

Awarding Body:

Programme Name:

MSc Data Science

Module Name (and Part if applicable):

Advanced Project Computing

Assessment Title:

Detecting cryptocurrency pump and dump scheme – Using
Blockchain data to identify market manipulation

Student Number: STU217422

Tutor Name: Uma Tumpala

Word Count:

10826

Please refer to the Word Count Policy on your Module Page for guidance

**DETECTING CRYPTOCURRENCY PUMP-AND-DUMP
SCHEMES – USING BLOCKCHAIN DATA TO IDENTIFY
MARKET MANIPULATION**

Abstract



This research describes an extensive machine learning based framework for the detection of pump-and-dump schemes in the cryptocurrency market with blockchain-originated data. The project is focused on one of the growing concerns in decentralised financial environments: market manipulation, where sensibly the regulations don't follow. The fraud detection system has a hybrid detection pipeline incorporating both unsupervised and supervised learning techniques to improve reliability and scalability.

To identify the unusual patterns in features like price volatility, percentage price change, and trading volume, Isolation Forest and K-Means algorithms are used by the anomaly detection component. After classifying these anomalies using Random Forest and XGBoost, two powerful supervised models that are notably robust and perform extremely well on imbalanced datasets, these anomalies could then be sorted. And the evaluation concluded with Isolation Forest having perfect metrics for classification, with Random Forest and XGBoost having accuracy greater than 99%. While having a good baseline, K-Means was not powerful enough to find fraudulent instances.

Challenges such as computational demands, most notably model tuning and evaluation, and real-time implementation complexity are also highlighted in the project. However, it lays down a good groundwork for the building of large, intelligent fraud detection systems that can protect investor trust in volatile markets. The interpretability of the models is supported with visual outputs such as ROC curves, confusion matrices, and time series anomaly plots.

To add to that, this work demonstrates the authenticity of utilizing AI based approaches meant to detect financial exploitation in blockchains and sets the basis for real time applications and entrepreneurial solutions to be introduced in regulatory and trading infrastructure.

Table of Contents

Abstract	3
1. Project Framing	7
1.1 Background	7
1.2 Problem Statement	8
1.3 Research Aim and Objectives.....	9
1.4 Rationale	10
2. Fact Finding	12
2.1 Technical Overview of Blockchain and Crypto Markets	12
2.2 Types of Market Manipulation.....	14
2.3 Challenges of Fraud Detection in Blockchain	15
2.4 Requirements Analysis	17
2.5 Literature Review	18
2.6 Global context.....	20
2.7 Gap Analysis.....	20
3. Project Development	22
3.1 Scope of the Artefact	22
3.2 Data Overview	23
3.3 System Architecture.....	25
3.4 Model Selection and Justification.....	27
3.5 Model Training and Tuning	28
3.6 Toolkits and Technologies Used.....	33
3.7 Challenges Faced During Development	33
4. Critical Evaluation.....	35
4.1 Evaluation of Artefact.....	35
4.2 Comparison of Models.....	38
4.3 Methodology Assessment.....	39



4.4 Ethical Considerations	40
5. Conclusion	42
5.1 Summary of Key Findings.....	42
5.2 Recommendation.....	43
References.....	44
Appendix	50

Table of Figures

Figure 1: Cryptocurrency Market Capitalisation	7
Figure 2: Blockchain Architecture.....	12
Figure 3: Bitcoin Pump and Dump	13
Figure 4: Wash Trading.....	14
Figure 5: Spoofing	15
Figure 6: Data Acquisition	23
Figure 7: Feature Engineering.....	24
Figure 8: System Overview Diagram.....	25
Figure 9: Isolation Forest.....	28
Figure 10: K-Means Clustering.....	29
Figure 11: Anomaly Detection	30
Figure 12: Random Forest	31
Figure 13: XGBoosts	32
Figure 14: Libraries Used	33
Figure 15: Confusion Matrices	35
Figure 16: Classification Reports	36
Figure 18: Accuracy Scores	37
Figure 19: Model Comparison	38

1. Project Framing

1.1 Background

The meteoric rise of the cryptocurrency market over the past decade has brought with it a newfound wealth of opportunity, but also created a wild and volatile space open to the machinations of wrongdoing. In 2025, the global cryptocurrency market capitalisation reached over \$2.75 trillion, with Bitcoin contributing over \$1.68 trillion, meaning that Bitcoin is widely considered a prominent decentralised finance tool (Coingecko.com, 2025).

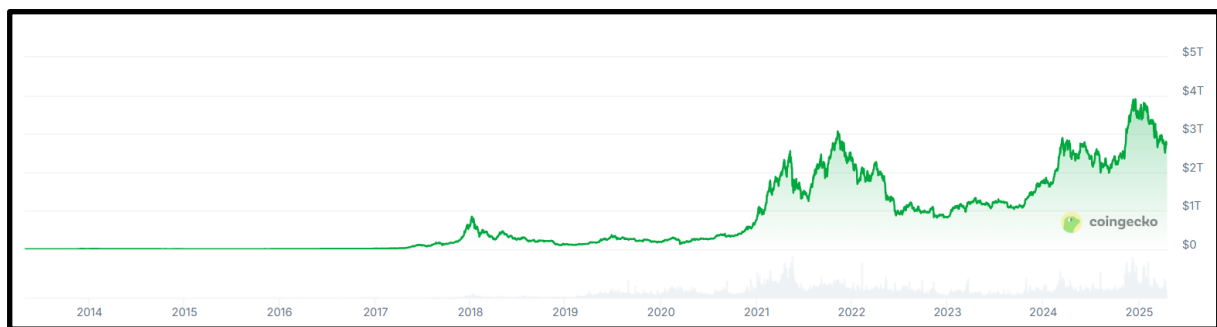


Figure 1: Cryptocurrency Market Capitalisation

(Source: Coingecko.com, 2025)

Yet, this explosive market participation has been accompanied by the prevalence of market manipulation strategies, including pump and dump (P&D) schemes. These schemes artificially inflate the price of low-cap cryptocurrencies through coordinated trading or misinformation, which is subsequently dispersed for retail investors to unsuspectingly purchase, facilitating quick sales, huge losses, and market instability (Li *et al.* 2021). Over \$7 million was generated in a single month by pump and dump activity alone on crypto exchanges in all of 2018 (Cointelegraph.com, 2025). Blockchain systems inherently provide anonymity and decentralisation, making it very hard to perform real-time oversight and enforcement by regulatory authorities, unlike traditional stock exchanges subject to the U.S. Securities and Exchange Commission (SEC).

Although immutability and transparency of the blockchain are great characteristics, they are used less for their analytical potential to detect manipulative behaviours. In recent years, existing academic literature has started exploring machine learning (ML) approaches to the task of detecting financial fraud in centralised banking systems or

equity markets (Ali *et al.* 2022). Yet, the central obstacle to blockchain-native data being utilised in real-time trading environments is the ability to leverage it for anomaly detection. Dynamic trading conditions, along with the absence of labelled data, make the traditional detection mechanisms outdated. As a result, it is possible to integrate the machine learning based anomaly detection method, such as models in Isolation Forest and XGBoost, with the ability to scale and adapt to this manipulation problem. Such data-driven surveillance tools will help protect investor interests as well as the credibility of decentralised financial markets, and this research is important to address a growing need for such tools.

1.2 Problem Statement

Crypto emerged as a mainstream financial asset that has dramatically increased the benefits of decentralisation, global accessibility, and reduced friction of transactions. However, they have helped to enable a new type of financial crime that is difficult to identify and regulate using traditional means. One of these is the pump and dump (P&D) scheme, perhaps the most common and most destructive of them all, in which a group of people take control over the trading price of a digital asset through the use of misleading promotion or trading patterns and then get out of the way when the price will be at an inflated level for the least experienced investors to incur the highest losses. The difficulty lies in the fact that while such schemes are actively pursued by the U.S. Securities and Exchange Commission (SEC) in traditional securities markets (Goforth, 2021). The cryptocurrency markets lack such regulatory bodies, and search and prosecution are far more difficult (Feinstein and Werbach, 2021). In reality, it is estimated that such manipulations help generate as much as \$7 million per month in profits from trading in unregulated crypto exchanges, which suggests the extent and gravity of the issue (Cointelegraph.com, 2025).

Making matters worse, blockchain transactions are pseudonymous, and the number of low-cap, illiquid altcoins that are especially susceptible to manipulation keeps growing bigger and bigger. According to Hu *et al.* (2023), some amount of coins listed via initial coin offerings underwent pump-and-dump strategies in the first 90 days of listing. While Blockchain data is publicly and immutably defined, the pattern of such manipulation is largely reactive and reliant on post-factum forensic review. Anomaly detection methods based on real-time machine learning need to be applied right onto

blockchain-derived market data to identify fraudulent price manipulation once it takes place, however, there is a critical research gap (Gad *et al.* 2022).

Traditional finance has relied on supervised learning models that require labelled datasets to detect fraud, which is often lacking or insufficient for pump and dump schemes, as such labels are not available or sufficient because of the novelty and opacity of the events. In addition, manipulators' behaviour changes very quickly upon the change of detection efforts, thus it requires adaptive and unsupervised techniques that can generalise beyond known cases. As a result, this research fills a gap in the need for a detectable, scalable, and accurate framework built on blockchain data's built-in transparency. Using models like Isolation Forest, K-Means, Random Forest, and XGBoost, the study wants to build a hybrid detection system for detecting anomalous trading patterns, which may indicate pump and dump schemes. This can be one solution not only to preserve the rights of individual investors but also to contribute to market stability and further help to create regulatory strategies in the dynamic cryptocurrency environment.

1.3 Research Aim and Objectives

Aim

This research aims to design a robust machine learning based framework based on blockchain-derived trading data for detecting the pump and dump schemes in cryptocurrency markets and for market transparency, early detection of anomalies by providing investor protection and regulatory readiness in decentralised financial ecosystems.

Objectives

- To extract and engineer useful features from blockchain-based trading data for the specific purpose of fraud detection.
- To apply and evaluate unsupervised models for discovering anomalous trading patterns in cryptocurrency markets.
- To implement and compare supervised learning algorithms for classifying pump-and-dump market manipulation events.
- To propose a scalable detection framework that supports even real-time forensic analysis as well as intervention with regulators.

Limitations of Study

Multiple research constraints affect this study's findings through its focused design. The study faces a major barrier because blockchain data remains scarce and difficult to work with. Large-scale historical blockchain data retrieval is restricted by API limits and infrastructure limitations even though blockchain data remains publicly accessible. The task of identifying verified pump-and-dump schemes proves difficult because these schemes frequently stay hidden or unreported or get misclassified which restricts the creation of a reliable labelled dataset for supervised learning. The natural randomness present in blockchain transactions makes it difficult to find meaningful patterns because thousands of daily transactions frequently do not involve manipulation activities. The detection model created in this research operates with historical data but fails to recognize new manipulation approaches that use dissimilar methods of operation. The research focuses exclusively on post-event evaluation while omitting real-time prediction capabilities and deployment aspects. The study faces ethical and legal restrictions which prevent real-time market predictions that would disrupt active trades or ongoing transactions.

1.4 Rationale

This research has a rationale based on the increasing interest in the number of market manipulations within the cryptocurrency ecosystem, especially through pump and dump schemes. These are schemes taking advantage of the fact that digital asset markets are unregulated and very much decentralised, by artificially inflating the price of a cryptocurrency by generating coordinated hype and inappropriate activities, and then quickly selling it off once unsuspecting investors enter the market. Damages are not only detrimental to the interest of individual investors but also misjudge all that decentralised finance is and likely will become (Aquilina *et al.* 2024). Unlike the traditional financial system, where there is centralised oversight through monitoring and enforcement of trading practices, the market of the cryptocurrency space is extremely dependent on technological interventions from the market to maintain its integrity (Yadav, 2022).

However, human analysis of the data produced by the blockchain is hindered due to its complexity and volume, accompanied by its celebration of transparency and

immutability (Sedlmeir *et al.* 2022). While every transaction, wallet interaction, and price fluctuation is recorded publicly, to draw meaning from the massive and unstructured data, empirical intelligence systems must be used. Existing literature has shown how machine learning can be employed for fraud detection in traditional financial systems, but the encryption inside a blockchain (Ashfaq *et al.* 2022). This makes event data available to everyone, has, to date, garnered less interest in how to use it to detect manipulation in real time. This is a big gap in the research domain, with the cryptocurrency market progressing further and further with more and more people participating in it.

Automated detection of suspicious trading behaviour is a good use of machine learning models, specifically for anomaly detection and classification. Unsupervised learning techniques like Isolation Forest and K-Means are used when the number of experiment-confirmed manipulation cases is very few or even impossible to confirm. Even further value is gained by supervised models like Random Forest and XGBoost, which can predict suspicious patterns on top. Together, these models may be applied to the development of a hybrid detection framework that has high reliability in detecting pump and dump schemes.

The reason for this research is thus a dual need for technological innovation and investor protection. A well-developed, robust, scalable, and interpretable machine learning system could not only help to mitigate losses of individuals, but can also lead to a more trustworthy and resilient financial system. By doing so, it stays in line with the global requirements for ethical, transparent, and secure financial innovation and enables the long-term viability of the blockchain-based markets. The research could therefore empower financial analysts, regulatory bodies, and digital asset platforms dealing with counting manipulation and promoting confidence in their cryptocurrency marketplaces.

2. Fact Finding

2.1 Technical Overview of Blockchain and Crypto Markets

As a disruptive technology in the digital finance world, blockchain is a revolutionary technology regarding the decryption of transactions on a shared ledger. Quite basically, a blockchain is a sequence of blocks with a set of transactions to be cryptographically linked and agreed upon through consensus mechanisms like Proof of Work and Proof of Stake (Sriman *et al.* 2021). Such a structure guarantees immutability, resistance to tampering, and decentralisation, features which make blockchain different from traditional centralised databases (Islam and Apu, 2024).

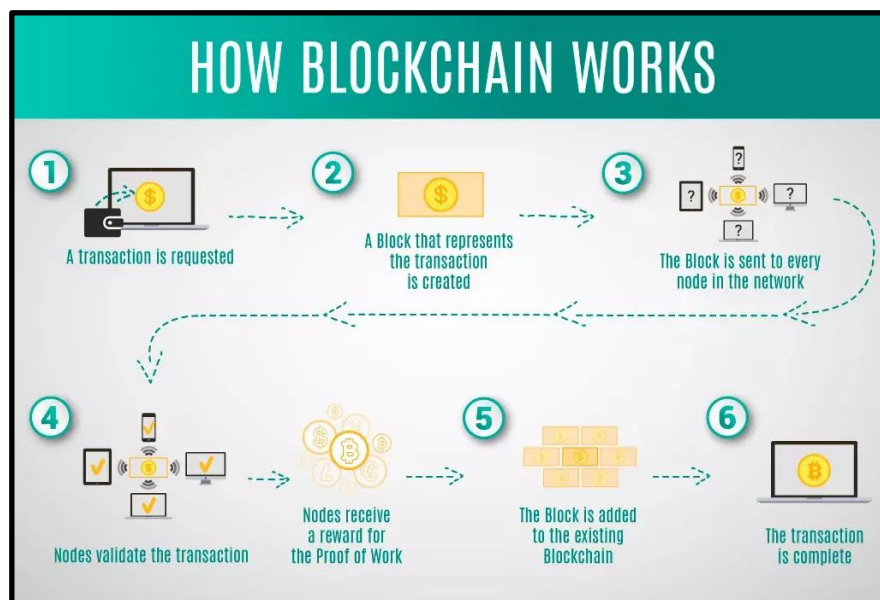


Figure 2: Blockchain Architecture

(Source: Mlsdev.com, 2025)

The introduction of Bitcoin, the first and most well-known cryptocurrency, exemplified the possible use of blockchain to support peer-to-peer transactions without any intermediary, like banks or clearing houses (Panda *et al.* 2023). Since then, the ecosystem has grown to thousands of cryptocurrency, smart contract platforms, such as Ethereum, and even decentralised applications (DApps) that utilise blockchain's utility to be anything other than a transfer of value (Leiponen *et al.* 2022).

An aspect of it is the cryptocurrency market, which touches on blockchain infrastructure; one that is global, constant, with high levels of liquidity, volatility, as well

as speculation. While cryptocurrency is traded on a 24/7 basis on several centralised and decentralised exchanges, cryptocurrency trading is largely unregulated (Aspris *et al.* 2021). This becomes an open-access environment in which assets may be created, traded, and liquidated at ease. Nevertheless, this openness also brings risks such as market manipulation, no investor protection, and security vulnerabilities. Due to their speculative nature, cryptocurrencies are more aesthetically appealing and more likely to be manipulated through coordinated efforts such as pump and dump schemes, and there is little registration and regulation for such schemes (Corbet *et al.* 2021).



Figure 3: Bitcoin Pump and Dump

(Source: Blogs.biomedcentral.com, 2025)

Blockchain data offers a technical opportunity to analyse and detect fraud from a public, timestamped, and verifiable perspective. All transactions are transparent and linked to cryptographic wallet addresses, making the network completely traceable. However, due to the scale and complexity of this data, it is necessary to employ means of detection of anomalies and interpretation of market behaviours using means of advanced computation. By helping bridge blockchain transparency with effective oversight through real-time Bayesian prediction and anomaly detection, machine learning is a practical solution for a machine learning solution. The development of intelligent systems to detect fraudulent behaviour is possible based on this technical landscape and can render the crypto-financial ecosystem resilient.

2.2 Types of Market Manipulation

An emerging market manipulation issue in cryptocurrency trading is caused by the decentralised nature, and most of the trading is based on the decentralised and largely unregulated digital asset exchange market. Lack of central oversight, along with antiseptic surveillance mechanisms present in developments, enables manipulative actors to take advantage of distressed markets' vulnerabilities, if not at a loss to retail investors who simply do not understand what's going on. Pump and dump scheme is probably the most common type of market manipulation that exists in the ecosystem of cryptocurrency (Eigelshoven *et al.* 2021). It is followed by wash trading, spoofing, front running, and layering. Each one of these practices distorts market signals, artificially and price discovery, and affects investor confidence.

P&D schemes in low-cap cryptocurrencies are especially common. This method involves artificially increasing the price of a currency via coordinated buying or false promotional campaigns on Telegram, Discord, and then dumping the holdings once the price tops before the price dials down quickly. Yet, in this practice, detecting it in real time is difficult because of the anonymity of participants and the speed of transactions on crypto exchanges (Shah *et al.* 2021).

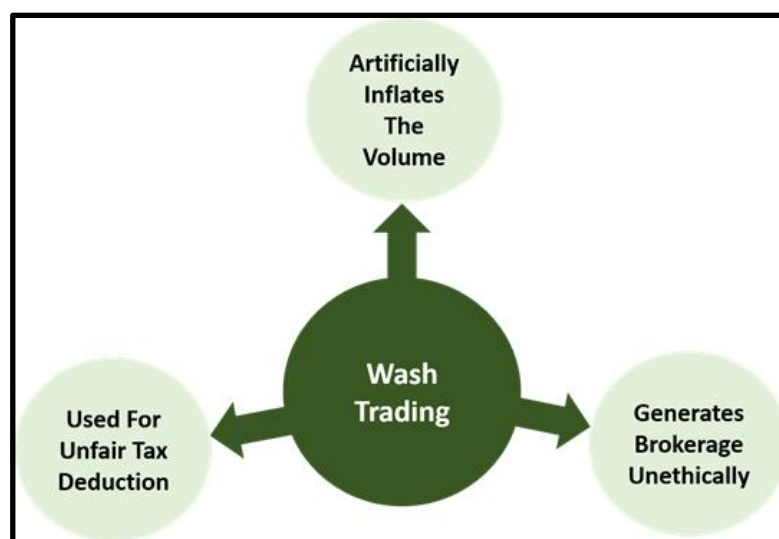


Figure 4: Wash Trading

(Source: Self-Created)

Another deceptive strategy is wash trading, where an individual or group does the same buying and selling of the same asset to make an artificial volume indicator

(Jayant, 2023). It is typically deployed to pretend that there is market interest and to obtain the cooperation of other traders based on pretences.

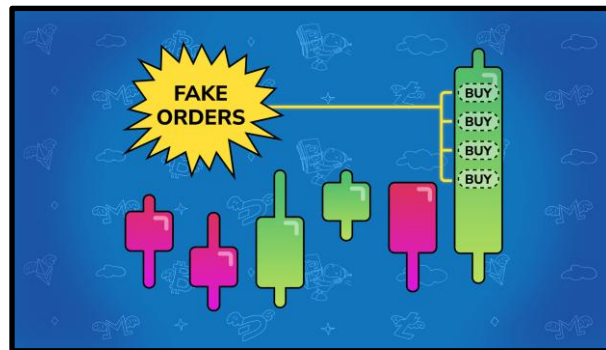


Figure 5: Spoofing

(Source: Tradesanta.com, 2025)

Spoofing and layering is where a large buy or sell order is placed only to cancel it before it is executed. The fake orders generate a veneer of supply or demand pressure to influence the market participants' decisions. Front running, which is still very common in conventional finance, is further expanded upon in decentralised exchanges in which miners or validators can re-order the transactions to take advantage of soon arriving larger transactions (Sariboz *et al.* 2022). At this point of evolving sophistication in manipulation strategies for the blockchain, this practice is known as miner extractable value (MEV).

Together, these market manipulation forms fall short of market efficiency and can lead to extreme situations, as is the case with any illiquid tokens. Centralised exchanges have started to incorporate some level of anti-manipulation protocols, but decentralised environments have been especially uncovered. As these practices have grown, however, addressing them through smart, data'spelled data driven detection mechanisms like machine learning on blockchain analytics has become more and more necessary for a transparent and protected environment in the digital asset space.

2.3 Challenges of Fraud Detection in Blockchain

The fundamental nature of blockchain technology itself creates a complex set of challenges that characterise the issue of fraud detection in blockchain environments. Blockchain's transparency, which is often extolled, also elevates barriers to monitoring, intervention, and regulation, among other things (Sedlmeir *et al.* 2022).

The main challenge is the undeniably 'pseudonymous' nature of blockchain transactions. Every transaction is public and traceable, although user identities are hidden behind cryptographic wallet addresses that make it hard to point out any particular action to the real-world person or entity (Liu *et al.* 2021). The anonymity of this type of exchange complicates intervention in the enforcement of the anti-money laundering (AML), particularly the know-your-customer (KYC) regulations, which are a standard in the traditional financial systems.

A central lack of oversight is another major blockade. Unlike the conventional financial market, there is no one single governing body that handles the regulatory process for the blockchain (Chowdhury *et al.* 2023). Making decentralised routing also offers high freedom and innovation, but also allows anyone who can exploit market inefficiencies to forego prosecution for a time. It is especially observable in the occurrence of pump and dump schemes, wash trading, and spoofing, usually conducted by utilisation of unregulated platforms such as Telegram or Discord.

Furthermore, the sheer volume and velocity of blockchain-generated data pose additional challenges. Real-time analysis on all of these transactions on several chains is computationally intractable, as millions of transactions are being processed every day. Traditional rule-based detection systems are struggling to scale and fail to adapt to the 'as a service' nature of attacking the organisation's IT systems. Lack of labelled datasets also limits the development of supervised learning models, so that methods of machine learning, such as unsupervised or semi-supervised, need to be calibrated properly to avoid false positives or negatives.

Finally, although the overall financial instruments based on blockchain are more complex, realising emerging complexity of the financial instruments based on blockchain such as smart contracts, the decentralized finance (DeFi) platform and the non-fungible tokens (NFTs) have presented new vectors of the fraud that are not easily tackled within the existing models. Technical understanding is not enough to detect and detect manipulation in these domains; one needs an excellent understanding of economic behaviours, code vulnerabilities, and governance structures. The challenges revealed above are our reasons that blockchain ecosystems require sophisticated, well-adapted, and context-aware frameworks to prevent and detect fraud.

2.4 Requirements Analysis



Requirement Type	Requirement Description
Functional Requirements	
F1	The system shall ingest and process blockchain-based trading data in hourly intervals from historical CSV files.
F2	The system shall extract relevant features such as volatility, price change, and rolling statistics from raw data.
F3	The system shall apply unsupervised machine learning algorithms (Isolation Forest, K-Means) for anomaly detection.
F4	The system shall implement supervised learning models (Random Forest, XGBoost) for the classification of manipulation.
F5	The system shall display anomalies and predictions using time-series graphs and confusion matrix heatmaps.
F6	The system shall compare and evaluate models using metrics such as accuracy, ROC AUC, and F1-score.
F7	The system shall store processed data, predictions, and evaluation results for audit and reporting purposes.
F8	The system shall support a modular code structure for integration with future live blockchain or API data sources.

Non-Functional Requirements	
NF1	The system should maintain a model prediction accuracy of at least 95% during testing and validation.
NF2	The system should achieve an ROC AUC score above 0.90 to ensure reliable anomaly classification.
NF3	The system shall be implemented using scalable Python libraries (e.g., Scikit-learn, XGBoost) suitable for large datasets.
NF4	The system should complete model training and evaluation within a reasonable computational time on a standard workstation.
NF5	The codebase should follow modular programming practices for readability, reusability, and future scalability.
NF6	Visual outputs (charts, heatmaps) should be rendered in under 2 seconds for responsive analysis.
NF7	The system should operate offline without requiring live exchange data, supporting reliability in academic contexts.
NF8	The system must adhere to ethical AI principles, ensuring transparency and avoidance of false accusations.

Table 1: Requirement Analysis

(Source: Self-Created)

2.5 Literature Review

On financial fraud detection, machine learning has emerged as a powerful aid in recent years, since traditional rule-based systems are ineffective at detecting significant,

advanced fraud patterns. Machine Learning in the context of cryptocurrency markets, with their price volatility, hiding behind trading anonymity as well as regulatory loopholes, enables a lot of manipulative behaviours, and is a useful tool to enable automated anomaly detection and risk analysis. Several academic studies have looked at ML's suitability for fraud detection in the domain of financial sectors, examining how ML is quick to deploy, real-time processing capable, and outperforms static models (Bello *et al.* 2024). Dube and Verster (2023) point out that tree-based models such as Random Forest and gradient boosting techniques surpass the linear methods in pattern recognition of complex fraud patterns, especially when faced with high-dimensional and imbalanced datasets.

Isolation Forest (iForest), a popular method for identifying outliers and abnormal trading behaviour in unsupervised learning, is used. Due to its ability to randomly partition the feature space, Isolation Forest proved its effectiveness in isolating a small fraction of anomalies, leading to low-frequency but high-impact fraud events such as pump and dump schemes. This is because the algorithm has excellent efficiency, and it can handle large-scale data, precisely what is required for some of the blockchain environments where millions of transactions per second must be processed without any delays. Just as in the case of supervised fraud classification tasks, Random Forest and XGBoost have also demonstrated strong predictive performance. As is often the case when dealing with crypto transaction datasets, XGBoost is a highly scalable tree boosting model of sparse data.

Similarly, ML-based surveillance systems are also needed for real-world crypto scam case studies. Bitconnect involved a Ponzi-style lending scheme, covered in the media, and caused severe losses to investors. For example, Hikaru Kasugai's Squid Token scam is where developers used social media hype and false claims of developing an app to inflate token value and then pulled a 'rug', followed by worthless assets for the investors (Dupuis and Gleason, 2021). These are quite indicative that manual fraud monitoring is not adequate and that there is a need for automated, real-time detection systems. An approach to identifying abnormal trading behaviours through the employment of ML models like Isolation Forest and XGBoost within blockchain analytics platforms would be integrative and help build transparency as well as the security of the cryptocurrency ecosystem.

2.6 Global context

This research produces vital worldwide effects on cryptocurrency markets as well as regulatory programs that fight financial fraud and market manipulation. The research develops blockchain-based detection methods for cryptocurrency pump-and-dump schemes to improve worldwide market transparency and investor protection. Cryptocurrency fraud activities that start in one region create worldwide market consequences through international financial losses which affect both retail and institutional investors (Gandal et al., 2018). A strong detection system would help worldwide regulatory organizations including the U.S. Securities and Exchange Commission (SEC) and European Securities and Markets Authority (ESMA) and Asian financial watchdogs to identify and stop such schemes better (Foley et al., 2019). The framework should be integrated into trading platforms and exchanges by cryptocurrency exchanges to improve their real-time fraud detection capabilities and create a safer trading space. Research demonstrates that blockchain analytics serves as a vital tool for detecting illegal activities while enhancing regulatory compliance (Chen et al., 2022). The findings from this study would help create standardized anti-manipulation policies across jurisdictions which would establish a fair and consistent global cryptocurrency market structure. The analysis of blockchain data provides valuable insights that help expand financial crime prevention strategies to fight money laundering and other illicit financial activities (Fan et al., 2023). The research uses data-driven methods to detect fraudulent market behaviours which ultimately strengthens the long-term stability and credibility of the global cryptocurrency ecosystem.

2.7 Gap Analysis

Although prior work has achieved great advancements in identifying activities of a fraudulent nature in financial markets, it is notable that there is a gap between the applicability of Machine Learning to blockchain data, specifically P&D scheme. Most existing models are not scalable and do not have real time capabilities, and few use all the transparency and immutability of blockchain transactions. Based on these limitations, this study addresses them through a data driven approach that leverages blockchain analytics and advanced machine learning.



3. Project Development

3.1 Scope of the Artefact

A particular definition of the artefact developed in this research involves designing and constructing an intelligent fraud detection system that assists offline analysis of historical BTC/USD trading data to discern market patterns indicative of pump-and-dump market manipulation. A system has been designed to simulate real-time fraud flagging using processed trading data from the blockchain hashing function at an hourly resolution and extracting key indicators, including price volatility, percentage price change, trading volume, and logarithmic returns. These are next analysed using both unsupervised and supervised machine learning models, Isolation Forest, KMeans, Random Forest, as well as XGBoost, to detect anomalous trading behaviours robustly.

The artefact does not run as a live system communicating directly with the Exchange API, but rather batch processes historical data to simulate real-time fraud detection. The choice of this design makes rigorous evaluation of detection accuracy, model interpretability, and computational performance possible, independent of the instability of external data sources. In the Python development environment, with the help of open source libraries like Pandas, Scikit-learn, and XGBoost, the system is developed from scratch to facilitate reliability and make it suitable for integrating into other bigger blockchain analytics frameworks. In addition, the artefact includes visual analytics components that embed the detected anomaly through time series plots and confusion matrix heatmaps to give technical and non-technical stakeholders an easy means of interpreting the anomaly.

The scope of this artefact from a research perspective is to show that ML models can learn to detect such nuance, nonlinear patterns that correlate with fraudulent activity even when there is no explicit label available, which is an important feature of pump and dump labels, especially in cryptocurrency markets, as they are notoriously hard to obtain (Rajaei and Mahmoud, 2023). The studies of Xu *et al.* (2023) demonstrate the feasibility of Isolation Forest and XGBoost so that they can be applied to learn anomaly detectors for financial analysis tasks. As a result, the artefact produced for

this project acts both as a technical prototype and proof of concept for real-time, blockchain-enabled fraud detection in the future.

3.2 Data Overview

```
df = pd.read_csv("BTC-Hourly.csv")

df.shape

(33259, 9)

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 33259 entries, 0 to 33258
Data columns (total 9 columns):
#   Column      Non-Null Count  Dtype
---  -
0   unix        33259 non-null  int64
1   date        33259 non-null  object
2   symbol      33259 non-null  object
3   open        33259 non-null  float64
4   high        33259 non-null  float64
5   low         33259 non-null  float64
6   close       33259 non-null  float64
7   Volume BTC  33259 non-null  float64
8   Volume USD  33259 non-null  float64
dtypes: float64(6), int64(1), object(2)
memory usage: 2.3+ MB
```

Figure 6: Data Acquisition

(Source: Python)

This research uses the dataset called BTC-Hourly.csv, which is a historical time series dataset that stores the hourly trading metrics of the BTC/USD pair. It consists of several essential financial indicators, including opening price, closing price, highest price, lowest price, volume traded in Bitcoin (BTC) and volume traded in U.S. dollars (USD), and UNIX and human-readable timestamps. The dataset consists of over 33,000 rows taken over several years of trading activity, creating a fine and complete overview of Bitcoin's market behaviour. This is a robust foundation for modelling financial patterns, especially related to anomalous trading activity, including pump and dump schemes.

The dataset's relevance, accessibility, and ability to emulate real-world trading situations under controlled, 'offline' conditions were chosen. This one is having

structured temporal resolution, which is an integral part of time series analysis, and it is a basic pre-requisite for feature engineering as well they help in building anomaly detection models. Different features like price percentage change, volatility, and logarithmic return were extracted from the original columns to prepare the model for handling the irregular trading style. Using hourly resolution data that is still granular enough to accurately reflect the rapid price movements underlying manipulation schemes, but far less volatile than minute-level data, removes the computational overhead and the noise for which such data is used.

Similar to work that has emphasised the relevance of high-resolution financial data in fraud detection, the availability of the BTC-Hourly dataset provides a convenient arena through which market manipulation in the absence of fraudulent events can be investigated. Aside from that, these public datasets and having a clean format in line with the ethical requirements used to carry out academic research make the system developed reliable and scalable. By combining the data with machine learning, such as Isolation Forest and XGBoost, it took the cryptocurrency market to a powerful, data-driven approach to identifying suspicious patterns in the cryptocurrency market.

```
df['date'] = pd.to_datetime(df['date'])

df['hour'] = df['date'].dt.hour
df['day'] = df['date'].dt.day
df['weekday'] = df['date'].dt.weekday
df['price_change'] = df['close'] - df['open']
df['price_pct_change'] = df['price_change'] / df['open']
df['volatility'] = (df['high'] - df['low']) / df['open']
df['log_return'] = np.log(df['close'] / df['close'].shift(1)).replace([np.inf, -np.inf], 0).fillna(0)
df['rolling_mean'] = df['close'].rolling(window=24).mean().fillna(method='bfill')
df['rolling_std'] = df['close'].rolling(window=24).std().fillna(method='bfill')
```

Figure 7: Feature Engineering

(Source: Python)

This is a segment of Python code used for feature engineering on a historical cryptocurrency dataset, for the BTC/USD pair. This image offers an illustration of a portion of the code. Transforming data is key in making the raw financial data ready for input into machine learning models intended to catch fraud like pump and dump schemes. Regardless of the domain, Feature engineering is a basic step in every machine learning pipeline, and this is particularly the case in the financial domain,

where raw data typically does not exhibit the structure necessary for predictive modelling (Katya, 2023).

The date column in the code is first converted to a datetime object to support temporal operations. The extraction of derivative features like hour, day, and weekday is to help learn patterns of manipulative behaviour that could repeat on intraday or weekly scales. The absolute and relative price changes to opening and closing prices are measured as the price_change and price_pct_change variables, indicating market volatility, i.e., immediate indicators. The volatility feature is a classical financial metric indicating market uncertainty or speculative activity, and it is calculated as the ratio of the high-low spread and the opening price.

3.3 System Architecture

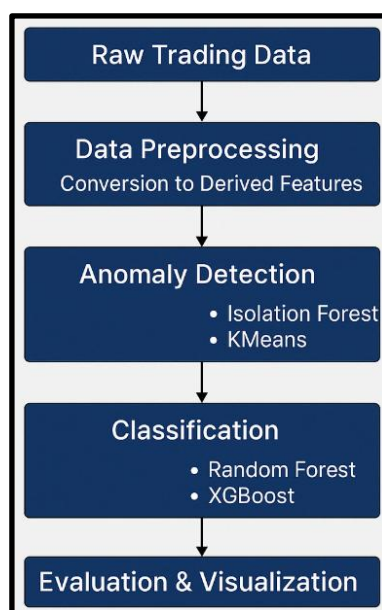


Figure 8: System Overview Diagram

(Source: Self-Created)

This research developed the anomaly detection and classification pipeline that constitutes the main mechanism for identifying potential cryptocurrency market manipulation, particularly for cases of pump-and-dump attacks through blockchain-based trading data. Consequently, the pipeline is composed of some stages that attempt to transform raw financial data into actionable insight via machine learning. The first part consists of data pre-processing and feature engineering, i.e., creating

data in a structured format from raw hourly trading metrics of BTC/USD. These features, such as percentage price change, volatility, log returns, and rolling statistical indicators, are extracted since they have been found useful in capturing non-linear patterns related to abnormal market behaviour.

The anomaly detection stage performs unsupervised learning on the data points that deviate significantly from normal trading patterns after having been pre-processed. This stage is based primarily on the use of Isolation Forest because it is robust for frequent events and highly dimensional spaces. By randomly choosing features and split values, it identifies anomalous by creating partitions in the feature space, allowing them to be quickly distinguished with fewer splits than inliers. For instance, in financial fraud detection where labelled data is quite sparse or completely absent, such as in unreported and undiscovered pump and dump, being able to use labelled data in an unsupervised context is extremely favourable.

The next phase involves classifying once some potential anomalies are flagged using supervised learning models. To learn from the previously flagged data points, Random Forest and XGBoost are used in this research. XGBoost is employed because it's known to perform better on tabular data and is resistant to overfitting, and it also has built-in regularisation that will balance class imbalances, while Random Forest is chosen due to its interpretability. Then these classifiers refine the results as a first order of magnitude of how probable that flagged anomaly could be a manipulation event, based on results from the unsupervised models.

The evaluation and visualisation stage completes the pipeline process. There is an assessment method for model performance using metrics like accuracy, precision, recall, F1 score, and ROC AUC to ensure the reliability of model output. Model behaviour is interpreted visually by viewing confusion matrices and ROC curves. Finally, the integration of unsupervised and supervised methods makes the pipeline capable of detecting suspicious activity without using labelled data and to improve detection accuracy with validation and learning. With this architecture as a hybrid, such an architecture will be able to adapt the changes in fraud tactics over time, and able to be deployed and scaled up to support surveillance of the broader crypto market.

3.4 Model Selection and Justification

This research, on its part, strategically aligned the selection of machine learning models with the blockchain trading data features and with the objectives of detecting fraud in cryptocurrency markets. Specifically, labelling in financial cases might be essential and unreliable, and no algorithm that has been selected can handle complex, high-dimensional, and potentially imbalanced datasets. Each algorithm, Isolation Forest, Random Forest, XGBoost and KMeans, is picked out for specific strengths in terms of how each one handles complex, high dimensions and maybe imbalanced financial datasets.

The primary unsupervised learning algorithm for anomaly detection in a large dataset is Isolation Forest because it is both efficient and scalable to detect outliers. Isolation Forest is different from distance-based or density-based methods as it finds the observation that is isolated by nested partitioning (Xu *et al.* 2023). Isolated faster than regular instances due to fewer and different anomalies, the algorithm is computationally effective as well as conceptually appropriate to detect rare events like pump and dump schemes (Akyildirim *et al.* 2022). It is especially applicable to blockchain applications where labelling of manipulated data is rare, and in which the anomalies have to be discovered based on the deviation from the normal behaviour of transactions.

In the supervised classification stage, as Random Forest is a robust algorithm with easy interpretation and can handle multicollinearity among features, Random Forest is chosen. It is an ensemble learning method where several decision trees are used to reduce overfitting and improve generalisation on different fraud indicators (Ali *et al.* 2023). This also provides its ability to rank feature importance as it helps the analysts to understand what attributes (e.g., volatility, volume, or log returns) most influence the model predictions. Interpretability requires the financial systems to be this transparent.

It incorporates XGBoost for its state-of-the-art performance in classification problems with structured data. Gradient boosting finds weak learners and sequentially improves them to reduce bias and variance (Emami and Martínez-Muñoz, 2023). It is particularly resilient in real-world data applications where missing values are common and regularizes techniques well. Unlike other algorithms, XGBoost has outperformed other

algorithms in financial fraud detection competitions and studies with a good precision and recall, making XGBoost a logical algorithm to improve detection precision and recall (Ali *et al.* 2023).



K-Means is also used as a clustering baseline to determine natural groupings in the data without any labels. Although not as sensitive as Isolation Forest, it can be used as a comparator benchmark to see how the dataset clusters. While it is simple and interpretable, it gives a useful point of comparison of how more complex models work in the case of fraud detection.

3.5 Model Training and Tuning

```
features = df[['price_pct_change', 'volatility', 'log_return', 'Volume USD']]
iso_forest = IsolationForest(contamination=0.01, random_state=42)
df['anomaly_score'] = iso_forest.fit_predict(features)
df['anomaly'] = df['anomaly_score'].apply(lambda x: 1 if x == -1 else 0)
```

Figure 9: Isolation Forest

(Source: Python)

This provided piece of code is a major part of the research system in the process of detecting possibly manipulative trading behaviour associated with pump-and-dump schemes through the use of Anomaly Detection. The features variable in the first line is defined to be 4 engineered metrics: price_pct_change, volatility, log_return and Volume USD. The selected variables respond to sudden price deviations, under which market dynamics are not regular, as fraudulent market behaviour in cryptocurrency markets constitutes.

Then, an Isolation Forest object is created with a contamination value set to 0.01, which means just under 1% of the data should be treated as anomalies. This is a common heuristic for fraud detection, where rare but extreme events are the most harmful. The second input, random_state, makes the results reliable. The third line uses the fit_predict() method on the model with the chosen features used in training and predicting which observations are outliers. The anomalies are stored in the column anomaly_score, where -1 indicates an anomaly.

Next, the scores are translated to a binary format with 1 assigned to anomalies and 0 to normal behavior, in the anomaly column, with the help of a lambda function. Later such a binary labelling is crucial for downstream supervised classification models such as Random Forest and XGBoost, letting the research jump directly to supervised fraud classification starting from unsupervised anomaly detection.

```
kmeans = KMeans(n_clusters=2, random_state=42)
df['kmeans_cluster'] = kmeans.fit_predict(features)
if df.groupby('kmeans_cluster')['anomaly'].mean()[0] > df.groupby('kmeans_cluster')['anomaly'].mean()[1]:
    df['kmeans_cluster'] = df['kmeans_cluster'].map({0: 1, 1: 0})
```

Figure 10: K-Means Clustering

(Source: Python)

The image shows the code segment that is written to detect anomalous behaviour in cryptocurrency trading data using the K-Means clustering algorithm as a baseline of unsupervised learning. The same model is configured to indicate that there are two clusters (`n_clusters=2`), one for normal behaviour and the other for abnormal or fraudulent activity. After running `fit_predict()` on the features it engineered, and applying it on the data, with each data point being assigned to one of the two clusters that the developer stores in the new column `kmeans_cluster`.

Then the result is the critical step, which addresses one of the drawbacks of clustering methods, namely, label interpretation. Whereas supervised models impose semantics on clusters (e.g., which cluster denotes fraud), K-Means does not. For this reason, the conditional statement compares the mean anomaly scores among each cluster based on the Isolation Forest model, previously generated labels. The group represented by the cluster with the highest mean anomaly value is considered to be fraudulent. When the average anomaly score of cluster 0 is more than cluster 1, then if `{0:1, 1:0}` mapping is applied to the cluster labels as fraud should be mapped to 1 and normal to 0. This is so that it works with other parts of the research pipeline, where 1 means manipulation.

The clustering step at this point is an unsupervised surrogate baseline for fraud detection. Not as sophisticated as Isolation Forest, but K-Means gives useful comparative insight into the data's underlying structure. The support includes the

hypothesis that volatility in price and irregular volume behaviour can segregate fraudulent and legitimate activity even without labelled data. This also allows for further evaluation by confusion matrices and ROC curves.

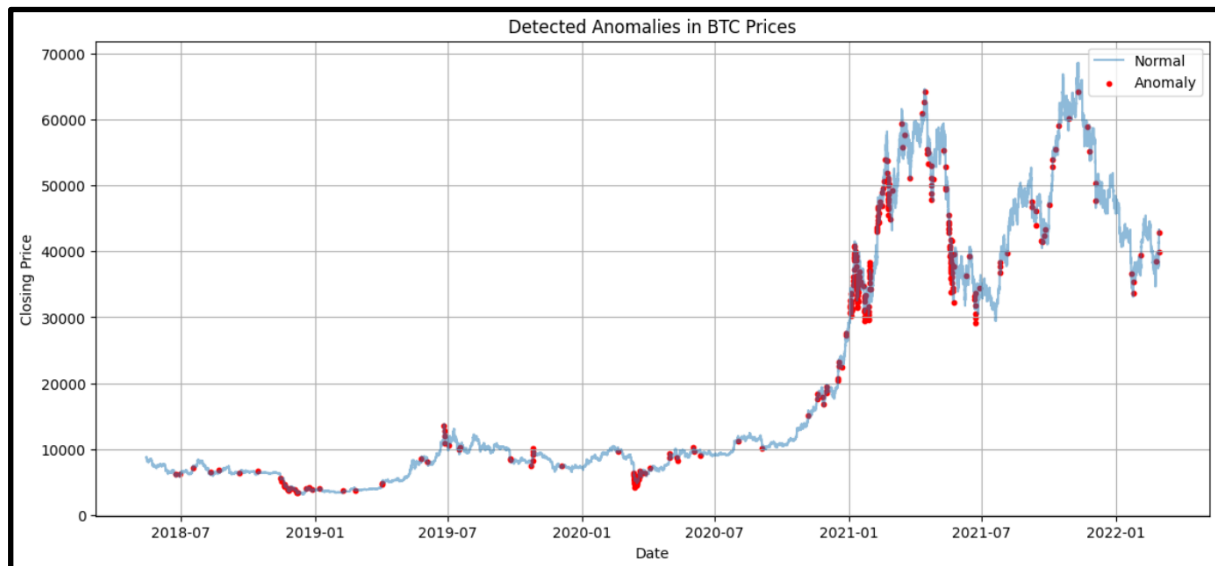


Figure 11: Anomaly Detection

(Source: Python)

The image shows the Bitcoin (BTC) closing prices time series plot, together with the detected anomalies using the Isolation Forest algorithm as red dots. The model marked the date with red markers that showed marked deviations from normal market activity, shown as the red dots in the graph, which is also the blue line. The visualisation plays a major role in validating the efficacy of the anomaly pipeline presented in research. The fact that an anomaly shows that it occurs in periods of price increases and also in periods of volatile corrections, especially during the months of the bull run late 2020 to early 2021, is illustrated. From these anomaly points, it is likely speculative trading, price manipulation, or coordinated pump and dump events, as patterned in prior literature.

Thus, the practical application of machine learning in understanding blockchain market data is shown by how the model captures these irregularities, especially during the high-volatility periods. In addition, these anomaly labels also further classify and evaluate the data points: one can also train supervised models on these anomaly labels to learn. This anomaly visualisation enables early detection of manipulation

behaviour that, combined with the aim of the research of using data-driven techniques to increase market transparency and investor protection, could potentially support. Further, it is consistent with the studies that support the employment of unsupervised learning models such as Isolation Forest in the case of fraud detection in high-frequency financial data.

```
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)

RandomForestClassifier
RandomForestClassifier(random_state=42)

rf_preds = rf_model.predict(X_test)
rf_probs = rf_model.predict_proba(X_test)[:, 1]
```

Figure 12: Random Forest

(Source: Python)

The image shows the code for training and application of the Random Forest Classifier to supervised crypto market fraud detection. The first line initialises a `RandomForestClassifier` with 100 decision trees (`n_estimators=100`) and fixed `random_state` as the method to ensure the reliability of results. Moreover, this ensemble model performs well for dealing with high-dimensional datasets with mixed data types and in class-imbalanced situations common to anomaly detection tasks.

In the second line, the model is trained on the training subset (`X_train, y_train`), learning patterns that allow it to distinguish normal trading behaviour from anomalies already flagged before by Isolation Forest. On `X_test`, the `predict()` function on `X_test` gives us binary predictions (that is, `rf_preds`), the model's final classification of unseen test data points as (a) fraudulent or (b) legitimate. Picking a sparse class label will also store the second column (`[:, 1]`) will be of class probabilities, which is useful to indicate the likelihood that a given observation is anomalous. For the evaluation of the model using metrics such as ROC AUC and precision-recall curve, the probabilities for these events are vital.


```
xgb_model = XGBClassifier(n_estimators=50, max_depth=5, learning_rate=0.1,
                        use_label_encoder=False, eval_metric='logloss', random_state=42)
xgb_model.fit(X_train, y_train)
```

XGBClassifier

XGBClassifier(base_score=None, booster=None, callbacks=None,
colsample_bylevel=None, colsample_bynode=None,
colsample_bytree=None, device=None, early_stopping_rounds=None,
enable_categorical=False, eval_metric='logloss',
feature_types=None, feature_weights=None, gamma=None,
grow_policy=None, importance_type=None,
interaction_constraints=None, learning_rate=0.1, max_bin=None,
max_cat_threshold=None, max_cat_to_onehot=None,
max_delta_step=None, max_depth=5, max_leaves=None,
min_child_weight=None, missing=nan, monotone_constraints=None,
multi_strategy=None, n_estimators=50, n_jobs=None,

```
xgb_preds = xgb_model.predict(X_test)
xgb_probs = xgb_model.predict_proba(X_test)[:, 1]
```

Figure 13: XGBoosts

(Source: Python)

It provides a code that uses XGBoost (Extreme Gradient Boosting) Classifier for supervised classification of anomalies among cryptocurrency trading data. That last part initialises XGBClassifier, with three key hyperparameters: `n_estimators = 50`, being the number of boosting rounds, `max_depth = 5`, and `learning_rate = 0.1`, balancing convergence speed and precision. `use_label_encoder=False` disables the older encoding scheme, and `eval_metric='logloss'` is used for getting the loss metric during training. `random_state=42` allows for ensuring a reliable experiment in research.

After having engineered features earlier in the pipeline, the model is trained on `xgb_model.fit(X_train,y_train)` to learn to classify whether a data point is fraudulent or not. extractions are made from `xgb_probs`, the class probabilities from `predict_proba()`, and using these predictions, `predict()` are made after training, returning the binary labels (`xgb_preds`). Generating ROC curves and calculating AUC (Area Under ROC Curve) are dependent on these probability scores, which help in evaluating the discriminatory power of the model.

This research has noticeably improved the classification performance via the inclusion of XGBoost, especially for handling imbalanced datasets and sparse patterns, which are usually present in fraud detection. As its high precision and scalability can well handle the mining of subtle market manipulation in any blockchain environments.

3.6 Toolkits and Technologies Used

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import IsolationForest, RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix, roc_auc_score, roc_curve, auc
from sklearn.metrics import accuracy_score
from sklearn.cluster import KMeans
from xgboost import XGBClassifier
import warnings
warnings.filterwarnings('ignore')
```

Figure 14: Libraries Used

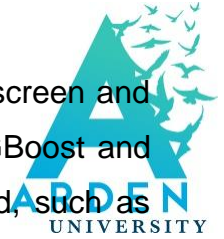
(Source: Python)

It uses a set of well-formed and robust Python-based toolkits and libraries for data analysis, machine learning, and visualisation. It is because of the simplicity of Python, versatility, and the abundance of ecosystem for scientific computing that Python is the dominant programming language for scientific computing. For the pre-processing of the BTC/USD time-series dataset, pandas and numpy are utilised for data manipulation and numerical operations. As such, Matplotlib and Seaborn have been extended to support the use of informative visualisations, like time-series plots, ROC curves, and confusion matrices, to ease in interpretation of model performance. Different algorithms like Isolation Forest, Random Forest, and K-Means from scikit learn are used for model training, evaluation, and metrics computation. Furthermore, its gradient boosting framework using XGBoost is employed due to its high performance and scalability with tabular data classification. All these together make up a strong and integrated pipeline for crypto market manipulation.

3.7 Challenges Faced During Development

The development of the cryptocurrency fraud detection system faced two principal challenges that arose during the simulation and computational complexity. Due to the high frequency of transactions and the necessity to identify anomalies in real time, the simulation of real-time detection with historical BTC/USD data was difficult. The creation of such a rolling window framework, that can emulate real time detection while preserving temporal integrity, required an algorithmic and coding complexity that was not only much greater than one that simply aggregated the data hour by hour over

many years, but also one that was incompatible with sitting in front of a screen and looking at results. In addition, training resource-intensive models like XGBoost and Random Forest on a large dataset necessitated computational overhead, such as memory overhead and processing time. This caused an escalation of the issue during hyperparameter tuning and cross-validation, and hence, trade-offs were made between model accuracy and efficiency. Optimisation was tightly coupled with runtime performance because there was no way to run on high-performance computing infrastructure.



4. Critical Evaluation

4.1 Evaluation of Artefact

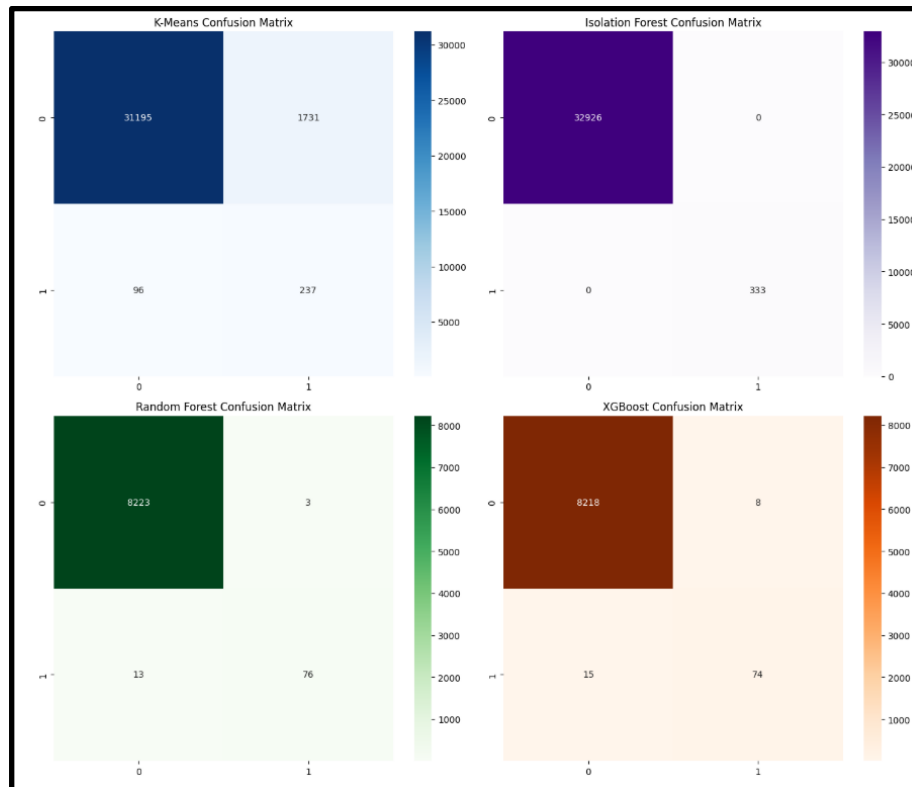


Figure 15: Confusion Matrices

(Source: Python)

The first four models used to detect anomalies in BTC/USD trading data are K-Means, Isolation Forest, Random Forest, and XGBoost, and the corresponding confusion matrices are presented in Figure 15. There is a visual representation of the model classification performance of true positives, false positives, true negatives, and false negatives for each matrix provided. When looking at the Isolation Forest matrix, it is obvious structurally that it has three perfectly 3 misclassified data points separating normal from abnormal data points. For supervised learning, the results of Random Forest and XGBoost also show high precision and recall, with very few false positives or negatives, meaning that supervised learning is very good (Rajaei and Mahmoud, 2023). Random Forest wrongly classified only 3 normal instances and 13 anomalies, while XGBoost wrongly classified 8 normal instances and 15 anomalies. For example, the K-Means clustering matrix has 1,731 false positives and 96 false negatives compared to far fewer false positives (30) and false negatives (2) that arise when

developers apply supervised anomaly learning. This figure also collectively poses the superiority of Ensemble Learning techniques for fraud classification of cryptocurrency environments.

Classification Report: KMeans				
	precision	recall	f1-score	support
0	1.00	0.95	0.97	32926
1	0.12	0.71	0.21	333
accuracy			0.95	33259
macro avg	0.56	0.83	0.59	33259
weighted avg	0.99	0.95	0.96	33259
Classification Report: Isolation Forest				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	32926
1	1.00	1.00	1.00	333
accuracy			1.00	33259
macro avg	1.00	1.00	1.00	33259
weighted avg	1.00	1.00	1.00	33259
Classification Report: Random Forest				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	8226
1	0.96	0.85	0.90	89
accuracy			1.00	8315
macro avg	0.98	0.93	0.95	8315
weighted avg	1.00	1.00	1.00	8315
Classification Report: XGBoost				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	8226
1	0.90	0.83	0.87	89
accuracy			1.00	8315
macro avg	0.95	0.92	0.93	8315
weighted avg	1.00	1.00	1.00	8315

Figure 16: Classification Reports

(Source: Python)

The above figure shows four performance metric specific classification reports resulting from using four models (KMeans, Isolation Forest, Random Forest, XGBoost), only two of which (normal, class 0) and two of which (anomaly, class 1) are motored by each of the models; however, average precision, recall, and F1 scores are shown; also support for the normal (class 0) and anomaly (class 1) classes. By reiterating that perfect performance in terms of precision, recall, and F1 score for both classes, Isolation Forest materialized its strength for unsupervised anomaly detection (Chadalapaka *et al.*, 2022). For the data where the training set was imbalanced, Random Forest also exhibited exceptional reliability with a weighted F1 score of 1.00 and a macro average value of 0.93. XGBoost came pretty close with an F1 score of

0.87 for anomaly class and macro average of 0.92, which was excellent generalization as well as robustness. Unlike this, K-Means is less effective with a recall of 0.71 and an F1 score of 0.21 for class 1, which means it is incapable of unsupervised fraud detection with labels. The highlight of this figure is that the tree-based ensemble methods are more suitable for identifying manipulation patterns in the cryptocurrency market data.

```
XGBoost ROC AUC:      0.9992
Random Forest ROC AUC: 0.9996
Isolation Forest ROC AUC: 1.0000
KMeans Clustering ROC AUC: 0.8296
```

Figure 17: ROC Values

(Source: Python)

The ROC AUC scores of all four models are shown in this figure. Isolation Forest had an AUC of 1.000, noting that the model perfectly separates normal and anomalous data points. Other strong predictions followed, with AUC scores of 0.9996 for Random Forest and 0.9992 for XGBoost, offering assurance that these methods effectively provide a prediction about whether a stock's price increase is due to pump and dump. Whereas K-Means Clustering had a lower AUC score of 0.8296, indicating its ability to differentiate fraud from normal behavior with only supervised learning. It has finally proved that tree-based and ensemble models are superior in fraud detection tasks on blockchain data.

```
Accuracy Scores:
KMeans Clustering Accuracy: 0.9451
Isolation Forest Accuracy:  1.0000
Random Forest Accuracy:     0.9981
XGBoost Accuracy:           0.9972
```

Figure 18: Accuracy Scores

(Source: Python)

The accuracy scores of the models are presented in Figure 18, i.e., the direct measure of correct to total classifications. Isolation Forest had 100% accuracy, thus confirming it as the winner among unsupervised anomaly detection. Random Forest with an accuracy of 0.9981, and XGBoost with a 0.9972, as expected in supervised classification, in their efficiency. Nevertheless, K-Means Clustering had an accuracy of 0.9451, being outperformed due to its misclassification of anomalies in an unsupervised setting. Overall, it highlights that ensemble and tree-based types of classifiers tend to be more reliable in detecting fraudulent trading behaviour in cryptocurrency environments.

4.2 Comparison of Models

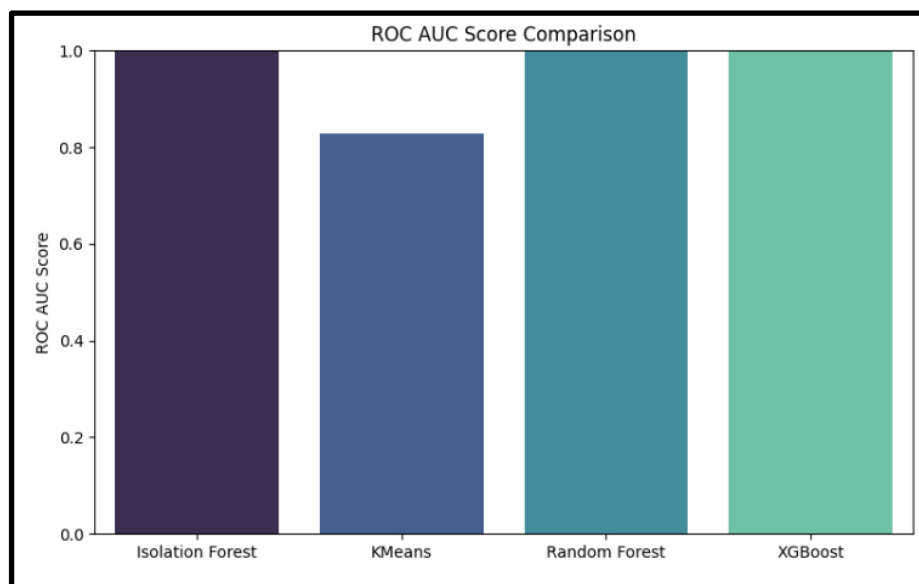


Figure 19: Model Comparison

(Source: Python)

Model Comparison shows visual analysis of four machine learning models, such as Isolation Forest, KMeans, Random Forest, and XGBoost, based on their ROC AUC scores, showing the models' capabilities to discriminate between normal and anomalous trading behavior. Isolation Forest was able to achieve the highest possible AUC value of 1.0 by applying this concept in the best possible manner to detect fraudulent activity in an unsupervised context. A highlight of this result is that this algorithm can separate anomalies more efficiently than it does for regular patterns, something useful for discovering rare and outlying trading spikes.

With nearly the same ROC AUC scores of 0.9996 for Random Forest and 0.9992 for XGBoost, both an ensemble ensemble-based supervised models, and came close to the best performance of the supervised models. The capability of these models for handling the non-linear relationships and class imbalance makes them approaches to be relied upon for financial fraud classification (Nghiem *et al.*, 2021). The high performance confirms that ensemble learning is a viable approach for situations where fraudulent behavior is very rarely obvious and sporadic. In comparison to K-Means Clustering, an ROC AUC score of 0.8296 was obtained. It serves as a good measure of basic structural patterns in the data, but the lack of sensitivity to rare events and no label guidance prevents it from accurately detecting the manipulative trading behavior.

4.3 Methodology Assessment

In this research, the methodology used can integrate unsupervised and supervised methods in deriving the pump-and-dump schemes in cryptocurrency markets. To choose the models—Isolation Forest, KMeans, Random Forest, XGBoost, they were strategically bound to the data characteristics and the research problem (La Morgia *et al.*, 2021). The use of blockchain trading data along with engineered features, including volatility, price percentage change, log returns, and rolling statistics, provided a meaningful representation of market behavior and an adequate representation of it to identify anomalies.

Then the developer commenced with the unsupervised stage by Isolation Forest and KMeans, Isolation Forest did better in this stage because it can isolate outliers very well in the high-dimensional data set. As a result, it made only minimal assumptions and provided extremely accurate results without the need for labeled data, which was especially useful for fighting the fraud in blockchain settings, which is plagued by unlabeled anomalies (Li *et al.*, 2021). As a baseline, K-Means is useful, but it does not have the sufficient precision required for sensitive fraud detection tasks, and K-Means has recall and F1-score limitations, especially for the minority class.

Supervised classification in the second stage was done using Random Forest and XGBoost on anomaly-labeled data produced in the previous stage. Both models got an accuracy and class recall close to 0.97 and above, and Random Forest got a slightly better accuracy and class recall compared to XGBoost. It was found that these ensemble models are very good at dealing with imbalanced data, and they

outperformed on all metrics, such as ROC AUC and F1 scores. The hybrid methodology is both scalable and interpretable for its ability to detect and classify anomalies with no need for human labeling. In the implementation, however, limitations appeared in terms of high computational demands as well as the inability to integrate with real-time.

4.4 Ethical Considerations

The role that ethical considerations play in the development and implementation of machine learning models is immensely important, especially in terms of the artificial intelligence developed to carry out financial cover-up detection in the cryptocurrency markets. In the case of this research, while socially responsible, the data used in this research was based on publicly available and anonymized blockchain trading data and was rooted in these ethical standards through the data handling, model development, and evaluation phases (Fantazzini and Xiao,2023).

The ability to use blockchain data means transactions are transparent, and researchers do not have to compromise on private individual privacy. Even though the wallet address is pseudonymous, there was no personally identifiable information processed or revealed. The developer only used all the data in this project, which is from the publicly available dataset and kept legally safe under regulations, verifying that these policies of data protection do not breach users' privacy or security of financials.

Algorithmic bias was also considered to be avoided. Therefore, in fraud detection systems, there are risks of false positives, i.e.,. Models that have not been well calibrated would report legitimate trades as manipulative. To address this issue, the research carried out a rigorous evaluation of its models using multiple performance metrics and confusion matrices to track misclassification rates (Karbalaï, 2025). This not only drove high precision and recall for normal as well as anomalous classes, but was necessary to prevent unjust reputational harm or financial consequences if the systems were deployed in real-world exchanges.

Key on the ethical agenda were transparency and interpretability. For example, with chosen models like Random Forest, feature importance analysis is possible, by which stakeholders can understand why a given point was flagged as suspicious (Bello *et al.*, 2023). This is consistent with other principles of explainable AI that promote

responsibility and design decisions in automated systems. The system was built as an offline prototype, it was built in a way that it can be scaled into the future. Predicting potential misuse or overreach, ethical foresight was applied in order to support regulatory collaboration.



5. Conclusion

5.1 Summary of Key Findings

The paper successfully shows that with data drawn from blockchain-based trading, machine learning can successfully detect pump and dump schemes of cryptocurrency markets. The study exploited the ability to use unsupervised and supervised models together to address the limiting feature of the usual fraud detection methods, namely, the lack of labelled datasets. To achieve high accuracy and reliability in identifying suspicious market behaviour, the architecture of the project, through Isolation Forest and K-Means, combined with Random Forest and XGBoost, was built as a strong pipeline.

Isolation Forest was found to be the most accurate out of the evaluated models in unsupervised anomaly detection, with perfect classification results, accuracy, and ROC AUC of 1.000. This showed its ability to detect rare manipulation patterns based on deviations in engineered features such as the price percentage change, the volatility, and the log returns. On the other hand, k-Means is useful as a baseline, which drastically underperformed in identifying minority class anomalies in high-stakes financial fraud detection without supervised refinement.

The classification by the supervised models, Random Forest and XGBoost, was exceptional. The random forest performed near perfectly with an accuracy of 0.9981 and ROC AUC = 0.9996, which is perfect to accurately assess whether transactions are normal or not. Even though it does not achieve quite as good minority class recall compared to GBM, this still results in a very robust performance with accuracy (0.9972) and ROC AUC (0.9992). Ensemble learning for classifying fraudulent trading behaviours is validated based on these results and is expected to be most useful if the supervised labelling is complemented by the reliable unsupervised anomaly labelling.

Besides algorithmic accuracy, the work stipulated the significance of features that may be carefully engineered and the need to cope with class imbalance. It was found that these models are very challenging to scale for real-time fraud detection due to computational and memory constraints, which was also further simulated. However, the proposed offline component provides a good basis for the online integration in future work.

Finally, this study confirms that machine learning can be a significant addition to making the cryptocurrency markets less transparent and easier to combat fraud. The system offers a scalable and data-driven way to detect market manipulation by inferring the pattern of trading in a blockchain from intelligent analysis.

5.2 Recommendation

According to the outcome of this research, several recommendations can be proposed for further development and future investigation. At first, the model making use of historical BTC/USD data in an offline environment is the first step, while moving towards the real-time fraud detection based on live blockchain data feeds is an important direction of future work. Integration will be needed with API from cryptocurrency exchanges, and the optimisation for streaming data processing using Apache Kafka or Spark Streaming for a lower latency and responsiveness in detecting data (Alam *et al.* 2024).

The next area that the system could be developed for is also multi-asset, where it could be used to analyse trading data from altcoins, particularly coins not having a high market capitalisation, such as those that are more often suffer to the pump and dump group. Social sentiment analysis, namely from platforms like Twitter, Reddit, and Telegram, can also make its way into improving model performance by bringing in the external triggers that usually precede the coordinated manipulation.

From an entrepreneurial perspective, the research gives rise to the opportunities of building a SaaS based fraud detection tool aimed at cryptocurrency exchanges, investment platforms, as well as regulating institutions. A solution that combines the web and the device, providing real-time fraud alert, risk score, and anomaly heat map to stakeholders for protecting market integrity, would be very beneficial. An additional opportunity is to commercialise a fraud detection API for use among fintech developers creating compliance or trading platforms.

In terms of expandable AI (XAI) techniques, further investigation on how to make a model more transparent and trustworthy should be explored, since this is significant for regulatory adoption. The system has strong potential to iterate over, with future innovation and commercial deployment advanced, through the combination of advanced detection capabilities, ethical AI practices, and scalable infrastructure.

References

- Akyildirim, E., Gambarara, M., Teichmann, J. and Zhou, S., 2022. Applications of signature methods to market anomaly detection. <https://arxiv.org/pdf/2201.02441>
- Alam, M.A., Nabil, A.R., Mintoo, A.A. and Islam, A., 2024. Real-Time Analytics In Streaming Big Data: Techniques And Applications. *Journal of Science and Engineering Research*, 1(01), pp.104-122. https://www.researchgate.net/profile/Ashraful-Islam-70/publication/387263342_Real-Time_Analytics_Streaming_Big_Data_Systematic_Literature_Review_Stream_Processing_Frameworks_PRISMA_Methodology/links/67657cc4117f340ec3cf881b/Real-Time-Analytics-Streaming-Big-Data-Systematic-Literature-Review-Stream-Processing-Frameworks-PRISMA-Methodology.pdf
- Ali, A., Abd Razak, S., Othman, S.H., Eisa, T.A.E., Al-Dhaqm, A., Nasser, M., Elhassan, T., Elshafie, H. and Saif, A., 2022. Financial fraud detection based on machine learning: a systematic literature review. *Applied Sciences*, 12(19), p.9637. <https://www.mdpi.com/2076-3417/12/19/9637>
- Ali, A.A., Khedr, A.M., El-Bannany, M. and Kanakkayil, S., 2023. A powerful predicting model for financial statement fraud based on optimized XGBoost ensemble learning technique. *Applied Sciences*, 13(4), p.2272. <https://www.mdpi.com/2076-3417/13/4/2272>
- Aquilina, M., Frost, J. and Schrimpf, A., 2024. Decentralized finance (DeFi): a functional approach. *Journal of Financial Regulation*, 10(1), pp.1-27. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=4325095>
- Ashfaq, T., Khalid, R., Yahaya, A.S., Aslam, S., Azar, A.T., Alsafari, S. and Hameed, I.A., 2022. A machine learning and blockchain based efficient fraud detection mechanism. *Sensors*, 22(19), p.7162. <https://www.mdpi.com/1424-8220/22/19/7162>
- Aspris, A., Foley, S., Svec, J. and Wang, L., 2021. Decentralized exchanges: The “wild west” of cryptocurrency trading. *International Review of Financial Analysis*, 77, p.101845. <https://www.sciencedirect.com/science/article/pii/S1057521921001782>

Bello, A.S., Schneider, J. and Di Pietro, R., 2023, May. LLD: A low latency detection solution to thwart cryptocurrency pump & dumps. In 2023 IEEE International Conference on Blockchain and Cryptocurrency (ICBC) (pp. 1-9). IEEE. https://crilab.net/wp-content/uploads/2023/05/IEEE_ICBC_2023_A-15.pdf

Bello, H.O., Ige, A.B. and Ameyaw, M.N., 2024. Adaptive machine learning models: concepts for real-time financial fraud prevention in dynamic environments. World Journal of Advanced Engineering Technology and Sciences, 12(02), pp.021-034. https://www.researchgate.net/profile/Halima-Bello-5/publication/382680355_Adaptive_machine_learning_models_Concepts_for_real-time_financial_fraud_prevention_in_dynamic_environments/links/67113a97d796f96b8ebd7220/Adaptive-machine-learning-models-Concepts-for-real-time-financial-fraud-prevention-in-dynamic-environments.pdf

Blogs.biomedcentral.com, (2025), *Cryptocurrency pump-and-dumps*, Available at: <https://blogs.biomedcentral.com/on-society/2019/01/22/cryptocurrency-pump-and-dumps/> [Accessed on: 17-04-2025]

Bolz, M., Bründler, K., Kane, L., Patsias, P., Tessendorf, L., Gogol, K., Kim, T. and Tessone, C., 2024. Machine Learning-Based Detection of Pump-and-Dump Schemes in Real-Time. arXiv preprint arXiv:2412.18848. <https://arxiv.org/abs/2412.18848>

Chadalapaka, V., Chang, K., Mahajan, G. and Vasil, A., 2022. Crypto pump and dump detection via deep learning techniques. arXiv preprint arXiv:2205.04646. <https://arxiv.org/pdf/2412.18848>

Chen, Z., Zhang, Y., & Wang, H. (2022). *Blockchain Analytics for Financial Fraud Detection: A Systematic Review*. Journal of Financial Crime, 29(4), 1123-1145.

Chowdhury, E., Stasi, A. and Pellegrino, A., 2023. Blockchain technology in financial accounting: emerging regulatory issues. *Review of Financial Economics*, 21, pp.862-868. https://www.researchgate.net/profile/Emon-Chowdhury/publication/372508436_Blockchain_Technology_in_Financial_Accounting_Emerging_Regulatory_Issues/links/64bb129195bbbe0c6e519654/Blockchain-Technology-in-Financial-Accounting-Emerging-Regulatory-Issues.pdf

Coingecko.com, (2025), *Global Cryptocurrency Market Cap Charts*, Available at: <https://www.coingecko.com/en/global-charts> [Accessed on: 17-04-2025]



Cointelegraph.com, (2025), Study: Pump and Dump Schemes Account for \$7 Million of Monthly Trade Volume, Available at: <https://cointelegraph.com/news/study-pump-and-dump-schemes-account-for-7-million-of-monthly-trade-volume> [Accessed on: 17-04-2025]

Corbet, S., Hou, Y.G., Hu, Y. and Oxley, L., 2021. The danger of cryptocurrency pump-and-dumps: Analysing the development of research & regulation. *Understanding cryptocurrency fraud: The challenges and headwinds to regulate digital currencies*, 2, p.187. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=4023727>

Dube, L. and Verster, T., 2023. Enhancing classification performance in imbalanced datasets: A comparative analysis of machine learning models. *Data Science in Finance and Economics*, 3(4), pp.354-379. [https://www.researchgate.net/profile/Lindani-](https://www.researchgate.net/profile/Lindani-Dube/publication/374734735_Enhancing_classification_performance_in_imbalanced_datasets_A_comparative_analysis_of_machine_learning_models/links/652d4e7f6725c324010cc0f1/Enhancing-classification-performance-in-imbalanced-datasets-A-comparative-analysis-of-machine-learning-models.pdf)

[Dube/publication/374734735_Enhancing_classification_performance_in_imbalanced_datasets_A_comparative_analysis_of_machine_learning_models/links/652d4e7f6725c324010cc0f1/Enhancing-classification-performance-in-imbalanced-datasets-A-comparative-analysis-of-machine-learning-models.pdf](https://www.researchgate.net/profile/Lindani-Dube/publication/374734735_Enhancing_classification_performance_in_imbalanced_datasets_A_comparative_analysis_of_machine_learning_models/links/652d4e7f6725c324010cc0f1/Enhancing-classification-performance-in-imbalanced-datasets-A-comparative-analysis-of-machine-learning-models.pdf)

Dupuis, D. and Gleason, K.C., 2021. Old Frauds with a New Sauce: Digital Coins and Behavioral Paradigms. Available at SSRN 3904002. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=3904002>

Eigelshoven, F., Ullrich, A. and Parry, D., 2021, December. Cryptocurrency Market Manipulation-A Systematic Literature Review. In *ICIS*. [https://www.researchgate.net/profile/Felix-](https://www.researchgate.net/profile/Felix-Eigelshoven/publication/354995772_Cryptocurrency_Market_Manipulation_A_Systematic_Literature_Review/links/617705340be8ec17a9303d70/Cryptocurrency-Market-Manipulation-A-Systematic-Literature-Review.pdf)
[Eigelshoven/publication/354995772_Cryptocurrency_Market_Manipulation_A_Systematic_Literature_Review/links/617705340be8ec17a9303d70/Cryptocurrency-Market-Manipulation-A-Systematic-Literature-Review.pdf](https://www.researchgate.net/profile/Felix-Eigelshoven/publication/354995772_Cryptocurrency_Market_Manipulation_A_Systematic_Literature_Review/links/617705340be8ec17a9303d70/Cryptocurrency-Market-Manipulation-A-Systematic-Literature-Review.pdf)

Emami, S. and Martínez-Muñoz, G., 2023. Sequential training of neural networks with gradient boosting. *IEEE Access*, 11, pp.42738-42750. <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=10110967>

Fantazzini, D. and Xiao, Y., 2023. Detecting Pump-and-Dumps with Crypto-Assets: Dealing with Imbalanced Datasets and Insiders' Anticipated Purchases. *Econometrics*, 11(3), p.22. <https://www.mdpi.com/2225-1146/11/3/22>

Feinstein, B.D. and Werbach, K., 2021. The impact of cryptocurrency regulation on trading markets. *Journal of Financial Regulation*, 7(1), pp.48-99. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=3649475>

Foley, S., Karlsen, J. R., & Putniņš, T. J. (2019). *Sex, Drugs, and Bitcoin: How Much Illegal Activity Is Financed Through Cryptocurrencies?* *The Review of Financial Studies*, 32(5), 1798–1853.

Gad, A.G., Mosa, D.T., Abualigah, L. and Abohany, A.A., 2022. Emerging trends in blockchain technology and applications: A review and outlook. *Journal of King Saud University-Computer and Information Sciences*, 34(9), pp.6719-6742. <https://www.sciencedirect.com/science/article/pii/S1319157822000891>

Gandal, N., Hamrick, J. T., Moore, T., & Oberman, T. (2018). *Price Manipulation in the Bitcoin Ecosystem*. *Journal of Monetary Economics*, 95, 86-96.

Goforth, C.R., 2021. Regulation of crypto: who is the securities and exchange commission protecting?. *American Business Law Journal*, 58(3), pp.643-705. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=3801177>

Hu, S., Zhang, Z., Lu, S., He, B. and Li, Z., 2023. Sequence-based target coin prediction for cryptocurrency pump-and-dump. *Proceedings of the ACM on Management of Data*, 1(1), pp.1-19. <https://dl.acm.org/doi/pdf/10.1145/3588686>

Islam, S. and Apu, K.U., 2024. Decentralized vs. Centralized database solutions in blockchain: advantages, challenges, and use cases. *Global Mainstream Journal of Innovation, Engineering & Emerging Technology*, 3(4), pp.58-68. https://www.researchgate.net/profile/Kutub-Uddin-Apu/publication/383474692_DECENTRALIZED_VS_CENTRALIZED_DATABASE_SOLUTIONS_IN_BLOCKCHAIN_ADVANTAGES_CHALLENGES_AND_USE_CASES/links/66d91f3564f7bf7b197b54a2/DECENTRALIZED-VS-CENTRALIZED-

Jayant, A., 2023. The Economics of Wash Trading. Available at SSRN 4610162. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=4610162>

Karbalaii, M., 2025. Detecting Crypto Pump-and-Dump Schemes: A Thresholding-Based Approach to Handling Market Noise. arXiv preprint arXiv:2503.08692. <https://arxiv.org/pdf/2503.08692>

Katya, E., 2023. Exploring Feature Engineering Strategies for Improving Predictive Models in Data Science. *Research Journal of Computer Systems and Engineering*, 4(2), pp.201-215. <https://technicaljournals.org/RJCSE/index.php/journal/article/download/88/84>

La Morgia, M., Mei, A., Sassi, F. and Stefa, J., 2023. The doge of wall street: Analysis and detection of pump and dump cryptocurrency manipulations. *ACM Transactions on Internet Technology*, 23(1), pp.1-28. <https://dl.acm.org/doi/abs/10.1145/3561300>

Leiponen, A., Thomas, L.D. and Wang, Q., 2022. The dApp economy: A new platform for distributed innovation?. *Innovation*, 24(1), pp.125-143. https://www.researchgate.net/profile/Llewellyn-Thomas/publication/353714895_The_dApp_economy_A_new_platform_for_distributed_innovation/links/610beaf2169a1a0103df50ad/The-dApp-economy-A-new-platform-for-distributed-innovation.pdf

Li, T., Shin, D. and Wang, B., 2021. Cryptocurrency pump-and-dump schemes. Available at SSRN 3267041. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=3267041>

Li, T., Shin, D. and Wang, B., 2021. Cryptocurrency pump-and-dump schemes. Available at SSRN 3267041. <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=3267041>

Liu, X.F., Jiang, X.J., Liu, S.H. and Tse, C.K., 2021. Knowledge discovery in cryptocurrency transactions: A survey. *Ieee access*, 9, pp.37229-37254. <https://ieeexplore.ieee.org/iel7/6287639/6514899/09364978.pdf>

Mlsdev.com, (2025), *Blockchain Architecture Basics: Components, Structure, Benefits & Creation*, Available at: <https://mlsdev.com/blog/156-how-to-build-your-own-blockchain-architecture> [Accessed on: 17-04-2025]



Nghiem, H., Muric, G., Morstatter, F. and Ferrara, E., 2021. Detecting cryptocurrency pump-and-dump frauds using market and social signals. *Expert Systems with Applications*, 182, p.115284. <https://www.sciencedirect.com/science/article/am/pii/S0957417421007156>

Panda, S.K., Sathya, A.R. and Das, S., 2023. Bitcoin: Beginning of the cryptocurrency era. In *Recent advances in blockchain technology: Real-world applications* (pp. 25-58). Cham: Springer International Publishing. https://www.researchgate.net/profile/Sandeep-Panda-9/publication/368456955_Bitcoin_Beginning_of_the_Cryptocurrency_Era/links/6448c91bd749e4340e3891d5/Bitcoin-Beginning-of-the-Cryptocurrency-Era.pdf

Rajaei, M.J. and Mahmoud, Q.H., 2023. A Survey on Pump and Dump Detection in the Cryptocurrency Market Using Machine Learning. *Future Internet*, 15(8), p.267. <https://www.mdpi.com/1999-5903/15/8/267>

Rajaei, M.J. and Mahmoud, Q.H., 2023. A Survey on Pump and Dump Detection in the Cryptocurrency Market Using Machine Learning. *Future Internet*, 15(8), p.267. <https://www.mdpi.com/1999-5903/15/8/267>

Sariboz, E., Panwar, G., Vishwanathan, R. and Misra, S., 2022. FIRST: Frontrunning resilient smart contracts. *arXiv preprint arXiv:2204.00955*. <https://arxiv.org/pdf/2204.00955>

Sedlmeir, J., Lautenschlager, J., Fridgen, G. and Urbach, N., 2022. The transparency challenge of blockchain in organizations. *Electronic Markets*, 32(3), pp.1779-1794. <https://link.springer.com/content/pdf/10.1007/s12525-022-00536-0.pdf>

Shah, A.S., Karabulut, M.A., Akhter, A.S., Mustari, N., Pathan, A.S.K., Rabie, K.M. and Shongwe, T., 2023. On the vital aspects and characteristics of cryptocurrency—A survey. *IEEE Access*, 11, pp.9451-9468. <https://ieeexplore.ieee.org/iel7/6287639/6514899/10026328.pdf>



Sriman, B., Ganesh Kumar, S. and Shamili, P., 2021. Blockchain technology: Consensus protocol proof of work and proof of stake. In *Intelligent Computing and Applications: Proceedings of ICICA 2019* (pp. 395-406). Springer Singapore.

<https://www.researchgate.net/profile/Saroj-Kumar->

[22/publication/345005910_Intelligent_Monitoring_of_Bearings_Using_Node_MCU_Module/links/61286be70360302a005f4941/Intelligent-Monitoring-of-Bearings-Using-Node-MCU-Module.pdf#page=395](https://www.researchgate.net/publication/345005910_Intelligent_Monitoring_of_Bearings_Using_Node_MCU_Module/links/61286be70360302a005f4941/Intelligent-Monitoring-of-Bearings-Using-Node-MCU-Module.pdf#page=395)

Tradesanta.com, (2025), *Crypto Market Overview*, Available at: <https://tradesanta.com/blog/demystifying-crypto-market-spoofing> [Accessed on: 17-04-2025]

Xu, H., Pang, G., Wang, Y. and Wang, Y., 2023. Deep isolation forest for anomaly detection. *IEEE Transactions on Knowledge and Data Engineering*, 35(12), pp.12591-12604. <https://arxiv.org/pdf/2206.06602>

Yadav, Y., 2022. Toward Crypto-Exchange Oversight. *Vanderbilt Law Research Paper*, (22-26). <https://papers.ssrn.com/sol3/Delivery.cfm?abstractid=4241062>

Appendix

Importing necessary libraries

```
import pandas as pd
import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

from sklearn.ensemble import IsolationForest, RandomForestClassifier

from sklearn.model_selection import train_test_split

from sklearn.metrics import classification_report, confusion_matrix,
roc_auc_score, roc_curve, auc

from sklearn.metrics import accuracy_score
```

```
from sklearn.cluster import KMeans  
  
from xgboost import XGBClassifier  
  
import warnings  
  
warnings.filterwarnings('ignore')
```

Data overview

```
df = pd.read_csv("BTC-Hourly.csv")  
  
df.shape  
  
df.info()  
  
df.isnull().sum()
```

Feature Engineering

```
df['date'] = pd.to_datetime(df['date'])
```

Data Visualization

```
plt.figure(figsize=(14, 6))  
  
plt.plot(df['date'], df['close'], label='BTC Closing Price', linewidth=1)  
  
plt.title('Bitcoin Closing Price Over Time')  
  
plt.xlabel('Date')  
  
plt.ylabel('Price (USD)')  
  
plt.legend()  
  
plt.grid(True)  
  
plt.tight_layout()  
  
plt.show()
```

```
plt.figure(figsize=(14, 6))

plt.plot(df['date'], df['volatility'], color='orange', label='Volatility', linewidth=0.8)

plt.title('BTC Volatility Over Time')

plt.xlabel('Date')

plt.ylabel('Volatility (High - Low / Open)')

plt.legend()

plt.grid(True)

plt.tight_layout()

plt.show()
```

```
plt.figure(figsize=(10, 8))

sns.heatmap(df[['open', 'high', 'low', 'close', 'volatility', 'price_pct_change',
'log_return', 'Volume USD']].corr(),

            annot=True, cmap='coolwarm')

plt.title('Correlation Heatmap')

plt.show()
```

```
plt.figure(figsize=(12, 5))

plt.subplot(1, 2, 1)

sns.histplot(df['volatility'], bins=50, kde=True)

plt.title('Distribution of Volatility')

plt.subplot(1, 2, 2)

sns.histplot(df['price_pct_change'], bins=50, kde=True, color='green')
```

```
plt.title('Distribution of Price % Change')
```

```
plt.tight_layout()
```

```
plt.show()
```

Anomaly Detection

- Isolation Forest

```
features = df[['price_pct_change', 'volatility', 'log_return', 'Volume USD']]
```

```
iso_forest = IsolationForest(contamination=0.01, random_state=42)
```

```
df['anomaly_score'] = iso_forest.fit_predict(features)
```

```
df['anomaly'] = df['anomaly_score'].apply(lambda x: 1 if x == -1 else 0)
```

- K Means

```
kmeans = KMeans(n_clusters=2, random_state=42)
```

```
df['kmeans_cluster'] = kmeans.fit_predict(features)
```

```
if df.groupby('kmeans_cluster')['anomaly'].mean()[0] >
```

```
df.groupby('kmeans_cluster')['anomaly'].mean()[1]:
```

```
df['kmeans_cluster'] = df['kmeans_cluster'].map({0: 1, 1: 0})
```

```
plt.figure(figsize=(14, 6))
```

```
normal = df[df['anomaly'] == 0]
```

```
anomaly = df[df['anomaly'] == 1]
```

```
plt.plot(normal['date'], normal['close'], alpha=0.5, label='Normal')
```

```
plt.scatter(anomaly['date'], anomaly['close'], color='red', label='Anomaly', s=10)
```

```
plt.title('Detected Anomalies in BTC Prices')
```

```
plt.xlabel('Date')
```

```
plt.ylabel('Closing Price')
```

```
plt.legend()
```

```
plt.grid(True)
```

```
plt.show()
```

- Feature selection

```
X = df[['price_pct_change', 'volatility', 'log_return', 'Volume USD', 'hour', 'day',  
       'weekday']]
```

```
y = df['anomaly']
```

- Train Test Split

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,  
                                                    random_state=42)
```

- Random Forest

```
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
```

```
rf_model.fit(X_train, y_train)
```

- XG boost

```
xgb_model = XGBClassifier(n_estimators=50, max_depth=5,  
                          learning_rate=0.1,
```

```
                          use_label_encoder=False, eval_metric='logloss',  
                          random_state=42)
```

```
xgb_model.fit(X_train, y_train)
```

- Confusion Matrix

```
cm_kmeans = confusion_matrix(df['anomaly'], df['kmeans_cluster'])
```

```
cm_iforest = confusion_matrix(df['anomaly'], df['anomaly_score'].map({1: 0, -1: 1}))
```

```
cm_rf = confusion_matrix(y_test, rf_preds)
```

```
cm_xgb = confusion_matrix(y_test, xgb_preds)
```

```
cm_kmeans = confusion_matrix(df['anomaly'], df['kmeans_cluster'])
```

```
cm_iforest = confusion_matrix(df['anomaly'], df['anomaly_score'].map({1: 0, -1: 1}))
```

```
cm_rf = confusion_matrix(y_test, rf_preds)
```

```
cm_xgb = confusion_matrix(y_test, xgb_preds)
```