same method is used to measure the same object repeatedly. Reliability indicators are mostly expressed by correlation coefficients, which can be roughly divided into three categories: stability coefficient (cross-time consistency), equivalence coefficient (cross-form consistency) and internal consistency coefficient (cross-project consistency). According to the characteristics and specific requirements of the survey, we finally use Cronbach's alpha to analyze the reliability. This method is suitable for the reliability analysis of attitude and opinion scale (questionnaire). The final calculated results are shown in the following table.

Table 1 Pre-investigation reliability test						
Variable name Measurement coefficient Cronbach's Alpha						
external factor	five	0.946				
internal factor	four	0.952	0.969			
Other factors	three	0.900				

By Table 1, it can be seen that the overall  $\alpha$  coefficient value (Cronbach's  $\alpha$  value) is 0.969, which is greater than 0.9, and the  $\alpha$  coefficient of a single question is also greater than 0.9, which is acceptable, so it shows that the questionnaire

has strong internal consistency and stability.

### 3.3.2 Validity analysis

Validity refers to the accuracy of the measurement tool, that is, the degree to which the measurement results can reflect the characteristics to be measured. The higher the validity, the better the purpose of the questionnaire test can be achieved. In measurement theory, validity is defined as the ratio of the true variance related to the purpose of measurement to the earned score variance in a series of measurements. The validity analysis results of the pre-survey data are shown in the following table (Table 2).

Ta	<b>Cable 2</b> KMO and Bartlett sphericity test in pre-investigation						
	KMO	O value	0.806				
		Approximate chi-square	501.684				
	Bartlett sphericity test	df	66				
		P value	0.00				

Secondly, the exploratory factor analysis of the questionnaire data shows that the overall KMO value is 0.806, which is acceptable, and the significance is less than 0.01. The variables are independent, so it is suitable for factor analysis and the questionnaire content is valid. Therefore, for the questions fed back in the pre-investigation, the content and structure of the questionnaire are further improved and optimized.

### 3.4 Sampling and Samples

#### 3.4.1 Determine the sample size

According to the actual situation, this survey mainly adopts stratified sampling to conduct questionnaire survey because consumers in various cities are the research objects. According to the sample size estimation formula of sampling survey, the sample size of this survey is calculated.

$$\mu_0 = \frac{Z^2 * P(1-P)}{E^2} \tag{1}$$

In general, the acceptable sampling limit error E in the above formula is 0.045, the confidence is 95%, and the estimated proportion P is 0.5. After calculation, the sample size is about 475.

According to the pre-survey distribution, we estimate that the efficiency of the questionnaire is about, so the questionnaire should be distributed. $n_0 = n_0/a \approx 495$ 

According to the general experience, it is estimated that the recovery rate of the questionnaire is 90% under the explanation of online researchers, so the sample size is adjusted to  $n_2 = n_1/b \approx 550$ 

For the convenience of statistical processing, we finally decided to distribute 550 questionnaires.

r

### 3.4.2 Sampling method

As the scope of this survey covers consumers in cities at various levels in China, and the designed population is large, stratified random sampling is conducted to ensure the effectiveness of sampling. The stratified standard is based on the data of the seventh national census, and the total resident population in eastern, central, western and northeastern regions is selected, and the sample size of each region is calculated by proportion, as shown in the following table (Table 3).

Table 3 Stratified Sampling Ratio by Region							
region Permanent population (10,000) Sampling proportion Sampling							
eastern region	56372	0.3993	220				
middle	36469	0.2583	142				
the west	38285	0.2712	149				
northeast	98515	0.0698	39				

### 3.5 Data Processing and Entry

A total of 550 questionnaires were distributed through online questionnaire survey and offline passers-by interview interception. After excluding invalid questionnaires such as logical contradictions, unanswered questions and short filling time, 457 valid questionnaires were finally obtained, with an effective rate of 83.09%, which met the statistical requirements. Input the questionnaire data through EXCEL, and analyze its reliability and validity through SPSS. The results are shown in the following table:

Table 4 Fo	rmal Investigation Reliab	ility Test Form
Variable name	Measurement coefficient	Cronbach's Alpha

			1
external factor	five	0.964	0.077
internal factor	four	0.965	0.977

Other factors	three	0.930	
			_

By Table 4, it can be seen that the Cronbach's alpha value of the result is 0.977, which means the reliability of the result is 97.7%. If it is greater than 0.8, it means that the reliability of the questionnaire is good, so we can continue statistical analysis.

The validity of the valid data of this formal investigation is analyzed, and KMO and Bartlett sphericity tests are used to get the results in the following table.

Table 5 Formal survey validity test form					
KMO	KMO value				
	Approximate chi-square	1265.588			
Bartlett sphericity test	df	91			
	P value	0			

By Table 5, it can be seen that the KMO value of the questionnaire data is 0.864, which is greater than 0.8; And the significance level is 0.000, less than 0.05. Then the synthesis can show that the questionnaire has passed the validity test and the expected results can be obtained.

### **3.6 Quality Control**

We mainly use online questionnaire survey to collect data, and there may be some phenomena such as random filling and random crossing, which will affect the overall data quality. Therefore, we mainly adopt the following methods for quality control:

Quality control of research planning: after determining the research topic, we know the scale and development status of the current Lipu taro deep-processing product market through online inquiry and market report, and consult relevant documents about the Lipu taro deep-processing industry, so that all members can think deeply about the research content and purpose, ensure that each member has a clear and unified idea about this research, and ensure the integrity and efficiency of the research work.

Quality control of questionnaire design: after consulting a large number of documents related to the research topic, communicate with the tutor to determine the questionnaire survey method as the main method of this formal research. When designing the questionnaire topic, make the questionnaire topic concise, logical and systematic. In addition, invalid answering questions are set on the "Tencent Questionnaire" platform to prevent the data from being unreliable due to the respondents' careless answers, so as not to affect the whole research work.

Quality control of questionnaire pre-investigation: before formal investigation, small-scale questionnaires are distributed for experiments, and whether the contents of the questionnaires are reasonable or not are checked through trial filling, and the problems existing in the questionnaires are corrected according to the feedback of the respondents, so as to improve the contents and order of the questionnaire topics again.

Quality control of formal investigation: In order to avoid data errors caused by respondents' cognitive errors, appropriate explanations should be given in time to improve the authenticity and accuracy of questionnaire answers. All questionnaires were screened, and the questionnaires with obviously inconsistent answers and logical contradictions were treated as invalid.

Quality control after investigation: after the investigation, the questionnaire will be reviewed twice. When more invalid questionnaires are found, more questionnaires need to be distributed to ensure the sample size, and the collected data will be sorted and summarized in time until the required sample size is met. Before the data analysis, the reliability and validity were tested to ensure that the questionnaire data were reliable and valid.

### 4 Descriptive Statistic

### 4.1 Analysis of the Basic Situation of the Respondents

In the gender distribution of the respondents, the proportion of men is 40.9%, and that of women is 59.1%. On the whole, the proportion of men and women is about 4:6, and the gender ratio is not much different. In the age distribution, the proportion of people aged 18-30 is 54.3%, and the proportion of people aged 31-45 is 35.2%. Most of the respondents are young, which is beneficial to our questionnaire survey of deep-processed products.



Figure 4 Gender and Age Distribution Map

In the regional distribution of the respondents (Figure 4), the provincial capital accounts for the largest proportion, accounting for 44.4% of the total, followed by districts and counties, accounting for 34.4% of the total, and rural areas account for 7%, which is the least number of all regions. Reflects the distribution of respondents in different regions. In the distribution of living expenses or wages, the number of people with expenses of 2000-4000 accounts for the most, accounting for 34.7%, followed by the number of people with expenses of 4000-6000, accounting for 26.7%, the number of people with expenses of more than 6000 accounts for 22.8%, and the number of people with expenses of less than 2000 accounts for 15.8%.



Figure 5 Distribution Map of Area and Living Expenses

In the distribution of working status, the number of employees is the largest (Figure 5), accounting for 64.3%, followed by students, accounting for 19%. The respondents are more suitable for people who know about deep processing. In the distribution of education level, undergraduate or junior college accounts for the largest proportion, up to 72.9%, followed by high school or technical secondary school. This group of people are willing to spend money on food and are more in line with the respondents (Figure 6).



Figure 6 Distribution Map of Working Status and Education Level



4.2.1 Respondents know the channel of Lipu taro



Figure 7 Respondents Know the Channel of Lipu taro

Such as Figure 7 As shown in the figure, it can be known that the channels for respondents to learn about Lipu taro are mainly through e-commerce platforms and social media, which shows that network communication plays a positive role in Lipu taro. Of course, through the recommendation of family and friends, TV advertisements and offline supermarkets, Lipu taro can also be spread to a certain extent, making the communication efficiency higher and letting more and more people know about this variety. In order to improve the market competitiveness of Lipu taro and promote the development of its deep-processed products, producers should make use of these or even more channels to publicize and improve their popularity.

### 4.2.2 The frequency of respondents buying deep-processed products of Lipu taro

According to the data in Figure 8 43.7% people buy the deep-processed products of Lipu taro 1-3 times a month, only 11.4% people buy the deep-processed products 1-3 times a week, and even 15.3% people never buy them. This means that most people are still willing to buy the deep-processed products of Lipu taro, but not all people like to buy them, and the time interval of purchase is determined by personal needs.



Figure 8 Frequency of Purchasing Deep-Processed Products by Respondents

### 4.2.3 Products that respondents prefer to buy

Such as Figure 9 as shown, the number of people who like to buy ready-to-eat snacks (such as taro strips, taro cake, etc.) and drinks (such as taro milk tea, taro plant milk, etc.) is the largest, with little difference, accounting for 27% and 29.2% respectively. However, fewer people buy condiments (such as taro sauce), only 6.6%. It can be found that, in fact, most people prefer dessert-like processed products, so positioning products on making rich and varied desserts can attract more customers, expand the market scale and promote the development of deep processing industry of taro in Lipu.



### 4.2.4 The purpose of the respondents to purchase the deep-processed products of Lipu taro

FromFigure 10, it can be seen that most people buy deep-processed products mainly for their daily consumption, and 33.9% people are willing to buy them to try new flavors, which shows that Lipu taro deep-processed products are more common in daily life. At the same time, 26.3% people choose deep-processed products as gifts, which shows that the deep-processed products of Lipu taro have a wide audience. For enterprises of this kind of products, different packaging can be designed according to customers' needs, and consumers can freely choose and buy goods according to the purpose of purchase when purchasing, and promote marketing according to the classified packaging strategy.



Figure 10 The Purpose of the Respondents to Buy Deep-Processed Products

### 4.3 Lipu Taro Deep-Processed Products Consumption Trend Analysis

### Have you ever heard of Lipu taro?

Regarding whether or not you have heard of Lipu taro, according to the data collected by the questionnaire, 81.4% of people have heard of Lipu taro, and only 18.6% have not, which shows that the publicity and promotion of Lipu taro are not bad. Enough people know about Lipu taro, which provides an opportunity for us to follow up on the development of its deep-processed products (Figure 11).



Figure 11 Have you Heard of Lipu Taro

For the respondents who have never bought the deep-processed products of Lipu taro.

### 4.3.1 The main factors of not buying deep-processed products

Such as Figure 12 as shown in the figure, there are three main reasons why the respondents didn't buy the deep-processed products of Lipu taro: they have doubts about the quality or safety, they don't know about such products and the price is too high, which account for 24.3%, 22.7% and 22.7% respectively, while 11% people prefer to buy fresh taro, which shows that it is very important to strengthen the promotion of its deep-processed products. Enterprises should introduce to the public that the deep-processed products still retain many nutrients and unique taste of Lipu taro, invite food bloggers to try and recommend products, and carry out promotional activities.



Figure 12 Main Factors of not Buying Deep-Processed Products

### 4.3.2 Factors that make respondents who have not bought try to buy

Such as Figure 13, it can be seen that 32.3% people are willing to buy if they try to eat offline; 26.9% people are willing to buy if there is a clear product efficacy description. In order to encourage more people to buy the deep-processed products of Lipu taro, improve the market competitiveness and expand the market scale, enterprises can choose suitable venues, invite others to try them on the premise of ensuring the quality of the tried products is consistent with the products sold, or distribute information listing the efficacy, ingredients and applicable people of the products.



Figure 13 Factors of Trying to Buy

In view of the respondents who have purchased the deep-processed products of Lipu taro.

### 4.3.3 External factors considered when purchasing deep-processed products

By Figure 14, it can be seen that consumers consider many factors when purchasing the deep-processed products of Lipu taro. First of all, whether there are promotional activities will greatly affect the public's choice. The preferential strength of promotional activities can attract customers who have never tried the product before, and can also give back to old customers, and customers can enjoy lower prices or more gifts. Secondly, the price of products is also a key factor for consumers to consider, and many consumers pursue high cost performance when buying products. The convenience of purchasing channels also enables consumers to obtain goods quickly and meet the current demand. In addition, the brand effect is also concerned by consumers, and having a good reputation will win the trust and love of consumers.



### 4.3.4 Internal factors considered when purchasing deep-processed products

By Figure 15, it can be seen that from the four aspects of taste, nutritional value, quality and shelf life, consumers should first consider the taste when purchasing the deep-processed products of Lipu taro. The taste is the first thing consumers feel when tasting food, which directly affects consumers' acceptance and love of the products. Secondly, quality and shelf life are the basic requirements of consumers for food safety and stability; Finally, the nutritional value is considered. Lipu taro itself is rich in various nutritional components, and it is more inclined to be used as a delicious snack or convenience food when purchasing its deep-processed products.

H. <mark>3</mark> % 21.1%	4!	5.4%	30.9%	
- 10.8%	30.4%	31.7%	23.7%	
- 9.8%	30.1%	35.1%	22.2%	

Figure 15 Social Factors Affecting the Purchase of Deep-Processed Products

### **5 MODEL APPLICATION**

# 5.1 Based on SEM Model, Explore the Influencing Factors of Consumption Willingness of Lipu Taro Deep-Processed Products

### 5.1.1 Correlation analysis of grey correlation degree

Question 16 (Would you like to recommend Lipu taro deep-processed products to your friends or family?) and question 18 (18. Are you willing to pay a higher price for the high-quality deep-processed products of Lipu taro?) and questions 13-15 are analyzed by grey correlation degree to explore the correlation between variables (Figure 16).

0 -	1.00	0.83	0.84	0.86	0.86	0.89	0.85	0.87	0.85	0.79	0.78	0.85	0.69	0.75	- 1.00
1 -	0.83	1.00	0.84	0.83	0.85	0.83	0.85	0.83	0.84	0.80	0.79	0.81	0.63	0.67	- 0.95
2 -	0.84	0.84	1.00	0.83	0.82	0.64	0.83	0.85	0.83	0.81	0.77	0.84	0.67	0.75	
3 -	0.86	0.83	0.83	1.00	0.86	0.86	0.86	0.87	0.87	0.79	0.78	0.85	0.66	0.71	- 0.90
4 -	0.86	0.85	0.82	0.86	1.00	0.85	0.84	0.87	0.86	0.79	0.79	0.65	0.70	0.74	- 0.85
5 -	0.89	0.83	0.84	0.86	0.88	1.00	0.85	0.90	0.88	0.77	0.76	0.87	0.71	0.74	
6 -	0.85	0.85	0.83	0.86	0.84	0.85	1.00	0.86	0.86	0.78	0.78	0.63	0.65	0.71	- 0.80
7 -	0.87	0.83	0.85	0.87	0.87	0.90	0.86	1.00	0.87	0.77	0.78	0.88	0.70	0.74	
8 -	0.85	0.84	0.83	0.87	0.86	0.88	0.86	0.87	1.00	0.79	0.79	0.85	0.67	0.71	- 0.75
9 -	0.79	0.80	0.81	0.79	0.79	0.77	0.78	0.77	0.79	1.00	0.85	0.81	0.59	0.68	- 0.70
10 -	0.78	0.79	0.77	0.78	0.79	0.76	0.78	0.78	0.79	0.85	1.00	0.79	0.56	0.66	
11 -	0.85	0.81	0.84	0.85	0.85	0.87	0.83	0.88	0.85	0.81	0.79	1.00	0.66	0.71	- 0.65
12 -	0.69	0.63	0.67	0.66	0.70	0.71	0.65	0.70	0.67	0.59	0.56	0.66	1.00	0.75	- 0.60
13 -	0.75	0.67	0.75	0.71	0.74	0.74	0.71	0.74	0.71	0.68	0.66	0.71	0.75	1.00	
	913_1 -	913_2 -	013_3 -	Q13_4 -	913_5 -	Q14_1 -	Q14_2 -	q14_3 -	Q14_4 -	- 1_210	q15_2 -	- E_210	- 910	Q18 -	

Figure 16 Heat Map of Correlation between Variables based on Grey Correlation Analysis

After dimensionless treatment of each variable by Python, the grey correlation values are calculated respectively, and the heat map as shown in the above figure is drawn[18]. According to the above figure, it can be found that there is a good correlation between the variables in this analysis.

### 5.1.2 Exploratory factor analysis

In order to explore the influencing factors of consumption willingness of deep-processed taro products in Lipu, 14 variables were introduced, namely question 13, question 14, question 15, question 16 and question 18.

Firstly, the exploratory factor analysis (EFA) method is used[19]. The principal component analysis of the selected sample data is carried out by SPSS software, and the factor load matrix is obtained by the maximum variance method, and the KMO and Bartlett test values are output to preliminarily determine whether this data is suitable for factor analysis, and to judge whether the factor structure determined in this paper is reasonable (Table 6).

Table 6 KMO and Bartlett S	phericity	Test in Ex	ploratory	/ Factor Ar	alysis
					~

KMO	metric	0.868
	Approximate chi-square	1097.925
Bartlett's sphericity test	df	55
	sig.	0.00

The test value of KMO and Bartlett is 0.868, which is higher than the threshold of 0.8, indicating that it is very suitable for factor analysis, and the value of sig. is 0.00, which is lower than the threshold of 0.5, so it can be concluded that the selected questionnaire sample data has a good modeling matching degree. The factors are extracted by dimensionality reduction. When four factors are extracted, the cumulative variance explanation rate is 90.7%, and the variance of each common factor is between 0.5 and 0.9, indicating that four factors can better explain the questions and variables in the questionnaire. Thus, the rotated component matrix is obtained, as shown in the following table (Table 7):

Table	7	Rotated	Componer	nt	Matrix
			1		

variable	ingredient			
variable	one	2	three	four
Q13 1: Price	0.831			
Q13 2: Packaging	0.803			
Q13_3: Brand Effect	0.761			
Q13 4: Channel Convenience	0.719			
Q13_5: Promotion Activities	0.698			
Q14_1: Taste		0.842		
Q14_2: Nutritional value		0.816		
Q14_3: Quality		0.779		
Q14_4: Shelf life		0.748		

66	RongJin Li	
Q15 1: Promotion of celebrity endorsement	0.885	
Q15_2:IP joint name	0.869	
Q15_3: Recommended by family and friends	0.818	
Q16: Would you like to recommend the deep-processed products of Lipu taro to your friends or family?	0.891	
Q18: Are you willing to pay a higher price for the high-quality deep-processed products of Lipu taro?	0.876	

Through the above analysis, and according to the item semantics and factor load size, we can know from the rotation component matrix table that we can get a four-factor model and name the common factors respectively:

Purchasing perception: Q13\_1, Q13\_2, Q13\_3, Q13\_4, Q13\_5. Product quality: Q14\_1, Q14\_2, Q14\_3 and Q14\_4. Community influence: Q15\_1, Q15\_2, Q15\_3. Consumption intention: Q16, Q18

### 5.2 Study the Path Hypothesis

According to the relevant literature research, this paper discusses the influencing factors of consumers' willingness to accept Lipu taro products combined with modern food. According to the analysis of influencing factors of product satisfaction by Xiong Wenzhen and Xu Jianxin, it is found that there is a significant positive effect on product price, packaging, taste and aroma when evaluating product purchase and satisfaction[20]. Therefore, combined with the results of exploratory factor analysis, this paper puts forward the following assumptions:

Hypothesis H1: Purchasing perception positively significantly affects consumers' willingness to consume Lipu taro products.

Suppose H2: the product quality positively and significantly affects consumers' willingness to consume Lipu taro deep-processed products.

In addition, this study also focuses on consumers' perception of community influence to promote their willingness to buy agricultural products. Zhang Qiyao and Li Na's empirical research based on the four-factor model of perceived value shows that perceived value has a significant positive impact on consumers' purchase intention of deep-processed agricultural products brands through the intermediary role of brand trust, while the negative adjustment of brand appeal of fresh agricultural products on the relationship between brand trust and purchase intention is not significant. Therefore, combined with the results of exploratory factor analysis, this paper puts forward the following assumptions:

Suppose H3: community influence positively and significantly affects consumers' willingness to consume Lipu taro deep-processed products.

### 5.3 Fit Test and Fitting Results of Equation Model

In this study, chi-square freedom ratio (CMIN/DF), approximate root mean square error (RMSEA), growth fitness index (IFI), Tucker-Lewis index (TLI) and comparative fitness index (CFI) were used to evaluate the fitting results. The fitting values of this study are as follows (Table 8):

Table 8 Test Results of Modified Model Fitness					
index	reference standard	Measured results	Fitness evaluation		
CMIN/DF	1/3 is excellent, $3/5$ is good.	2.624	excellent		
RMSEA	< 0.05 is excellent, $< 0.08$ is good.	0.065	good		
IFI	> 0.9 is excellent, $> 0.8$ is good.	0.904	excellent		
TLI	> 0.9 is excellent, $> 0.8$ is good.	0.854	good		
CFI	> 0.9 is excellent, $> 0.8$ is good.	0.901	excellent		

According to the fitting results of the model constructed in this study, except RMSEA (root mean square error) and CMIN/df (chi-square freedom ratio) are acceptable standards, the other indicators GFI, RMR, IFI and CFI all meet the ideal standards, indicating that the model constructed in this study has a good overall fitting degree and can be further analyzed.

Table 9 Model Path Fitting Result					
Path relation	Standard path coefficient	S.E.	C.R.	Р	
Consumer Willingness <-Purchase Perception	0.730	0.082	8.923	***	
Consumer willingness <-product quality	0.795	0.113	7.044	***	
Consumption intention <-community influence	-0.24	0.103	-1.761	0.078	
According to Table 9, it can be seen that the product quality positively affects the consumption intention ( $\beta$ =0. 795, p<0.05) in the path hypothesis relationship test of this study, so the hypothesis H1 holds; Purchase perception positively affects consumption intention ( $\beta$ =0.730, p<0.05), so it is assumed that H2 holds; Community influence negatively affects consumption intention ( $\beta$ =-0.24, p>0.05), so it is assumed that H3 is not valid.

## 5.4 Test Results and Summary

Table 10 Hypothetical Test Results	
research hypothesis	Research conclusion
H1: Purchasing perception positively and significantly affects consumers' willingness to consume the deep-processed products of Lipu taro.	found
H2: The product quality has a positive and significant impact on consumers' willingness to consume Lipu taro deep-processed products.	found
H3: The community influence positively and significantly affects consumers' willingness to consume the deep-processed products of Lipu taro.	false

(1) Product quality has a significant positive effect on customers' consumption intention, and the path coefficient is 1.18. Taste and nutritional value are the factors that consumers pay more attention to the quality of products. Taste directly enhances the consumption experience, and nutritional value meets the health needs. In addition, quality and shelf life enhance trust and ensure the safety and quality of products. Quality is the cornerstone of brand, and safety is the bottom line of products. Turning the bottom line of safety into the high line of market competition will eventually form a strong cognitive bond of "quality is brand" in consumers' minds (Table 10).

(2) Purchasing perception has a significant positive effect on consumers' willingness to spend, and the path coefficient is 0.96. Price, packaging, brand effect and channel convenience jointly drive consumption decision. It shows that consumers are influenced by purchase perception factors when they buy deep-processed taro products in Lipu, Guangxi. Increasing the cost performance of products, establishing a good brand image and enhancing the convenience of purchase have become important factors affecting consumers' purchase behavior, which should be paid attention to.

(3) The influence of community has no significant influence on customers' willingness to consume, which shows that the influence of community has no obvious influence on the purchase and premium willingness of products when consumers buy, so the influencing factors of purchase perception can be ignored in the subsequent promotion and marketing of products.

## 5.5 K-Means Clustering

#### 5.5.1 K-means clustering algorithm

K-means clustering analysis is a common unsupervised learning algorithm, which is used to divide the data set into clusters, so that the data points in the same cluster have high similarity, while the data points in different clusters are quite different. Its goal is to minimize the total distance between the data points in the cluster and the center (centroid) of the cluster, so the K-means clustering algorithm is described as follows[21]:

## (1) initialization

Let the total sample set be a set of n sample combinations, and the number of clusters is, divide the sample set into classes at will, mark it as, calculate the corresponding initial clustering centers, mark it as, and calculate:  $G = \{w_j, j = 1, 2, ..., n\}C(2 \le C \le n)GCG_1, G_2, ..., G_cCm_1, m_2, ..., m_cJ_e$ 

$$Y_e = \sum_{i=1}^{C} \sum_{w \in G_i} ||w - m_i||^2$$
(2)

Among them, the smallest cluster is the optimal result under the criterion of sum of squares of errors.  $J_e$ 

#### (2) Iteration of clustering centers

According to the principle of minimum distance, the samples are clustered, namely:  $G_i = \emptyset(i = 1, 2, ..., C)w_j = (j = 1, 2, ..., n)$ 

if

$$d(w_j, G_k) = \min_{1 \le i \le C} d(w_j, m_i)$$
<sup>(3)</sup>

Then, and recalculate the cluster center:  $w_j \in G_k$ ,  $G_k = G_k \cup \{w_j\}$ , j = 1, 2, ..., n

$$m_i = \frac{1}{n_i} \sum_{w_j \in G_i} w_j, i = 1, 2, ..., C$$
(4)

Where is the number of samples in the current class and recalculate. $n_i G_i J_e$ 

(3) If the two iterations are unchanged, the algorithm terminates, otherwise the algorithm goes to  $(2)J_e$ 

# 5.5.2 Selection of the number of clusters

In this paper, the consumer's gender, age, living area, salary level, working status, education level, willingness to buy products and payment premium are used as the basis for classification, and K-means clustering is used to classify consumer behavior. Python is used to analyze the data, and the optimal K value is selected by the Elbow Method. With

#### Volume 2, Issue 2, Pp 51-72, 2025

the increase of cluster number K, the total sum of squares (SSE) of the clustering model gradually decreases, but when K increases to a certain critical value, the SSE value decreases slowly, forming an "inflection point" or "elbow". The k value corresponding to this inflection point is usually considered as the optimal number of clusters. As shown in the figure below, we choose the value of k as 3 (Figure 17).



## 5.5.3 Cluster result analysis

According to the clustering results, the sample data are divided into three categories based on K-means model, and the proportion of people in each category is as follows Figure 18 as shown.



Figure 18 Percentage of Consumers in Each Cluster Center

According to the situation of each cluster feature, the radar map of cluster feature is made to observe the situation of each cluster consumer group more clearly (Figure 19).



Figure 19 Cluster Feature Radar Map

Combined with the questionnaire data, the whole consumer is pictured first, and then the whole consumer is subdivided into three categories by integrating the K-means clustering proportion diagram and the characteristic radar diagram, and the in-depth consumer portrait and difference analysis are carried out. The analysis results are as follows:

Overall portrait of consumers: According to the respondents' information, the age group of 18-30 years old accounts for 54.3%, among which students and employees are the main consumers. There are more highly educated people, accounting for 72.9%, which shows the preference of highly educated people for Lipu taro deep-processed products. Most of these people live in cities with high pressure and moderate consumption level, ranging from 2000 to 4000 yuan. Such consumers generally recognize the deep-processed products of Lipu taro and are willing to pay a certain premium for the deep-processed products of Lipu taro.

The first cluster: balanced consumers. This group presents a balanced feature of "no outstanding shortcomings and no obvious preferences", with a large age span of 25-45 years old. The monthly consumption level is not high, around 2000-4000, and the consumption intention is at a medium level. The demand for deep-processed products of Lipu taro is concentrated in daily high-frequency scenes (such as family breakfast and refreshments). Consumer decision-making relies on public reputation and regular promotion, and prefers basic products with high cost performance (such as instant taro paste and taro cakes). Widely distributed in the region, it has low sensitivity to the cultural connotation or health attributes of products, but it has basic requirements for the safety of raw materials. It is suggested that through omni-channel distribution and mass marketing, the label of "daily nutrition" should be strengthened, and large-capacity products for family packaging should be developed to improve the repurchase rate.

The second cluster: key consumers. The core characteristics of this group are "regional concentration and high income stability", and the areas are mostly distributed in provincial capitals and districts and counties, with stable jobs and high income, and strong acceptance of the payment premium for the deep-processed products of Lipu taro. Prefer joint products with regional characteristics (such as taro moon cakes with local brands), and the consumption scene has social attributes. The decision-making path is easily driven by promotional activities, and at the same time, it pays attention to the brand and tends to make products with brand effect. It is suggested that products should be customized according to the culture of the target area, accurately promoted through the local life platform, and the attribute of "quality socialization" should be strengthened.

The third cluster: potential consumers. This group is marked by "high education, high income and high health concern". It is a youth group aged 28-35, and most of them live in provincial capitals. Pursue the healthy attributes of Lipu taro deep-processed products, and is willing to raise the consumer price for the technological innovation and cultural connotation of the products. The consumption scene is biased towards fitness meal replacement and high-end catering customization, and the decision-making relies on professional evaluation and knowledge platform information, and the willingness to share socially is strong (such as spreading a healthy lifestyle through Xiaohongshu). For this kind of consumer groups, it is suggested to introduce high-end products with organic certification and nutrition visualization (such as dietary fiber meal replacement powder), build a "health+quality" content matrix with nutritionists and fitness KOL, and tap the historical and cultural value of Lipu taro to match its pursuit of quality of life, so as to cultivate and improve the consumption willingness of this kind of groups.

# 6 RESEARCH CONCLUSIONS AND SUGGESTIONS

## 6.1 Conclusion

Starting from the development of rural revitalization product "Lipu Taro" and guided by the popular trend of modern food, this study conducted a research on the combination mode of Guangxi Lipu Taro and modern food, and made an inquiry and analysis on the market analysis and development trend, market opportunity mining and marketing

suggestions. Among them, the questionnaire survey was distributed by means of the Questionnaires platform, and the questionnaire structure was adjusted after the pre-survey. Finally, 550 questionnaires were collected and 93 invalid questionnaires were eliminated, and the number of valid questionnaires was 475, with an effective rate of about 83.09%. The quality of the questionnaire is controlled in terms of survey planning, questionnaire design, questionnaire pre-survey and formal survey. The results show that the reliability of the questionnaire is 0.977 and the validity is 0.864, indicating that the questionnaire data has good reliability and validity. In this paper, Python web crawler, structural equation model (SEM) and K-means clustering model are mainly used for research and modeling analysis. Finally, through the analysis of the questionnaire data, the following conclusions are drawn:

# 6.1.1 Integration and innovation: a new era of Lipu taro and modern food industry

In the current wave of healthy food consumption, the combination of the traditional value of Lipu taro and modern food processing technology has opened a new chapter in the deep processing of characteristic agricultural products. By developing innovative products such as taro paste pre-products, instant taro chips and taro baking raw materials, it not only retains the dense and sweet characteristics of Lipu taro, but also meets the diversified needs of the market for convenience foods, healthy snacks and catering raw materials. This innovation breakthrough is not only reflected in the expansion of product matrix, but also runs through the whole process of supply chain optimization, scene marketing and regional brand building. With the improvement of consumers' awareness of the quality of regional agricultural products and the emphasis on nutritional functions such as dietary fiber, the economic value and cultural connotation of Lipu taro have been deeply developed. The integration of traditional ingredients and modern food industry marks a new stage of improving quality and increasing efficiency in the field of deep processing of agricultural products.

## 6.1.2 Potential release: multidimensional value reconstruction of characteristic agricultural products

As a national geographical indication product, Guangxi Lipu taro is showing amazing development potential. The research obtained the consumption data of e-commerce platform and social media through Python crawler, and found that the interaction volume of topics related to "Lipu Taro" increased. The hot words of consumer concern generated by WordCloud show that the characteristics of "regional scenery", "taste evaluation" and "food collocation" are the most concerned. SnowNLP emotional analysis shows that consumers' satisfaction with taro deep-processed products is 42.46%, and negative emotions account for 33.82%, indicating that they are more positive about Lipu taro. Semantic network analysis reveals that the market demand has been upgraded from primary agricultural products to ready-to-eat and functional directions, and consumers' acceptance of the brand premium of "Lipu Geographical Indications" has increased.

## 6.1.3 Value identification: the core logic of building consumption stickiness

In the highly competitive health food market, it is very important to establish consumers' continuous recognition of regional agricultural products. Through the structural equation model (SEM) analysis, it is found that three factors, purchase perception (0.96), product quality (1.18) and community influence (-0.30), have significant influence on repurchase intention, among which the concept of "food modernization" associated with product quality factor has the greatest influence weight. This shows that consumers' experience of purchasing products directly affects their willingness to pay premium. Combined with IPA-KANO mixed model analysis, the research and development should focus on two high-value dimensions of "instant convenience" and "nutrition visualization". K-means-based consumer clustering shows that young customers pay more attention to product innovation, while family customers pay attention to the traceability of raw materials. Therefore, building a three-dimensional operation system based on quality, cultural empowerment and value resonance is the key strategy to enhance customer stickiness.

## 6.2 Suggestions

Based on the above conclusions, in order to promote the development of deep-processed products of Lipu taro, the following are specific suggestions:

# 6.2.1 Deepen market research and accurately locate consumer demand

Enterprises should invest more resources in market research to understand the preferences and needs of consumers of different ages, genders, occupations and regions for Lipu taro and modern food products. The research content can include consumers' preferences on taste, sweetness, packaging, shelf life and price of Lipu taro deep-processed products, as well as their cognition and acceptance of Lipu taro culture. By accurately positioning consumer demand, enterprises can develop innovative products that are more in line with market demand and improve the market competitiveness of products.

## 6.2.2 Strengthen product innovation and create differentiated competitive advantages

Product innovation is the key for enterprises to gain sustainable competitive advantage. Enterprises can rely on its unique taste and high dietary fiber characteristics to develop diversified product matrix through modern food technology. Research and develop ready-to-use frozen taro paste and seasoning taro powder for the catering market, introduce non-fried freeze-dried taro chips for the healthy snack track, and innovate cross-border fusion products, such as taro-flavored vegetable protein drinks and flowing taro cake. The shelf life and taste of deep-processed products of Lipu taro are improved by modern technology. The two-wheel drive mode of "traditional flavor+modern technology" not only retains the regional characteristics of Lipu taro, but also forms a differentiated product system with technical barriers.

#### 6.2.3 Pay attention to brand building and enhance brand influence

Brand is an important asset of an enterprise and plays an important role in enhancing the added value and market competitiveness of products. Enterprises should pay attention to brand building and improve consumers' awareness and loyalty to products by creating unique brand image and communication strategy. For example, we can enhance the brand's popularity and reputation by holding cultural activities, developing public welfare undertakings and cooperating with well-known brands. At the same time, enterprises should also strengthen the protection of intellectual property rights, avoid the occurrence of infringement and safeguard the legitimate rights and interests of brands.

## 6.2.4 Strengthen marketing and expand market share

Marketing promotion is an important means for enterprises to expand market share and improve sales performance. Enterprises should formulate scientific and reasonable marketing strategies and make full use of various marketing channels and tools for promotion. For example, brand promotion and product promotion can be carried out through social media, online advertisements and offline activities; Can cooperate with e-commerce platform to carry out online sales; We can cooperate with catering enterprises to introduce products into more consumption scenarios. In addition, enterprises can continuously optimize marketing strategies and improve marketing effects according to market feedback and consumer demand.

## 6.2.5 Strengthen industrial chain cooperation and achieve coordinated development

Industrial chain cooperation is an effective way to share resources, reduce costs and enhance competitiveness. Enterprises should establish close cooperative relations with upstream and downstream enterprises to jointly promote the industrial development of Lipu taro and modern food. For example, we can establish a long-term and stable cooperative relationship with taro planting base to ensure the quality and supply stability of raw materials; Can cooperate with logistics companies to optimize the logistics distribution system and reduce transportation costs; We can cooperate with scientific research institutions to jointly develop new varieties, new processes and new technologies to enhance the scientific and technological content and added value of products.

# 6.2.6 Inherit Lipu taro culture and enhance brand value

As a geographical indication product in China, Lipu taro bears a long planting history and regional culture. Enterprises can hold taro culture festival, explore the historical story of "royal tribute", and combine modern design to integrate Zhuang township patterns and ancient cellar-keeping techniques into product packaging and experience scenes. Developing taro carvings to create peripheral and joint-name non-legacy foods, and implanting cultural IP into prefabricated vegetables and deep-processed products not only retain the traditional charm, but also enhance consumers' value recognition of "Lipu Taro" and promote the upgrading of regional agricultural products to cultural brands.

## 6.2.7 Social responsibility and sustainable development

Adhere to the concept of environmental protection, pay attention to the application of environmental protection concept in product production and packaging design, and reduce the pollution and damage to the environment. Adopt recyclable materials and environmentally friendly packaging to reduce the burden of products on the environment. Strengthen the development of rural revitalization strategy, actively participate in rural revitalization, promote local economic development and increase farmers' income by developing Lipu taro industry, and cooperate with local government to carry out poverty alleviation projects or support farmers to plant high-quality taro to give back to society and achieve sustainable development. Establish corporate culture, establish correct corporate values and cultural concepts, pay attention to employee care and welfare benefits, and enhance employee satisfaction and loyalty. At the same time, actively fulfill corporate social responsibility, participate in public welfare undertakings and charitable activities, and establish a good corporate image and social reputation.

To sum up, in order to highlight the characteristics of Lipu taro in the modern food market and attract more consumers, enterprises need to deeply understand the target market, innovate products, strengthen brand building, and expand market share through diversified marketing strategies. At the same time, work closely with all parties in the industrial chain to jointly promote industrial development. On this basis, we should attach importance to cultural inheritance and environmental protection, actively assume social responsibilities and realize sustainable development. These comprehensive strategies will help enterprises to enhance their competitiveness, win the trust of consumers and promote the prosperity of Lipu taro industry.

## **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

# REFERENCES

- Liu Jian, Li Yonggui, Zhou Junyuan. The 7th Lipu taro Culture Festival continued to polish the "golden signboard". Guilin Daily, 2023(005).
- [2] Zhou Wenqiong. "Guilin Courtesy Tourism Commodities": Welcome visitors from all directions with supreme courtesy. Guilin Daily, 2024(007).
- [3] Ping Shaohua, Geng Lan. Lipu taro with wings. China Fruit and Vegetable, 2016, 36(05): 62-64.
- [4] He Pingjiang, Zhou Junyuan, Xie Qizhu. Where did sugar orange go? Guilin Daily, 2023(002).

- [5] Liang Zhongchao, Zhou Junyuan. Lipu, Guangxi: The "new" power brings new agricultural atmosphere. China Food News, 2024(001).
- [6] Li Yuanyuan, Zhong Xiangjing, Zhang Haijun. Research on Agricultural Products Processing and Food Innovation under the Integration of Agronomy and Food Science. High-tech and Industrialization, 2024, 30(06): 47-49.
- [7] Lei Yuliang, Xiao Lin. Analysis on the Impact of Agricultural Products Processing Industry Agglomeration on Agricultural Modernization from the Perspective of Industrial Relevance. Scientific Decision-making, 2024(04): 139-152.
- [8] Guan Yonghua, Zhou Qiuhong. Research on the Optimization Path of Guangxi Agricultural Products Processing Industry from the Perspective of Chinese Modernization. Cereals, Oils and Feeds Science and Technology, 2023(01): 73-75.
- [9] Bai Jian, Hong Xiaojuan. Online public opinion text mining and sentiment analysis based on barrage. Software Engineering, 2022, 25(11): 44-48.
- [10] Ding Liuhua, Ye Xinliang, Sun Ruihong. Research on Perceptual Dimensions and Semantic Networks of Negative Online Comments on Cruise Tourism // China tourism academy. Proceedings of the 2022 China Tourism Science Annual Conference: Tourism Talents Construction and Young Talents Training. Shanghai University of Engineering Science, 2022: 13.
- [11] Shi Xiaochen, Li Mingzi, Li Jinyuan. Study on the Influencing Factors of Consumers' Purchase Intention of Low-alcohol Liquor-Based on SEM Structural Equation Model. National Circulation Economy, 2023(14): 24-27.
- [12] Cui Hongcheng, Chen Qingguo. Research on the Willingness to Use Fitness Intelligent Wearable Devices Continuously. Journal of Hebei Institute of Physical Education, 2024: 1-10.
- [13] Qingliang Meng, Bian Lingling, Lin He, et al. Detection method of express service quality integrating Kano model and IPA analysis. Industrial Engineering and Management, 2014, 19(02): 75-80+88.
- [14] Du Jili, Li Dongsheng, Yan Yafei. Evaluation of users' satisfaction of ruins park based on KANO-IPA analysis —— Taking Xiyuan Park in Luoyang as an example. Agriculture and Technology, 2023, 43(05): 130-134.
- [15] Cai Shaolin, Wu Liwen, Zheng Dongrong. Portrait Construction and Precision Marketing Strategy of Agricultural Products Consumer Groups Based on K-means Cluster Analysis. Rural Economy and Technology, 2022, 33(22): 251-254+262.
- [16] Shi Lemeng, LinYu Cheng, Hei Minnan. Modern sugar control "revolution": sweet, not sugar! -factors influencing the purchase behavior of sugar-reduced products by young consumers. Translation of Economic Data, 2022(02): 34-44.
- [17] Wei Zongbian. On the development status and countermeasures of farmers' professional cooperatives in Lipu County. Southern Horticulture, 2015, 26(05): 29-31.
- [18] Chen Miao, Lin Zheng. Study on the effectiveness simulation of the linkage between green tax system and carbon tax system-based on grey relational analysis and CGE model. Ecological Economy, 2025, 41(03): 13-23+74.
- [19] Deng Xiaohua, Yang Bo. Construction and Evaluation of Digital Literacy Index System for Open University Teachers — Based on Structural Equation Model. Journal of Hubei Open University, 2024, 44(06): 28-35.
- [20] Shi Jingjuan, Bao Yuting, Qiao Nan, et al. Study on the influencing factors of farmers' willingness to moderate scale management based on structural equation model-taking Yuzhong County of Lanzhou City as an example. Tropical Agricultural Engineering, 2025, 49(01): 107-113.
- [21] Nahoujy R M .Applying a K-means model to TSD data to find categories for the structural assessment of flexible pavements[J].Transportation Engineering,2025,20100342-100342.

# MACHINE LEARNING APPROACHES FOR ACCURATE DEMAND FORECASTING IN SUPPLY CHAIN MANAGEMENT

Liu Zhen, Yang Lin\*

School of Computer Science, Southeast University, Nanjing 210000, Jiangsu, China. Corresponding Author: Yang Lin, Email: Ylin394021@seu.edu.cn

**Abstract:** Accurate demand forecasting is a cornerstone of effective supply chain management, enabling companies to align production, inventory, and distribution with market needs. Traditional statistical models often fail to capture the nonlinear and complex patterns in consumer demand, particularly in the presence of seasonal shifts, promotional events, and external shocks. In recent years, machine learning (ML) has emerged as a powerful tool for enhancing demand forecasting accuracy by leveraging large-scale historical and real-time data. This paper reviews the core machine learning techniques applied to demand forecasting, including supervised learning, time series forecasting models, and ensemble methods. We develop and evaluate a hybrid forecasting framework that integrates Long Short-Term Memory (LSTM) neural networks with gradient boosting to capture both sequential patterns and feature-based dependencies. The proposed approach is validated using a retail demand dataset, and its performance is benchmarked against traditional models. The results demonstrate that ML-based methods significantly outperform classical forecasting techniques, offering improvements in forecast precision, robustness to noise, and responsiveness to dynamic market signals.

**Keywords:** Demand forecasting; Supply chain management; Machine learning; LSTM; Gradient boosting; Time series prediction; Forecast accuracy; Retail analytics

# **1 INTRODUCTION**

Demand forecasting is one of the most critical components of supply chain management, influencing decisions related to procurement, production planning, inventory control, logistics scheduling, and customer service[1]. Inaccurate forecasts can lead to a range of operational inefficiencies, including stockouts, overstocking, excess holding costs, and loss of customer loyalty[2]. Therefore, organizations have a strong incentive to develop accurate and timely demand forecasting systems that can adapt to rapidly changing consumer behavior and market dynamics[3].

Historically, demand forecasting has relied heavily on classical statistical techniques such as exponential smoothing, moving averages, and autoregressive integrated moving average (ARIMA) models[4]. While these methods are computationally efficient and interpretable, they often struggle with high-dimensional data, nonlinear relationships, seasonality, and sudden changes in demand due to exogenous factors such as promotions, weather, or economic shocks[5]. Moreover, traditional models require significant manual intervention and domain knowledge for feature engineering and parameter tuning, limiting their scalability and adaptability[6].

In contrast, machine learning (ML) offers a data-driven alternative that is capable of learning complex, nonlinear, and hierarchical relationships from historical data without relying on rigid parametric assumptions [7]. ML models can be trained to automatically capture patterns across large datasets, incorporate a wide array of features including categorical and temporal variables, and adapt continuously as new data becomes available[8]. Recent advancements in deep learning, particularly recurrent neural networks (RNNs) and their variants such as Long Short-Term Memory (LSTM) networks, have further pushed the frontier of demand forecasting by modeling long-term dependencies in time series data [9].

This paper explores the application of various machine learning approaches to demand forecasting in retail supply chains. We propose a hybrid framework that combines LSTM networks for temporal sequence learning with gradient boosting for incorporating static and contextual features. The proposed methodology is evaluated using a real-world retail dataset and benchmarked against classical forecasting models. We aim to demonstrate that machine learning not only improves forecasting accuracy but also enhances robustness and adaptability in complex supply chain environments.

# **2** LITERATURE REVIEW

The field of demand forecasting has undergone a significant transformation with the advent of machine learning, evolving from conventional statistical approaches to data-driven algorithms capable of modeling complex temporal and cross-sectional patterns[10]. In traditional supply chain operations, statistical techniques such as exponential smoothing, ARIMA, and seasonal decomposition were widely used due to their interpretability and ease of implementation[11]. However, these models are inherently limited in their ability to capture nonlinear dependencies, interaction effects among multiple variables, and sudden regime shifts caused by promotional events, competitor actions, or macroeconomic changes[12].

Machine learning has emerged as a promising alternative, offering flexible models that can learn directly from data without requiring strong assumptions about underlying distributions or data-generating processes[13]. Supervised learning models, such as decision trees, random forests, support vector machines, and gradient boosting machines, have demonstrated superior performance in scenarios where a rich set of features is available[14]. These algorithms can incorporate exogenous variables such as holidays, weather, regional economic indicators, and product-level metadata, making them highly suitable for retail forecasting tasks[15].

In addition to classical machine learning algorithms, deep learning has gained substantial attention due to its ability to automatically learn hierarchical representations from raw input data[16]. Specifically, RNNs and their gated variants like LSTM and GRU (Gated Recurrent Unit) have shown notable success in capturing long-term dependencies in time series data[17]. These models can maintain internal memory states that are updated dynamically based on the input sequence, allowing them to model temporal lags, seasonality, and abrupt changes in demand patterns[18].

Another strand of literature has focused on hybrid modeling strategies that combine the strengths of different machine learning methods[19]. For instance, some frameworks use gradient boosting to process static features such as store location, product category, and promotion types, while relying on LSTM networks to model the temporal dynamics[20]. This fusion of static and dynamic modeling components enhances both short-term responsiveness and long-term trend recognition[21].

Furthermore, the rise of big data technologies has facilitated the use of high-frequency data sources such as clickstream logs, customer transaction histories, and point-of-sale information[22]. These data sources provide granular insights into consumer behavior, enabling forecasting models to move beyond simple SKU-level predictions and incorporate customer segmentation, behavioral clustering, and personalized demand forecasting[23]. The literature also highlights the importance of feature engineering, model interpretability, and forecast explainability, especially in business contexts where decisions based on forecasts have significant operational implications[24].

Despite the advancements, several challenges persist. One major issue is the lack of transparency and interpretability in black-box models, which hinders their adoption in risk-averse industries[25]. Moreover, machine learning models are sensitive to data quality, missing values, and outliers, which are common in real-world supply chain datasets. Another concern is the requirement for continuous retraining and validation to ensure sustained performance over time[26].

The literature reveals a growing consensus that while no single model universally outperforms others across all forecasting scenarios, machine learning models—particularly when customized and hybridized—offer substantial improvements in forecast accuracy, adaptability, and decision support. These findings form the foundation for our proposed framework, which aims to integrate advanced ML architectures with real-time retail data streams to enable dynamic and granular demand forecasting in modern supply chains.

## **3** METHODOLOGY

This section outlines the methodological approach employed to implement ML models for demand forecasting in supply chain management. The methodology consists of four key stages: data acquisition and preprocessing, feature engineering and selection, model development and training, and evaluation and validation. Each stage is crucial in ensuring the accuracy, robustness, and applicability of the forecasting models in real-world supply chain operations.

## 3.1 Data Acquisition and Preprocessing

The foundation of any forecasting model lies in high-quality, relevant data. In this study, historical sales data, inventory records, pricing logs, promotional calendars, weather indicators, and macroeconomic variables were collected from multiple retail outlets over a three-year period. The raw data exhibited typical challenges such as missing values, seasonality, outliers, and inconsistent temporal resolution.

To address these issues, we employed a multi-step preprocessing strategy as in Figure 1. First, time series were resampled to a uniform daily granularity using forward filling and linear interpolation. Outliers were detected using seasonal hybrid extreme studentized deviate (S-H-ESD) tests and were capped to maintain statistical integrity. Missing categorical values were imputed using mode substitution, while continuous variables used k-nearest neighbor (KNN) imputation based on correlated features.



## Figure 1 Multi-Step Preprocessing Strategy

#### **3.2 Feature Engineering and Selection**

Effective forecasting requires the transformation of raw data into meaningful, compact representations. We constructed a feature matrix consisting of temporal attributes (e.g., day of week, month, holiday indicators), lagged demand values, rolling statistics (e.g., moving averages, standard deviations), product attributes (e.g., brand, category, shelf life), and exogenous signals such as promotions or weather conditions.

To identify the most informative predictors, we applied SHAP (SHapley Additive exPlanations) value analysis on a baseline XGBoost model as in Figure 2. This allowed us to rank features by their marginal contribution to model predictions. Temporal proximity (e.g., lagged 1-day sales), promotion indicators, and rolling demand volatility emerged as the most predictive variables across multiple retail SKUs.



Figure 2 Feature Importance Based on SHAP Values

Furthermore, feature reduction was performed via recursive feature elimination (RFE) and principal component analysis (PCA) for models sensitive to multicollinearity. These techniques enhanced computational efficiency and helped prevent overfitting in high-dimensional models.

## 3.3 Model Development and Training

We developed and compared multiple ML algorithms, including Random Forest (RF), Gradient Boosting Machines (GBM), LSTM, and Temporal Fusion Transformers (TFT). Each model was tailored for time series forecasting, using rolling-window cross-validation to preserve temporal dependencies.

Hyperparameter tuning was conducted using Bayesian Optimization with cross-validated mean absolute percentage error (MAPE) as the objective function. Early stopping was used to mitigate overfitting, particularly for deep learning models. In addition, each model was trained on a distributed computing cluster using GPU acceleration where applicable, to handle high-volume data and model complexity efficiently.



Figure 3 Forcasting Error Metrics by Model

Throughout the training process, performance metrics such as root mean square error (RMSE), MAPE, and symmetric mean absolute percentage error (sMAPE) were recorded for model comparison as in Figure 3. Ensemble methods were also evaluated by aggregating outputs from multiple models to enhance prediction robustness.

## 4 RESULTS AND DISCUSSION

The experimental evaluation of machine learning approaches for demand forecasting was conducted on a real-world retail dataset, encompassing daily sales data across multiple product categories, regions, and seasons. The models evaluated included Random Forest, Gradient Boosting, LSTM, TFT, and a final ensemble method combining the strengths of the best-performing models.

The performance of each model was assessed using two key error metrics: MAPE and RMSE. These metrics provide complementary insights—MAPE reflects the relative prediction error, which is crucial for inventory decisions across products with different volume scales, while RMSE penalizes larger errors more heavily, highlighting forecasting robustness in high-variance scenarios.

The results demonstrated that classical tree-based models such as Random Forest and Gradient Boosting provided reasonable accuracy, with MAPE values around 11–12%. However, deep learning models significantly outperformed them. The LSTM model achieved a MAPE of 10.3% due to its ability to model temporal dependencies, while the TFT model pushed the error down further by incorporating attention mechanisms and covariate information. The ensemble model, which combined predictions from LSTM and TFT with weighted averaging based on validation set performance, yielded the best results overall, achieving a MAPE of 8.7% and RMSE of 18.5.



Figure 4 Model Performance Comparison

These findings underscore the importance of choosing forecasting methods that are sensitive to both seasonality and covariate shifts, especially in volatile retail environments. Deep learning models, particularly those designed to process sequential and multivariate data, are well-suited for such settings. Moreover, the ensemble strategy mitigates the weaknesses of individual models and stabilizes forecasts, making it highly applicable in operational supply chains.

In addition to quantitative metrics, qualitative analysis of demand forecast plots revealed that advanced models more effectively captured demand surges (e.g., promotional spikes) and long-tail seasonal trends, whereas traditional methods often lagged or oversmoothed the response. This has direct implications for inventory management, as underestimating demand during such periods can lead to costly stockouts and lost revenue, while overestimating leads to increased holding costs.

Overall, these results validate the application of modern ML architectures in real-world forecasting pipelines and provide a compelling argument for integrating these approaches into retail demand planning systems.

# 5 CONCLUSION

Accurate demand forecasting remains a cornerstone of effective supply chain management, directly influencing inventory control, logistics coordination, and customer satisfaction. This study explored and compared various machine learning approaches, including traditional models and advanced deep learning architectures, for forecasting demand in dynamic retail environments. Our findings highlight the clear performance advantages of deep learning models—particularly LSTM and Temporal Fusion Transformer—over classical methods such as Random Forest and Gradient Boosting. By leveraging temporal dependencies and exogenous features, these models deliver more nuanced, accurate, and adaptable forecasts.

The ensemble strategy further improved forecasting robustness by combining the strengths of individual models, leading to a significant reduction in both MAPE and RMSE. These quantitative improvements translated into practical benefits for supply chain operations, such as reduced stockouts, optimized safety stock levels, and better alignment between supply and demand.

Moreover, our analysis shows that the use of machine learning not only enhances forecasting accuracy but also supports more proactive and strategic supply chain decisions. Rather than reacting to demand fluctuations, companies can preemptively adjust procurement, production, and logistics plans based on high-quality predictive insights.

Future work can extend this research by incorporating external data sources such as economic indicators, competitor pricing, and weather conditions to further enhance forecast accuracy. Additionally, integrating explainability frameworks such as SHAP values or attention weight analysis can improve model transparency and support human decision-makers in validating and adjusting forecasts.

In conclusion, the integration of advanced machine learning models into demand forecasting workflows represents a critical step toward building agile, resilient, and data-driven supply chains capable of thriving in an increasingly complex market landscape.

# **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

- [1] Tadayonrad Y, Ndiaye A B. A new key performance indicator model for demand forecasting in inventory management considering supply chain reliability and seasonality. Supply Chain Analytics, 2023, 3: 100026.
- [2] Jean G. Inventory Management Strategies: Balancing Cost, Efficiency, and Customer Satisfaction. 2024.
- [3] Wu B, Shi Q, Liu W. Addressing Sensor Data Heterogeneity and Sample Imbalance: A Transformer-Based Approach for Battery Degradation Prediction in Electric Vehicles. Sensors, 2025.
- [4] Van Chau D, He J. Machine Learning Innovations for Proactive Customer Behavior Prediction: A Strategic Tool for Dynamic Market Adaptation. 2024.
- [5] Wang Y, Xing S. AI-Driven CPU Resource Management in Cloud Operating Systems. Journal of Computer and Communications, 2025.
- [6] Leon J. Forecasting Imported Fruit Prices in the United States Using Neural Networks (Doctoral dissertation, National University). 2024.
- [7] Zhang Q, Chen S, Liu W. Balanced Knowledge Transfer in MTTL-ClinicalBERT: A Symmetrical Multi-Task Learning Framework for Clinical Text Classification. Symmetry, 2025, 17(6): 823.
- [8] Ravishankar S, Battineni G. A Survey on Recent Advancements in Auto-Machine Learning with a Focus on Feature Engineering. Journal of Computational and Cognitive Engineering, 2025, 4(1): 56-63.
- [9] Strielkowski W, Vlasov A, Selivanov K, et al. Prospects and challenges of the machine learning and data-driven methods for the predictive analysis of power systems: A review. Energies, 2023, 16(10): 4025.
- [10] Wilson A, Anwar M R. The Future of Adaptive Machine Learning Algorithms in High-Dimensional Data Processing. International Transactions on Artificial Intelligence, 2024, 3(1): 97-107.
- [11] Jin J, Xing S, Ji E, et al. XGate: Explainable Reinforcement Learning for Transparent and Trustworthy API Traffic Management in IoT Sensor Networks. Sensors (Basel, Switzerland), 2025, 25(7): 2183.
- [12] Mienye I D, Swart T G, Obaido G. Recurrent neural networks: A comprehensive review of architectures, variants, and applications. Information, 2024, 15(9): 517.
- [13] Giannopoulos P G, Dasaklis T K, Tsantilis J, et al. Machine learning algorithms in intermittent and lumpy demand forecasting: A review. Available at SSRN 5231788, 2025.
- [14] Borucka A. Seasonal methods of demand forecasting in the supply chain as support for the company's sustainable growth. Sustainability, 2023, 15(9): 7399.
- [15] Tan Y, Wu B, Cao J, et al. LLaMA-UTP: Knowledge-Guided Expert Mixture for Analyzing Uncertain Tax Positions. IEEE Access, 2025.
- [16] Benigno G, Foerster A, Otrok C, et al. Estimating macroeconomic models of financial crises: An endogenous regime-switching approach (No. w26935). National Bureau of Economic Research, 2020.
- [17] Pichler M, Hartig F. Machine learning and deep learning—A review for ecologists. Methods in Ecology and Evolution, 2023, 14(4): 994-1016.
- [18] Choudhury A, Mondal A, Sarkar S. Searches for the BSM scenarios at the LHC using decision tree-based machine learning algorithms: a comparative study and review of random forest, AdaBoost, XGBoost and LightGBM frameworks. The European Physical Journal Special Topics, 2024, 233(15): 2425-2463.
- [19] Tiainen M. Forecasting seasonal demand at the product level in grocery retail. 2021.
- [20] Dargan S, Kumar M, Ayyagari M R, et al. A survey of deep learning and its applications: a new paradigm to machine learning. Archives of computational methods in engineering, 2020, 27: 1071-1092.
- [21] Wang J, Tan Y, Jiang B, Wu B, Liu W. Dynamic Marketing Uplift Modeling: A Symmetry-Preserving Framework Integrating Causal Forests with Deep Reinforcement Learning for Personalized Intervention Strategies. Symmetry, 2025, 17(4): 610.
- [22] Waqas M, Humphries U W. A critical review of RNN and LSTM variants in hydrological time series predictions. MethodsX, 2024: 102946.
- [23] Lindemann B, Müller T, Vietz H, et al. A survey on long short-term memory networks for time series prediction. Procedia Cirp, 2021, 99: 650-655.
- [24] Yang J, Li P, Cui Y, et al. Multi-Sensor Temporal Fusion Transformer for Stock Performance Prediction: An Adaptive Sharpe Ratio Approach. Sensors, 2024, 25(3): 976.
- [25] Joshi R, Iyer R, Chopra R, et al. Enhancing E-Commerce Demand Prediction Using Long Short-Term Memory Networks and Gradient Boosting Machines. Innovative AI Research Journal, 2021, 10(2).
- [26] Olayinka O H. Big data integration and real-time analytics for enhancing operational efficiency and market responsiveness. Int J Sci Res Arch, 2021, 4(1): 280-96.

- [27] Feliu Juan M. Enhancing demand forecasting: an analysis of factors impacting sales and implementation of improved methodologies for accurate prediction (Bachelor's thesis, Universitat Politècnica de Catalunya). 2024.
- [28] Eskandari H, Saadatmand H, Ramzan M, et al. Innovative framework for accurate and transparent forecasting of energy consumption: A fusion of feature selection and interpretable machine learning. Applied Energy, 2024, 366: 123314.
- [29] Han X, Yang Y, Chen J, et al. Symmetry-Aware Credit Risk Modeling: A Deep Learning Framework Exploiting Financial Data Balance and Invariance. Symmetry (20738994), 2025, 17(3).
- [30] Ghodake S P, Malkar V R, Santosh K, et al. Enhancing Supply Chain Management Efficiency: A Data-Driven Approach using Predictive Analytics and Machine Learning Algorithms. International Journal of Advanced Computer Science & Applications, 2024, 15(4).
- [31] Guo L, Hu X, Liu W, et al. Zero-Shot Detection of Visual Food Safety Hazards via Knowledge-Enhanced Feature Synthesis. Applied Sciences, 2025.

