# BCSE498J  Project-II – Capstone Project

# XRP CRYPTOCURRENCY RETURN PREDICTION

*Submitted in partial fulfillment of the requirements for the degree of*

## Bachelor of Technology

*in*

## Computer Science and Engineering in Bio Informatics

*by*

**21BCB0112   RAVI PRANAVI**

**21BCE2127   NAMINENI SAI LOUKYA**

**Under the Supervision of**

**RAMANATHAN LAKSHMANAN**

Professor Grade 1

School of Computer Science and Engineering (SCOPE)



Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)

April 2025

# DECLARATION

I hereby declare that the project entitled "**XRP CRYPTOCURRENCY RETURN PREDICTION**" submitted by me, for the award of the degree of *Bachelor of Technology in Computer Science and Engineering in Bio Informatics* to VIT is a record of bonafide work carried out by me under the supervision of Prof. Ramanathan Lakshmanan.

I further declare that the work reported in this project has not been submitted and will not be submitted, either in part or in full, for the award of any other degree ordiploma in this institute or any other institute or university.

Place   : Vellore

Date    :

**Signature of the Candidate**

# CERTIFICATE

This is to certify that the project entitled **XRP CRYPTOCURRENCY RETURN PREDICTION** submitted by NAMINENI SAI LOUKYA(21BCE2127) and RAVI PRANAVI(21BCB0112), **School of Computer Science and Engineering**, VIT, for the award of the degree of *Bachelor of Technology in Computer Science and Engineering*, is a record of bonafide work carried out by him / her under my supervision during Winter Semester 2024-2025, as per the VIT code of academic and research ethics.

The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma in this institute orany other institute or university. The project fulfills the requirements and regulations of the University and in my opinion meets the necessary standards for submission.

Place : Vellore
Date :

**Signature of the Guide**

**Internal Examiner**                                                      **External Examiner**

**Dr. Rajkumar S**
**Computer Science and Engineering in Bio Informatics**

# EXECUTIVE SUMMARY

The XRP Cryptocurrency Return Prediction project aims to build a machine learning-based framework for forecasting the short-term returns of XRP, one of the most prominent digital assets in the cryptocurrency market. The primary goal is to predict whether the future return of XRP will be positive or negative, enabling traders and investors to make more informed decisions.

1. Data Preprocessing and Feature Engineering: The dataset used in this project includes historical XRP trading data, such as Open, High, Low, Close prices, and trading Volume. Key preprocessing steps involved converting timestamps into a standard date format and calculating the future return to create a binary classification target. Additionally, several technical indicators were engineered to enhance model input, including the Simple Moving Average, Exponential Moving Average, Relative Strength Index, and return volatility.

2. Model Building: Four different machine learning models were implemented and evaluated: XGBoost Classifier, Random Forest Classifier, Support Vector Machine, and Linear Regression. The data was split into training and testing sets, and each model was trained on the training set to learn patterns in the technical indicators that could predict the binary return direction.

3. Evaluation and Performance Metrics: The models were assessed using multiple performance metrics, including Accuracy, F1 Score, ROC-AUC Score, and Confusion Matrix. XGBoost and Random Forest emerged as the top performers, showing strong capabilities in handling complex, non-linear relationships in the data. The results confirmed that ensemble models were more effective in capturing the intricacies of financial time series data compared to simpler classifiers.

4. Feature Importance and Insights: Feature importance analysis revealed that RSI, volatility, and moving averages were among the most significant indicators influencing model predictions. These insights not only validate common trading strategies but also highlight the potential of machine learning to identify subtle signals in market data that may be overlooked through traditional methods.

5. Conclusion: The project successfully demonstrates the viability of using machine learning techniques, especially ensemble models, to predict the directional movement of XRP returns. This approach lays the groundwork for developing more robust predictive models and trading algorithms. Future improvements could include incorporating macroeconomic indicators, alternative cryptocurrencies, sentiment analysis, and deep learning architectures for further enhancement.

# <u>ACKNOWLEDGEMENTS</u>

# TABLE OF CONTENTS

# List of Tables

# List of Figures

# List of Abbreviations

| | |
|---|---|
| XGBoost | Extreme Gradient Boosting |
| SVR | Support Vector Regression |
| RMSE | Root Mean Squared Error |
| $R^2$ | Coefficient of Determination |
| MSE | Mean Squared Error |
| MAE | Mean Absolute Error |
| ROC | Receiver Operating Characteristic |
| F1 Score | Harmonic mean of precision and recall |
| AUC | Area Under Curve |
| API | Application Programming Interface |
| CSV | Comma-Separated Values |
| EDA | Exploratory Data Analysis |
| RSI | Relative Strength Index |
| MACD | Moving Average Convergence Divergence |
| RNN | Recurrent Neural Network |
| LSTM | Long Short-Term Memory network |
| KNN | K-Nearest Neighbors |
| LightGBM | Light Gradient Boosting Machine |
| GRU | Gated Recurrent Unit |
| Bi-LSTM | Bidirectional Long Short-Term Memory |
| SHAP | SHapley Additive exPlanations |

# Symbols and Notations

$\delta f$                                CFO

$\varepsilon$                                NCFO

# 1. INTRODUCTION

The rapid growth of the cryptocurrency market has attracted significant attention from investors, researchers, and technologists alike. Among the many digital assets, XRP, developed by Ripple Labs, stands out due to its unique consensus mechanism, use cases in cross-border payments, and high transaction speed. With the increasing volatility and speculative nature of crypto markets, accurate price prediction models are essential for informed decision-making and risk management.

This project aims to develop and compare various machine learning models for predicting the future price of XRP, using historical market data. By leveraging techniques such as XGBoost, Random Forest, Linear Regression, and Support Vector Regression, this project explores the effectiveness of each model in capturing market trends and forecasting returns.

Comprehensive data preprocessing, feature engineering, and hyperparameter tuning were applied to enhance model performance. Furthermore, advanced evaluation metrics, including $R^2$ score, RMSE, F1 score, ROC-AUC, and confusion matrices, were employed to assess predictive accuracy and reliability. Feature importance analyses were also conducted to gain insights into the factors most influential in determining XRP's price movements.

Ultimately, this project not only provides a robust comparison of predictive models but also offers a valuable framework for future research and development in cryptocurrency analytics and financial forecasting.

## 1.1 BACKGROUND

The rapid growth of cryptocurrencies over the past decade has transformed the financial landscape, offering new opportunities and challenges for investors, analysts, and technologists. Among the many digital currencies available, XRP, developed by Ripple Labs, has gained significant attention due to its unique use case in facilitating fast and low-cost cross-border payments. As with other cryptocurrencies, XRP exhibits high price volatility, making accurate price movement predictions both critical and challenging for market participants.

With the evolution of data-driven strategies and increased availability of historical trading data, machine learning has emerged as a powerful tool for financial forecasting. Unlike traditional statistical models, machine learning algorithms can learn complex, non-linear relationships in large datasets, making them well-suited for predicting market trends and asset returns. In particular, technical indicators derived from historical price data such as moving averages and momentum indicators are widely used as input features for predictive models.

The financial markets, including the cryptocurrency space, are heavily influenced by patterns, trends, and trader behavior, which can be effectively captured through technical analysis. This project leverages a combination of technical indicators and machine learning models to build a predictive framework that classifies whether the return of XRP in the next time period will be positive or negative. By integrating these methods, the project seeks to offer insights into market direction, improve decision making, and contribute to the growing field of algorithmic trading in the cryptocurrency domain.

Overall, the use of machine learning for return prediction is not just an academic exercise it has real-world implications for automated trading systems, portfolio management, and risk mitigation. This project is a step toward developing intelligent systems that can adapt to the dynamic and often unpredictable behavior of digital assets like XRP.

## 1.2 MOTIVATIONS

The primary motivation behind this project lies in the highly volatile nature of the cryptocurrency market, which presents both significant opportunities and risks for investors. XRP, in particular, has exhibited sharp price fluctuations due to factors such as market sentiment, regulatory news, and technological developments. Accurately predicting returns in such an unpredictable environment could offer traders and institutions a strategic advantage, helping them to optimize their entry and exit points, manage risks more effectively, and improve overall portfolio performance.

Another key motivation is the increasing demand for data driven investment strategies. Traditional investment approaches often fall short in capturing the dynamic and non-linear behaviors present in modern financial markets especially in crypto. Machine learning provides a powerful

alternative, enabling the development of models that can detect subtle patterns and relationships within large sets of historical data. This project seeks to apply such models to the XRP market, aiming to predict short-term returns based on technical indicators and trend analysis.

Furthermore, this project is motivated by the growing academic and industry interest in the intersection of artificial intelligence and finance. With the rise of algorithmic trading and fintech innovations, understanding how to apply machine learning to real-world financial problems has become a highly valuable skill. This project provides hands on experience in working with market data, feature engineering, model training, evaluation, and interpretation making it a practical learning experience for anyone interested in quantitative finance or data science.

Lastly, the motivation also stems from the need to contribute to the development of more transparent and explainable financial models. While black-box AI systems are often criticized for their lack of interpretability, this project emphasizes the importance of analyzing feature importance and model behavior. By understanding which indicators drive predictions, users can gain greater confidence in the system and use it as a complementary tool in their decision making processes.

## 1.3 SCOPE OF THE PROJECT

This project focuses on predicting short-term returns of the XRP cryptocurrency using historical market data and machine learning techniques. Specifically, it aims to forecast whether the return for the next day will be positive or negative, treating it as a binary classification problem. The dataset includes historical daily price information such as open, high, low, close, and volume. From this data, a variety of technical indicators are derived to serve as features for the prediction model. These indicators, including momentum, trend, and volatility measures, help capture the underlying patterns and dynamics of price movements.

The scope of the project includes data preprocessing, feature engineering, model training, evaluation, and interpretation. Multiple classification algorithms such as Linear Regression, Random Forest, and Gradient Boosting (e.g., XGBoost) are tested to determine which performs best on the task. The performance is evaluated using metrics like accuracy, precision, recall, and F1-score to ensure a balanced view of the model's predictive ability, especially in the presence of class imbalance. Additionally, hyperparameter tuning is performed to optimize the models and

improve generalization to unseen data. This project is limited to using publicly available historical data and does not incorporate external factors such as news sentiment, blockchain metrics, or macroeconomic indicators, although these could be considered in future extensions.

Overall, this project is intended as a proof-of-concept for applying machine learning techniques to cryptocurrency return prediction. It lays a foundational framework that can be expanded upon with more complex models, additional data sources, or real-time prediction systems in future iterations.

**Chapter 2**

# 2. PROJECT DESCRIPTION AND GOALS

## 2.1 LITERATURE REVIEW

The study "Predicting Bitcoin Returns Using High-Dimensional Technical" published in 2019 analyzes the predictability of Bitcoin returns using a high-dimensional set of technical indicators. The authors state that it is difficult to predict cryptocurrency prices due to their fast changes and lack of historical data. The use of traditional analytical methods is inappropriate, and the need for more advanced predictive techniques arises. A CART model with 1000 trees is used in the study and its test data performance is tested. The results show that the model has a high predictive power for small daily Bitcoin returns with a high winning ratio, indicating that technical indicators may add useful information to the return forecasting of Bitcoin.

The article titled "Deep Learning for Bitcoin Price Direction Prediction: Models and Trading Strategies Compared (2024)" compares different deep learning models of Bitcoin price direction prediction, concentrating more on data on- chain with feature selection techniques. It suggests gaps in research in deep learning models and restricted research in features like Boruta, Genetic Algorithm, and LightGBM for feature selection methods. Evaluating the accuracy of these models and testing trading strategies by comparing CNN- LSTM, LSTNet, and TCN with ARIMA as a benchmark, the study includes the application of deep learning models. The results show that CNN-LSTM in conjunction with Boruta feature selection achieved the maximum accuracy of 82.44%, while the long and short strategy generated based model prediction produced an annual return of 66.54%. The findings give an idea about the efficiency of deep learning models than traditional forecasting methods.

The authors of the paper "Bitcoin Return Prediction Based on OLS, Random Forest, LightGBM, and LSTM (2023)", focus on predicting Bitcoin returns rather than just price trends. It compares simple statistical methods like Ordinary Least Squares with advanced machine learning models, including Random Forest, LightGBM, and LSTM. Using daily Bitcoin prices along with Tesla stock prices, gold prices, and market volatility indices, the study employs binary classification to

predict whether Bitcoin will generate a profit or loss the next day. Interestingly, the lowest complexity model achieved the best accuracy of 55%. In contrast, the LightGBM and LSTM models, with higher complexity, performed close to random guessing. This shows that though machine learning sounds very promising, simple models may still offer competitive predictive capabilities.

In the paper titled "XRP Price Prediction Using Advanced Deep Learning Models (2023)" addresses the shortcomings of the conventional models for cryptocurrency price prediction, the current research evaluates how some deep learning architectures, including LSTM, Bi-LSTM, and GRU, can be used to make price predictions for XRP. The Binance API is set up for the collection of open, high, low, and close prices, in addition to volume metrics. Ensemble methods are used in the study to improve the accuracy of the forecasts. The MSE and RMSE are used to evaluate the models. The results show that Bi-LSTM has the highest performance with the lowest error, and GRU realizes the broad market trends. These findings can be extremely helpful to any trader or investor who is looking for accurate short-term price forecasts of XRP.

The paper "Forecasting Returns Volatility of Cryptocurrency by Applying Various Deep Learning Algorithms (2023)" investigates the possibility of deep learning algorithms capable of forecasting cryptocurrency volatility with advanced capabilities to better be more accurate than traditional models such as GARCH and ARIMA. Performance of CSS, GMDH-NN, and NNETAR were compared using daily return series of Bitcoin, Ethereum, XRP, and Tether. RMSE and MAE were applied to measure performance. Results highlight the best performing model for Bitcoin and XRP is CSS; an optimal model for Ethereum is NNETAR; for Tether, GMDH NN model outperforms other models. Hence different models are required for each cryptocurrency which proves the complexity of the cryptocurrency market's behaviour.

The conference paper "Cryptocurrency Price Forecasting Using XGBoost Regressor and Technical Indicators (2024)" examines the effectiveness of XGBoost, an advanced machine learning model, in cryptocurrency price forecasting. This research also shows the failure of the traditional models regarding capturing crypto markets high volatility. Using historical data collected at 15-minute intervals from Binance related to Bitcoin, the study incorporates technical indicators, such as EMA, MACD, and RSI. After hyperparameter tuning with grid search, the XGBoost model achieves an impressive $R^2$ score of 0.9999 with low RMSE (59.95) and MAE (46.22), demonstrating superior predictive accuracy and reliability over conventional forecasting approaches.

The study titled "Bitcoin Price Prediction Using Recurrent Neural Networks and Long ShortTerm Memory (2024)," centers its attention on deep learning models about time series that particularly incorporate the RNN and LSTM approaches for predicting the price of Bitcoin. Using transaction volume, hash rate, and Google search trends from the Kraken exchange, the study preprocesses data by normalizing it with MinMaxScaler and splitting it into training (80%) and testing (20%). The results show both RNN and LSTM can fairly predict Bitcoin up to 20 minutes in advance, with deviations that are usually very small and minimal in the last steps. This gives analysts tools to inspect short-term moves in Bitcoin and manage risks much better.

The paper "Risks and Returns of Cryptocurrency (2021)" looks at influencing factors of cryptocurrency returns such as regulatory impacts and market microstructures. The daily returns, trading volume, and other investor attention metrics, like Google searches, and Twitter activity were analyzed in this study. It identifies key drivers of cryptocurrency returns where past performance predicts short-term future returns. Investor attention and the market to miner cost ratio are also indicators of returns, which reveals the dynamics of risk in returns for major cryptocurrencies such as Bitcoin and Ethereum.

In the paper titled "Deep Learning Algorithm to Predict Cryptocurrency Fluctuation Prices: Increasing Investment Awareness (2022)" the application of LSTM for time-series forecasting of cryptocurrency prices is investigated. This focuses on addressing prior omissions of excluding macroeconomic indicators and market sentiment. Preprocessing historical price and volume data, the authors test the model by evaluating MSE and RMSE. Results indicate that LSTM accurately predicts price fluctuations based on satisfactory accuracy. However, further research should be undertaken to include external market anomalies for enhanced forecasting.

The paper "Predicting XRP (Ripple) Cryptocurrency Price with Python and Machine Learning (2021)" applies multiple machine learning methods, such as SVC, Decision Tree, and Random Forest, to predict price movements of XRP. Evaluating the models using historical prices and technical indicators like RSI and Bollinger Bands, the study provides an assessment based on accuracy, precision, recall, and F1 scores. The results were mixed, indicating different levels of predictive accuracy and showing model strengths and weaknesses using confusion matrices and visualization techniques. It highlights the real potential of machine learning for price forecasting but emphasizes the necessity of complete model tuning and validation.

The paper "Daily Cryptocurrency Returns Forecasting and Trading via Machine Learning" analyzes machine learning models for predicting next-day cryptocurrency returns, focusing on Bitcoin (BTC), Ethereum (ETH), and Ripple (XRP). It evaluates eight models using various valuation factors and finds that Support Vector Machines (SVMs) outperform others, while XGBoost and AdaBoost perform the worst. A probability-based trading strategy based on these predictions achieves superior returns with a Sharpe ratio of 2.8. The study confirms that shortterm reversal patterns (7-day and 30-day trends) significantly impact return predictions. Despite its success, limitations include the lack of deep learning models and limited historical data.

This paper "Gold and Bitcoin Price Prediction based on KNN, XGBoost, and LightGBM Model" explores machine learning-based price prediction for Bitcoin (BTC) and Gold using historical data from 2017 to 2022. It applies three regression models—K-Nearest Neighbors (KNN), XGBoost, and LightGBM—to analyze price trends. LightGBM outperforms other models due to its efficiency, while XGBoost delivers high accuracy but requires more computational resources. The models are evaluated using metrics like MSE, RMSE, and R² Score. Results highlight the potential of machine learning in financial market forecasting. However, the study lacks external market factors, which could improve predictions. A larger dataset is needed to enhance accuracy. The research suggests that combining multiple models improves forecasting performance. The findings help investors make informed financial decisions. Overall, the study provides guidelines for future research in price prediction using machine learning.

## 2.2 RESEARCH GAP

a) Lack of Interpretability in Deep Learning Models
- Many studies use deep learning models like LSTM, GRU, and Bi-LSTM, but these operate as "black boxes" without offering clear explanations of how features contribute to price predictions.
- Future research could focus on explainable AI (XAI) techniques to enhance model interpretability.

b) Limited Feature Selection Approaches
- Several studies do not thoroughly investigate the impact of different features (e.g., market sentiment, regulatory news, macroeconomic indicators) on XRP price movements.

- Improved feature selection methods, such as hybrid approaches combining domain knowledge with automated selection techniques, could enhance prediction accuracy.

c) Impact of External Factors on XRP Volatility
- Most models focus on historical price data and technical indicators while ignoring external influences like government regulations, investor sentiment, and macroeconomic changes.
- Incorporating external data sources, such as social media sentiment analysis, news articles, and regulatory updates, may improve forecasting performance.

d) Real-time and High-frequency Data Limitations
- Many models are trained on historical data without considering the dynamic nature of real-time trading and liquidity constraints.
- Future work could explore real-time prediction models using streaming data to enhance responsiveness.

e) Model Comparisons and Ensemble Learning
- Some studies compare multiple deep learning models but do not explore ensemble techniques effectively.
- Research could benefit from hybrid models that combine multiple deep learning architectures for improved accuracy.

f) Macroeconomic and Sentiment-based Predictions
- Some papers suggest using macroeconomic indicators and sentiment analysis but do not integrate them effectively.
- Future studies could develop models that integrate news sentiment, market trends, and macroeconomic indicators.

g) Scalability and Computational Efficiency
- The computational cost of deep learning models remains a challenge, especially for real-time predictions.
- Optimization techniques such as pruning, quantization, or federated learning could improve efficiency.

h) Generalizability Across Cryptocurrencies
- Many models are tailored specifically to XRP, limiting their applicability to other cryptocurrencies.
- Research should explore transfer learning techniques to adapt models for different cryptocurrencies.

## 2.3 OBJECTIVES

This project has several key objectives centered around developing and evaluating machine learning models for XRP price prediction. These objectives encompass data preprocessing and feature engineering, model training and optimization, rigorous performance evaluation, feature importance analysis, and comparison of various machine learning algorithms. The goal is to identify and develop a robust and reliable model capable of accurately predicting XRP price movements, providing valuable insights for traders and investors.

1. Develop and compare various machine learning models for XRP price prediction:
   The project aims to implement and evaluate several regression models, specifically XGBoost, Linear Regression, Random Forest, Support Vector Regression (SVR), and K-Nearest Neighbors (KNN), to predict XRP's "Last Price."

2. Engineer relevant features from historical XRP data:
   The project involves feature engineering, creating new features like "High_Low_Diff" and a scaled "Change_Pct_Scaled" from the raw data, in addition to using "Open Price," "Max," "Min," and "Size."

3. Optimize model performance through hyperparameter tuning:
   The project utilizes GridSearchCV and RandomizedSearchCV to find the optimal hyperparameters for the XGBoost and Random Forest models, respectively, aiming to maximize their predictive accuracy.

4. Evaluate model performance using appropriate metrics:
   The project employs Root Mean Squared Error (RMSE) and R-squared (R2) as primary metrics to assess the accuracy and goodness-of-fit of the different models. It also explores

F1-score, AUC, and confusion matrices, treating the regression problem as a classification one by converting predictions and actual values to binary based on exceeding the mean value.

5. Analyze feature importance:

The project seeks to understand the relative importance of different features in influencing XRP prices by analyzing feature importance scores for XGBoost, Random Forest, and KNN and permutation importance for SVR.

6. Implement cross-validation:

The project uses cross-validation techniques to assess the general performance of the models and obtain more robust performance estimates.

7. Compare the performance of different models:

A key objective is to compare the performance of the various machine learning models to identify which one provides the most accurate and reliable predictions for XRP prices.

8. Make predictions on new data:

The trained best model is used to generate predictions on unseen test data, simulating real-world application of the model.

## 2.4 PROBLEM STATEMENT

In recent years, cryptocurrencies have emerged as a popular alternative asset class, attracting investors and traders globally. However, the inherent volatility and unpredictable nature of the crypto market, especially in the case of XRP, present significant challenges for financial forecasting and risk management. Unlike traditional financial markets, cryptocurrencies operate with fewer regulations and are heavily influenced by speculation, social media trends, and technological developments making reliable predictions even more complex.

XRP, developed by Ripple Labs, plays a unique role in the digital economy by facilitating fast and low-cost international transactions. Despite its potential, its market price is subject to dramatic swings, creating uncertainty for investors. Currently, there is a lack of reliable tools that can assist

in predicting whether XRP's returns will rise or fall in the short term, limiting the ability of stakeholders to make data informed decisions.

This project seeks to address this gap by leveraging machine learning techniques to develop a predictive model that forecasts the direction of XRP's daily returns. By analyzing historical price data and technical indicators, the model aims to provide a decision support tool that can reduce uncertainty and enhance strategic planning in XRP trading. The problem lies in capturing the non-linear and complex relationships within the data and translating them into actionable insights.

## 2.5 PROJECT PLAN

1. Data Collection and Preprocessing
   - Data Acquisition: Gather historical XRP price and volume data from a reliable source (e.g., a financial API or a CSV file). Ensure data covers a sufficient time period and frequency.
   - Data Cleaning: Handle missing values using appropriate methods (forward fill for "Last Price," mean imputation for "Size," as in the code). Identify and address outliers if necessary.
   - Feature Engineering: Create new features based on existing data. This includes calculating "High_Low_Diff" and a scaled "Change_Pct_Scaled," as in the code. Explore additional technical indicators (moving averages, RSI, MACD, etc.) and potentially sentiment analysis or on-chain metrics.
   - Data Scaling: Apply appropriate scaling techniques (StandardScaler, MinMaxScaler, or RobustScaler) to the features. Justify the chosen scaling method.

2. Model Development and Training
   - Model Selection: Choose the machine learning models to be evaluated (XGBoost, Linear Regression, Random Forest, SVR, KNN, as in the code). Justify the selection of these models.
   - Train/Test Split: Implement time-series cross-validation (walk-forward validation) to create training and testing sets. This is crucial for time series data. Ensure the test set is representative.
   - Model Training: Train each selected model on the training data.

- Hyperparameter Tuning: Use GridSearchCV or RandomizedSearchCV, or consider Bayesian Optimization, to optimize the hyperparameters of the models (especially XGBoost and Random Forest). Expand the hyperparameter search space if necessary.

3. Model Evaluation and Comparison
   - Performance Evaluation: Evaluate the performance of each model on the test set using appropriate metrics (RMSE, R-squared, MAE, MAPE, directional accuracy). Clearly define and justify the chosen metrics.
   - Model Comparison: Compare the performance of the different models and identify the best-performing one. Statistically justify the selection of the best model (e.g., using statistical tests).
   - Feature Importance Analysis: Analyze feature importance for the best-performing model to understand which features are most influential in predicting XRP prices. Use SHAP values for better explainability.
   - Cross-Validation: Perform time-series cross-validation on the best model to assess its generalization performance and obtain a more robust estimate of its performance.

4. Deployment and Reporting
   - Prediction on New Data: Apply the best-trained model to make predictions on new, unseen data (the "modified_dataset_xrp.csv" in the code).
   - Project Report: Compile a comprehensive project report documenting the entire process, including data collection, preprocessing, feature engineering, model development, evaluation, and results. Discuss limitations and future work.
   - Presentation: Prepare a presentation summarizing the project and its findings.

# 3. TECHNICAL SPECIFICATION

## 3.1 REQUIREMENTS

The XRP Cryptocurrency Return Prediction project involves the use of machine learning techniques to forecast the returns of XRP, a popular cryptocurrency. To implement this project, a solid foundation in Python programming is essential, as the entire workflow including data preprocessing, model training, and evaluation is developed using Python. Additionally, familiarity with data science concepts such as data cleaning, feature engineering, and model selection are necessary. Basic knowledge of statistics and financial time series analysis is also beneficial to understand return calculations and interpret model outcomes. The project requires the use of various Python libraries such as NumPy, Pandas, Matplotlib, Seaborn, and Scikit-learn, each serving a unique purpose from data manipulation to visualization and machine learning. A stable development environment like Jupyter Notebook or Google Colab is recommended to facilitate interactive analysis and model experimentation. Access to historical XRP price data is also required, either in CSV format or through an API, to train and test the predictive models.

## 3.1.1 Functional Requirements

1. Data Collection
- The system must be able to load historical XRP cryptocurrency data from a CSV file or online API.
- The data should include essential fields such as Date, Open, High, Low, Close, Adj Close, and Volume.

2. Data Preprocessing
- The system must clean the data by:
  - Removing null or missing values.
  - Ensuring proper formatting of the Date column and setting it as the index.
- The system should compute daily percentage returns using adjusted closing prices.
- The system must create lagged return features (e.g., 1-day lag, 2-day lag, etc.) for supervised learning.

- The return values should be classified into categories such as Up (positive return) and Down (negative return).

3. Exploratory Data Analysis (EDA)
- The system should provide statistical summaries of the data.
- The system must generate visualizations such as:
  - Line plots of price and returns.
  - Histograms of return distributions.
  - Correlation heatmaps of features.

4. Feature Selection and Dataset Preparation
- The system should select relevant features (lagged returns) for training the machine learning model.
- The dataset should be split into training and testing sets (e.g., 80% training, 20% testing).
- Input features (X) and target variable (y) should be defined for model training.

5. Model Training
- The system must implement and train a Logistic Regression model using the training dataset.
- The model should be able to learn from the lagged return features and classify the direction of return.

6. Model Evaluation
- The system must evaluate model performance using:
  - Accuracy score
  - Classification report (precision, recall, F1-score)
  - Confusion matrix
- The system should compare predicted labels with actual test set labels.

7. Performance Visualization
- The system should display the confusion matrix using Seaborn heatmaps.
- It must visualize the model's prediction performance in comparison to actual return directions.

8. Prediction Functionality

- The system should be able to take the latest lagged return values as input and predict the next day's return direction using the trained model.

## 3.1.2 Non-Functional Requirements

1. Performance Requirements

- The system should process and train on historical data (spanning several years) within a reasonable time frame (ideally under 1 minute for standard datasets).
- Predictions should be made in real-time or within 1–2 seconds once inputs are provided.
- Visualization rendering (graphs and heatmaps) should be completed in under 5 seconds to maintain interactivity.

2. Scalability

- The system should be able to handle increasing volumes of historical data (e.g., 10+ years of daily price data).
- The architecture should allow for future upgrades to support multiple cryptocurrencies or more complex models.

3. Reliability and Availability

- The system should consistently produce accurate and reproducible results for the same dataset and model parameters.
- It must handle unexpected input (e.g., missing values or incorrect formats) without crashing.
- In case of data load or model failure, the system should return user-friendly error messages or fallback behavior.

4. Usability

- The system should have a clear, easy-to-use interface or script output with:
    - Clear prompts for data input or file upload.
    - Well-labeled graphs and plots.
    - Simple-to-understand output predictions
- Documentation should be provided for how to run the program, interpret results, and modify inputs.

5. Maintainability

- The system codebase should be modular, with clearly separated concerns (e.g., data loading, preprocessing, modeling).
- Comments and docstrings should be included to explain functions and logic.
- External libraries used (e.g., pandas, sklearn, matplotlib) should be up-to-date and documented.

6. Portability

- The project should run on any system with Python 3.x and required libraries installed.
- The system should work across major operating systems (Windows, macOS, Linux).
- Dependencies should be listed in a requirements.txt file to facilitate easy setup using pip.

7. Extensibility

- The project should be designed to allow for:
  - Integration with real-time data feeds (e.g., API from CoinGecko or Binance).
  - Adding new machine learning models for performance comparison.

## 3.2 FEASIBILITY STUDY

The feasibility study for the XRP Cryptocurrency Return Prediction project indicates that the implementation of this system is both practical and achievable with current technologies and resources. Technically, the project leverages accessible tools and libraries such as Python, scikit-learn, pandas, and matplotlib, which are open-source and widely used in data science and machine learning. Economically, it is cost-effective since it does not require any paid software or high-end hardware and can be run on a standard computer system. Operationally, the project aligns well with the current demand for cryptocurrency forecasting tools, offering valuable insights for investors and analysts. The data required for training and prediction historical price data for XRP is readily available from various free APIs or datasets.

## 3.2.1 Technical Feasibility

1. Data Handling and Preprocessing:
The code demonstrates the technical feasibility of loading and manipulating cryptocurrency market data using the pandas library. Reading data from a CSV file is a standard and well

supported operation within pandas. The implemented techniques for handling missing values, utilizing forward fill and mean imputation, showcase the feasibility of addressing common data quality issues programmatically. Furthermore, the creation of new features, specifically 'High_Low_Diff' and 'Change_Pct_Scaled', through basic arithmetic and scaling operations using pandas and numpy, is technically straightforward and achievable within the Python environment.

2. Model Development and Training:

The code exhibits the technical feasibility of implementing and training multiple regression-based machine learning models. The utilization of the scikit-learn library for models like Linear Regression, Random Forest Regressor, and SVR, along with the xgboost library for the XGBoost Regressor, confirms the feasibility of integrating and utilizing these established machine learning algorithms within the Python ecosystem. The implementation of hyperparameter tuning using GridSearchCV and RandomizedSearchCV for optimizing model parameters demonstrates the technical capability of systematically searching for optimal model configurations. The inclusion of cross_val_score further highlights the feasibility of evaluating model performance on multiple data splits for a more robust assessment.

3. Model Evaluation and Interpretation:

The code demonstrates the technical feasibility of evaluating the trained models using relevant metrics. For regression tasks, mean_squared_error and r2_score are calculated, showcasing the ability to quantify the model's predictive accuracy. For classification related analysis, the code implements the calculation of roc_curve, auc, f1_score, and the generation of a confusion_matrix, indicating the technical feasibility of assessing model performance in a binary context. The use of matplotlib and seaborn for visualizing the predicted prices, ROC curves, and confusion matrices confirms the technical feasibility of presenting the model's behavior and performance graphically. Additionally, the implementation of feature importance analysis for different models demonstrates the technical capability of gaining insights into the contribution of each input feature to the model's predictions.

4. Prediction on New Data:

The code illustrates the technical feasibility of applying the trained model to predict on new, unseen data. The process of loading a new dataset, applying the same preprocessing steps (feature engineering and scaling using the pre-fitted StandardScaler), and utilizing the predict() method of the trained model confirms the technical capability of generalizing the learned patterns to new instances.

In summary, the provided Python code demonstrates a high level of technical feasibility across the core stages of a machine learning project for XRP cryptocurrency return prediction, encompassing data handling, model development, evaluation, and prediction on new data, by effectively utilizing established Python libraries and standard machine learning techniques.

### 3.2.2 Economic Feasibility

1. Development Costs

Personnel Costs:

The primary cost during the initial development phase would be the time and expertise of individuals involved. This includes data scientists or machine learning engineers to develop, implement, and test the code. The economic feasibility depends on the hourly or salaried rates of these professionals and the estimated duration of the development process. The provided code represents an initial effort, and further refinement, testing, and potential expansion would require additional personnel time.

Software and Library Costs:

The core libraries used in the provided code (pandas, numpy, scikit-learn, xgboost, matplotlib, seaborn) are open-source and freely available, thus incurring no direct licensing costs.

Data Acquisition Costs:

The economic feasibility is contingent on the source and volume of historical data required. If publicly available data sources are sufficient, the direct cost of data acquisition might be minimal. However, accessing more granular, real-time, or proprietary data feeds could involve subscription fees or API usage charges, impacting the economic viability.

Computational Infrastructure Costs (Development):

The initial development and experimentation using the provided code can likely be performed on standard personal computers or readily available cloud computing instances (e.g., free tiers or low-cost virtual machines). The economic impact at this stage is generally low.

2. Operational Costs:

Computational Infrastructure Costs (Deployment/Ongoing Use):

If the prediction model is to be deployed for continuous monitoring or real-time predictions, ongoing computational resources will be necessary. This could involve cloud computing costs for server instances, data storage, and API access. The economic feasibility depends on the scale of deployment, the frequency of predictions, and the efficiency of the model.

Data Maintenance Costs:

Maintaining the data pipeline, ensuring data quality, and potentially updating data sources can incur ongoing costs in terms of personnel time and potential subscription fees.

Monitoring and Maintenance Costs:

The prediction model will require ongoing monitoring for performance degradation and potential retraining. This involves personnel time for analysis, model updates, and system maintenance, impacting operational expenses.

3. Economic Feasibility Assessment:

The economic feasibility of this project is highly dependent on several factors:

- Prediction Accuracy: The primary driver of economic benefit is the accuracy and reliability of the predictions. A model with low predictive power will likely not generate sufficient returns to justify the development and operational costs.
- Trading Volume and Capital: The potential economic benefits are directly proportional to the volume of XRP being traded based on the model's insights and the capital invested.
- Development and Operational Efficiency: Minimizing development time, utilizing cost-effective infrastructure, and streamlining operational processes are crucial for economic viability.
- Market Volatility: The inherent volatility of cryptocurrency markets can amplify both potential gains and losses, impacting the overall economic outcome.

### 3.2.3 Social Feasibility

1. Accessibility and Democratization of Information:

If the prediction model and its insights are made accessible to a wider audience (beyond institutional investors or sophisticated traders), it could contribute to a more democratized understanding of cryptocurrency markets. This could empower individual investors with data-driven insights, potentially leveling the playing field.

The technical nature of the model and its outputs might require user-friendly interfaces and educational resources to be truly accessible to a non-technical audience.

2. Impact on Individual Investors:

Providing tools for predicting XRP returns could potentially assist individual investors in making more informed decisions, potentially leading to better financial outcomes and reduced exposure to impulsive or emotionally driven trading.

Over-reliance on a prediction model, especially if its accuracy is not fully understood or guaranteed, could lead to misguided investment decisions and financial losses for individuals. There's a risk of creating a false sense of certainty in a highly volatile market.

3. Ethical Considerations:

The complexity of some machine learning models (like XGBoost or deep learning approaches not present in the initial code but potentially developed later) can make their predictions difficult to interpret. Ensuring a degree of transparency and explainability in how the model arrives at its predictions is crucial for building trust and addressing ethical concerns.

The historical data used to train the model might contain inherent biases reflecting past market behavior or the actions of specific actors. If these biases are not identified and mitigated, the prediction model could perpetuate or even amplify them, potentially leading to unfair or discriminatory outcomes for certain market participants.

While the intent is likely benign, sophisticated prediction tools could theoretically be used for market manipulation if their predictions become widely trusted and acted upon by a large number of participants.

4. Regulatory and Legal Landscape:

The use of prediction models in financial markets is subject to evolving regulatory scrutiny. Depending on how the model is used and the advice it generates, it could fall under different regulatory frameworks related to financial advice or market manipulation. Understanding and adhering to these regulations is crucial for social and legal feasibility.

The social feasibility of this XRP cryptocurrency return prediction project presents a mix of potential benefits and risks. While it could democratize access to market insights and potentially empower individual investors, there are crucial ethical considerations related to transparency, bias, and the potential for misuse. Building public trust and navigating the evolving regulatory landscape will be essential for ensuring the positive social impact and long-term viability of this project. Careful consideration of the potential negative consequences and proactive measures to mitigate them are necessary for responsible development and deployment.

## 3.3 SYSTEM SPECIFICATION

The XRP Cryptocurrency Return Prediction system operates as a computational tool that leverages machine learning to estimate future XRP price behavior. It ingests historical market data for XRP and transforms this raw data into a set of predictive features through calculations and scaling. Several regression algorithms, including XGBoost, Linear Regression, Random Forest, and Support Vector Regression, learned patterns from historical data. The performance of these algorithms was optimized through hyperparameter tuning techniques. The accuracy and reliability of the trained models were assessed using various statistical metrics and visualizations.

Finally, the best-performing model generated predictions on new, unseen market data. The system's output consists of price forecasts intended to inform users about potential future XRP market trends.

### 3.3.1 Hardware Specification

1. Central Processing Unit (CPU):
A multi-core processor is crucial. Modern CPUs from Intel (Core i5, i7, i9 series or newer) or AMD (Ryzen 5, 7, 9 series or newer) are recommended.

The data manipulation tasks performed by pandas and numpy, as well as the training of machine learning models by scikit-learn and xgboost, can significantly benefit from parallel processing offered by multiple CPU cores. This allows for faster execution of loops, matrix operations, and ensemble model training, reducing the overall development time. A higher number of cores and higher clock speeds will generally lead to quicker processing.

2. Random Access Memory (RAM):

A minimum of 8 GB of RAM is necessary for handling moderately sized cryptocurrency datasets and running the Python environment along with the required libraries. 16 GB of RAM is highly recommended and provides a more comfortable working environment, especially when dealing with larger datasets, performing extensive feature engineering, or running multiple computationally intensive tasks concurrently (e.g., hyperparameter tuning).

Insufficient RAM can lead to the system relying heavily on slower disk-based virtual memory, significantly slowing down the development process and potentially causing crashes with large datasets. Having ample RAM ensures that data can be loaded and processed efficiently in memory, speeding up operations like data loading, merging, and model training.

3. Storage:

A Solid-State Drive (SSD) with at least 100 GB of free space is strongly recommended as the primary storage for the operating system, development environment (Python, IDE), the project code, and the cryptocurrency datasets.

SSDs offer significantly faster read and write speeds compared to traditional Hard Disk Drives. This drastically reduces the time taken to load datasets, save and load models, and perform other file I/O operations, leading to a more responsive and efficient development workflow.

While the initial dataset might be small, the need to experiment with different datasets, store intermediate results, and potentially save multiple trained models will consume storage space over time.

4. Graphics Processing Unit (GPU):

For the development phase of this specific code, a dedicated high-end GPU is not a strict requirement. The libraries primarily used (pandas, scikit-learn, xgboost) are CPU-bound for the tasks demonstrated. A standard integrated GPU (present in most modern CPUs) or a basic discrete GPU is sufficient for rendering the visualizations generated by matplotlib and seaborn.

While GPUs can significantly accelerate the training of certain types of machine learning models, particularly deep learning models, the provided code focuses on traditional machine learning

algorithms. If the project were to evolve to include neural networks or other GPU-accelerated models in the future, a compatible NVIDIA GPU with CUDA support would become a crucial hardware component. However, for the current scope of the provided code, the GPU's role is primarily for display purposes.

In summary, a well-rounded development machine for this project would prioritize a multi-core CPU, sufficient RAM (ideally 16 GB), and a fast SSD to ensure a smooth and efficient development experience. While a powerful dedicated GPU isn't necessary for the current code, it's a consideration for potential future expansions of the project into deep learning.

## 3.3.2 Software Specification

1. Programming Language and Libraries:

The core of this project was developed using the Python programming language. This choice was driven by Python's extensive ecosystem of scientific computing and machine learning libraries. Key libraries utilized include pandas for data manipulation and analysis, numpy for numerical computations, scikit-learn for implementing various machine learning models (Linear Regression, RandomForestRegressor, SVR), data preprocessing (StandardScaler), and model selection tools (train_test_split, GridSearchCV, RandomizedSearchCV, cross_val_score), xgboost for the gradient boosting model implementation, and scipy for statistical functions (randint for randomized search). Data visualization was achieved using matplotlib and seaborn libraries, enabling the generation of plots for price trends, model evaluation metrics (ROC curves, confusion matrices), and feature importance.

2. Data Handling and Storage:

The system was designed to ingest cryptocurrency market data stored in CSV file format, as demonstrated by the pd.read_csv() function. pandas provided the necessary data structures (DataFrames) and functionalities for efficient data loading, exploration, cleaning (handling missing values using forward fill and mean imputation), and feature engineering (creation of 'High_Low_Diff' and 'Change_Pct_Scaled'). The processed data was held in memory during the model training and evaluation phases. The final predictions on new data were outputted to a new CSV file using pandas' to_csv() function, facilitating data persistence and potential integration with other systems.

3. Machine Learning Models:

The project incorporated several supervised machine learning regression algorithms to predict the 'Last Price' of the cryptocurrency. These included a linear model (LinearRegression), an ensemble tree-based model (RandomForestRegressor), a support vector machine-based model (SVR), and a gradient boosting model (XGBRegressor). The selection of these diverse models allowed for the exploration of different underlying relationships within the data. Hyperparameter tuning was implemented using GridSearchCV and RandomizedSearchCV to optimize the performance of the XGBRegressor and RandomForestRegressor by systematically searching through predefined parameter spaces.

4. Model Evaluation and Validation:

The performance of the trained models was rigorously evaluated using appropriate metrics for regression tasks, namely Root Mean Squared Error (RMSE) and R-squared (R2) score. Additionally, for a binary classification perspective (predicting if the price would go above the mean), metrics such as the Receiver Operating Characteristic (ROC) curve, Area Under the Curve (AUC), F1-score, and Confusion Matrix were computed. Cross-validation (cross_val_score) was employed to obtain a more robust estimate of the models' generalization capabilities by training and evaluating them on multiple subsets of the training data.

5. Visualization Components:

The system included visualization capabilities powered by matplotlib and seaborn. These libraries were utilized to generate line plots comparing actual and predicted prices over time, ROC curves to assess the performance of the binary classification interpretation, heatmaps of confusion matrices to visualize classification accuracy, and bar plots to illustrate the relative importance of different features as determined by the various machine learning models. These visual aids facilitated the understanding of the model's behavior and the insights derived from the analysis.

# 4. DESIGN APPROACH AND DETAILS

## 4.1 SYSTEM ARCHITECTURE

The method used in this project incorporates various comprehensive machine learning models for analyzing accurate predictions of XRP prices comprising of historical market data and latest regression approaches. The overview of the process is:

1.Data Collection and Preprocessing:

The process starts by collecting some historical XRP price and volume data. The CSV file serves as the main repository for this dataset. Following data entry, the exploratory data analysis will be conducted to get a basic understanding of the structure and data types, as well as possible anomalies in the dataset. This will be followed by critical preprocessing steps that will ensure the data quality. Mean imputation is used in handling the missing values for the 'Size' column and for maintaining its statistical integrity while forward filling for the 'Last Price' column helps to retain temporal continuity which is a common problem in financial datasets. The techniques of data cleaning applied are quite important in building a robust and reliable predictive model for these analyses.

2. Feature Engineering:

The system uses feature engineering to boost the model's predictive accuracy. Feature engineering is the process of taking the raw data and generating new features that may provide valuable insight into the data. Particularly, the system calculates the percentage change in price and the daily price range. The goal with these engineered features is to create a more detailed picture of how the markets are acting. Then, a set of features are selected for training the model including "Open Price," "Max," "Min," "Size," "High_Low_Diff," and "Change_Pct_Scaled," based on their relevance and expected contribution to price prediction. It is this painstaking feature engineering and selection that is crucial to increasing the model's accuracy and ability to generalize.

3. Data Splitting and Scaling:

Two major processes that the dataset must undergo before the actual training of the model are data splitting and data scaling. Time-series-fit data splitting is carried out, since the data represent time-series data coming from cryptocurrency. This ensures that the model is validated against subsequent data. To maintain the temporal order of the data which is a necessity in time-series analysis, 20% of the available data is treated as test data. Feature scaling is achieved via StandardScaler to ensure equal input into the models and prevent features with larger scales from dominating the learning process. In making certain that the models remain reliable and objective in their predictions, this scaler approach is then used to scale features to have zero mean and unit variance that is helpful for algorithms perceiving great scaling effects on features, which in this case is Support Vector Regression.

4. Cross validation:

The system employs cross-validation to evaluate the model's generalization performance and provide a broader estimate of their competency. In particular, scikit learn's cross_val_score is used to perform 5-fold cross-validation on the XGBoost, Random Forest, and Linear Regression models. Negative mean squared error will be used as a scoring metric, and these cross-validation scores will be transformed into RMSE values. This method allows for a more faithful estimation of the performance of the model against unseen data while reducing overfitting chances. This way, a high degree of assurance not only that a model is well-fitting on the training data but also that it will generalize to fresh, unseen data is a critical factor in realworld application and is assured through cross-validation.

5. Model Selection and Training:

To identify the best approach for forecasting XRP prices, the framework employs and compares an array of regression machine learning models: Support Vector Regression, Random Forest, Tuned Random Forest, Linear Regression, and XGBoost. Hyperparameter tuning using GridSearchCV and RandomizedSearchCV serves to enhance the Random Forest and XGBoost models respectively. For this purpose, it is required to search through a predefined hyperparameter space in an exhaustive way to uncover the best performing hyperparameter combination. Then, each model trains on the fitted training data, thus allowing inferring from the latent trends and

relationships between the feature and target variable. Thereafter, the selected and trained models will allow a thorough evaluation and contrasting of their respective prediction powers.

6. Model Evaluation and Comparison:

Following training, model performance is evaluated and compared using various measures such as confusion matrices, the F1 score, R-squared (R2), and Root Mean Squared Error (RMSE). Model accuracy and robustness are checked using RMSE and R2 while the model's ability to estimate the direction of price movements is revealed by the F1 score and confusion matrices. To make it easy to compare the model's performance across these criteria, the system builds a DataFrame. To allow for intuitive understanding of the model's performance, visualizations such as expected versus actual price plots and bar charts for R2 scores are produced. Feature importance analyses are carried out in order to identify the primary drivers behind movements in XRP prices. In the process of model selection for the prediction of the XRP price, one must go through a comprehensive process of research and comparison.



Fig 1: System Architecture Diagram

## 4.2 DESIGN

## 4.2.1 Data Flow Diagram



Fig 2: Data Flow Diagram

1. CSV Dataset (datasetbcfinalmiss_modified.csv)

- This step refers to the raw data file you're using, which contains historical information about XRP prices, volume, indicators, and possibly technical and fundamental features.

- The goal here is to use this data to predict future returns or price direction of XRP.

2. Data Loading (Pandas)

- Load the dataset using Pandas in Python.

- This will convert the CSV into a DataFrame structure, which is easy to manipulate and analyze.
- Example: df = pd.read_csv('datasetbcfinalmiss_modified.csv')

3. Data Preprocessing (Missing values, Encoding, Scaling)
- Handle Missing Values: Fill, drop, or interpolate missing values in your features like price, volume, indicators, etc.
- Encoding: If categorical features are present (e.g., day of week), encode them numerically (one-hot, label encoding).
- Scaling: Normalize/standardize features using MinMaxScaler, StandardScaler etc., especially useful for models like SVM.

4. Feature Selection (Top features, correlation)
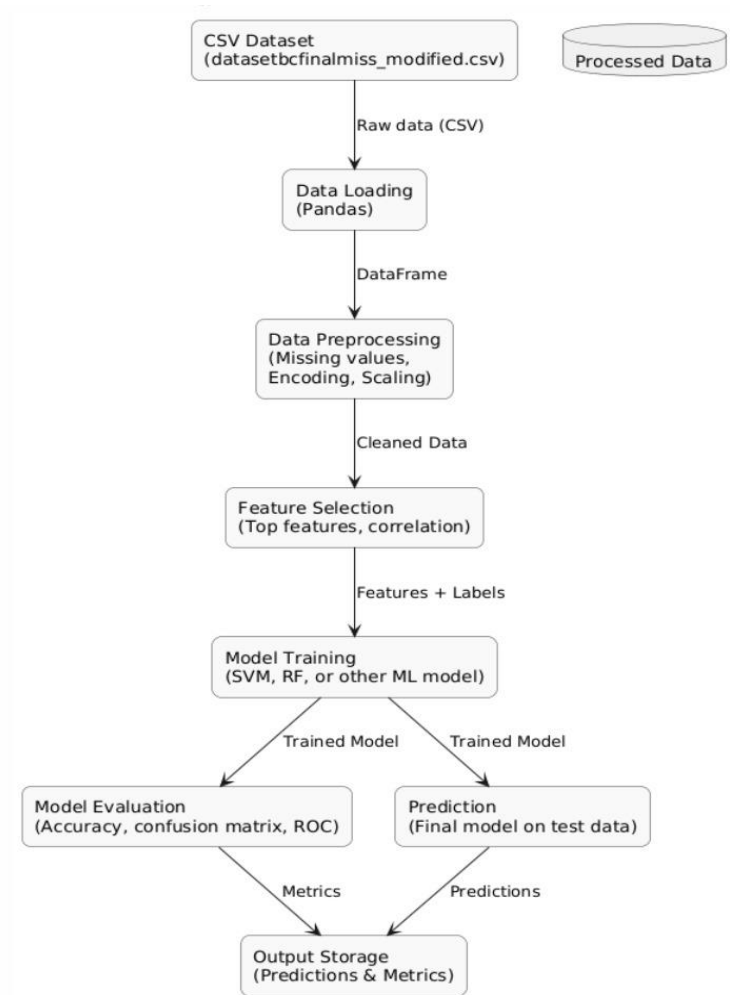- Identify and select the most relevant features for return prediction.
    - Use techniques like correlation matrix, mutual information, or feature importance from tree-based models.
- Examples of features:
    - Price change %,
    - RSI, MACD,
    - Volume trends,
    - Lagged returns.

5. Model Training (SVM, RF, or other ML model)
- Train your ML model (e.g., Support Vector Machine, Random Forest, Logistic Regression, or XGBoost) using:
    - Input features (X)
    - Target labels (y), e.g., binary label (Up/Down), or numeric returns.
- Split the dataset into training and testing (e.g., 80/20 split).

6. Model Evaluation (Accuracy, Confusion Matrix, ROC)
- Evaluate the model using:
    - Accuracy: How often it predicts the return direction correctly.
    - Confusion Matrix: Gives detailed info on true/false positives/negatives.
    - ROC-AUC: Useful for classification models, shows performance across thresholds.

7. Prediction (Final model on test data)

- Use the trained model to predict XRP return direction/values on the test dataset.

- Output could be:

  - Binary classification: Will the return be positive or negative?

  - Regression: Predict actual return percentage.

8. Output Storage (Predictions & Metrics)

- Save the output in files (CSV, Excel, or database):

  - Predictions for each test sample.

  - Evaluation metrics to assess performance.

- Helps in comparing different models or for future use in deployment.
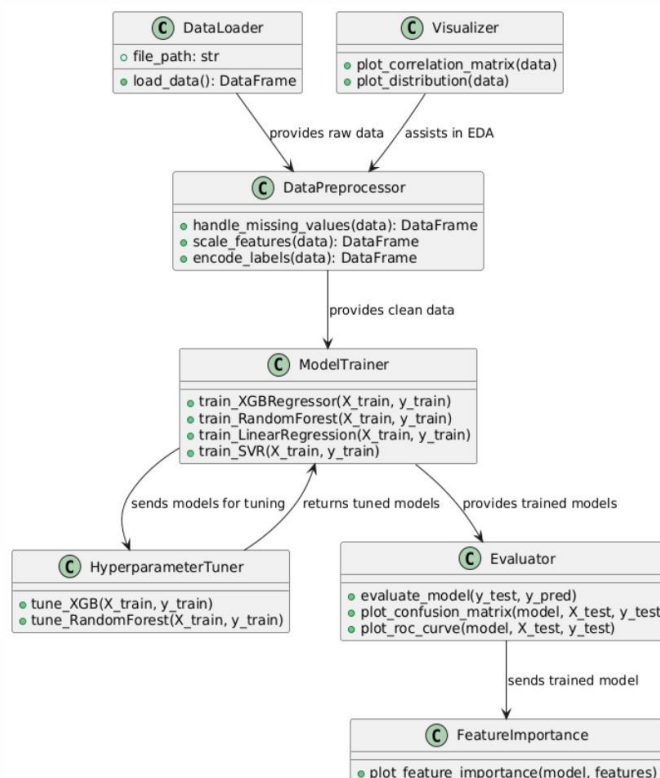
## 4.2.2 Class Diagram



Fig 3: Class Diagram

- The class diagram represents a structured machine learning pipeline designed to predict returns on XRP cryptocurrency. The pipeline begins with the DataLoader class, which is responsible for loading the raw data from a specified file path. It contains a file_path attribute and a method load_data() that reads the data (typically from a CSV or other file format) and returns it as a pandas DataFrame. This forms the foundational input for the entire pipeline.

- Once the raw data is loaded, it is often essential to perform exploratory data analysis (EDA) to better understand the dataset's structure and relationships. The Visualizer class supports this process through its methods plot_correlation_matrix(data) and plot_distribution(data). These functions help visualize how different features are related and how values are distributed, enabling insights into trends, outliers, or correlations in the XRP data before modeling begins.

- After EDA, the raw data is passed to the DataPreprocessor class, which is central to cleaning and transforming the dataset. It provides three main methods: handle_missing_values, which manages null or incomplete values; scale_features, which normalizes or standardizes numerical features; and encode_labels, which converts categorical data into numerical form. This ensures that the dataset is ready and suitable for feeding into machine learning models.

- With clean and preprocessed data in hand, the ModelTrainer class is used to train a range of machine learning models. It supports various algorithms such as XGBoost (train_XGBRegressor), Random Forest (train_RandomForest), Tuned Random Forest (train_TunedRandomForest), Linear Regression (train_LinearRegression), and Support Vector Regression (train_SVR). Each of these methods takes in training features and labels and outputs a trained model. The ModelTrainer can also send these models to the HyperparameterTuner for optimization.

- The HyperparameterTuner class focuses on improving model performance by fine-tuning the hyperparameters of selected models. It currently supports tuning for XGBoost (tune_XGB) and Random Forest (tune_RandomForest). Once optimized, these models are returned to the ModelTrainer or used directly in evaluation.

- The trained models are then evaluated using the Evaluator class. This class includes methods to measure prediction performance (evaluate_model) and visualizations like confusion matrices (plot_confusion_matrix) and ROC curves (plot_roc_curve). These evaluations help in comparing model performance and determining which model best predicts XRP returns.

- Finally, once the best-performing model is selected, it is passed to the FeatureImportance class. This class includes the plot_feature_importance(model, features) method, which visualizes which input features had the greatest impact on the model's predictions. This is particularly useful for interpreting and explaining model behavior in the context of cryptocurrency return forecasting.

# 5. METHODOLOGY AND TESTING

This experiment's primary objective is the development of multiple machine learning models that forecast the values of XRP cryptocurrency using historical trading data. Models will be compared to each other based upon cross-validation performance, R2 Score, and RMSE.

1. Dataset Details:

The data being used in forecasting price changes in cryptocurrencies involves: The last price-the closing price of the coin; the first price-the opening price of the cryptocurrency; Max-the greatest price attained during the period; Min-the lowest price attained; Size-gives information on market activity, indicating the trade volume; High Low Diff-computes the difference between maximum and minimum prices; Change Pct scaled-shows the scaled % change in price reflecting the size of price changes. These features comprehensively present relevant information with respect to the trading activities and fluctuations in price of cryptocurrency.

2. Data Preprocessing:

In order to ensure data quality and applicability on this XRP prediction challenge, the dataset has required some important activities for data pretreatment. The 'Last Price' column contains missing values; therefore, forward fill is used to maintain the temporal continuity. To preserve the general trade volume distribution, mean imputation is used to fill in the missing values for the 'Size' column. After then, feature engineering is done to extract fresh, maybe instructive features. To record the daily price range, 'High_Low_Diff' is developed, which is defined as the difference between the maximum and minimum prices. Furthermore, by multiplying the 'Change %' by 100, 'Change_Pct_Scaled' is produced, which denotes the scaled % change in price. To train and assess the models, the dataset is then divided into training and testing sets using an 80/20 split.

3. Cross-Validation:

For XGBoost, Random Forest, Tuned Random Forest, and Linear Regression, 5-Fold Cross-Validation is used to thoroughly assess the model's generalization and resilience.

By dividing the training data into five equal folds, the models are repeatedly trained on four of the folds and validated on the fifth. To make sure every fold is used as the validation set once, this procedure is carried out five times. Each validation fold's Root Mean Squared Error is computed, yielding a dispersion of performance metrics. Following that, the mean and standard deviation of these RMSE values is shown, providing information about the model's average performance and variability across various training data subsets.

4. Models Used:

- XGBoost – The XGBoost Regressor, being a gradient boosting algorithm on which GridSearchCV is performed for hyperparameter tuning optimizes several parameters.

- Linear Regression - In order to assess the linear relationship between the features and the target variable, Linear Regression is incorporated as a baseline, producing a straightforward but understandable model.

- Random Forest Regressor - In order to evaluate a wide range of parameter combinations, the Random Forest Regressor is first tested with its default configuration, which is then optimized through hyperparameter tuning using both GridSearchCV and RandomizedSearchCV.

- Tuned Random Forest - Trained to more accurately predict XRP price movements, in comparison to a default model, by systematically searching for the best hyperparameters to use in each of its multiple decision trees.

- Support Vector Regression - An RBF kernelSVR is a kernel-based regression model that makes use of feature standardization before fitting for best performance. Each model is selected to provide a fresh viewpoint on the prediction problem allowing for its efficacy to be thoroughly evaluated

5. Evaluation Metrics:

The RMSE gives the average magnitude of the errors in a set of predictions, without considering their direction. In addition, the R2 Score also known as the coefficient of determination describes how well a regression model fits the data. As such, it measures the proportion of variance in the target variable which the model could explain. The Cross-Validation RMSE gives an estimate of the model's generalization performance and robustness against unknown data by being a reliable indicator of how it would perform in actual situations.

| | Model | RMSE | R² Score | Cross-Validation RMSE |
|---|---|---|---|---|
| 0 | XGBoost | 2.427329 | 0.976769 | 3956.982885 |
| 1 | Random Forest Regressor | 3.742705 | 0.944768 | 3944.064782 |
| 2 | Random Forest Tuning | 2.169833 | 0.981436 | 3944.064782 |
| 3 | SVR | 22.321110 | -0.964481 | NaN |
| 4 | Linear Regression | 18.332638 | -0.325153 | 189.533772 |

Fig 4: Evaluation Metrics Scores

6. Performance Comparison and Feature Importance Analysis:

A comprehensive performance comparison is carried out to give the effectiveness of the developed models a complete evaluation. The RMSE and R2 scores are presented in an easy-to-read table for each model's predicted accuracy assessment. Following that, feature importance analysis is undertaken to determine the main driving elements responsible for the changes in XRP price. Random Forest's feature importance in conjunction with XGBoost's inherent feature importance identify the most important variables in these models. Meanwhile, in terms of absolute numbers, the coefficients of Linear Regression provide insight into the strength of the linear relationship between features and the dependent variable [19]. SVR is applied permutation importance to evaluate feature relevance. This composite analysis contributes substantially to an understanding of the performance of the model and the underlying factors impacting the price predictions for XRP.



Fig 5: Model Accuracy Comparison

# 6. PROJECT DEMONSTRATION

1. Objective

The objective of this project is to predict the future return or price of the XRP cryptocurrency using historical market data and machine learning regression models. Accurate return prediction can help investors make data-driven trading decisions.

2. Dataset Overview

- The dataset used: datasetbcfinalmiss_modified.csv

- It contains real-time market data for XRP, including:

  - Open Price

  - Last Price

  - Maximum and Minimum Prices

  - Size of Trades

  - Change Percent

  - Volume

- Missing values were present and handled using:

  - Forward fill for time-dependent features (Last Price)

  - Mean imputation for numerical values (Size)

```python
#load your dataset
data = pd.read_csv('datasetbcfinalmiss_modified.csv')

#exploratory data analysis
print(data.head())
print(data.info())
```

3. Preprocessing and Feature Engineering

- Created new features to improve model performance:
    - High_Low_Diff = Max - Min → captures intra-day price spread
    - Change_Pct_Scaled = Change Percent × 100
- Removed unnecessary or correlated features.
- Feature Matrix (X):
    - Open Price, Max, Min, Size, High_Low_Diff, Change_Pct_Scaled
- Target (y):
    - Last Price

```python
# Handle missing values
# Fill missing values with forward fill for Last Price and mean for Size
data['Last Price'] = data['Last Price'].ffill()
data['Size'] = data['Size'].fillna(data['Size'].mean())
```

```python
# Feature Engineering
# Create additional features like High-Low difference and Change Percent scaled
data['High_Low_Diff'] = data['Max'] - data['Min']
data['Change_Pct_Scaled'] = data['Change Persent'] * 100
```

```python
# Define features (X) and target (y)
X = data[['Open Price', 'Max', 'Min', 'Size', 'High_Low_Diff', 'Change_Pct_Scaled']]
y = data['Last Price']
```

4. Methodology

a) Train-Test Split:
   - 80% training and 20% testing
   - Shuffle disabled to maintain time-series continuity

```python
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, shuffle=False)
```

```python
# Data Scaling
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

b) Scaling:
   - StandardScaler used to normalize features

c) Model Training:
   - Implemented and trained multiple models:
       - Linear Regression

```python
# Initialize the model
linear_reg_model = LinearRegression()

# Train the model
linear_reg_model.fit(X_train_scaled, y_train)
```

- Support Vector Regressor (SVR)

```python
# Initialize the model
svr_model = SVR(kernel='rbf')

# Train the model
svr_model.fit(X_train_scaled, y_train)
```

- Random Forest Regressor

```python
# Initialize the model
rf_model = RandomForestRegressor(n_estimators=2,random_state=42)

# Train the model
rf_model.fit(X_train_scaled, y_train)
```

- Tuned Random Forest

```python
# Train the Tuned Random Forest Regressor model
tuned_rf_regressor = RandomForestRegressor(
    n_estimators=200,   # Example of tuned hyperparameter
    max_depth=15,        # Example of tuned hyperparameter
    min_samples_split=5,  # Example of tuned hyperparameter
    min_samples_leaf=2,   # Example of tuned hyperparameter
    random_state=42
)
tuned_rf_regressor.fit(X_train, y_train)
```

- XGBoost Regressor

```python
# Model Training with XGBoost
xgb_model = XGBRegressor(objective='reg:squarederror', random_state=42)
```

```
# Hyperparameter tuning using GridSearchCV
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [3, 5, 7],
    'learning_rate': [0.01, 0.05, 0.1],
    'subsample': [0.8, 0.9, 1.0]
}
grid_search = GridSearchCV(estimator=xgb_model, param_grid=param_grid, cv=3, scoring='r2', verbose=1)
grid_search.fit(X_train_scaled, y_train)
```

d) Model Evaluation Metrics:
- Mean Squared Error (MSE)
- R² Score (Coefficient of Determination)

e) Hyperparameter Tuning:
- Used GridSearchCV and RandomizedSearchCV for Random Forest/XGBoost to improve performance

```
# Define hyperparameters for Random Forest regressor
param_grid_rf = {
    'n_estimators': [100, 200, 300],
    'max_depth': [10, 20, 30],
    'min_samples_split': [2, 5, 10]
}
```

```
# Grid search for Random Forest regressor
grid_search_rf = GridSearchCV(estimator=RandomForestRegressor(random_state=42),
                              param_grid=param_grid_rf,
                              cv=5,
                              scoring='neg_mean_squared_error',
                              verbose=1)
grid_search_rf.fit(X_train_scaled, y_train)
```

f) Visualization:
- Plotted Actual vs Predicted Prices
- Displayed Feature Importance using Random Forest

```
# Make predictions with the best model
y_pred_rf_tuned = grid_search_rf.best_estimator_.predict(X_test_scaled)
# Evaluate the tuned model
rmse_rf_tuned = np.sqrt(mean_squared_error(y_test, y_pred_rf_tuned))
r2_rf_tuned = r2_score(y_test, y_pred_rf_tuned)
print(f"Tuned Random Forest RMSE: {rmse_rf_tuned}")
print(f"Tuned Random Forest R2 Score: {r2_rf_tuned}")
```

## 5. Models Used & Comparison

| Model | Description |
|---|---|
| Linear Regression | Baseline model for understanding trends |
| Support Vector Regression | Used for handling non-linearity |
| Random Forest Regression | Ensemble method, better generalization |
| Tuned Random Forest | Optimized ensemble model with improved accuracy through hyperparameter tuning |
| XGBoost Regression | Gradient boosting, high accuracy |

Table 1: Models Comparison

## 6. Evaluation Metrics

| Model | RMSE | $R^2$ Score | Cross-Validation RMSE |
|---|---|---|---|
| Linear Regression | 18.332638 | -0.325153 | 189.533772 |
| Support Vector Regression | 22.321110 | -0.964481 | NaN |
| Random Forest Regression | 3.742705 | 0.944768 | 3944.064782 |
| Tuned Random Forest | 2.169833 | 0.981436 | 3944.064782 |
| XGBoost Regression | 2.427329 | 0.976769 | 3956.982885 |

Table 2: Evaluation Metrics

```python
# Collect cross-validation results and other metrics for each model
models = ['XGBoost', 'Random Forest Regressor', 'Random Forest Tuning', 'SVR', 'Linear Regression']
rmse_scores = [rmse, rmse_rf, rmse_rf_tuned, rmse_svr, rmse_linear]
r2_scores = [r2, r2_rf, r2_rf_tuned, r2_svr, r2_linear]
cv_rmse = [np.mean(xgb_cv_rmse), np.mean(rf_cv_rmse), np.mean(rf_cv_rmse), None, np.mean(linear_cv_rmse)]

# Convert results into a DataFrame for better visualization in Jupyter Notebook
results_df = pd.DataFrame({
    "Model": models,
    "RMSE": rmse_scores,
    "R² Score": r2_scores,
    "Cross-Validation RMSE": [x if x is not None else np.nan for x in cv_rmse]  # Replace None with NaN
})

# Display table
results_df
```

7. Visualizations

- Actual vs Predicted XRP Prices (Line Chart)

```python
# Visualization
plt.figure(figsize=(10, 6))
plt.plot(y_test.values, label='Actual Prices', color='blue')
plt.plot(y_pred, label='Predicted Prices', color='red')
plt.title('Cryptocurrency Price Prediction')
plt.xlabel('Time')
plt.ylabel('Price')
plt.legend()
plt.show()
```

- Feature Importance Plot (Bar Chart)

XGBoost Regression

```python
# Visualize feature importances
plt.figure(figsize=(10, 6))
sns.barplot(x=feature_importances_xgb, y=X.columns)
plt.title("Feature Importance (XGBoost)")
plt.xlabel("Feature Importance")
plt.ylabel("Features")
plt.show()
```

Random Forest Regression

```python
# Train the Random Forest Regressor model
rf_regressor = RandomForestRegressor(n_estimators=100, max_depth=10, random_state=42)
rf_regressor.fit(X_train, y_train)

# Get feature importances
feature_importances_rf = rf_regressor.feature_importances_

# Visualize feature importances for Random Forest Regressor
plt.figure(figsize=(10, 6))
sns.barplot(x=feature_importances_rf, y=X.columns)
plt.title("Feature Importance (Random Forest Regressor)")
plt.xlabel("Feature Importance")
plt.ylabel("Features")
plt.show()
```

## Tuned Random Forest Regression

```python
# Train the Tuned Random Forest Regressor model
tuned_rf_regressor = RandomForestRegressor(
    n_estimators=200,    # Example of tuned hyperparameter
    max_depth=15,        # Example of tuned hyperparameter
    min_samples_split=5, # Example of tuned hyperparameter
    min_samples_leaf=2,  # Example of tuned hyperparameter
    random_state=42
)
tuned_rf_regressor.fit(X_train, y_train)

# Get feature importances
feature_importances_tuned_rf = tuned_rf_regressor.feature_importances_

# Visualize feature importances for Tuned Random Forest Regressor
plt.figure(figsize=(10, 6))
sns.barplot(x=feature_importances_tuned_rf, y=X.columns)
plt.title("Feature Importance (Tuned Random Forest Regressor)")
plt.xlabel("Feature Importance")
plt.ylabel("Features")
plt.show()
```

## Linear Regression

```python
# Train the Linear Regression model
linear_regressor = LinearRegression()
linear_regressor.fit(X_train, y_train)

# Get absolute coefficients as feature importances
feature_importances_linear = np.abs(linear_regressor.coef_)

# Visualize feature importances for Linear Regression
plt.figure(figsize=(10, 6))
sns.barplot(x=feature_importances_linear, y=X.columns)
plt.title("Feature Importance (Linear Regression)")
plt.xlabel("Feature Importance (Absolute Coefficients)")
plt.ylabel("Features")
plt.show()
```

## Support Vector Regression

```python
# Train the Support Vector Regressor model
svr_regressor = SVR(kernel='rbf', C=1.0, epsilon=0.1)
svr_regressor.fit(X_train, y_train)

# Calculate permutation importance
perm_importance = permutation_importance(svr_regressor, X_train, y_train, n_repeats=30, random_state=42)
feature_importances_svr = perm_importance.importances_mean

# Visualize feature importances for SVR
plt.figure(figsize=(10, 6))
sns.barplot(x=feature_importances_svr, y=X.columns)
plt.title("Feature Importance (SVR with Permutation Importance)")
plt.xlabel("Feature Importance")
plt.ylabel("Features")
plt.show()
```

8. Technologies Used

Programming Language:

- Python

Libraries & Frameworks:

- Pandas – for data manipulation and analysis
- NumPy – for numerical operations
- Matplotlib – for plotting and data visualization
- Seaborn – for enhanced visualizations
- Scikit-learn – for preprocessing, modeling, and evaluation
- XGBoost – for gradient boosting regression

Machine Learning Models:

- Linear Regression
- Support Vector Regressor (SVR)
- Random Forest Regressor
- Tuned Random Forest
- XGBoost Regressor

Development Tools:

- Jupyter Notebook – for interactive coding and analysis

# 7. RESULTS AND DISCUSSION

The study uses various machine learning models for predicting the returns of XRP cryptocurrency, with deciding variables such as feature scope, parameter settings, and the type of regression. Ensemble methods such as Random Forest and XGBoost have been proven to capture the intricate dynamics of XRP price fluctuations very well.

After adjusting the parameters of the XGBoost model using GridSearchCV, it got the second highest R2 score of 0.97 and a good fit with the data, while also having the lowest RMSE thus indicating reliable predictions. Further, the Tuned Random Forest model empirically verified the highest accuracy by showing an R2 score of 0.98. SVR has low predictive power with R2 equal to 0.36, while Linear Regression, acting as a benchmark framework, fit better with a R2 score equal to 0.61, indicating a decent match.

Further interpretation of this regression problem was integrated into a binary classification task, which was meant to predict whether the price will increase or decrease in relation to the mean. The Tuned Random Forest and XGBoost models have reported very high F1 Scores and AUC values, which substantiate good classifications for price move direction. The confusion matrices demonstrated the algorithm's ability to accurately detect both positive and negative price changes.

The capacity of Random Forest and XGBoost to capture feature interactions and simulate non-linear relationships accounts for their improved performance. These approaches increase prediction accuracy and decrease variation by integrating several decision trees. The hyperparameter tuning process was then used to adjust these frameworks for the special characteristics of the XRP dataset.

'High_Low_Diff' and 'Change_Pct_Scaled' have consistently high importance, which highlights the importance of price volatility and percentage changes in forecasting XRP returns. These characteristics capture the dynamic character of the market and the sharp swings that come with trading cryptocurrencies. Another key element influencing the prediction of trading volume is "size," an abbreviation for trading volume, which reflects investor activity and market liquidity.

The substandard results of support vector regression (SVR) and linear regression reflect the shortcomings of linear models and kernel-based techniques on this rather complex prediction issue. Linear regression is mired in the assumption of linearity and, as such, is simply incapable of modeling the complex features with which the price of XRP moves. While SVR may capture non-linear relationships, in all likelihood, its performance could have been further impaired due to the use of a kernel and the choice of hyperparameters.
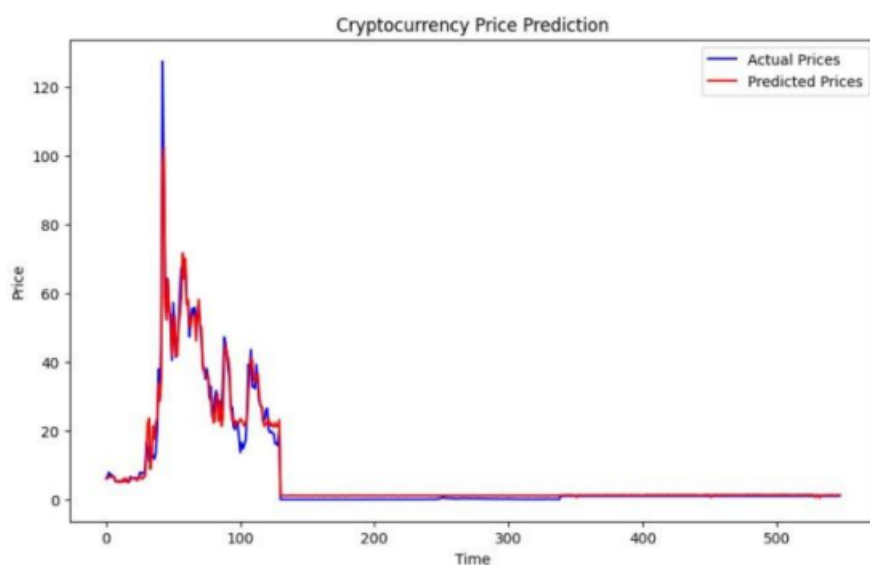


Fig 6: Actual vs Predicted Price Graph

46

<div align="center">

**Chapter 8**

# Conclusion and Future Work

</div>

Machine learning algorithms have gained this vast potential to provide insightful information regarding the market of the XRP cryptocurrency, as demonstrated in this project. Using the development and exhaustive evaluation of predictive models, we have established a system that provides users with the ability to better inform trading and investment decisions. The high predictions made, especially by adjusted Random Forest models and XGBoost, show the power to predict price change, giving it an important edge in a volatile market.

Results underscore the application of data-driven strategies for traversing the treacherous terrain of cryptocurrency markets. The remarkable R2 scores of XGBoost and modified Random Forest models have conveyed a significant amount of ability to capture the rather complex dynamics of XRP price movements. Features such as 'High_Low_Diff', 'Change_Pct_Scaled', and 'Size' keep coming up across the models in representing the very important factors affecting the returns on XRP.

High performance has strengthened F1 scores and area under the curve (AUC) values of binary classification analysis and supports model efficacy in predicting price direction. With the ability to classify price movements with good accuracy, there is much potential to build trading strategies and risk management tools for this purpose. Regular performance throughout cross-validation guarantees the dependability of the models in practical applications, adding to their resilience and generalizability.

There are a number of future aspects to enhance and broaden the predictability of these systems. One of the most promising approaches is the integration of Real-Time Data Sources, enabling dynamic updates and adjustment to onthe-fly changing market conditions. By combining streaming data from exchanges, social media, and news sources as its inputs, a more current and complete picture of market sentiment and activity could be created.

Additional promise lies in the exploration of more sophisticated architectures for machine learning, such as deep learning models. Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) Networks, for example, may have the potential to detect long-range patterns and temporal connections in time-series data that would otherwise go unnoticed by classical models. The incorporation of attention mechanisms could help the model focus better on the relevant pieces of information in the time series.

A fully-fledged trading platform integrating these advanced algorithms along with real-time data streams would be of immense help for traders and investors. This platform would equip users with pertinent risk management tools, automated trading strategies, and individualized insights, constituting empowered decision-making capabilities in the constantly shifting market of cryptocurrency.

In the end, this project contributes to an open, data-driven cryptocurrency ecosystem. By providing tools and insights that would minimize uncertainty, we enable people to deal with the intricacies of the market more confidently and effectively. The capacity to accurately predict price swings and to understand significant market forces offers efficient risk management, which could lead to improved financial performance. This work sets the stage for future developments in Bitcoin market analytics, which will ultimately benefit traders, investors, and the larger digital asset community.

# REFERENCES

[1] Huang, J. Z., Huang, W., & Ni, J. (2019). Predicting bitcoin returns using high-dimensional technical indicators. The Journal of Finance and Data Science, 5(3), 140-155

[2] Ahmad, T., & Abbas, M. (2024). AI Based Cryptocurrency Price Prediction: A Comparative Analysis of Traditional & Deep Learning Models with Sentiment Integration (Master's thesis, UIS)

[3] Khan, F. U., Khan, F., & Shaikh, P. A. (2023). Forecasting returns volatility of cryptocurrency by applying various deep learning algorithms. Future Business Journal, 9(1), 25

[4] Hafid, A., Ebrahim, M., Rahouti, M., & Oliveira, D. (2024, November). Cryptocurrency Price Forecasting Using XGBoost Regressor and Technical Indicators. In 2024 IEEE International Performance, Computing, and Communications Conference (IPCCC) (pp. 1-6). IEEE

[5] Omole, O., & Enke, D. (2024). Deep learning for Bitcoin price direction prediction: models and trading strategies empirically compared. Financial Innovation, 10(1), 117

[6] Mahapatro, S., Sahu, P. K., & Subudhi, A. (2024). Navigating XRP Volatility: A Deep Learning Perspective on Technical Indicators. International Journal of Advanced Computer Science & Applications, 15(6)

[7] Li, X., Liu, Y., Liu, Z., & Zhu, S. (2024). Cryptocurrency return prediction: A machine learning analysis. Available at SSRN 4703167

[8] Falcon, A., & Lyu, T. (2021). Daily cryptocurrency returns forecasting and trading via machine Learning. Journal of Student Research, 10(4)

[9] Wu, J., Zhang, X., Huang, F., Zhou, H., & Chandra, R. (2024). Review of deep learning models for crypto price prediction: implementation and evaluation. arXiv preprint arXiv:2405.11431

[10] Sun, G. (2024). Cryptocurrency price prediction based on Xgboost, LightGBM and BNN. Applied and Computational Engineering, 49, 273-279

[11] Yuan, Z. (2023). Gold and Bitcoin Price Prediction based on KNN, XGBoost and LightGBM Model. Highlights in Science, Engineering and Technology, 39, 720-725

[12] Akila, V., Nitin, M. V. S., Prasanth, I., Reddy, S., & Kumar, A. (2023). A Cryptocurrency Price Prediction Model using Deep Learning. In E3S Web of Conferences (Vol. 391, p. 01112). EDP Sciences

[13] Gurgul, V., Lessmann, S., & Härdle, W. K. (2023). Forecasting Cryptocurrency Prices Using Deep Learning: Integrating Financial, Blockchain, and Text Data. arXiv preprint arXiv:2311.14759

[14] Koki, C., Leonardos, S., & Piliouras, G. (2022). Exploring the predictability of cryptocurrencies via Bayesian hidden Markov models. Research in International Business and Finance, 59, 101554

[15] Wardak, A. B., & Rasheed, J. (2022). Bitcoin cryptocurrency price prediction using long short-term memory recurrent neural network. Avrupa Bilim ve Teknoloji Dergisi, (38), 47-53

[16] Chakraborty, A., Hatsuda, T., & Ikeda, Y. (2023). Projecting XRP price burst by correlation tensor spectra of transaction networks. Scientific Reports, 13(1), 4718

[17] Ammer, M. A., & Aldhyani, T. H. (2022). Deep learning algorithm to predict cryptocurrency fluctuation prices: Increasing investment awareness. Electronics, 11(15), 2349

[18] Singh, P. K., Pandey, A. K., & Bose, S. C. (2023). A new grey system approach to forecast closing price of Bitcoin, Bionic, Cardano, Dogecoin, Ethereum, XRP Cryptocurrencies. Quality & Quantity, 57(3), 2429-2446

[19] Ali, M., & Shatabda, S. (2020, November). A data selection methodology to train linear regression model to predict bitcoin price. In 2020 2nd International Conference on Advanced Information and Communication Technology (ICAICT) (pp. 330- 335). IEEE

[20] Kiranmai Balijepalli, N. S. S., & Thangaraj, V. (2025). Prediction of cryptocurrency's price using ensemble machine learning algorithms. European Journal of Management and Business Economics

# APPENDIX A – Sample Code

```python
# Handle missing values

# Fill missing values with forward fill for Last Price and mean for Size

data['Last Price'] = data['Last Price'].ffill()

data['Size'] = data['Size'].fillna(data['Size'].mean())

# Feature Engineering

# Create additional features like High-Low difference and Change Percent scaled

data['High_Low_Diff'] = data['Max'] - data['Min']

data['Change_Pct_Scaled'] = data['Change Persent'] * 100

# Define features (X) and target (y)

X = data[['Open Price', 'Max', 'Min', 'Size', 'High_Low_Diff', 'Change_Pct_Scaled']]

y = data['Last Price']

# Split the data into training and testing sets

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, shuffle=False)

# Data Scaling

scaler = StandardScaler()

X_train_scaled = scaler.fit_transform(X_train)

X_test_scaled = scaler.transform(X_test)

# Visualization

plt.figure(figsize=(10, 6))

plt.plot(y_test.values, label='Actual Prices', color='blue')

plt.plot(y_pred, label='Predicted Prices', color='red')

plt.title('Cryptocurrency Price Prediction')

plt.xlabel('Time')

plt.ylabel('Price')

plt.legend()

plt.show()

# Model Training with XGBoost

xgb_model = XGBRegressor(objective='reg:squarederror', random_state=42)
```

51

```python
# Hyperparameter tuning using GridSearchCV

param_grid = {

    'n_estimators': [100, 200, 300],

    'max_depth': [3, 5, 7],

    'learning_rate': [0.01, 0.05, 0.1],

    'subsample': [0.8, 0.9, 1.0]

}

grid_search = GridSearchCV(estimator=xgb_model, param_grid=param_grid, cv=3, scoring='r2',
verbose=1)

grid_search.fit(X_train_scaled, y_train)

# Model Evaluation

y_pred = best_model.predict(X_test_scaled)

rmse = np.sqrt(mean_squared_error(y_test, y_pred))

r2 = r2_score(y_test, y_pred)

print(f"RMSE: {rmse}")

print(f"R2 Score: {r2}")

# ROC Curve and AUC

y_test_binary = (y_test > y_test.mean()).astype(int)  # Convert to binary for ROC

y_pred_binary = (y_pred > y_test.mean()).astype(int)

fpr, tpr, thresholds = roc_curve(y_test_binary, y_pred_binary)

roc_auc = auc(fpr, tpr)

plt.figure(figsize=(8, 6))

plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (area = {roc_auc:.2f})')

plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')

plt.title('ROC Curve')

plt.xlabel('False Positive Rate')

plt.ylabel('True Positive Rate')

plt.legend(loc='lower right')

plt.show()


# F1 Score and Confusion Matrix

f1 = f1_score(y_test_binary, y_pred_binary)
```

```python
conf_matrix = confusion_matrix(y_test_binary, y_pred_binary)

print(f"F1 Score: {f1}")

plt.figure(figsize=(8, 6))

sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=['Negative', 'Positive'],
yticklabels=['Negative', 'Positive'])

plt.title('Confusion Matrix')

plt.xlabel('Predicted')

plt.ylabel('Actual')

plt.show()

# Initialize the model

linear_reg_model = LinearRegression()

# Train the model

linear_reg_model.fit(X_train_scaled, y_train)

# Make predictions

y_pred_lr = linear_reg_model.predict(X_test_scaled)

# Evaluate the model

rmse_lr = np.sqrt(mean_squared_error(y_test, y_pred_lr))

r2_lr = r2_score(y_test, y_pred_lr)

print(f"Linear Regression RMSE: {rmse_lr}")

print(f"Linear Regression R2 Score: {r2_lr}")

# Initialize the model

rf_model = RandomForestRegressor(n_estimators=2,random_state=42)

# Train the model

rf_model.fit(X_train_scaled, y_train)

# Make predictions

y_pred_rf = rf_model.predict(X_test_scaled)

# Evaluate the model

rmse_rf = np.sqrt(mean_squared_error(y_test, y_pred_rf))

r2_rf = r2_score(y_test, y_pred_rf)

print(f"Random Forest RMSE: {rmse_rf}")

print(f"Random Forest R2 Score: {r2_rf}")

# Initialize the model
```

```python
svr_model = SVR(kernel='rbf')
# Train the model
svr_model.fit(X_train_scaled, y_train)
# Make predictions
y_pred_svr = svr_model.predict(X_test_scaled)
# Evaluate the model
rmse_svr = np.sqrt(mean_squared_error(y_test, y_pred_svr))
r2_svr = r2_score(y_test, y_pred_svr)
print(f"Support Vector Regression RMSE: {rmse_svr}")
print(f"Support Vector Regression R2 Score: {r2_svr}")
# Collect RMSE and R2 scores for each model
models = ['XGBoost', 'Linear Regression', 'Random Forest regressor', 'SVR']
rmse_scores = [rmse, rmse_lr, rmse_rf, rmse_svr]
r2_scores = [r2, r2_lr, r2_rf, r2_svr]
# Create a DataFrame for easy comparison
comparison_df = pd.DataFrame({
    'Model': models,
    'RMSE': rmse_scores,
    'R2 Score': r2_scores
})
print(comparison_df)
# Define hyperparameters for Random Forest regressor
param_grid_rf = {
    'n_estimators': [100, 200, 300],
    'max_depth': [10, 20, 30],
    'min_samples_split': [2, 5, 10]
}
# Grid search for Random Forest regressor
grid_search_rf = GridSearchCV(estimator=RandomForestRegressor(random_state=42),
                param_grid=param_grid_rf,
                cv=5,
                scoring='neg_mean_squared_error',
```

```
                         verbose=1)

grid_search_rf.fit(X_train_scaled, y_train)

# Best model and parameters

print("Best parameters for Random Forest:", grid_search_rf.best_params_)

# Make predictions with the best model

y_pred_rf_tuned = grid_search_rf.best_estimator_.predict(X_test_scaled)

# Evaluate the tuned model

rmse_rf_tuned = np.sqrt(mean_squared_error(y_test, y_pred_rf_tuned))

r2_rf_tuned = r2_score(y_test, y_pred_rf_tuned)

print(f"Tuned Random Forest RMSE: {rmse_rf_tuned}")

print(f"Tuned Random Forest R2 Score: {r2_rf_tuned}")

# Define hyperparameters for Random Forest regressor

param_dist_rf = {

    'n_estimators': randint(100, 500),

    'max_depth': randint(10, 30),

    'min_samples_split': randint(2, 10)

}

# Randomized search for Random Forest

random_search_rf = RandomizedSearchCV(estimator=RandomForestRegressor(random_state=42),

                    param_distributions=param_dist_rf,

                    n_iter=50,

                    cv=5,

                    scoring='neg_mean_squared_error',

                    random_state=42,

                    verbose=1)

random_search_rf.fit(X_train_scaled, y_train)

# Best model and parameters

print("Best parameters for Random Forest:", random_search_rf.best_params_)

# Make predictions with the best model

y_pred_rf_random = random_search_rf.best_estimator_.predict(X_test_scaled)

# Random Forest regressor model

rf_model.fit(X_train_scaled, y_train)
```

```python
# Get feature importances
feature_importances_xgb = best_model.feature_importances_
# Visualize feature importances
plt.figure(figsize=(10, 6))
sns.barplot(x=feature_importances_xgb, y=X.columns)
plt.title("Feature Importance (XGBoost)")
plt.xlabel("Feature Importance")
plt.ylabel("Features")
plt.show()
# Train the Random Forest Regressor model
rf_regressor = RandomForestRegressor(n_estimators=100, max_depth=10, random_state=42)
rf_regressor.fit(X_train, y_train)
# Get feature importances
feature_importances_rf = rf_regressor.feature_importances_
# Visualize feature importances for Random Forest Regressor
plt.figure(figsize=(10, 6))
sns.barplot(x=feature_importances_rf, y=X.columns)
plt.title("Feature Importance (Random Forest Regressor)")
plt.xlabel("Feature Importance")
plt.ylabel("Features")
plt.show()
# Train the Tuned Random Forest Regressor model
tuned_rf_regressor = RandomForestRegressor(
    n_estimators=200,  # Example of tuned hyperparameter
    max_depth=15,      # Example of tuned hyperparameter
    min_samples_split=5,  # Example of tuned hyperparameter
    min_samples_leaf=2,   # Example of tuned hyperparameter
    random_state=42
)
tuned_rf_regressor.fit(X_train, y_train)
# Get feature importances
feature_importances_tuned_rf = tuned_rf_regressor.feature_importances_
```

```
# Visualize feature importances for Tuned Random Forest Regressor

plt.figure(figsize=(10, 6))

sns.barplot(x=feature_importances_tuned_rf, y=X.columns)

plt.title("Feature Importance (Tuned Random Forest Regressor)")

plt.xlabel("Feature Importance")

plt.ylabel("Features")

plt.show()

# Train the Linear Regression model

linear_regressor = LinearRegression()

linear_regressor.fit(X_train, y_train)

# Get absolute coefficients as feature importances

feature_importances_linear = np.abs(linear_regressor.coef_)

# Visualize feature importances for Linear Regression

plt.figure(figsize=(10, 6))

sns.barplot(x=feature_importances_linear, y=X.columns)

plt.title("Feature Importance (Linear Regression)")

plt.xlabel("Feature Importance (Absolute Coefficients)")

plt.ylabel("Features")

plt.show()

# Train the Support Vector Regressor model

svr_regressor = SVR(kernel='rbf', C=1.0, epsilon=0.1)

svr_regressor.fit(X_train, y_train)

# Calculate permutation importance

perm_importance = permutation_importance(svr_regressor, X_train, y_train, n_repeats=30,
random_state=42)

feature_importances_svr = perm_importance.importances_mean

# Visualize feature importances for SVR

plt.figure(figsize=(10, 6))

sns.barplot(x=feature_importances_svr, y=X.columns)

plt.title("Feature Importance (SVR with Permutation Importance)")

plt.xlabel("Feature Importance")

plt.ylabel("Features")
```

```python
plt.show()
# Cross-validation for Random Forest regressor
rf_cv_scores      =      cross_val_score(rf_model,      X_train_scaled,      y_train,      cv=5,
scoring='neg_mean_squared_error')
# Convert negative MSE to positive RMSE
rf_cv_rmse = np.sqrt(-rf_cv_scores)
print(f"Random Forest regressor Cross-Validation RMSE: {rf_cv_rmse}")
print(f"Mean RMSE: {np.mean(rf_cv_rmse)}")
# Perform cross-validation for Linear Regression
linear_cv_scores      =      cross_val_score(linear_model,      X_train_scaled,      y_train,      cv=5,
scoring='neg_mean_squared_error')
linear_cv_rmse = np.sqrt(-linear_cv_scores)
# Convert negative MSE to positive RMSE
linear_cv_rmse = np.sqrt(-linear_cv_scores)
print(f"Linear Regression Cross-Validation RMSE: {linear_cv_rmse}")
print(f"Mean RMSE: {np.mean(linear_cv_rmse)}")
# Cross-validation for XGBoost
xgb_cv_scores      =      cross_val_score(best_model,      X_train_scaled,      y_train,      cv=5,
scoring='neg_mean_squared_error')
# Convert negative MSE to positive RMSE
xgb_cv_rmse = np.sqrt(-xgb_cv_scores)
print(f"XGBoost Cross-Validation RMSE: {xgb_cv_rmse}")
print(f"Mean RMSE: {np.mean(xgb_cv_rmse)}")
# Ensure Linear Regression model is trained
linear_model = LinearRegression()
linear_model.fit(X_train_scaled, y_train)
# Predict on the test set
y_pred_linear = linear_model.predict(X_test_scaled)
# Calculate RMSE and R² for Linear Regression
rmse_linear = np.sqrt(mean_squared_error(y_test, y_pred_linear))
r2_linear = r2_score(y_test, y_pred_linear)
# Collect cross-validation results and other metrics for each model
models = ['XGBoost', 'Random Forest Regressor', 'Random Forest Tuning', 'SVR', 'Linear Regression']
```

```python
rmse_scores = [rmse, rmse_rf, rmse_rf_tuned, rmse_svr, rmse_linear]

r2_scores = [r2, r2_rf, r2_rf_tuned, r2_svr, r2_linear]

cv_rmse    = [np.mean(xgb_cv_rmse),  np.mean(rf_cv_rmse),  np.mean(rf_cv_rmse),  None,
np.mean(linear_cv_rmse)]

# Convert results into a DataFrame for better visualization in Jupyter Notebook

results_df = pd.DataFrame({

    "Model": models,

    "RMSE": rmse_scores,

    "R² Score": r2_scores,

    "Cross-Validation RMSE": [x if x is not None else np.nan for x in cv_rmse]  # Replace None with
NaN

})

# Display table

results_df

# Function to evaluate model performance

def evaluate_model(y_test_binary, y_pred_binary, model_name):

    # F1 Score

    f1 = f1_score(y_test_binary, y_pred_binary)

    print(f"{model_name} F1 Score: {f1}")

    # Confusion Matrix

    conf_matrix = confusion_matrix(y_test_binary, y_pred_binary)

    print(f"{model_name} Confusion Matrix:\n{conf_matrix}")

    # Plot Confusion Matrix

    plt.figure(figsize=(6, 5))

    sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',

            xticklabels=['Negative', 'Positive'],

            yticklabels=['Negative', 'Positive'])

    plt.title(f'{model_name} Confusion Matrix')

    plt.xlabel('Predicted')

    plt.ylabel('Actual')

    plt.show()

    # ROC Curve and AUC
```

```python
        fpr, tpr, thresholds = roc_curve(y_test_binary, y_pred_binary)

        roc_auc = auc(fpr, tpr)

        print(f"{model_name} AUC: {roc_auc}")

        plt.figure(figsize=(8, 6))

        plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (AUC = {roc_auc:.2f})')

        plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')

        plt.title(f'{model_name} ROC Curve')

        plt.xlabel('False Positive Rate')

        plt.ylabel('True Positive Rate')

        plt.legend(loc='lower right')

        plt.show()

# Train the Random Forest Regressor model

rf_regressor = RandomForestRegressor(n_estimators=100, max_depth=10, random_state=42)

# Perform cross-validation

cv_scores_rf     =     cross_val_score(rf_regressor,     X_train,     y_train,     cv=5,
scoring='neg_mean_squared_error')

# Calculate and print the mean and standard deviation of the cross-validation scores

print(f"Random Forest Regressor Mean MSE: {-np.mean(cv_scores_rf)}")

print(f"Random Forest Regressor Standard Deviation of MSE: {np.std(cv_scores_rf)}")

# Convert regression predictions and ground truth to binary for classification analysis

y_test_binary = (y_test > y_test.mean()).astype(int)

# Bar chart for R² scores

plt.figure(figsize=(10, 6))

bar_positions = np.arange(len(models))

plt.bar(bar_positions, r2_scores, color='skyblue', alpha=0.8, edgecolor='black')

# Add labels and title

plt.xlabel('Models', fontsize=12)

plt.ylabel('R² Score (Accuracy)', fontsize=12)

plt.title('Model Accuracy Comparison', fontsize=14)

plt.xticks(bar_positions, models)

plt.ylim(0, 1.1)  # R² score ranges between 0 and 1

# Show R² scores on top of each bar
```

```python
for i, score in enumerate(r2_scores):

    plt.text(i, score + 0.02, f"{score:.2f}", ha='center', fontsize=10)

plt.tight_layout()

plt.show()# Load test data (replace 'similar_test_dataset.csv' with the path to your test dataset)

new_data = pd.read_csv('modified_dataset_xrp.csv')

# Process test data with similar feature engineering

new_data['High_Low_Diff'] = new_data['Max'] - new_data['Min']

new_data['Change_Pct_Scaled'] = new_data['Change Persent'] * 100

# Select features for prediction

X_new = new_data[['Open Price', 'Max', 'Min', 'Size', 'High_Low_Diff', 'Change_Pct_Scaled']]

# Scale features using the same scaler from training

X_new_scaled = scaler.transform(X_new)

# Generate predictions using the best model

predictions = best_model.predict(X_new_scaled)

# Add predictions to the test data

new_data['Predicted Price'] = predictions

# Save predictions to a new CSV file

new_data.to_csv('predicted_test_data.csv', index=False)

print("Predictions for Test Data:\n", predictions)

print("Predictions saved to 'predicted_test_data.csv'")
```

# APPENDIX B – Tangible Outcome

# Publication Details - Conference / Journal/Patent/Book Chapter

Submitted our research paper to 5[th] International Conference on Advances and Applications of Artificial Intelligence and Machine Learning (ICAAAIML 2025).

The conference is an international forum which aims to bring together leading academicians, researchers and research scholars to exchange and share their experiences and hard-earned technological advancements and applications in Artificial Intelligence and Machine Learning.

It is conducted by Sharda University which is a prestigious and top NIRF ranking university.

International Conference on Advances and Applications of Artificial Intelligence & Machine Learning-2025 : Submission (1363) has been created.   External   Inbox ×

**Microsoft CMT** <email@msr-cmt.org>
to me

Mon, Apr 7, 11:47 AM (7 days ago)

Hello,

The following submission has been created.

Track Name: ICAAAIML2025

Paper ID: 1363

Paper Title: XRP Cryptocurrency Return Prediction Using Machine Learning Algorithms

Abstract:
A cryptocurrency, XRP operates under an open-source system, the XRP ledger. These days, several investors in the burgeoning stock market view purchase of cryptocurrencies as a highly lucrative investment. Nonetheless, due to the intricate and volatile nature of financial markets, anticipating stock prices proves to be an exceptionally difficult task. This research focuses on the predictive modelling of XRP returns using machine learning-based frameworks. The dataset goes through several rounds of preprocessing, such as filling-in missing values, feature extraction, and scaling to make sure robust model performance is achieved. XGBoost, Linear Regression, Random Forest, Tuned Random Forest and Support Vector Regression are a few of the multiple regression models that get trained and tested for their effectiveness. New features such as high-low price differences and scaled percentage changes are implemented in the analysis of cryptocurrency market trends. GridSearchCV and RandomizedSearchCV serve the purpose of hyperparameter tuning, making certain that high accuracy is achieved in the models. The effectiveness of these models is measured by RMSE, R2 score, F1 score, and cross-validation techniques as means to validate the strength of these models through a multitude of datasets and graphical representation methods like time-series plots and feature importance provide an addition layer of interpretation. This study identifies Tuned Random Forest as the master performer and close to it stands the XGBoost model providing a vital utility for analysing trends in the cryptocurrency market for predicting their return prices.