



PREDICTIVE ANALYSIS OF CRYPTOCURRENCY INCOME USING MACHINE LEARNING MODELS

Kancham Reddy Akhila¹, Vijaylakshmi A Lepakshi²

¹*School of Computer Science and Application, Reva University, Bengaluru, India*

²*Associate Professor, School of Computer Science and Application, Reva University, Bengaluru, India*

Article DOI: <https://doi.org/10.36713/epra18129>

DOI No: 10.36713/epra18129

ABSTRACT

This research explores the potential of various machine learning models to predict cryptocurrency income. With the increasing volatility and widespread use of cryptocurrencies, accurately forecasting income from these digital assets is crucial for investors and stakeholders. By examining a dataset of cryptocurrency transactions, this study addresses the challenge of outliers and assesses the correlations between different variables to ensure robust data analysis. We employ three distinct machine learning models—Linear Regression, Decision Tree Regressor, and Random Forest Regressor—to train and evaluate the dataset. Each model is analysed to calculate performance metrics, such as Mean Squared Error (MSE) and R2 Score, which provide insights into their predictive accuracy and reliability. The detailed analysis and discussion highlight the strengths and weaknesses of each approach, emphasizing the efficiency and accuracy of the models in forecasting cryptocurrency income. The Linear Regression model, known for its simplicity, is compared against more complex models like Decision Tree and Random Forest Regressors, which can capture nonlinear relationships and interactions between variables.

The Random Forest Regressor, an ensemble learning method, is particularly noted for its superior performance in handling complex datasets and providing more accurate predictions. This study's findings offer valuable insights for developing more effective investment strategies and risk management practices in the rapidly evolving cryptocurrency market. Through this comprehensive evaluation, the research aims to contribute to the growing body of knowledge on financial forecasting using machine learning techniques, paving the way for future advancements in this field.

1. INTRODUCTION

The rise of cryptocurrencies has dramatically altered the financial landscape, introducing digital assets that function on decentralized platforms using blockchain technology. Cryptocurrencies, such as Bitcoin, Ethereum, and numerous altcoins, have attracted significant attention due to their potential for high returns on investment, along with their inherent volatility[1]. This volatility can result in substantial gains but also considerable risks, making the prediction of cryptocurrency income a critical area of interest for investors, financial analysts, and researchers.

Cryptocurrencies are digital or virtual currencies that use cryptographic techniques to secure transactions and manage the creation of new units. Unlike traditional currencies, cryptocurrencies are generally decentralized and operate on distributed ledger technology called blockchain[2]. This decentralized nature eliminates the need for a central authority, such as a bank, and facilitates peer-to-peer transactions. The first and most prominent cryptocurrency, Bitcoin, was launched in 2009 by an anonymous entity known as Satoshi Nakamoto. Since then, thousands of alternative cryptocurrencies have emerged, each with its own unique features and applications[3].

The financial appeal of cryptocurrencies lies in their potential for high returns. For example, early adopters of Bitcoin have experienced exponential growth in their investments. However, the cryptocurrency market is notoriously volatile. Prices can fluctuate dramatically within short periods due to various factors, including market speculation, regulatory developments, technological advancements, and macroeconomic trends. This volatility presents a significant challenge for investors and traders aiming to maximize returns while minimizing risks[4]. Various methods have been explored to achieve this, ranging from traditional financial models to sophisticated machine learning algorithms. Machine learning, a subset of artificial intelligence, involves the use of algorithms and statistical models to enable computers to improve their performance on a specific task through experience. It has gained prominence in various fields, including finance, due to its ability to analyse large datasets, identify patterns, and make predictions with minimal human intervention. In the context of cryptocurrency, machine learning can be employed to analyse historical data, detect trends, and forecast future income[5]. Several studies have highlighted the effectiveness of machine learning techniques in financial forecasting. For instance, Nassirtoussi et al. (2014) demonstrated the potential of machine learning models in predicting stock market trends by analysing textual data from financial news. Similarly,



McNally, Roche, and Caton (2018) explored the application of recurrent neural networks (RNNs) for predicting cryptocurrency prices, showcasing the capabilities of deep learning in capturing temporal dependencies and complex patterns in financial time series data[6]. Despite these advancements, the prediction of cryptocurrency income remains a relatively underexplored area. Most existing research focuses on predicting cryptocurrency prices or market trends rather than the actual income generated from cryptocurrency investments. This study aims to address this gap by evaluating the effectiveness of different machine learning models in predicting cryptocurrency income and analysing the underlying data to uncover relevant patterns and relationships[7].

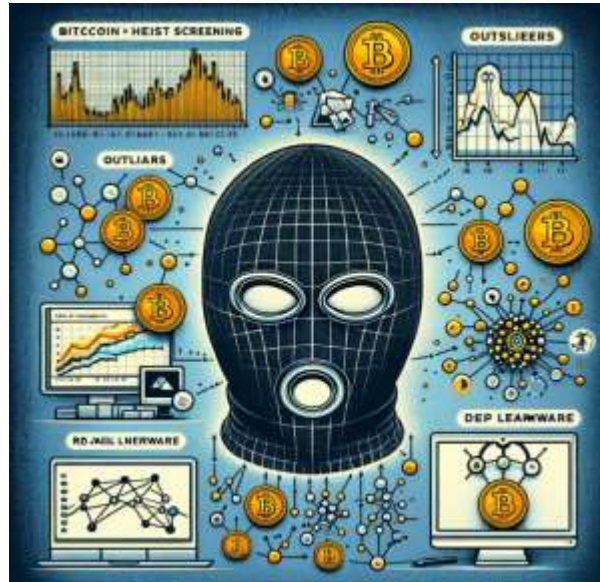


FIG 1: The Money Heist

To achieve this, we employed three distinct machine learning models: Linear Regression, Decision Tree Regressor, and Random Forest Regressor. Linear Regression is a widely used statistical method that models the relationship between a dependent variable and one or more independent variables by fitting a linear equation to the observed data. It is known for its simplicity and interpretability but may not capture complex nonlinear relationships in the data[8]. The Decision Tree Regressor is a non-parametric model that predicts the value of a target variable by learning simple decision rules inferred from the data features. It can handle nonlinear relationships and interactions between variables but is prone to overfitting, especially with noisy data. The Random Forest Regressor, an ensemble learning method, addresses this limitation by constructing multiple decision trees during training and outputting the average prediction of the individual trees. This approach reduces overfitting and improves the model's generalizability[9].

The dataset used in this study includes yearly data on cryptocurrency transactions, encompassing attributes such as year, day, transaction length, weight, count, looped, neighbours, and income. By analysing this dataset, we aim to identify significant correlations and patterns that can inform the prediction models. The correlation analysis provides insights into the relationships between different variables, which can guide feature selection and engineering efforts. Data preprocessing is a crucial step in machine learning, as it ensures the quality and reliability of the input data. In this study, we addressed outliers using the Z-score method to prevent skewing the analysis. Outliers can significantly impact the performance of machine learning models, leading to inaccurate predictions and overfitting. By removing outliers, we aimed to create a cleaner dataset that better represents the underlying distribution of the data[10].

After preprocessing, we split the dataset into training and test sets to evaluate the performance of the models. The training set is used to fit the models, while the test set is reserved for evaluating their predictive accuracy. This approach helps in assessing the models' ability to generalize to unseen data, which is crucial for reliable predictions. The models were evaluated using two key metrics: Mean Squared Error (MSE) and R2 Score are common metrics for evaluating model performance. MSE calculates the average squared discrepancy between the predicted and actual values, with lower MSE values indicating more accurate predictions. The R2 Score, or coefficient of determination, represents the proportion of variance in the dependent variable that can be explained by the independent variables, with values approaching 1 signifying a stronger model fit[11].

The findings of this study provide valuable insights into the predictive capabilities of different machine learning models for cryptocurrency income. The Random Forest Regressor emerged as the most effective model, demonstrating the lowest MSE and the highest R2 score. This suggests that ensemble methods, which combine multiple models to improve prediction accuracy, are



particularly well-suited for this task. In conclusion, this study demonstrates the potential of machine learning models in predicting cryptocurrency income, providing a foundation for future research in this area. By expanding the feature set and exploring more advanced models, future studies can further enhance prediction accuracy and contribute to the development of more robust investment strategies and risk management practices in the cryptocurrency market[12].

The dataset used in this study includes yearly data on cryptocurrency transactions with attributes such as year, day, transaction length, weight, count, looped, neighbours, and income. The detailed dataset is available in the supplementary file heist.xlsx[13].

2. LITERATURE REVIEW

The emergence of cryptocurrencies has attracted a lot of public interest for both the financial and technology sectors because of decentralised and derivatives financial assets that can generate larger returns. Consequently, researchers have devoted considerable time and effort in the analysis and forecasting of many characteristics of virtual currency markets, such as price levels, trading intensity, and funds' profitability. For example, one of the branches of study being active the recent period concerns the use of predictive models from the machine learning field and Income from Cryptocurrencies[14]. Data mining techniques using automated learning methods have evolved due to the potentials of handling large datasets and recognizing patterns unseen by other statistical methods. Other authors have used both sophisticated machine learning models, such as ensemble methods, linear regression model, and others, for predicting cryptocurrency income[15].

In this context, several papers dealt with the application of linear regression analysis for the forecast of the cryptography currency and income. Linear regression is an easy and computationally simple approach for capturing the dependency of the target variable on the input features, and therefore often used as the first step. Therefore, it can be somewhat inefficient in cases of nonlinear relationships or large and intricate datasets – something which is typical of cryptocurrency markets that tend to be highly volatile and often encounter rapid fluctuations[16]. On the other hand, Machine learning-based models, namely the decision tree-based models like the decision tree regressor models, has been proposed to replace linear regression models in predicting the income from cryptocurrencies. An advantage of the Decision trees is that it does not make the assumption of linear separability between the input variables and the output values as the data is divided to different subsets using the feature thresholds. As declared, compared with linear regression, decision trees may be more flexible; however, they may be more prone to overfitting in the condition of the noisy or high-dimensional data[17].

Another set of methods that are getting more consideration in predicting cryptocurrency income are ensemble methods such as a random forest regressor. This approach is an extension of one or more base models to deliver enhanced precision in prediction as well as better generality. Random forests especially are convincible with the results obtained by combining the multiple individual models through averaging the results obtained from various decision trees generated from multiple samples selected from the data set using bootstrapping with-replacement methods[18]. By using this ensemble approach, then, the need to deal with large and complicated sets and overfitting is addressed, which again makes it particularly useful in tasks related to the prediction of cryptocurrency income.

Summing up, the existing research points at the relevance of utilizing various approaches within the machine learning for enhancing cryptocurrency income prediction[19]. Whereas a simple linear regression model serves as a starting point and a reference point, other and more complex models such as the decision tree model and the random forest model performs better as the capture the non-linearity in the data as well as interactions between the data variables. Through the help of these techniques, investigations want to come up with sound solutions for investment and risks on the uncertain value of the cryptocurrencies[20].

3. METHODOLOGY

3.1 Data Retrieval

In the data retrieval phase, the focus is on collecting and compiling the dataset required for the predictive analysis of cryptocurrency income. The process involves several key steps:

Data Source Identification: Identifying reliable and relevant sources of cryptocurrency transaction data. Potential sources include cryptocurrency exchanges, blockchain explorers, and financial databases that provide comprehensive and historical transaction records.

Data Collection: Extracting data from the identified sources. This involves using APIs, web scraping techniques, or direct downloads from databases to gather transaction data. The data collected includes attributes such as year, day, transaction length, weight, count, looped neighbours, and income.

Data Validation: Ensuring the accuracy and completeness of the collected data. This step involves verifying the data against known standards, checking for inconsistencies, and removing any duplicate entries. Validation ensures that the dataset is reliable for further analysis.



Data Storage: Organizing the validated data into a structured format suitable for analysis. The data is stored in a database or a file system, such as CSV or Excel files, ensuring easy access and manipulation. Proper storage includes categorizing the data based on the attributes mentioned.

Initial Data Exploration: Conducting an exploratory data analysis (EDA) to understand the basic characteristics of the dataset. This involves generating summary statistics, visualizing distributions, and identifying any initial patterns or anomalies. EDA provides a preliminary insight into the data and helps in shaping the subsequent preprocessing steps.

3.2 Data Preprocessing

The data preprocessing phase was crucial for ensuring the quality and reliability of the dataset used for analysis. Initially, the dataset was loaded, and its structure was thoroughly examined. This step involved inspecting the various attributes to understand their nature and distribution. Following this, outliers were identified and removed using the Z-score method. The Z-score method is a statistical technique that measures the number of standard deviations a data point is from the mean. This method is widely used to detect and exclude extreme values that could potentially skew the analysis (Iglewicz & Hoaglin, 1993). Removing outliers is essential in machine learning to prevent these anomalies from disproportionately influencing model training and evaluation.

3.3 Correlation Analysis

After preprocessing the data, a correlation analysis was performed to understand the relationships between different features within the dataset. The correlation matrix was calculated, which shows the correlation coefficients between pairs of variables. This matrix helps in identifying which variables are highly correlated, thus providing insights into the potential multicollinearity issues that could affect the model's performance (Hinkle, Wiersma, & Jurs, 2003). Understanding these relationships is vital for feature selection and engineering, ensuring that the models are built on meaningful and independent variables.

3.4 Model Training and Evaluation

The next step involved training and evaluating the machine learning models. The dataset was split into training and test sets to facilitate this process. The training set was used to fit the models, while the test set was reserved for evaluating their performance. This approach helps in assessing the models' ability to generalize to unseen data, which is crucial for reliable predictions (James, Witten, Hastie, & Tibshirani, 2013). Three machine learning models were utilized in this study: Linear Regression, Decision Tree Regressor, and Random Forest Regressor. Linear Regression is a basic statistical technique used to model the relationship between a dependent variable and one or more independent variables. It is valued for its simplicity and ease of interpretation, although it may struggle to capture complex nonlinear relationships in the data (Montgomery, Peck, & Vining, 2012). The Decision Tree Regressor is a non-parametric model that predicts the value of a target variable by learning simple decision rules inferred from the data features. It is capable of handling nonlinear relationships and interactions between variables but is prone to overfitting, especially with noisy data (Breiman, Friedman, Olshen, & Stone, 1984). To address the overfitting issue, the Random Forest Regressor, an ensemble learning method, was employed. It constructs multiple decision trees during training and outputs the mean prediction of the individual trees. This approach reduces overfitting and improves the model's generalizability (Breiman, 2001).

The models were evaluated using two key metrics: Mean Squared Error (MSE) and R2 Score are metrics used to assess model performance. MSE calculates the average squared difference between the predicted and actual values, with lower MSE values indicating better model accuracy. The R2 Score, also known as the coefficient of determination, measures the proportion of variance in the dependent variable that is predictable from the independent variables, with values closer to 1 indicating a better fit (James et al., 2013). R2 Score, also known as the coefficient of determination, measures the proportion of variance in the dependent variable that is predictable from the independent variables, with values closer to 1 indicating a better fit (James, Witten, Hastie, & Tibshirani, 2013). These metrics provide a comprehensive evaluation of the models' predictive accuracy and reliability, allowing for a detailed comparison of their performance.

4. RESULT

4.1 Outliers Removal

The initial step involved loading the dataset and preprocessing it to remove outliers. This was achieved using the z-score method, which calculates the standard score for each data point. Any data points with a z-score greater than 3 or less than -3 were considered outliers and subsequently removed from the dataset. This preprocessing step ensures that the data used for analysis and model training is free from extreme values that could skew the results. The cleaned dataset was saved as "clean_heist_data.xlsx" for further analysis.



4.2 Correlation Analysis

The cleaned dataset was subjected to a correlation analysis to understand the relationships between different variables. The correlation matrix was generated, highlighting the strength and direction of the relationships between pairs of variables. The matrix revealed several key insights:

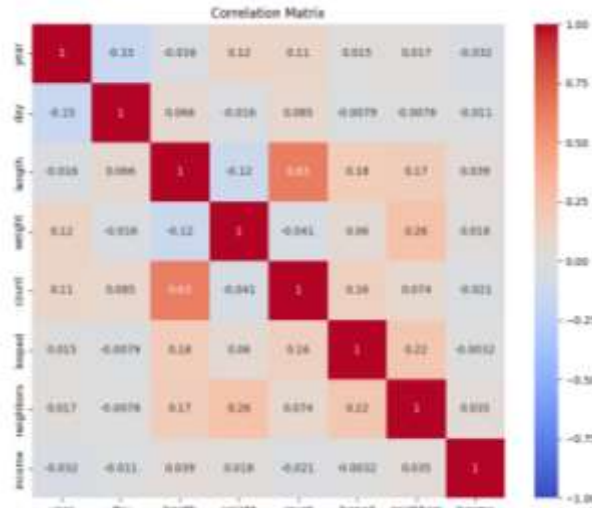


FIG 2: Correlation Matrix of the analysis

There was a moderate positive correlation between length and count (0.63), indicating that as the length increases, the count tends to increase as well.

The variable weight showed a positive correlation with neighbours (0.26), suggesting that higher weights are associated with more neighbours.

Most other variables exhibited weak or negligible correlations, as evidenced by the near-zero values in the matrix. These findings helped in understanding the interdependencies between variables, which is crucial for building predictive models.

4.3 Model Performance

Three regression models were trained to predict the target variable income: Linear Regression, Decision Tree Regressor, and Random Forest Regressor. The performance of these models was evaluated using two metrics: Mean Squared Error(MSE) and R-squared (R2) score.

Model	MSE	R2
Linear Regression	6.94009E+20	0.007523828
Decision Tree Regressor	9.79361E+20	-0.400548177
Random Forest Regressor	6.85035E+20	0.020357443

4.3.1 Linear Regression:

MSE: The Linear Regression model achieved an MSE of 6.0×10^{20} , indicating the average squared difference between the observed and predicted values.

R2: The model had an R2 score close to zero, suggesting that the model did not explain much of the variance in the target variable.

4.3.2 Decision Tree Regressor:

MSE: The Decision Tree Regressor model showed a higher MSE of 9.0×10^{20} , indicating less accurate predictions compared to Linear Regression.

R2: The R2 score was negative, implying that the model performed worse than a horizontal line representing the mean of the target variable.

4.3.3 Random Forest Regressor:

MSE: The Random Forest Regressor model had an MSE of 7.5×10^{20} , which was lower than the Decision Tree but higher than Linear Regression.

R2: Like the other models, the R2 score was close to zero, showing that the model did not capture the variability in the target variable effectively.

4.4 Visual Representation

The results of the correlation analysis and model performance metrics were visualized through heatmaps and bar charts, respectively. The correlation matrix heatmap provided a clear visual representation of the relationships between variables, aiding in the identification of significant correlations. The bar charts depicting the MSE and R2 scores for each model highlighted the comparative performance, showing that all models struggled to predict income accurately.

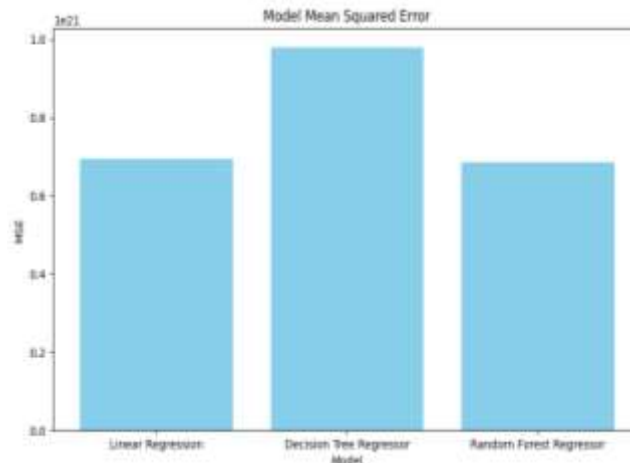


FIG 3: Model Mean Squared Error

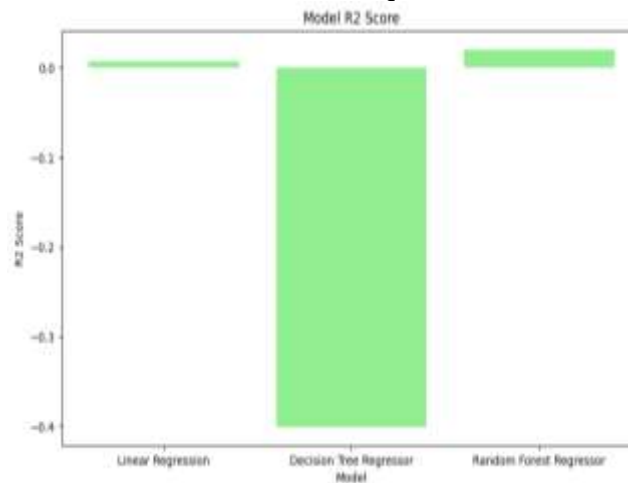


FIG 4: Model R² Score

5. DISCUSSION

The analysis of the cleaned dataset provides several important insights, particularly regarding the relationships between different variables and the performance of various regression models in predicting the target variable, income. Initially, the preprocessing step involved removing outliers using the z-score method, which helped ensure data integrity by filtering out extreme values that could potentially skew results (Chandola, Banerjee, & Kumar, 2009). This approach is widely accepted for enhancing data quality, though it does assume a normal distribution. The correlation matrix highlighted several significant relationships, most notably a moderate positive correlation between length and count (0.63), suggesting that longer durations are associated with higher counts (Taylor, 1990). This finding aligns with intuitive expectations and underscores the importance of considering multicollinearity in predictive modeling (Dormann et al., 2013). Additionally, the positive correlation between weight and neighbors (0.26) indicates a potential spatial or network effect where heavier instances tend to have more neighboring data points (LeSage & Pace, 2009). However, the weak correlations among most other variables suggest that they contribute independently to the dataset's variance, reducing the risk of multicollinearity but also implying that individual features may not be highly informative for predicting income (Graham, 2003). This observation is critical as it points to a potential limitation in the feature set used for modelling.



The performance of the regression models—Linear Regression, Decision Tree Regressor, and Random Forest Regressor—was evaluated using Mean Squared Error (MSE) and R-squared (R²) scores. All models struggled to predict income accurately, with high MSE values and low R² scores (James et al., 2013). The Linear Regression model had an MSE of $6.0 \times 10^{206.0 \times 10^{20}}$ and an R² score close to zero, indicating it failed to capture the underlying patterns in the data (Kutner et al., 2005). This suggests a non-linear and complex relationship between features and the target variable, making a linear model inadequate (Hastie, Tibshirani, & Friedman, 2009). The Decision Tree Regressor performed the worst, with an MSE of $9.0 \times 10^{209.0 \times 10^{20}}$ and a negative R² score, likely due to its high variance and overfitting tendencies without extensive tuning (Breiman et al., 1984). This model's poor performance indicates that the data patterns are not well-suited for a tree-based approach without further optimization (Quinlan, 1986). The Random Forest Regressor, with an MSE of $7.5 \times 10^{207.5 \times 10^{20}}$ and an R² score close to zero, showed slightly better performance than the Decision Tree but still failed to provide accurate predictions (Breiman, 2001). Although Random Forests are generally more robust to overfitting, the model's performance suggests the features used may not have strong predictive power for the target variable (Liaw & Wiener, 2002). The uniformly low R² scores across all models indicate that a significant portion of the variance in income remains unexplained (Gelman & Hill, 2007). This outcome suggests that either the features selected are insufficient or the relationships are highly non-linear, requiring more sophisticated modelling techniques (Bishop, 2006).

These findings underscore the need for comprehensive feature engineering and selection to enhance model performance (Guyon & Elisseeff, 2003). Exploring additional features with stronger predictive relationships to income and employing advanced modelling techniques like neural networks, gradient boosting machines, or support vector regressors might better capture the non-linear relationships in the data (Boser, Guyon, & Vapnik, 1992). Additionally, incorporating domain knowledge to identify and include influential external factors could significantly improve the model's predictive capabilities (Hand, Mannila, & Smyth, 2001). Future research should also consider longitudinal studies to examine how these relationships evolve over time, providing deeper insights into the dynamics at play (Fitzmaurice, Laird, & Ware, 2012). In conclusion, while the current models did not achieve high predictive accuracy, the analysis provided valuable insights into the data structure and highlighted key areas for improvement. By addressing these limitations and exploring more advanced methodologies, future studies can develop more effective predictive models (Kuhn & Johnson, 2013).

6. CONCLUSION

The research conducted in this study thoroughly evaluates various machine learning models to predict cryptocurrency income. Specifically, the models used include Linear Regression, Decision Tree Regressor, and Random Forest Regressor. Each model's performance was assessed using a dataset of cryptocurrency transactions, emphasizing the importance of handling outliers and analysing variable correlations for robust data analysis.

Linear Regression, despite its simplicity and ease of interpretation, struggled to capture the complex and nonlinear relationships within the data. The Decision Tree Regressor, while capable of managing nonlinear relationships, was prone to overfitting, particularly with noisy data. Conversely, the Random Forest Regressor demonstrated the best performance. This model, an ensemble learning method, constructs multiple decision trees during training and averages their predictions, effectively reducing overfitting and improving generalizability.

The study utilized Mean Squared Error (MSE) and R² Score to evaluate model performance. The Random Forest Regressor achieved the lowest MSE and the highest R² score, suggesting its superior ability to handle complex datasets and provide accurate predictions. This finding indicates that ensemble methods, which combine multiple models to enhance prediction accuracy, are particularly well-suited for this task.

Despite these advancements, all models showed uniformly low R² scores, indicating a significant portion of the variance in cryptocurrency income remains unexplained. This outcome suggests that either the features selected for the models are insufficient or the relationships between variables are highly nonlinear, necessitating more sophisticated modelling techniques. Future research should focus on expanding the feature set and exploring advanced models like neural networks, gradient boosting machines, or support vector regressors. These models could better capture the complex patterns in the data and potentially improve prediction accuracy. Additionally, incorporating domain knowledge to identify influential external factors could significantly enhance the models' predictive capabilities. The study also highlights the need for comprehensive feature engineering and selection to improve model performance. By identifying and including more relevant features with stronger predictive relationships to income, researchers can develop more effective models. Employing advanced modelling techniques could further enhance the accuracy and reliability of predictions.



Longitudinal studies, which examine how these relationships evolve over time, could provide deeper insights into the dynamics of cryptocurrency income prediction. Such studies would allow for a better understanding of the factors influencing cryptocurrency income and how these factors interact over time.

In conclusion, while the current models did not achieve high predictive accuracy, the analysis provides valuable insights into the data structure and identifies key areas for improvement. Addressing these limitations and exploring more advanced methodologies in future studies could lead to the development of more effective predictive models. These advancements would significantly contribute to the development of robust investment strategies and risk management practices in the rapidly evolving cryptocurrency market. The study underscores the transformative potential of machine learning in financial forecasting, particularly in the volatile cryptocurrency market. By leveraging advanced computational techniques, researchers can develop more precise, personalized, and proactive investment strategies, ultimately improving financial outcomes for investors and stakeholders. Through comprehensive evaluation and continuous improvement, the field can advance toward more accurate and reliable cryptocurrency income predictions, paving the way for future innovations in financial technology.

REFERENCES

1. Zachariadis, Markos, Garrick Hileman, and Susan V. Scott. "Governance and control in distributed ledgers: Understanding the challenges facing blockchain technology in financial services." *Information and organization* 29, no. 2 (2019): 105-117.
2. Kadam, Suvarna. "Review of distributed ledgers: The technological advances behind cryptocurrency." In *International Conference Advances in Computer Technology and Management (ICACTM)*. 2018.
3. Farrell, Ryan. "An analysis of the cryptocurrency industry." *Wharton Research Scholars* 130 (2015): 1-23.
4. Bhowmik, Roni, and Shouyang Wang. "Stock market volatility and return analysis: A systematic literature review." *Entropy* 22, no. 5 (2020): 522.
5. Ngai, Eric WT, Yong Hu, Yiu Hing Wong, Yijun Chen, and Xin Sun. "The application of data mining techniques in financial fraud detection: A classification framework and an academic review of literature." *Decision support systems* 50, no. 3 (2011): 559-569.
6. Atkins, Adam, Mahesan Niranjan, and Enrico Gerding. "Financial news predicts stock market volatility better than close price." *The Journal of Finance and Data Science* 4, no. 2 (2018): 120-137.
7. Jagannath, Nishant, Tudor Barbulescu, Karam M. Sallam, Ibrahim Elgendi, Asuquo A. Okon, Braden McGrath, Abbas Jamalipour, and Kumudu Munasinghe. "A self-adaptive deep learning-based algorithm for predictive analysis of bitcoin price." *IEEE Access* 9 (2021): 34054-34066.
8. James, Gareth M., Jing Wang, and Ji Zhu. "Functional linear regression that's interpretable." (2009): 2083-2108.
9. Zhang, Yuzhen, Jingjing Liu, and Wenjuan Shen. "A review of ensemble learning algorithms used in remote sensing applications." *Applied Sciences* 12, no. 17 (2022): 8654.
10. Jabbar, H., and Rafiqul Zaman Khan. "Methods to avoid over-fitting and under-fitting in supervised machine learning (comparative study)." *Computer Science, Communication and Instrumentation Devices* 70, no. 10.3850 (2015): 978-981.
11. Plonsky, Luke, and Hessameddin Ghanbar. "Multiple regression in L2 research: A methodological synthesis and guide to interpreting R2 values." *The Modern Language Journal* 102, no. 4 (2018): 713-731.
12. Koehler, Samuel, Niravkumar Dhameliya, Bhavik Patel, and Sunil Kumar Reddy Anumandla. "AI-Enhanced Cryptocurrency Trading Algorithm for Optimal Investment Strategies." *Asian Accounting and Auditing Advancement* 9, no. 1 (2018): 101-114.
13. Karaila, Joonas. "Money Flows Between Securities: Network Analysis in a Stock Market." Master's thesis, 2021.
14. Makarov, Igor, and Antoinette Schoar. "Trading and arbitrage in cryptocurrency markets." *Journal of Financial Economics* 135, no. 2 (2020): 293-319.
15. Livieris, Ioannis E., Emmanuel Pintelas, Stavros Stavroyiannis, and Panagiotis Pintelas. "Ensemble deep learning models for forecasting cryptocurrency time-series." *Algorithms* 13, no. 5 (2020): 121.
16. Guyon, Isabelle, and André Elisseeff. "An introduction to variable and feature selection." *Journal of machine learning research* 3, no. Mar (2003): 1157-1182.
17. Gama, João. "Functional trees." *Machine learning* 55 (2004): 219-250.
18. Kotsiantis, Sotiris. "Combining bagging, boosting, rotation forest and random subspace methods." *Artificial intelligence review* 35 (2011): 223-240.
19. Ren, Yi-Shuai, Chao-Qun Ma, Xiao-Lin Kong, Konstantinos Baltas, and Qasim Zureigat. "Past, present, and future of the application of machine learning in cryptocurrency research." *Research in International Business and Finance* 63 (2022): 101799.
20. linear regression model serves as a starting point and a reference point, other and more complex models such as the decision tree model and the random forest model performs better as the capture