

A Work Project, presented as part of the requirements for the Award of a Master's degree in  
Finance from the Nova School of Business and Economics.

## ALGORITHMIC TRADING WITH CRYPTOCURRENCIES

-

Does Twitter Sentiment impact short-term price fluctuations in Bitcoin?

Tim Schnülle

Work project carried out under the supervision of:

Leid Zejnilovic

17-12-2021

## **Abstract**

Since its inception in 2009, Bitcoin has gained popularity and importance in financial markets. The Bitcoin price is highly volatile entailing high risk and chances of high returns for traders. This work is part of a work project, which performs a holistic approach to build an intraday Bitcoin trading algorithm based on predictive analysis of Machine Learning models. This part performs a Sentiment Analysis on Twitter data, showing a Granger causal relationship between the extracted Sentiment and the Bitcoin price.

## **Keywords**

Forecasting, Business Analytics, Cryptocurrency, Bitcoin, Social Media influencer, Price Prediction, Algorithmic Trading, Granger causality, Vader, Twitter,

## **Acknowledgements**

I gratefully acknowledge the help from supervisor Leid Zejnilovic as main representative of the Nova Data Science Knowledge Center and Nova SBE faculty.

This work used infrastructure and resources funded by Fundação para a Ciência e a Tecnologia (UID/ECO/00124/2013, UID/ECO/00124/2019 and Social Sciences DataLab, Project 22209), POR Lisboa (LISBOA-01-0145-FEDER-007722 and Social Sciences DataLab, Project 22209) and POR Norte (Social Sciences DataLab, Project 22209).

## List of Abbreviations

ADA	Cardano
ADF	Augmented Dickey Fuller
ARIMA	Auto-Regressive Integrated Moving Average
AVAX	Avalanche
BNB	Binance Coin
BUSD	Binance USD
DL	Deep Learning
DOGE	Dogecoin
DOT	Polkadot
ETH	Ethereum
FI	Feature Importance
GRU	Gated Recurring Unit
LSTM	Long Short-Term Memory
LTC	Litecoin
LUNA	Luna Coin
MC	Multicollinearity
ML	Machine Learning
RF	Random Forest
RMSE	Root mean square error
RNN	Recurrent Neural Network
SMBO	Sequential model-based optimization
SOL	Solana
TPE	Tree-structured Parzen Estimator
UNI	Uniswap

USDC

USD Coin

VADER

Valence Aware Dictionary and Sentiment Reasoner

XGB

Extreme Gradient Boosted trees

XRP

Ripple

# Table of Contents

<b>1</b>	<b>INTRODUCTION</b> .....	<b>2</b>
<b>2</b>	<b>LITERATURE REVIEW</b> .....	<b>6</b>
<b>3</b>	<b>METHODOLOGY</b> .....	<b>8</b>
3.1	PROBLEM STATEMENT .....	8
3.2	DATA .....	11
3.2.1	<i>Data Collection</i> .....	11
3.2.2	<i>Exploratory Data Analysis</i> .....	14
3.2.3	<i>Data Transformation</i> .....	18
3.2.3.1	Data Preprocessing .....	18
3.2.3.2	Feature Transformation .....	19
3.2.3.3	Technical Analysis .....	21
3.2.3.4	Study I: Sentiment Analysis .....	24
3.2.3.4.1	Introduction .....	24
3.2.3.4.2	Related Work .....	25
3.2.3.4.3	Twitter .....	27
3.2.3.4.4	Data .....	27
3.2.3.4.5	Methodology .....	31
3.2.3.4.6	VIP Sentiment .....	34
3.2.3.4.7	Results and Discussion .....	35
3.2.3.4.8	Conclusion and Future Work .....	39
<b>4</b>	<b>RESULTS AND DISCUSSION</b> .....	<b>40</b>
<b>5</b>	<b>CONCLUSION AND FUTURE WORK</b> .....	<b>45</b>
	<b>LIST OF APPENDICES</b> .....	<b>46</b>
	<b>LIST OF TABLES</b> .....	<b>46</b>
	<b>LIST OF FIGURES</b> .....	<b>47</b>
	<b>APPENDIX</b> .....	<b>48</b>
	<b>REFERENCES</b> .....	<b>51</b>

# 1 Introduction

Since its inception in 2009, Bitcoin has gained popularity and importance in the international financial landscape, attracting media coverage, the attention of regulators, government institutions, investors, academia, and the public (Sebastião and Godinho 2021). Following Bitcoin, other cryptocurrencies were introduced over the past decade. As of 14<sup>th</sup> of October 2021, there are 6,590 different cryptocurrencies on the market, amounting to a market capitalization of around \$2.4 trillion for the entire cryptocurrency market (CoinMarketCap 2021). The opportunities to own and trade cryptocurrencies have increased significantly in recent years. With the rise of online wallet companies, trading is made easier and accessible to the public, which is reflected in higher trading volume and an increase in the number of wallets (Blockchain.com 2021).

On 14<sup>th</sup> of October 2021, the 24-hour trading volume of the entire cryptocurrency market amounted to \$92 billion (CoinMarketCap 2021). On that day, the trading volume of Bitcoin was \$43 billion compared a trading volume of \$10 billion for Apple and \$10 billion for Tesla as commonly known stocks (Wall Street Journal 2021). Bitcoin is the most relevant cryptocurrency on the market with a market capitalization of around \$1 trillion (October 14<sup>th</sup>), accounting for around 46% (followed by Ethereum accounting for around 19%) of the total market capitalization of the cryptocurrency market (CoinMarketCap 2021). In addition to its high market capitalization and trading volume, the cryptocurrency market is characterized by high price fluctuations, i.e., high volatility. High volatility results in high risk and return because volatility is considered as an alternative measure for risk and risk has a positive and significant relation to returns (Bali and Lin, 2006). Cryptocurrencies do not follow the development of major financial asset classes but are driven by behavioral factors like for example herding factors where traders follow other people instead of relying on their own analysis (Sebastião and Godinho 2021). Machine Learning (ML) algorithms discover patterns and drivers for the

financial development of an asset, enabling to develop a model that predicts future price movements, and generates returns superior to its benchmark if executed in the market (Tao, et al. 2021). Prior research has been conducted to analyze the applicability of different ML algorithms to predict the development of cryptocurrencies. We identified 73 papers that discuss the prediction of cryptocurrency prices (cf. Table 1). 63 of these papers (86%) analyze data from 2019 and previous years. 64 papers (88%) translate the prediction problem into either a Regression (42 paper, 58%) or Classification (22 paper, 30%) analysis. 53 papers (73%) use either Statistical or ML algorithms but do not compare both. The Feature Selection is characterized by endogen cryptocurrency features (69 paper, 95%). 21 papers (29%) consider a trading strategy to evaluate the model.

Due to recent developments in the Bitcoin market, patterns of the Bitcoin movement have changed. Before 2019, the highest trading volume (\$120 billion) per week was accorded in the first calendar week in 2018. After 2019, the week with the highest trading volume grew by 538% to a total of \$765 billion in calendar week eight in 2021 (Finance 2021). The state of research is limited as 63 papers do not incorporate data from 2019 on and dismiss the current pattern in the development of Bitcoin. Algorithms need to be trained with recent data. Research on the implementation of a real-time trading algorithm for Bitcoin is not covered by any paper. Limited research has been conducted regarding a holistic approach for the development of a trading algorithm, including both Regression and Classification problems, comparing several algorithms, considering endogen (Supply & Demand) as well as exogen features (Crypto market, Macro Financial, Political and Sentiment) and including a trading strategy for final evaluation. A holistic approach for algorithmic trading brings scientific novelty and can discover new insights for data-driven trading. Based on the identified gaps in the current state of research this work seeks to answer three major questions: (Study I) Does Twitter Sentiment impact short-term price fluctuations in Bitcoin? (Study II) What is the optimal modelling design

for Bitcoin price and trend prediction? (Study III) How to translate multiple model predictions into an algorithmic trading strategy?

The main objective of this work is to build a Bitcoin trading algorithm based on the predictive analysis of ML models. Past research is analyzed and aggregated to develop a holistic approach, from Data Collection, Feature Engineering, Feature Selection, Model Implementation and Model Selection to the definition of a trading strategy. Using a simulation setting for real-time trading during a test period, the final evaluation is conducted on economic performance measures of trading strategies that combine multiple model predictions. The results are compared to benchmark strategies. The findings indicate that a trading algorithm derived from ML model predictions is able to generate positive returns and to outperform its benchmark strategies. Ensemble trading strategies that combine predictions of multiple Long Short-Term memory (LSTM) Regression models have the highest overall performance.

This work is organized in 5 major sections. In section 2, we review the related work. In section 3 we outline the methodology of this work and describe Problem Definition, Data, Modelling, and Trading Strategy. The individual studies (Study I, Study II, Study III) mentioned above represent self-contained analyses and are included in this section. In section 4 we report and discuss the results. Section 5 concludes this work and gives an outlook for future research opportunities.



**Table 1: Coverage of major topics in the defined 73 papers**

	<b>Details</b>	<b>Number</b>
Observation period	till 2019	63
	2019 - now	10
	Total	73
Prediction	Regression	42
	Classification	22
	Both	9
	Total	73
Algorithm	Statistical	20
	ML	33
	Both	20
	Total	73
Features	Supply & Demand	69
	Crypto market	0
	Macro-financial	12
	Political	6
	Sentiment	11
	Total (Higher due to duplication)	98
Trade strategy	Yes	21
	No	52
	Total	73

## 2 Literature Review

Systematic search of the literature ensures qualitative scientific work (Timmins and Mccabe 2005). We applied a forward and backward search process to identify relevant papers and used filter criteria to ensure a state-of-the-art literature base, as shown in figure 1.



Figure 1: Literature research approach

We used the EBSCO library, a collection of scientific databases, for our initial research search. EBSCO is a leading provider of research databases and includes paper, e-journals, magazines, and e-books (Williams and Foster 2011). EBSCO displays the peer review status of papers to ensure academic scientific quality. During forward search, we used the keywords in table 2 to find papers focusing on Bitcoin and cryptocurrency prediction or volatility.

Table 2: Combinations of the keyword search query

Keyword 1	Keyword 2
Bitcoin	Prediction
Bitcoin	Forecasting
Bitcoin	Volatility
Cryptocurrencies	Prediction
Cryptocurrencies	Forecasting
Cryptocurrencies	Volatility

We gathered 61 papers during the first step. We selected papers that focus on prediction or forecasting in the second step to reduce the literature base to 23 papers. We applied backward search, scanning related work for additional relevant paper. The scope totaled to 73 papers. We filtered the papers by content for ML algorithms and only included paper that were published in 2019 or later to ensure state-of-the-art. The remaining five papers represent the focus papers of our work, shown in table 3.

The focus papers provide guidance for our work in presenting the latest research results as well as represent a benchmark to compare our work. All papers that are included in this work are peer-reviewed to ensure high academic quality.

**Table 3: Overview of focus paper**

<b>Author (Year)</b>	<b>Prediction</b>	<b>Forecast</b>	<b>Trade strategy</b>	<b>Supply &amp; Demand</b>	<b>Crypto market</b>	<b>Macro-financial</b>	<b>Political</b>	<b>Sentiment</b>
Chen, Li and Sun (2020)	Classification	5 minutes & 1 day	N/A	x	N/A	x	N/A	x
Cocco, Tonelli and Marchesi (2021)	Regression	1 day	N/A	x	N/A	N/A	N/A	N/A
Dutta, Kumar and Basu (2020)	Regression	1 day	x	x	N/A	x	x	x
Mudassir, et al. (2020)	Both	1 day, 1 week, 1 month	N/A	x	N/A	N/A	N/A	N/A
Sebastião and Godinho (2021)	Both	1 day	x	x	x	x	N/A	N/A

Focus papers are categorized according to the ML problem they are analyzing into Regression, Classification, and a combination of both learning problems. Prediction in general is defined as estimating the output for unseen data. Forecasting is a part of prediction and is concerned with time-series data (Matsuo 2003). In this work, we use the general wording “prediction”. The focus papers leverage different features which can be aggregated into five feature categories: Supply & Demand, Crypto market, Macro Financial, Political and Sentiment. The features categories will be explained in depth in section 2.2.

Chen, Li and Sun (2020) use high dimensional features of Supply & Demand, Macro-financial and Sentiment on a five-minute interval basis to predict the Bitcoin price trend in five minutes and for the next day. They highlight the importance of sample granularity and feature dimensions on ML model performance (Chen, Li and Sun 2020). Cocco, Tonelli and Marchesi (2021) and Dutta, Kumar and Basu (2020) predict the daily closing Bitcoin price. Cocco,

Tonelli and Marchesi (2021) compare several ML frameworks to predict the prices of Bitcoin and Ethereum. They use five technical indicators that are calculated from the cryptocurrency price and provide insights how to build efficient trading frameworks (Cocco, Tonelli and Marchesi 2021). Dutta, Kumar and Basu (2020) investigate Feature Engineering of twenty features from Supply & Demand, Crypto Market, Macro-financial, Political and Sentiment for ML algorithms. They implement a simple trading strategy and demonstrate the possibility of financial gain through algorithmic cryptocurrency trading (Dutta, Kumar and Basu 2020). Mudassir, et al. (2020) and Sebastião and Godinho (2021) consider a Regression and Classification problem, predicting price and price trend. Mudassir, et al. (2020) predict Bitcoin volatility on a daily, weekly, and monthly base. They use 700 features based on technical indicators and show that it is possible to predict the daily Bitcoin price with low error rates (Mudassir, et al. 2020). Sebastião and Godinho (2020) predict the daily price of Bitcoin, Ethereum and Litecoin and implement trading strategies. They consider Supply & Demand, Crypto Market and Macro Financial features and find that ML is a good technique to predict cryptocurrency prices and price trends, enabling profitable algorithmic trading of cryptocurrencies.

None of the identified papers combines Regression and Classification, the usage of all five feature categories (Supply & Demand, Crypto market, Macro Financial, Political and Sentiment) and the implementation of a trading strategy.

### **3 Methodology**

#### **3.1 Problem Statement**

The goal of this work is to develop a trading algorithm that generates superior returns to its benchmarks with automated trading of the cryptocurrency Bitcoin. Financial time-series data is challenging to analyze due to the dynamic, non-linear, non-stationary, highly volatile, and chaotic nature of financial markets. ML algorithms can be used to analyze large amounts of

seemingly chaotic data, to discover patterns in the data and to predict future data. Additionally, an automated trading bot can react much faster to developments in the market than any human (Borges and Neves 2020). To build a trading algorithm that is based on a ML, we must convert the problem of profitable trading into a ML problem defining an output that can be generated by a ML model. In this paper, we conduct price and trend prediction. Price prediction represents a Regression problem and trend prediction a Classification problem. All learning problems represent a Supervised Learning Problem because the target variable, i.e., the Bitcoin price or the trend calculated from the Bitcoin price, is given, and can be tested against.

To evaluate our trading algorithm, we create a simulation design that is aligned to the prerequisites a real-time trading algorithm requires. The architecture of our work is represented in figure 2. Data is collected through API connections. The features are categorized in five feature categories: Supply & Demand, Cryptocurrency Market, Political, Macro Financial and Sentiment. The collected data is pre-processed, and Feature Transformation, Technical Analysis and Sentiment Analysis are applied to enrich the data and build the final dataset for modelling purposes. We implement Regression and Classification algorithms. Finally, trading strategies are developed that translate model outputs into trading actions.

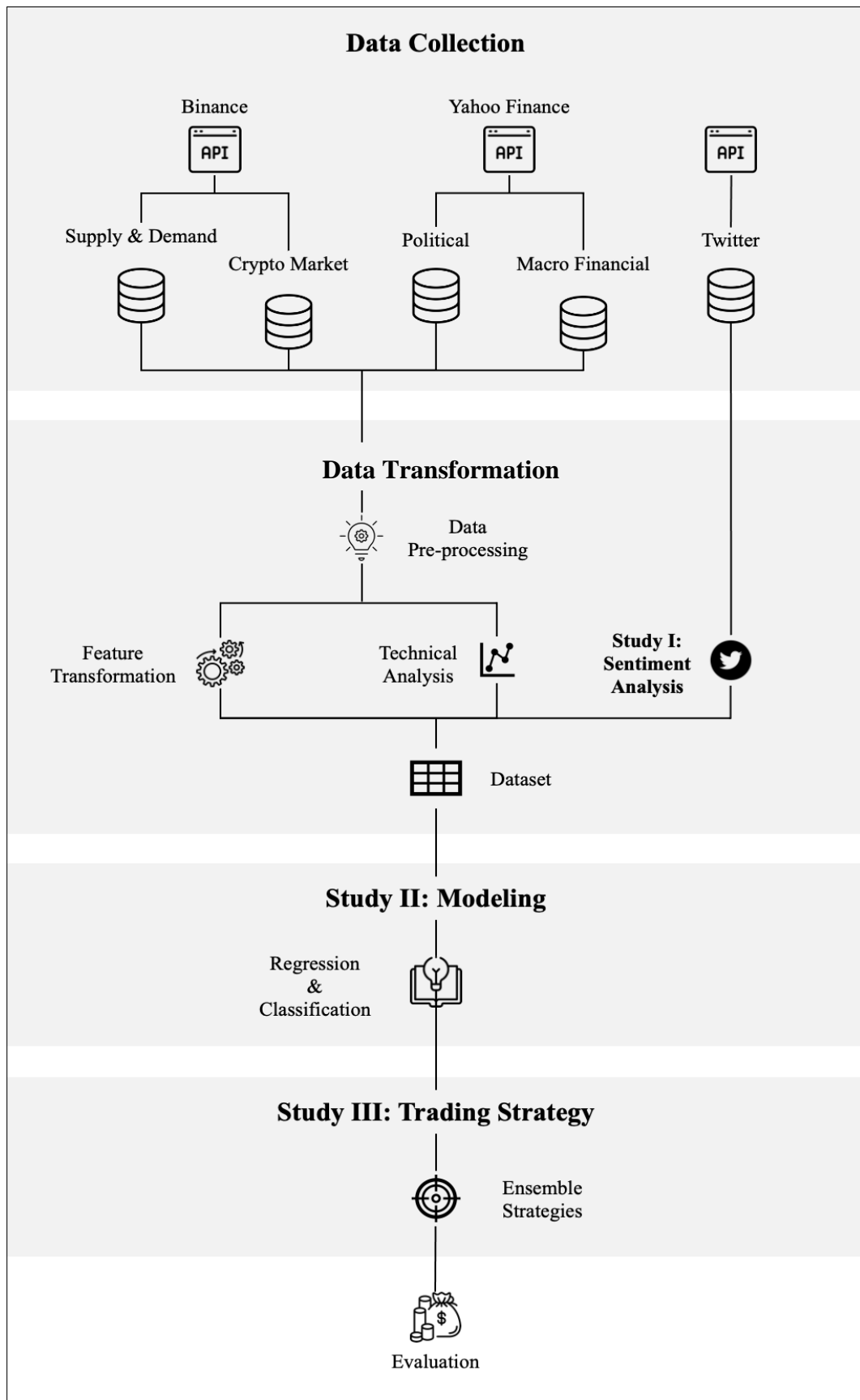


Figure 2: Architecture of the work project

We use Python 3.9 for this work. The training, validation, and evaluation of our ML models was executed through cloud computing provided by Genesis Cloud. We used two NVIDIA GPU GEDorce GTX 1080Ti cores. We used the computing power on demand for a two-week period. The main used Pytoch libraries are PyTorch 1.6 for DL algorithms, scikit-learn 1.0.1 for ML algorithms and Optuna 2.10.0 for Hyperparameter Tuning. The host operating system was Linux.

## **3.2 Data**

### **3.2.1 Data Collection**

A solid and valid data base is the prerequisite for any data analysis and the cornerstone of this project (Gupta, et al. 2021). Data Collection is a challenging process for algorithmic trading. Past data is needed to train the algorithm and real-time data must constantly be fed into the algorithm follow market fluctuations and adjust predictions accordingly. Data Collection requires time, restricting the selection of data sources.

We performed a detailed analysis of the Data Collection process introduced in our focus paper and identified two major characteristics. First, data is collected with different techniques and in different formats and second, features collected for the prediction of cryptocurrencies differ between papers. We observe three different techniques for the Data Collection process: Data retrieval from external files, e.g., CSV (Ahmed and Mafrachi, 2021), data retrieval using web scraping (Kim, et al. 2021), and Application Programming Interface (API) (Chen, Li and Sun 2020). Feature selection is performed differently in terms of categorization and number of features. Several studies examine for example the influence of S&P 500, Gold, other Cryptocurrencies and Sentiment on Bitcoin price fluctuations (Bouri, et al. 2017, Abraham, et al. 2018, Mallqui and Fernandes 2019). Factors influencing the Bitcoin price can be categorized in endogen and exogen features (Bouri, et al. 2017). We follow the approach of Sovbetov (2018) as the work provides the most comprehensive collection of factors influencing Bitcoin

fluctuations. Sovbetov (2018) divides the features for cryptocurrency predictions into four categories: Supply & Demand, Cryptocurrency Market, Political and Macro Financial. In our work, we will extend the collection of Sovbetov (2018) by a fifth category: Sentiment. Considering the increasing importance of social media in recent years, prior research investigates a causal relationship between Online Sentiment and Bitcoin price fluctuations (Kraaijeveld and De Smedt 2020, Pano and Kashef 2020). Table 4 summarizes the five categories and gives examples for each. Endogen features are connected to supply and demand of cryptocurrencies. Exogen features are not directly connected to the observed cryptocurrency but measurements of other influencing factors.

**Table 4: Overview of feature categories**

<b>Categories</b>	<b>Influence</b>	<b>Example Features</b>
Supply & Demand	Endogen	Exponential Moving Average etc.
Cryptocurrency Market	Exogen	Ethereum, Solana, Cardano, Dogecoin etc.
Political	Exogen	CBOE Volatility
Macro Financial	Exogen	S&P 500, CAC40, DAX40, Nikkei 225 etc.
Sentiment	Exogen	Twitter

Within the focus papers (cf. table 3) Data Collection is mostly performed for daily data, only one paper analyzes intraday data. The access to past intraday data is limited compared to the access of daily data for most data sources. From the Yahoo Finance API, past intraday data can only be retrieved for a maximum period of 60 days (Aroussi 2021). There are data sources which can provide data for a longer period, which will imply further costs. Therefore, we decided to use the Yahoo Finance API. Four focus paper analyze data from 2019. No paper is based on data from 2021.



**Table 5: Overview of Data Collection in the focus paper**

<b>Authors and Year</b>	<b>End Time</b>	<b>Frequency</b>	<b>Sources</b>	<b>Observations</b>
Chen, Li and Sun (2020)	02.2019	5 minutes	1	50.000
Cocco Tonelli and Marchesi (2021)	04.2020	Daily	1	1.216
Dutta, Kumar and Basu (2020)	06.2019	Daily	9	3.469
Mudassir et al. (2020)	12.2019	Daily	N/A	N/A
Sebastiano and Godinho (2020)	03.2019	Daily	2	1.297

Binance provides past intraday data for cryptocurrencies nearly without limitations since the opening of the trading platform in 2017 (Binance 2021). The collection of Twitter data is different and depends on the arranged tweet limit. During our academic research the limit was set to 10.000.000 Tweets (Twitter, Developer Platform 2021a).

In our analysis, we exclusively use data sources that offer an API connection. APIs have an advantage over the other two Data Collection techniques because data is retrieved in a concise time. Downloading and reading CSV files or web scraping consumes more time. Uploading multiple files and scraping multiple websites increases the numbers of sources that need to be monitored, increasing the risk for changes in data formats which would interrupt the automated trading algorithm. Aiming to minimize the risk of unwanted changes in data formats, we use a minimum number of APIs that provide high data quality and ensure a data base that includes features from all five feature categories.

To determine the endogenous factors about the Bitcoin Price, the following features in the corresponding time interval were extracted. To examine the influence of other cryptocurrencies on the Bitcoin price, the 15 of the largest following Cryptocurrencies, based on market capitalization in October (coinmarketcap 2021) were included in the dataset. Along with the Bitcoin movement, 84,7% (coinmarketcap 2021) of the whole market capitalization is being tracked and analyzed.

As there is a relationship between Macro Financial Movement and the Price Development of Bitcoin (Walther 2019). The ten most important countries sorted by GDP were selected for further analysis, for each country the primary equity index and the currency for the respective country were extracted (Silver 2020). For countries with the same currency or a currency that is already included in the dataset the value has been skipped.

Furthermore, the most actively traded commodities according to Futures Industry Association were also added to the analysis (FIA 2021).

Table 6 gives an overview of the sum of the total number of features extracted. For example, if a feature is extracted like “Ethereum” it will be counted as one feature which will contain some more sub features like close price and trading volume. The Cryptocurrency Market (13 Features) and Macro-Financial (28 Features) category contain the most features. For Supply & Demand we extracted 12 features through technical analysis.

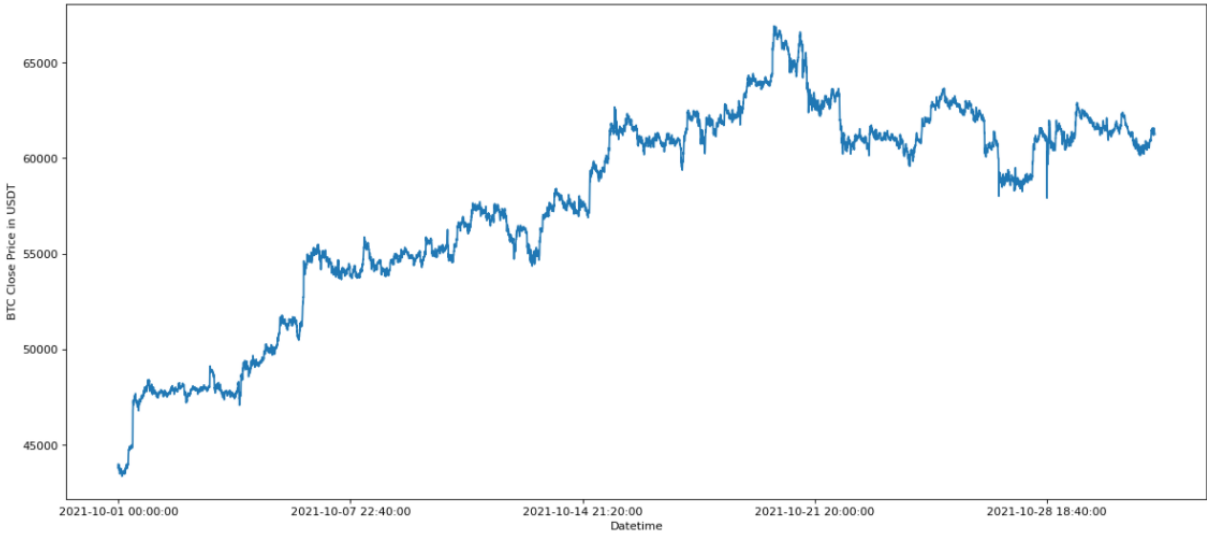
**Table 6: Overview of collected features**

<b>Feature Category</b>	<b>Feature</b>	<b>Sub Features</b>	<b>Engineered Features</b>
Supply and Demand	1	5	12
Cryptocurrency Market	13	78	0
Political	1	1	0
Macro-Financial	28	28	7
Sentiment	1	2	0
	<b>Total</b>	115	19

### 3.2.2 Exploratory Data Analysis

Figure 3 provides an overview of the observed data and gives a detailed picture of the development of the target variable. The data analyzed in this paper contains 8,929 observations, representing a time-period of 31 days from the 01.10.2021 00:05 to the 31.10.2021 23:55. The target variable is the closing price of Bitcoin measured in Tether (BTCUSDT), visualized in figure 3. Tether (USDT) is a cryptocurrency whose value is linked to the US dollar and called

a stable coin. At the start of the observation period the price of Bitcoin is 43,981 USDT and reaches a price of 61,243 USDT at the last observation. For the analyzed period the average Bitcoin price was 57,653 USDT, while the median was 59,516 USDT. The price was subject to strong fluctuations and ranged from 43,361 USDT as a minimum to a highest price of 66,908 USDT within the period. The standard deviation was 5,285.



*Figure 3: Development of Bitcoin price*

We included 115 features which are divided in 5 categories as described in the data collection part. The current memory usage is 16,1 MB. In addition to the collected data, we add 19 features that are calculated using the collected data and which will be described in Feature Engineering part. An overview of the features is provided in Appendix 1.

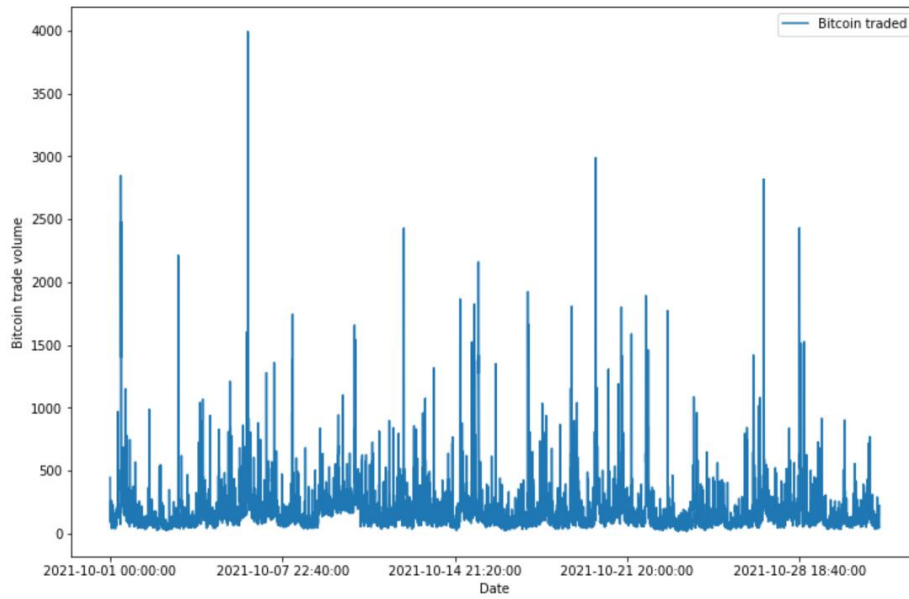


Figure 4: Bitcoin trade volume in 5-minute intervals

Figure 4 shows the trading volume of Bitcoin within a 5-minute time interval. On average, 175 Bitcoins are traded every 5 minutes on Binance, making the Binance platform a liquid trading venue.

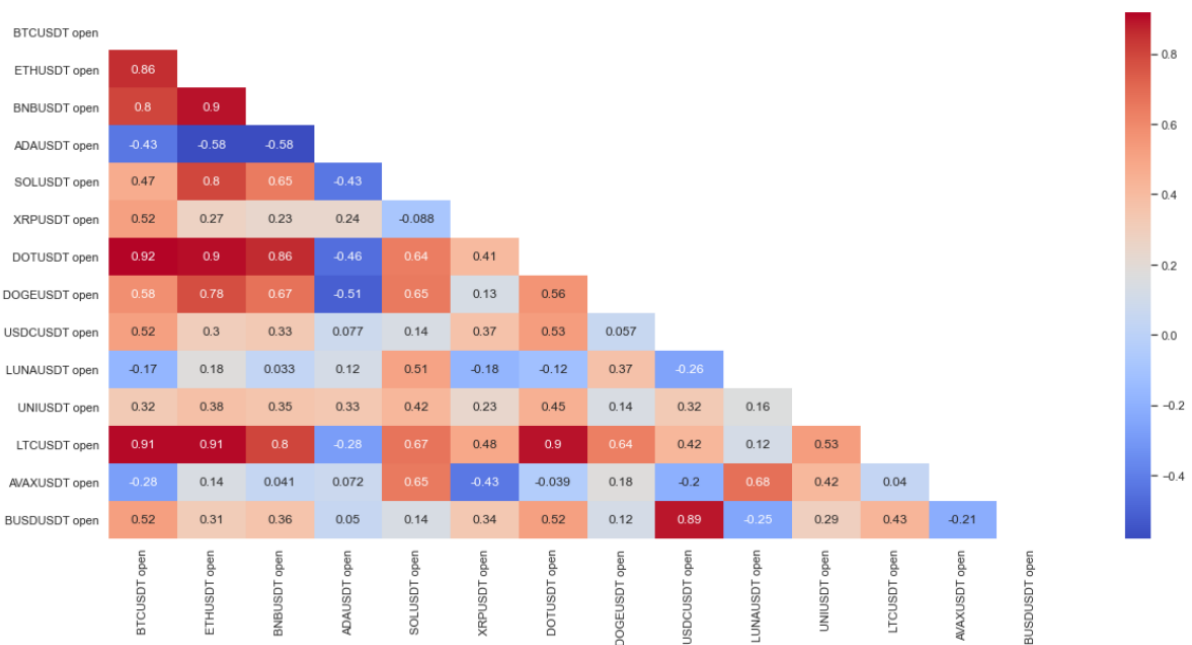


Figure 5: Correlation matrix of cryptocurrency market features

Figure 5 shows a correlation analysis of the cryptocurrencies observed in this work. The correlations range from -0.51 to 0.91. The differing correlations between cryptocurrencies

indicate that the cryptocurrency market is not always moving in the same direction. The highest positive correlation with Bitcoin has Polkadot (DOT), i.e., 0.92, and the highest negative correlation Cardano (ADA), i.e., -0.43. In figure 6 the normalized price developments of the observed cryptocurrencies are visualized. The different developments support the findings of the correlation matrix. DOT experiences the highest increase in the observed period, ADA experiences the worst development.



*Figure 6: Normalized price development of observed cryptocurrencies*

Cryptocurrency trading is not limited to opening hours of stock exchanges but is possible 24 hours, seven days a week. Data for cryptocurrencies exists for every 5-minute timestamp of the observed period as shown in figure 6. Data of the commodity market is received for opening hours of the market. This is visualized in figure 7.



*Figure 7: Normalized price development of commodity and Bitcoin prices*

### **3.2.3 Data Transformation**

#### **3.2.3.1 Data Preprocessing**

Different sources were used to retrieve data and the data was merged to build a comprehensive base for analysis and modelling purposes. As described in the previous section, we observe missing values in 98 of our 115 input features. As we are dealing with time-series data from the cryptocurrency as well as the general stock market data, missing values in our data have specific characteristics. While cryptocurrencies can be traded 24 hours a day, seven days a week, the trading of stocks, commodities and other securities is bound to opening hours of stock market exchange providers (Dutta, Kumar and Basu 2020). The number of missing values can differ between features, because of after-hours trading and differences in opening hours between Stock Exchanges. After-hours trading occurs after regular market hours. Due to after-hours trading and after-hours volatility, the opening price for a stock on the following day can differ quite extensively from the price at which it closed the previous day (Barclay and Hendershott

2003). Missing values need to be accounted for by deleting respective observations or features, or by imputation, because most of the existing ML algorithms don't work well with missing values.

Other researchers that have used general stock market data to predict cryptocurrency prices have used imputation methods to fill missing values. One of the simplest imputation methods that has been widely used is the Last Observation Carried Forward (LOCF) time-series imputation method (Vo and Yost-Bremm 2018). Therefore, it is assumed that stock, commodity, and other security prices do not change after closing hours, i.e., after-market trading is ignored (Dutta, Kumar and Basu 2020).

We use the LOCV imputation method in combination with Next Observation Carried Backward (NOCB) method to account for missing values that arise within time-horizons when the general stock market is closed. As for the LOCV, missing values are imputed as the previously observed value, i.e., the last observation is carried forward. In case there is no previous value the NOCV method is used. Thus, the follow-up value is used to impute the previous value. The combination of the observed and imputed data is then analyzed as there were no missing data.

### **3.2.3.2 Feature Transformation**

The architecture of our model requires a transformation of our target variable, i.e., Bitcoin price. Sebastião and Godinho (2021) found out that model assembling enables profitable trading strategies. We formulate the analysis as a Regression and Classification problem with three different prediction horizons  $ph = \{ \text{"one hour"}: 1h, \text{"two hours"}: 2h, \text{"three hours"}: 3h \}$ . The horizons are selected, aiming to take advantage of intraday trading and avoid exaggerated transaction costs. The prediction is evaluated each five minutes as it is the most granular level, we can collect a comprehensive number of features. In addition to a Regression problem, we are dealing with a Classification problem, and we need to transform the Bitcoin price into a categorical variable, i.e., trend variable. For the Classification problem, we create a categorical

variable that is equal to 0, 1 or 2. The variable is set equal to 0 if the Bitcoin price is decreasing and set to 2 if the Bitcoin price is increasing. In all other cases the variable is set to 1. We use a threshold of 0.5% to calculate the trend variable for the Classification problem. The distribution of the trend variables for each time lag is shown in figure 8. For the Regression problem, we create lagged variables for all prediction horizons 1h, 2h, and 3h. Each lagged variable represents the target variable for one prediction horizon and is used to validate and test the respective model.

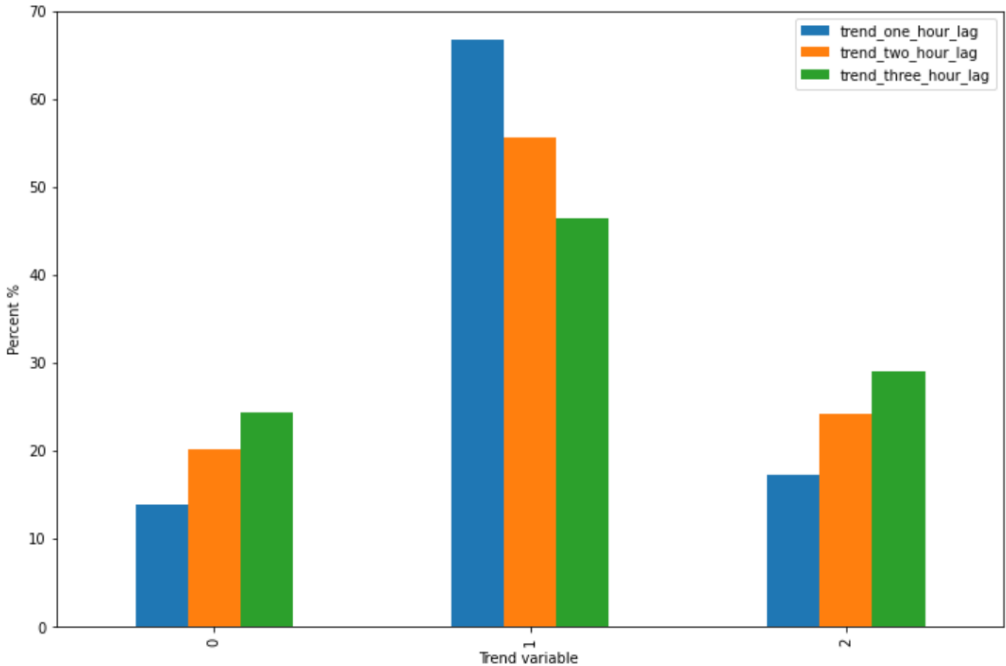


Figure 8: Distribution of trend variables

Aiming to improve modelling performance we add new features by transforming existing data. Lagged variables are common features to be included for predictive analysis of cryptocurrency prices and price trends (Sebastião and Godinho 2021). We include past lagged trend variables with time lags equal to those of the trend target variables, i.e., one, two and three hours. They show the trend of the Bitcoin price compared to the previous instance according to the defined time lags.

The day-of-the-week effect represents a well-known phenomenon in the study of financial markets where differing returns between days of the week are observed in a persistent way.



Indications for this anomaly are observed for many products on the financial market (Aharona and Qadan 2019). The price fluctuations of cryptocurrencies, especially Bitcoin, seem to depend on the day of the week (Sebastião and Godinho 2021). In this paper, we created daily dummy variables for each weekday to capture effects that are related to certain days of the week.

### **3.2.3.3 Technical Analysis**

Technical analysis is the study of historical prices and price movements in the market to get an estimation of the price or its trend in the future (Borges and Neves 2020). Technical analysis intends to identify specific rules like price trends, market cycles, momentum, volatility, or price chart patterns, under the assumption that prices move in trends and historic movements repeat themselves (Huang, Huang and Ni 2019). Extensive research regarding the impact of technical analysis of stocks has been conducted, constituting a high importance of technical features on future price predictions and trends (Fang, et al. 2020). For the cryptocurrency market, prior researchers have also used technical analysis for price prediction and have also concluded that technical features are an important factor to predict price movements (Kristjanpoller and Minutolo 2018) (Nakano, Takahashi and Takahashi 2018) (Huang, Huang and Ni 2019), (Abbad, Fardousi and Abbad 2014).

A vast amount of different technical indicators has been used in prior research with the intention to improve prediction of future price movements. Technical indicators differ in their purpose and can be divided in different categories like overlap study indicators, momentum indicators, cycle indicators, volatility indicators, and pattern recognition indicators. Prior research on the predictability of Bitcoin prices using a large set of 124 technical indicators has been conducted. (Huang, Huang and Ni 2019). Other researchers focus on a small number of technical indicators that are widely accepted. The most represented technical indicators that we identified in our research are Exponential Moving Average (EMA), Moving Average Convergence–Divergence

(MACD), Relative Strength Index (RSI), On Balance Volume (OBV) (Borges and Neves 2020) (Vo and Yost-Bremm 2018) (Nakano, Takahashi and Takahashi 2018).

In this paper, we calculated and included the following technical indicators: Exponential Moving Average (EMA), Moving Average Convergence–Divergence (MACD), Relative Strength Index (RSI), On Balance Volume (OBV) and Stochastic Oscillator. These technical indicators have comprehensively been used in prior research and are widely accepted by traders on the market. All indicators depend only on past Bitcoin prices. In the following paragraphs, we will provide more detailed information about the technical indicators that are used in this paper. Further elaborations and explanations of technical analysis and each technical feature can be found in (Murphy 1999).

### **Exponential Moving Average (EMA)**

A moving average (MA) is a technical indicator that helps to smooth out the price data by dampen the effects of short-term oscillations, through a constantly updated average price. The MA is a trend following indicator that reacts to the market by announcing a trend that has already begun. An EMA is a variation of the MA that assigns more weight and significance the most recent data points, having the ability to react faster to recent price variations (Borges and Neves 2020). As this paper intends to analyze short-term predictions of the highly volatile Bitcoin price, we use the EMA. The EMA is calculated using the following equation:

$$EMA_t = EMA_{t-1} + \frac{\text{smoothing factor}}{n+1} + [Price_t - EMA_{t-1}] \quad (1)$$

In equation (1),  $t$  refers to the current period,  $n$  refers to the number of time periods the EMA is calculated on, and the *smoothing factor* represents a smoothing parameter that is set to the most common value of two for all calculations. For this study, we used 6 different time periods corresponding to  $n = \{12, 24, 48, 96, 288, 576\}$ , representing time periods of one, two, four and eight hours, as well as one and two days. The first  $n$  values of EMA are set to an initial average of the first  $n$  time periods for each of the time periods, respectively.

### **Moving Average Convergence–Divergence (MACD)**

The MACD is calculated using the difference between two trend following indicators, EMAs, of different time periods. As a trend-following momentum indicator it combines the purpose of trend-following and momentum (Borges and Neves 2020). The MACD is popular among traders thanks to its simplicity and effectiveness. It shows how the two EMAs converge and diverge and helps to understand whether the bullish or bearish movement in the price is increasing or decreasing (Vo and Yost-Bremm 2018). Traditionally, a 26- period EMA is subtracted from a 12-period EMA to calculate the MACD (Borges and Neves 2020). We use a 24- period EMA and a 12-period EMA representing a period of  $1h$  and  $2h$ , respectively. The equation to calculate the MACD is the following:

$$MACD_t = EMA_t \times 12 - EMA_t \times 24 \quad (2)$$

In equation (2),  $t$  refers to the current period and  $12$  and  $24$  refer to the EMA with the respective period  $n$  at time  $t$ .

### **Relative Strength Index (RSI)**

The RSI measures the magnitude of recent price changes and is used to identify general price trends (Vo and Yost-Bremm 2018). It is a momentum oscillator that is used to evaluate whether a market is overbought or oversold. It represents a line that moves between two extremes and can take a value between 0 and 100 (Borges and Neves 2020). The RSI is calculated using the following equation:

$$RSI_t = 100 - \frac{100}{(1 + [\textit{average gain}_{t-14} / \textit{average loss}_{t-14,t}])} \quad (3)$$

In equation (3),  $t$  refers to the current period. *Average gain* and *average loss* are calculated using the gains and losses of the past 14 observations, while losses are set to zero to calculate the *average gain* and gains are set to zero to calculate the *average loss*.

## On Balance Volume (OBV)

While the previous indicators utilize prices and price movements, OBV is a technical momentum indicator that focuses on volume flow to predict price changes. OBV is built on the idea that volume movement precedes price movement and is a key factor behind markets. An increase of OBV signals a price move up while a decrease of OBV signals a decrease (Vo and Yost-Bremm 2018). The RSI is calculated using the following equation:

$$OBV_t = \begin{cases} OBV_{p-1} + Volume_p, & \text{if } Price_t > Price_{t-1} \\ OBV_{p-1} - Volume_p, & \text{if } Price_t < Price_{t-1} \\ OBV_{p-1}, & \text{if } Price_t = Price_{t-1} \end{cases} \quad (4)$$

In equation (4),  $t$  refers to the current period and  $Volume$  refers to the amount of trading volume in the past 5 minutes prior to  $t$ .

### 3.2.3.4 Study I: Sentiment Analysis

#### 3.2.3.4.1 Introduction

Bitcoin is a speculative asset and its price is highly volatile (Valencia, Gómez-Espinosa and Valdés-Aguirre 2019). Trading Bitcoin is associated with high risk but also the chance for extraordinary profits. As a decentralized cryptocurrency, Bitcoin is trading 24 hours, 7 days a week, which makes it an exciting target for price speculations and predictions (Kraaijeveld and De Smedt 2020). Bitcoin does not have an inherent value and is not backed by any government (Bugár and Somogyvári 2020). Being detached from the characteristics of traditional assets, the value drivers of Bitcoin are continuously discussed and researched (Woebeking 2021). The long-term development of Bitcoin represents a major discussion, where opinions range from Bitcoin being a bubble to being a solid investment (Valencia, Gómez-Espinosa and Valdés-Aguirre 2019). Therefore, the Bitcoin price is driven by behavioral factors where people do not follow their own analysis but, e.g., the opinion of a majority (Sebastião and Godinho 2021). Social media platforms are predominantly used to exchange information on Bitcoin and social

media sentiment proved to have an influence on cryptocurrency prices (Kraaijeveld and De Smedt 2020). Additionally, statements from influencers like Elon Musk have a significant influence on the Bitcoin price, leading to abnormal returns (Ante 2021).

The integration of high frequency features derived from Sentiment Analysis in algorithmic trading is barely represented in past academic literature. The importance of sentiment for algorithmic trading of Bitcoin is analyzed by answering the following research question:

*Does Twitter Sentiment impact short-term price fluctuations in Bitcoin?*

In this work project, I perform a Sentiment Analysis for intraday sentiment extracted from Twitter data and investigate the impact on short-term price fluctuations of Bitcoin. I define a group of high influential accounts, i.e., VIP accounts, and analyze the difference between the impact of *Overall Sentiment* and *VIP Sentiment* on the Bitcoin price.

#### 3.2.3.4.2 Related Work

Previous research observed the relationship between Twitter sentiment and the development of cryptocurrency prices. Kraaijeveld and Smedt (2020) observe *tweets* on a daily interval and find evidence that the observed sentiment has a Granger causal relation to Bitcoin returns. They formulate the hypothesis that an intraday analysis of Twitter sentiment further improves the results. The hypothesis is supported by findings of Pano and Kashef (2020), who identify a higher correlation for shorter time horizons if an optimal text pre-processing strategy is applied. Text pre-processing is a technique to reduce noise and is commonly used for Sentiment Analysis, especially for social media content. Multiple pre-processing methods for *tweets* are observed and compared by Pano and Kashef (2020) to extract the correct sentiment and categorize *tweets* in positive and negative. Sentiment can be extracted by two approaches. The first approach is lexicon-based and does not require labelled data and rely in predefined lexicons (Turner, Labille and Gauch 2021). The second approach includes ML and requires manually labelled data to create a training set. Where humans have to manually categorize the sentiment

of a fraction of the *tweets* (Pant, et al. 2018). The VADER lexicon-based approach is the most popular method for Sentiment Analysis in social media, as described by Kraaijeveld and De Smedt (2020) and Stenqvist and Lönnö (2017). Dr. Rajab and Jahjah (2020) find a positive correlation between Twitter sentiment and the Bitcoin price of the next day. To further analyze a correlation for a “causal” relationship, the Granger causality test is a common procedure. The Granger causality test is used in prior research to investigate the relationship between stock data and investor sentiment (Chu, Wu and Qiu 2015) and the relationship between cryptocurrency data and social media sentiment (Kraaijeveld and De Smedt 2020) (Dastgir, et al. 2019).

The findings of Pant, et al. (2018) indicate that there is also a relationship of specific Twitter cryptocurrency news accounts and the Bitcoin price.

This work project performs an intraday Sentiment Analysis of tweets and investigates a Granger causal relationship of the observed sentiment and the Bitcoin price. Furthermore the observed news accounts from Pant, et al. (2018) will be extended and also tested for a Granger causal relationship.

After an introduction of Twitter, the project is divided into five different steps: (1) Data Collection, (2) Data Pre-processing, (3) Sentiment Analysis, (4) Filtering and a (5) Granger causality test. Further visualized in Figure 9:

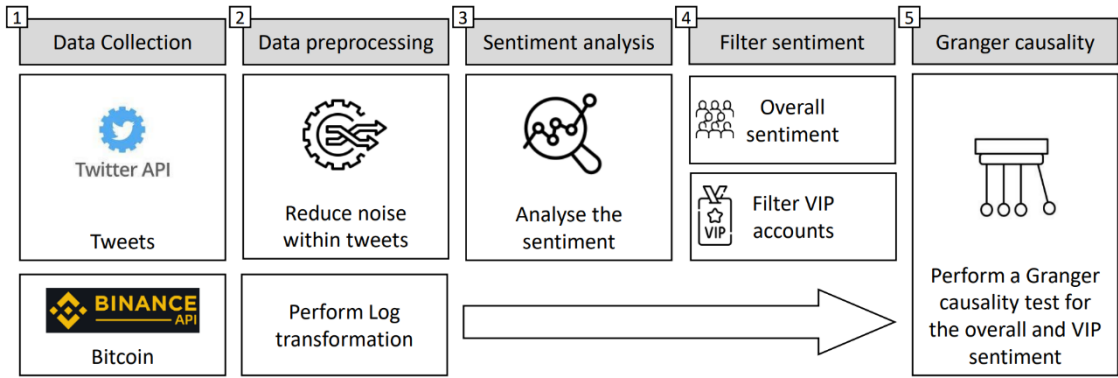


Figure 9: Structure of the Sentiment Analysis

#### 3.2.3.4.3 Twitter

"Twitter is a global platform for public self-expression and conversation in real time. Twitter allows people to consume, create, distribute, and discover content and has democratized content creation and distribution." (Twitter, Fiscal Year 2020 Annual Report 2021, 6). Twitter has developed a rapid growth in users and popularity, becoming one of the most used social media networks globally (Wollams 2021). In 2020, average daily active user amounted to 192 million, increasing 27% compared to last year (Twitter, Fiscal Year 2020 Annual Report 2021). In 2021, Twitter was the fourth most visited website and the second most visited social network and online community globally (Neufeld 2021). It offers the users to share text, images, or short videos (Antonakak, Fragopoulou and Ioannidis 2021). A post on Twitter is called a *tweet*. It represents a short message addressed to a wide variety of receivers and is limited to a maximum of 280 characters (Twitter, Developer Platform 2021d). Therefore, users must be precise with their statements.

Due to the limitation in characters, the high number of users, the substantial increase of popularity, the publication of opinions and trends makes the Twitter database helpful in analyzing public opinions on specific topics (Kraaijeveld and De Smedt 2020).

#### 3.2.3.4.4 Data

##### 3.2.3.4.4.1 Data Collection

For this work, data was collected from 01.10.2021 up to 31.10.2021. Bitcoin prices are collected through the Binance API in a 5-minute frequency summing up to 8,893 observations in total (Binance 2021). *Tweets* are collected through the Twitter API (v2: Early access), launched in June 2020 (Twitter, Developer Platform 2021a). Counting 4,81 million *tweets* in English, queried for the keyword: "Bitcoin". Queries in Twitter are case-insensitive for all characters (Twitter, Developer Platform 2021b). Which, for example, would also include the following notations: "bitcoin", "BITCOIN", or "BiTcOiN". For each *tweet*, the following data is received:

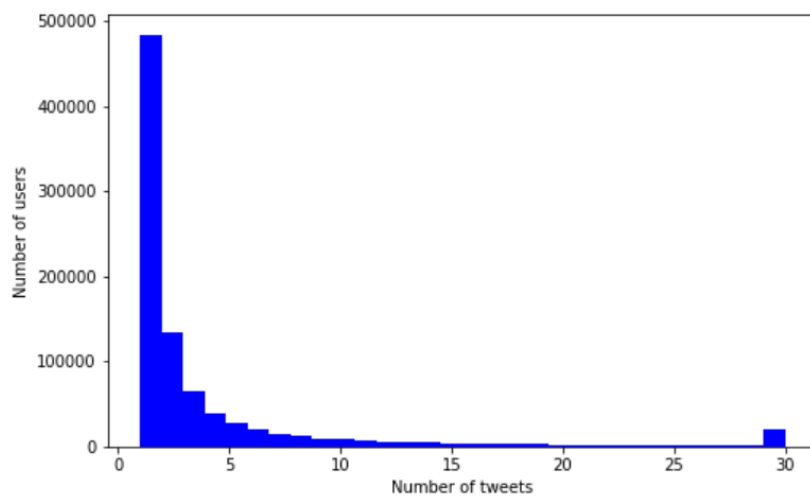
**Table 7: Received Twitter information**

Information	Detail
Created_at	Timestamp in UTC when the post was published
Author_id	Unique author id from which the post was published
Public_metrics	Retweet, reply, like, and quote count
Source	The operating system was the post was made on
Text	The content of the tweet

If a tweet is retweeted, it will contain the same message as the original *tweet* but starts with: RT @username. Public metrics are measured at the point in time when the tweet was pulled. *author\_id* is a unique number given to each account.

#### 3.2.3.4.4.2 Data Descriptive statistics before data pre-processing

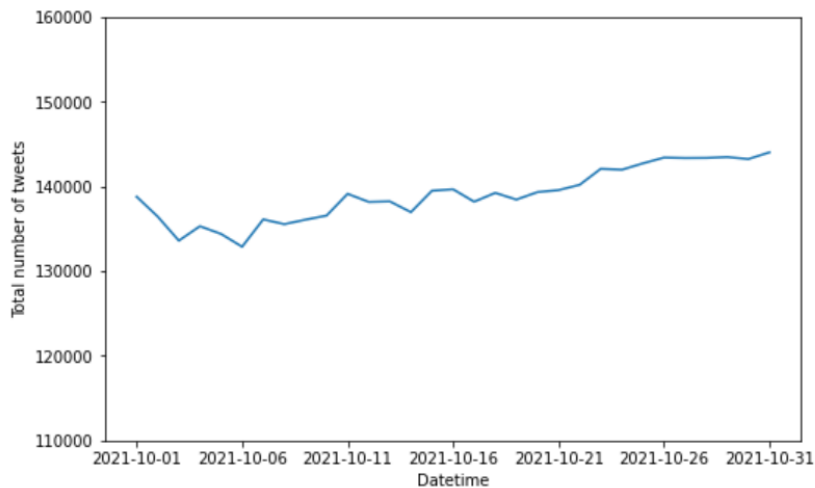
A total number of 864,968 unique *author\_ids* were identified.



*Figure 10: Distribution of tweets each user*

Figure 10 shows how many times a unique *author\_id* tweeted during October 2021. If the account tweeted more than 30 times, it was counted towards bin 30. Most of the accounts (478,000) only *tweet* once a month. As some accounts are tweeted 30 or more times indicate that the population is split into active tweeters and heavily active tweeters.





*Figure 11: Development of total number of tweets*

Figure 11 shows the number of total *tweets* observed per day, roughly 139,000 *tweets*. Towards the end of the month the number of *tweets* increases to 144,000 *tweets* a day.

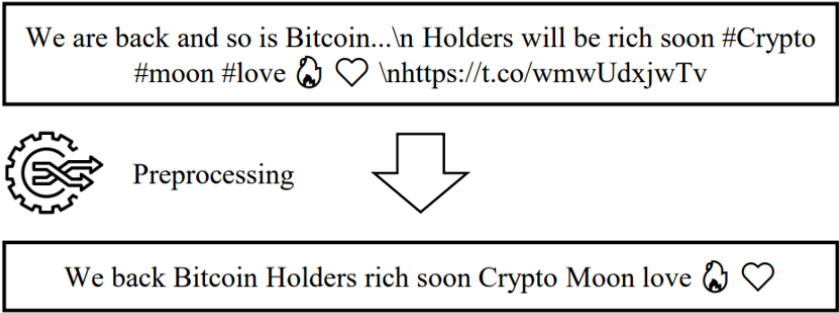
#### 3.2.3.4.4.3 Data Pre-processing

A *tweet* has a different structure than a newspaper article or a literary text from a book. The publication takes place without a correction. *Tweets* contain significantly more slang, emojis, or spelling mistakes (Pano and Kashef 2020). Hashtags (#), retweets (RT), mentions (@), or URLs, which are frequently used in *tweets*, add noise to the text (Antonakak, Fragopoulou and Ioannidis 2021). For example, the tagged username @love or @nice will influence the sentiment score of a sentence although it just represents a username without any polarity. A hashtag connected to a word #like or #love could also distort the analysis. It is recommended to clean *tweets* before performing the Sentiment Analysis (Elbagir and Yang 2019). To reduce noise in *tweets*, I performed the following data preparation steps, following the documentation proposed by (Pano and Kashef 2020). Table 8 shows all preprocessing steps performed.

**Table 8: Preprocessing steps**

Action	Example
Delete hashtags	#
Delete tags of users	@sampleusername
Delete URLs	<a href="https://www.sampledomain/">https://www.sampledomain/</a>
Delete HTML entities	&amp;, \n, etc.
Delete stopwords	the, is, a, etc.
Delete numbers	1-9

Figure 12 visualizes the step of Preprocessing for a sample tweet.



*Figure 12: Preprocessing visualization*

3.2.3.4.4.4 Descriptive Statistics

After cleaning with the above-mentioned techniques, a common technique to visualize insights about textual data is a Word Cloud (Heimerl, et al. 2014). The most used words in connection with Bitcoin are identified and first interpretations about important topics connected to Bitcoin can be made. The more giant and bold a word is represented in the graphical output, the higher is the frequency of the observed text. Figure 13 presents the Word Cloud for the search conducted in this work. Next to Bitcoin, the most frequently used words are: Crypto, BTC, gift, Price, and Ethereum. This shows that people that tweet about Bitcoin also tweet about other cryptocurrencies and the price. This is a first indication that the price is linked to the content of the observed *tweets*.



#### 3.2.3.4.5.2 Machine Learning based Tools

The process of Sentiment Analysis, with ML tools, can be interpreted as a sentiment Classification. A given ML algorithm is trained on a data set with a label/class for each data point (Liu, Bi, and Fan 2017). The model is tested on the entire data set if its accuracy on a labeled validation data set is high enough (Jain and Jain 2019).

A problem with this kind of method is the need for labeled data. As the current dataset contains 4.81 million *tweets* it would take a huge amount of time and people to label a train set manually. Furthermore, it could lead to unwanted biases in the results if the data has been classified by only a few people. ML based tools are also less interpretable because, respective of the given model, the tool acts as a kind of black box on the data (Carvalho, Pereira and Cardoso 2019). This problem arises for traditional ML algorithms and increases for Deep Learning algorithms.

#### 3.2.3.4.5.3 Lexicon-based Tools

Lexicon-based approach is based on lexicons of features – in this case: words – with their respective polarity. These lexicons were already filled and labeled. Texts are analyzed by querying the words and aggregating the produced polarity values (Turner, Labille and Gauch 2021). A challenge for lexicons is that some words in English do not possess a single polarity score but influence the polarity of other words, intensifiers like “very” or “most” or downtowners such as “slightly” or “somewhat” (Taboada, et al. 2011).

Taboada, et al. (2011) also show that negators and intensifiers (boosters) are words that drastically change the sentiment of a sentence without having a polarity value on their own such as "not", "don't" or “nobody”. However, they do change the sentiment of a sentence or a word they are linked to.

Furthermore, in social media posts, emotions are often expressed in emoticons rather than in words which also need to be measured to extract the sentiment. In the following paragraph two lexicon-based sentiment analysis tools are presented and compared.

#### 3.2.3.4.5.3.1 SentiArt

SentiArt is a lexicon-based Sentiment Analysis tool that has produced good results on literary texts. It was developed by (A. M. Jacobs 2019). The tool consists of lexicons from three different languages (English, German and Dutch). These lexicons have different feature columns that express the computed score for six emotions such as “fear” or “joy”. The feature which compared to the compound score of VADER is the AAP (Affective-Aesthetic Potential) score (Jacobs, et al. 2020). While SentiArt proved to produce good results on sentiment of novels and literary texts, it does not possess the ability to analyze the sentiment of emojis or punctuation (A. M. Jacobs 2019).

#### 3.2.3.4.5.3.2 VADER

The lexicon-based tool Valence Aware Dictionary and Sentiment Reasoner (VADER) was first introduced by (Hutto and Gilbert 2014). VADER consists of a lexicon of words that have been developed by researchers using qualitative and quantitative methods (Elbagir and Yang 2019). The lexicon’s features were established using a wisdom-of-the-crowd approach where different people assess a word's sentiment, and all assessments are aggregated to form a precise assessment (Hutto and Gilbert 2014). The developed lexicon is then combined with grammatical and syntactical rules to establish sentences' estimates in different combinations. (Hutto and Gilbert 2014). VADER has demonstrated to improve social media data results compared to traditional sentiment lexicons, due to the following reasons (Hutto and Gilbert 2014): A major part of the VADER model is its ability to take symbols into account. For instance, emoticons play a significant role on Twitter as they convey a user’s sentiment (Eisner, et al. 2016). VADER possesses an emoji dictionary that maps an emoticon to a fitting description. For example, the emoji 😊, which conveys a positive sentiment, has the description “smiling face with smiling eyes”, such that VADER evaluates the sentiment of the emoji using only its description.

Another part of *tweets* are punctuation marks. Using many punctuation marks, users convey boosted opinions and sentiments. The sentiment of “I really like this” differs from the sentiment of “I really like this!!!”. The second sentence is much more intense than the first. Similarly, the use of ALL-CAPS modifies a sentence's sentiment (Hutto and Gilbert 2014). Using the example from above, “I REALLY like this” conveys more intensity than the original sentence. These problems are tackled in VADER by using an empirically derived mean sentiment intensity rating increase and multiplying this scalar value with the given sentence or word as established in (Hutto and Gilbert 2014). Additionally, VADER is accurate in catching the negations mentioned above (Swarnkar 2020). VADER returns a compound score ranging from -1 (negative) to +1 (positive), which is further defined as the Sentiment score (Boldt and Borg 2020). Figure 14 visualizes the Vader Sentiment Analysis of the above cleaned tweet:

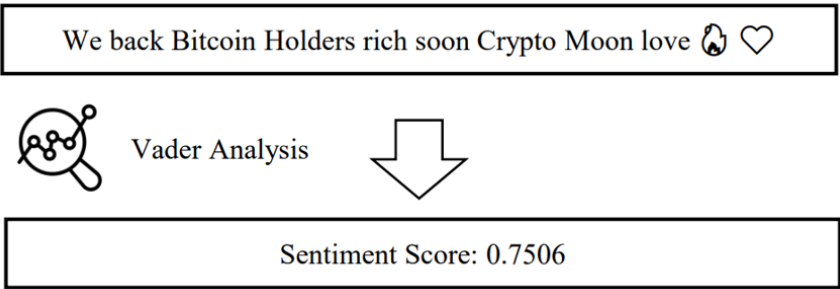


Figure 14: VADER Sentiment Analysis procedure

3.2.3.4.6 VIP Sentiment

A post from someone who is an expert or has a high number of followers is more valuable than a post from an average person, which is supported by the study from (Pant, et al. 2018) (Pant, et al. 2018).

To give these accounts a different weight, high influential Twitter accounts from people and organizations are identified, considering four different influential groups. First, the seven Bitcoin news accounts identified by Pant, et al. (2018) are included. The second group consists of the ten most followed news accounts on Twitter, according to (McCabe 2019). The third group is represented by the account of Elon Musk, chosen due to the observed abnormal returns

in the work performed by Ante (2021). Finally, the list is completed by the most influential people on Twitter regarding cryptocurrencies and Bitcoin, according to Coinbound (2021). Appendix 2 shows all accounts that were added and analyzed within the VIP Sentiment.

#### 3.2.3.4.7 Results and Discussion

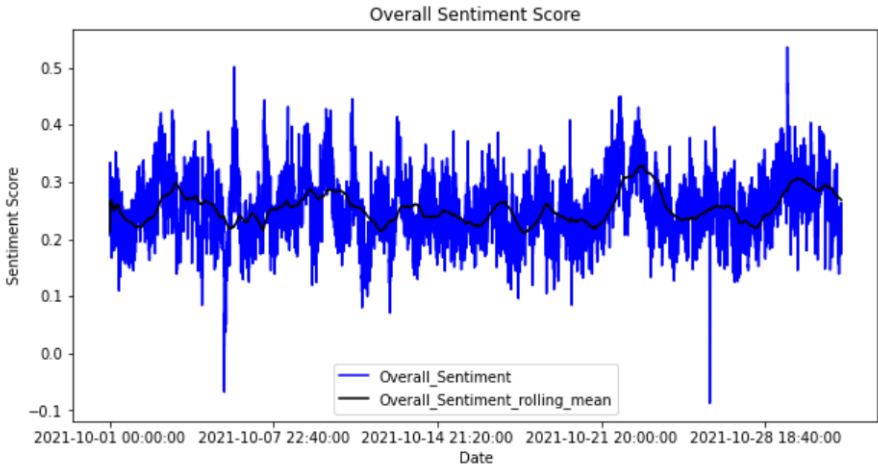
##### 3.2.3.4.7.1 Sentiment Results

In this section, I present the results of the performed VADER Sentiment Analysis and compare the differences between the *Overall Sentiment* and the *VIP Sentiment*.

**Table 9: Results of the Sentiment Analysis**

Measure	Overall Sentiment	VIP Sentiment
Mean	0.254	0.166
Standard deviation	0.055	0.355
Minimum	-0.087	-0.908
Maximum	0.536	0.917

Table 9 describes the *Overall Sentiment* results of the Sentiment Analysis. The mean is higher for the *Overall Sentiment* (0.254) than the *VIP Sentiment* (0.166), but both are positive. The standard deviation is significantly higher for the *VIP Sentiment* (0.355) compared to the *Overall Sentiment* (0.055). In line with the higher standard deviation, minimum and maximum values are more extreme for the *VIP Sentiment*. The minimum value of the *Overall Sentiment* is particularly striking because it is only (-0.087), which indicates that the average Sentiment nearly never gets negative.



*Figure 15: Development Overall Sentiment score*



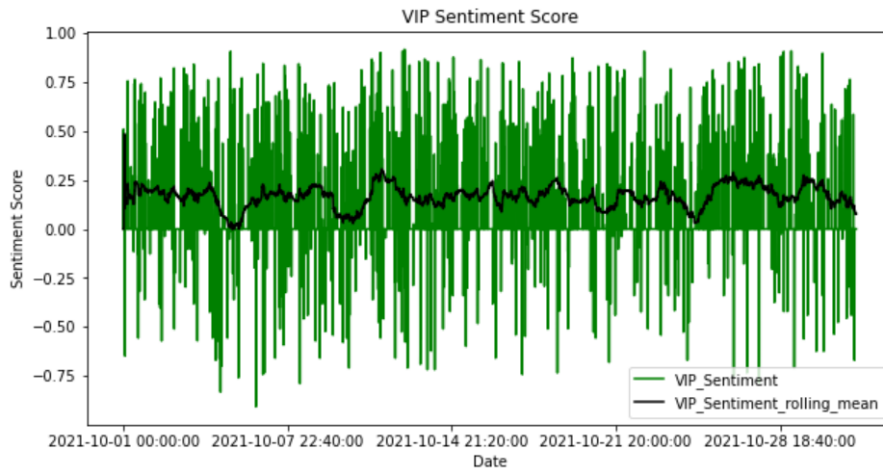


Figure 16: Development of the VIP Sentiment score

Figures 15 and 16 visualize the development of the *Overall Sentiment* and the *VIP Sentiment* during October 2021. The black line is the daily rolling average for each respective Sentiment which is also compared in one graph in figure 17. The *Overall Sentiment* has four outstanding peaks during October, with two highs and two lows.

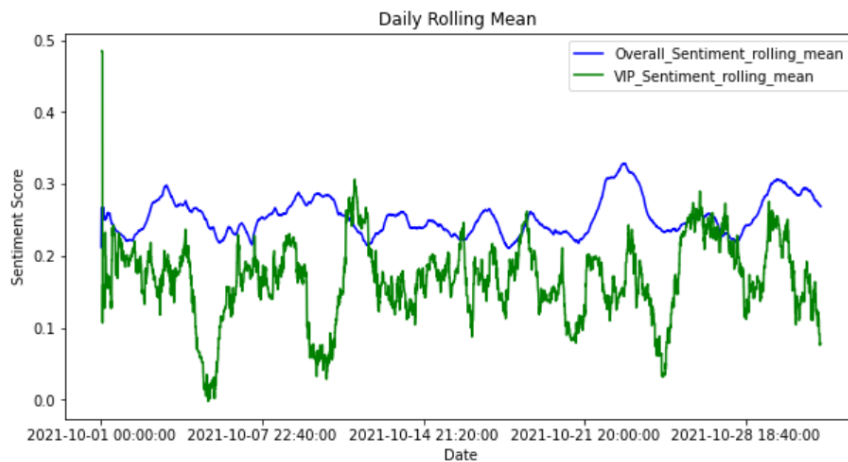


Figure 17: Comparison between Overall Sentiment and the VIP Sentiment

### 3.2.3.4.7.2 Granger Causality

“For a strictly stationary bivariate process  $\{(X_t, Y_t)\}$ ,  $\{X_t\}$  is a Granger cause of  $\{Y_t\}$  if past and current values of  $X$  contain additional information on future values of  $Y$  that is not contained in past and current  $Y$ -values alone.” (Diks and Pancheko 2006). A variable is helpful if the included variable  $X$  reduces a prediction error of the target variable  $Y$  (Clower 2021). Testing for Granger causality within a Sentiment Analysis is critical to check whether the

sentiment Granger causes the Bitcoin price or if the Bitcoin price Granger causes a sentiment (Kraaijeveld and De Smedt 2020). It is the most common “causality” test for Sentiment Analysis and is widely used in a broad range of papers for cryptocurrencies and stocks (Kraaijeveld and De Smedt 2020) (Chu, Wu and Qiu 2015) (Behrendt and Schmid 2018).

One of the requirements for testing for Granger causality is that all variables need to be stationary (Granger 1969), which can be tested with the Augmented Dickey Fuller (ADF) Test (Dickey and Fuller 1981) that is commonly used in prior research (Li 2020).

If the p-value is  $\leq 0.05$ , the null hypothesis, that the series is nonstationary, is rejected.

**Table 10: Augmented Dickey Fuller Test**

<b>Data</b>	<b>ADF Statistics</b>	<b>p-value</b>
Bitcoin Price	-2.734	0.068
<i>Overall Sentiment</i>	-9.236	0.000
<i>VIP Sentiment</i>	-37.173	0.000

For the ADF Test, a p-value  $> 0.05$  is observed for the Bitcoin price. The null hypothesis cannot be rejected. The *Overall Sentiment* and the *VIP Sentiment* both have a p-value of 0.00 and do not need to be further transformed. A common technique to transform skewed data is taking log values (Li 2020). After calculating the log values and reperforming the ADF Test the p-value is less than 0.05.

**Table 11: ADF Test after transformation**

<b>Data</b>	<b>ADF Statistics</b>	<b>P-Value</b>
Log Bitcoin Price	-3.088	0.0274

After transforming, the Granger causality test can be performed, testing the null hypothesis:

$$H_0 : \{X_t\} \text{ is not Granger causing } \{Y_t\}$$

The graph below plots the p-values for each time lag  $t$  from 1 to 40 in 5-minute steps for the respective sentiment. The null hypothesis is rejected when a p-value smaller  $\leq 0.05$  is observed, visualized as the red line in figure 18.

The blue and green highlighted areas represent the time intervals for which the null hypothesis can be rejected. Blue highlights the *Overall Sentiment* and green the *VIP Sentiment*.

For the *Overall Sentiment* the null hypothesis can be rejected with a time lag of one interval (5 minutes). The *VIP Sentiment* has a significant p-value given a time lag of nineteen intervals (95 minutes).

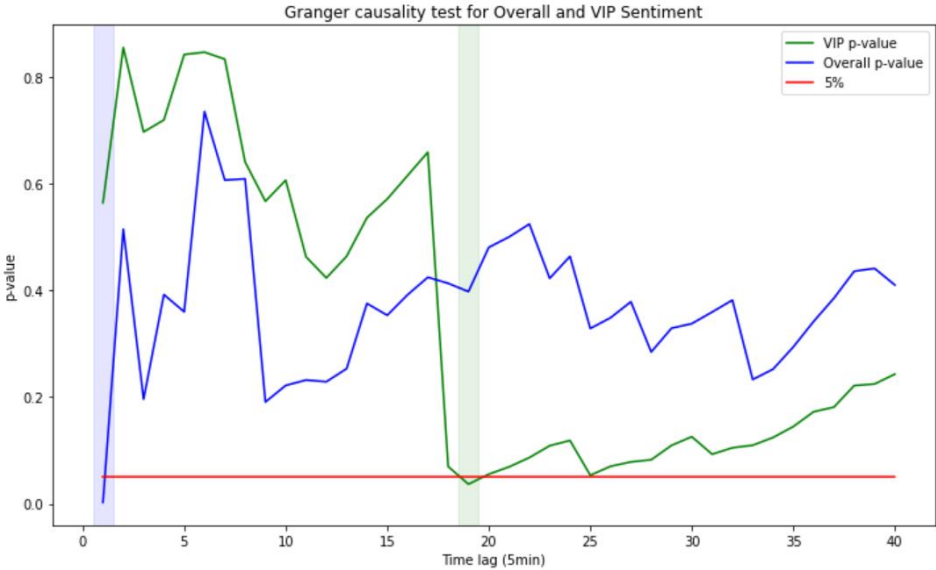


Figure 18: Granger causality test

### 3.2.3.4.8 Conclusion and Future Work

The relationship of *Overall* and *VIP Sentiment* to the intraday Bitcoin log price was analyzed using a Granger causality test. I find a Granger causal relationship of both Sentiments to the Bitcoin price. The *Overall Sentiment* demonstrates predictive influence in 5 minutes. The results of the *VIP Sentiment* show a Granger causal relationship to the Bitcoin log price in 95 minutes. The results show that Twitter Sentiment can reduce the predictive error. The effect of the Twitter *VIP Sentiment* is lagged against the *Overall Sentiment*. *Overall* and *VIP Sentiment* will be included as a feature for algorithmic trading.

Further work can focus on analyzing the interactions between *Overall* and *VIP Sentiment* and reveal the interrelationship. There could exist a group of accounts which are Granger causing the Overall Sentiment. In addition, the number of followers and the number of likes could be integrated into the analysis.

## 4 Results and Discussion

In this section, we describe the findings from the simulation of a real time Bitcoin trading algorithm. First, we provide an overview of the generated insights regarding the influence of twitter sentiment on the Bitcoin price (Study I), the optimal modeling design for Bitcoin price and trend prediction (Study II) and the best performing trading strategy derived from the predictive analysis (Study III). Second, we evaluate the profitability of the trading strategy including costs associated with running a real time trading algorithm. Finally, we discuss feasibility of implementation and limitations of financial evaluation.

We introduce a real time Bitcoin trading algorithm that covers the entire process from Data Collection to the translation of trading signals, analyzing the influence of Twitter, the modeling design and trading strategies. We investigate that *Overall Sentiment* and *VIP Sentiment* for Twitter have a Granger causal relationship to the Bitcoin price. We find that LSTM yields the best price prediction performance for Bitcoin price prediction in *1h*, *2h* and *3h*. DL outperforms RF and XGBoost for Classification. GRU, LSTM and RNN provide the best trend prediction of the Bitcoin price for *1h*, *2h* and *3h* trend predictions, respectively. We find evidence that algorithmic trading of Bitcoin using predictive analysis of ML algorithms can earn positive returns. Strategies derived from regression models have a higher financial performance than strategies derived from Classification models. We identify the ensemble strategy *reg\_consensus* that combines the predictions of LSTM regression models with three different prediction horizons to be the best performing strategy. The *reg\_consensus* strategy generates a ROI of

7.32% for the test period, which is superior to the ROI of the *buy\_and\_hold* strategy (0.13%). The ARIMA represents the benchmark for our Regression models. The *ARIMA\_consensus* strategy has the best overall results among the ARIMA models but is not able to trigger profitable trading decisions. The PV development of the *reg\_consensus* strategy is visualized and compared to *buy\_and\_hold* and *ARIMA\_consensus* strategy in figure 23.



Figure 19: PV development of *buy\_and\_hold*, *ARIMA\_consensus* and *reg\_consensus*

### Profitability

The calculation of the ROI that is performed in this work, includes trading associated costs based on the cost-settings of Binance. Costs that arise for the development, implementation and deployment of the ML models are not factored into the calculation. For a final evaluation of the trading algorithm, costs for Data Collection and Computing Power need to be included. Several platforms offer the integration of algorithmic trading strategies, providing API access, computing power, back-testing analysis, and other services (Fang, et al. 2020). The usage of these services presents multiple opportunities to structure costs but this work project intends to provide a stand-alone approach for algorithmic trading. The Data collection of this work is aimed to incur a minimum of costs. While the Yahoo Finance and Binance API are available

free of charge, a paid API is required to collect Twitter data. For academic research like our work, the Twitter Developer API is freely available after application review and signing a non-commercial use agreement. The implementation of a trading algorithm requires a commercial Twitter API that costs \$2.499 per month and allows to retrieve a maximum of five million tweets per month (Twitter, Developer Platform 2021c). The costs for the commercial Twitter API during the period tested in this work, would amount to \$280. Due to the calculational complexity of Hyperparameter Tuning for DL algorithms, GPU computation is required (Cocco, Tonelli and Marchesi 2021). External computing power needs to be purchased to train and deploy the developed algorithms. A GeForce GTX 1080Ti with two GPU cores is required to train the LSTM Regression model for predictions in *1h*, *2h* and *3h*. Monthly costs total \$876. The costs for the computing power during the period tested in this work, amount to \$100 (Genesis, 2021). Additional costs include the development and monitoring of the algorithm as well as the connection to the platform API for real-time trading. These costs are difficult to quantify, and we do not include them in the following evaluation. Gains from cryptocurrency trading are not subject to taxes in Portugal and therefore not included in the calculation (Cointaxlist 2021). The ROI and the total profit of *buy\_and\_hold*, the *ARIMA\_consensus* strategy as the best performing strategy among the ARIMA strategies, and the strategy *reg\_consensus*, are shown with and without the inclusion of costs for Twitter API and Computing Power in table 19. The performance metrics are calculated for the test period from 28.10.2021 09:50 to 31.10.2021 21:05. During this period our trading algorithm, based on the *reg\_consensus* strategy, earns a total profit after costs of 352 USDT, equal to a ROI of 3.52%. *ARIMA\_consensus* and *buy\_and\_hold* only achieve a ROI of 0.13% and -1.56%, respectively. Considering costs for Data Collection and Computing Power our trading algorithm outperforms its benchmark strategies.

**Table 12: ROI of *reg\_consensus* strategy and benchmark strategies including costs**

<b>All Amounts in USDT</b>	<b>Buy_and_hold</b>	<b>ARIMA_consensus</b>	<b>reg_consensus</b>
Start Balance	10,000	10,000	10,000
Final Balance	10,013	9,834	10,732
+ Profit	13	-156	732
- Twitter API	0	0	280
- Computing Power	0	0	100
Profit (after costs)	13	-156	352
ROI (after costs)	0.13%	-1.56%	3.52%

### Feasibility

Automated trading based on the developed trading algorithm requires considerations for feasibility of a real-time implementation. While previous research fails to address components of real-time implementation, we discuss limitations of the trading algorithm when it comes to real-time trading. Training of ML algorithms requires computing time dependent on the provided computing power.

LSTM algorithms have a high computational complexity (Cocco, Tonelli, and Marchesi 2021). Using computing power from *Genesis* cloud it took around eight hours for each LSTM algorithm to train, limiting the frequency of applying a newly trained model when Computing Power is used as described in this work. An alternative to the presented modeling design, i.e., batch learning, is online learning. Online learning is a promising technique for learning from continuous streams of data but requires a different modeling architecture. For online learning, the algorithm takes the current model and subsequently uses new observations to further adjust the weights of each parameter. Online learning is faster to train but more difficult to maintain as the algorithms rely on a constant flow of data points (Hoi, et al. 2021). Constantly collecting data in the same format is challenging. The data for this work project is entirely collected using APIs. Although APIs are specifically designed to support Data Collection, APIs are subject to changes. For example, Twitter updated its API in June 2020 with an Early Access for the v2 *API*. In November 2021, the usage was published for all developers, causing changes in the

Data Collection process. The four main variations are: Endpoint URLs, app and project requirements, response data format, and request parameters (Twitter, Twitter Developer 2021). Changes in API structure need to be monitored and code needs to be adjusted to prevent prediction errors. Trading signals triggered by the trading algorithm can either notify the trader or directly execute a trading action. An automated execution of trading signals requires the implementation of a real-time connection between the trading algorithm and a trading platform, e.g., Binance. While this work presents the foundation for a real-time implementation of the trading strategy, the connection to the Binance platform is not in the scope of this work. Considering these feasibility issues, we find that our developed trading algorithm has the prerequisites to be used for real-time trading.



## 5 Conclusion and Future Work

We define a holistic approach to build an intraday Bitcoin trading algorithm derived from predictive analysis of ML models and test the developed trading algorithm in a simulation setting. Special focus is placed on the impact of Twitter Sentiment (Study I), the Modeling Design (Study II), and the Trading Strategy (Study III). Finally, we evaluate profitability and feasibility of the trading algorithm in real-time implementation, which previous research fails to address.

We combine Regression and Classification models, with features from five feature categories (Supply & Demand, Crypto market, Macro Financial, Political and Sentiment). Study I identifies a Granger causal relationship between the overall and VIP Twitter Sentiment on the Bitcoin price. Study II concludes that LSTM models yield the best prediction performance for Bitcoin price prediction and GRU, LSTM and RNN generate the best Bitcoin trend predictions in *1h*, *2h* and *3h*, respectively. Study III finds superior profitability of ensemble trading strategies over individual trading strategies and identifies a Regression ensemble strategy to achieve the best overall results. Combining the findings of Study I, II, and III, we provide a holistic design of a trading algorithm. Finally, we evaluate the profitability of our trading algorithm for a real-time implementation considering costs for Data Collection and Computing Power and evaluate feasibility concerns. Our findings indicate that our intraday trading algorithm can be implemented for real-time trading and generates positive returns that exceed the returns of benchmark strategies.

Direct extensions to this work can investigate the real-world implementation of the presented design for an intraday Bitcoin trading algorithm. Special focus should be placed on the deployment of online learning for continuous model development. Further work can elaborate on developing additional business plans for monetarizing the presented trading algorithm.

## List of Appendices

Appendix 1: Feature Overview .....	48
Appendix 2: VIP List .....	49

## List of tables

Table 1: Coverage of major topics in the defined 73 papers.....	5
Table 2: Combinations of the keyword search query.....	6
Table 3: Overview of focus paper .....	7
Table 4: Overview of feature categories .....	12
Table 5: Overview of Data Collection in the focus paper.....	13
Table 6: Overview of collected features .....	14
Table 7: Received Twitter information .....	28
Table 8: Preprocessing steps .....	30
Table 9: Results of the Sentiment Analysis .....	36
Table 10: Augmented Dickey Fuller Test.....	38
Table 11: ADF Test after transformation.....	38
Table 12: ROI of <i>reg_consensus</i> strategy and benchmark strategies including costs .....	42

## List of figures

Figure 1: Literature research approach.....	6
Figure 2: Architecture of the work project.....	10
Figure 3: Development of Bitcoin price.....	15
Figure 4: Bitcoin trade volume in 5-minute intervals .....	16
Figure 5: Correlation matrix of cryptocurrency market features .....	16
Figure 6: Normalized price development of observed cryptocurrencies .....	17
Figure 7: Normalized price development of commodity and Bitcoin prices .....	18
Figure 8: Distribution of trend variables .....	20
Figure 9: Structure of the Sentiment Analysis .....	26
Figure 10: Distribution of tweets each user .....	28
Figure 11: Development of total number of tweets.....	29
Figure 12: Preprocessing visualization .....	30
Figure 13: Tweet Word Cloud .....	31
Figure 14: VADER Sentiment Analysis procedure .....	34
Figure 15: Development Overall Sentiment score .....	36
Figure 16: Development of the VIP Sentiment score .....	37
Figure 17: Comparison between Overall Sentiment and the VIP Sentiment.....	37
Figure 18: Granger causality test .....	39
Figure 19: PV development of buy_and_hold, ARIMA_consensus and reg_consensus.....	41

## Appendix

### Appendix 1: Feature Overview

Feature	Category	Interval	Sources	Start	End
BTC	Supply & Demand	5 min	Binance API	01.10.2021	31.10.2021
ETH	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
BNB	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
ADA	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
SOL	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
XRP	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
DOT	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
DOGE	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
USDC	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
LUNA	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
LTC	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
AVAX	Cryptocurrency Market	5 min	Binance API	01.10.2021	31.10.2021
S&P 500	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
SSE Composite	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Nikkei 225	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Dax 40	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
BSE Senex	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
FTSE 100	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
CAC 40	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
BOVESPA	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
FTSE MIB	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
TSX Composite	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
CNY	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
JPY	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
EUR	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
INR	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
GBP	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
BRL	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
CAD	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
CBOE Volatility	Political	5 min	Yahoo API	01.10.2021	31.10.2021
Brent	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021

Natural Gas	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Soybeans	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Corn	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Gold	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Copper	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Silver	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
WTI	Macro Financial	5 min	Yahoo API	01.10.2021	31.10.2021
Twitter	Sentiment	5 min	Twitter API	01.10.2021	31.10.2021

## Appendix 2: VIP List

Name	Twitter	ID
CNN Breaking News	@cnnbrk	428333
The New York Times	@nytimes	807095
CNN	@cnn	759251
BBC Breaking News	@bbcbreaking	5402612
BBC World	@bbcworld	742143
The Economist	@theeconomist	5988062
Reuters Top News	@reuters	1652541
The Wall Street Journal	@wsj	3108351
Time	@time	14293310
BitcoinNews	@BTCTN	3367334171
CryptoCurrency	@cryptocurrency	216304017
CryptoYoda	@CryptoYoda1338	852256178021294080
BitcoinMagazine	@BitcoinMagazine	361289499
CoinDesk	@coindesk	1333467482
Roger Ver	@rogerkver	176758255
Erik Voorhees	@ErikVoorhees	61417559
Ty Smith	@TyDanielSmith	961971412528517120
Tone Vays	@ToneVays	2577886615
CryptoCobain	@CryptoCobain	2259434528
Tyler Winklevoss	@Tyler	24222556
Vitalik Buterin	@VitalikButerin	295218901
CryptoWendyO	@CryptoWendyO	935742315389444096
StackingUSD	@StackingUSD	431243238
Girl Gone Crypto	@Girlgone_Crypto	1150790822813560833
David Gokhshtein	@davidgokhshtein	170049408
Hailey Lennon	@HaileyLennonBTC	3740778132
Justin Sun	@justinsuntron	902839045356744704

Ivan on Tech	@IvanOnTech	390627208
Kenn Bosak	@KennethBosak	4693571508
Scott Melker	@ScottMelker	17351167
TheCryptoDog	@TheCryptoDog	887748030304329728
BitBoy Crypto	@Bitboy_Crypto	954005112174862336
Dan Held	@DanHeld	1598709350
LayahHeilpern	@LayahHeilpern	455937214
Elon Musk	@elonmusk	44196397

---

## References

- Abbad, Jumah, Bashar Fardousi, and Muneer Abbad. 2014. "Advantages of Using Technical Analysis to Predict Future Prices on the Amman Stock Exchange." *International Journal of Business and Management* 1-16.
- Abraham, Jethin, Daniel Higdon, John Nelson, and Juan Ibarra. 2018. "Cryptocurrency Price Prediction Using Tweet Volumes and Sentiment Analysis." *SMU Data Science Review: Vol. 1 : No. 3 , Article 1 22*.
- Abramovich, Felix, and Claudia Angelini. 2006. "Testing in mixed-effects FANOVA models." *Journal of Statistical Planning and Inference* 136(12): 4326-4348.
- Aharona, David Yechia, and Mahmoud Qadan. 2019. "Bitcoin and the day-of-the-week effect." *Finance Research Letters* 31 415-424.
- Ahmed, Walid M.A., and Mustafa AL Mafrachi. 2021. "Do higher-order realized moments matter for cryptocurrency returns?" *International Review of Economics and Finance* 72 483-499.
- Ahn, Yongkil, and Dongyeon Kim. 2021. "Emotional trading in the cryptocurrency market." *Finance Research Letters* 42 101912.
- Ante, Lennart. 2021. "How Elon Musk's Twitter activity moves cryptocurrency markets." *BRL Working Paper Series No. 16 13*.
- Antonakak, Despoina, Paraskevi Fragopoulou, and Sotiris Ioannidis. 2021. "A survey of Twitter research: Data model, graph structure, sentiment analysis and attacks." *Expert Systems with Applications. Elsevier BV Volume 164 25*.
- Aroussi, Ran. 2021. *yfinance 0.1.67*. Accessed December 07, 2021. <https://pypi.org/project/yfinance/>.

- Atsalakis, George S., Ioanna G. Atsalaki, Fotios Pasiouras, and Constantin Zopounidis. 2019. "Bitcoin price forecasting with neuro-fuzzy techniques." *European Journal of Operational Research* 276 770-780.
- Bali, Turan G., and Peng Lin. 2006. "Is there a risk-return trade-off? Evidence from high-frequency data." *Journal of Applied Econometrics* 21 1169-1198.
- Barclay, Michael J., and Terrence Hendershott. 2003. "Price Discovery and Trading After Hours." *The Review of Financial Studies* 1041-1073.
- Behrendt, Simon, and Alexander Schmid. 2018. "The Twitter myth revisited: Intraday investor sentiment, Twitter activity and individual-level stock return volatility." *Journal of Banking and Finance* Volume 96 355–367.
- Bergstra, James, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. 2011. "Algorithms for hyper-parameter optimization." *Proceedings of the 24th International Conference on Neural Information Processing*. Granada, Spain.
- Bernal, Armando, Sam Fok, and Rohit Pidakparthi. 2012. "Market Time Series Prediction with Recurrent Neural Networks."
- Binance. 2021. *API Documentation*. December 03. Accessed December 07, 2021. <https://binance-docs.github.io/apidocs/futures/en/#change-log>.
- Binance. 2021. *Binance*. December 7. Accessed December 7, 2021. <https://www.binance.com/en/fee/schedule>.
- Blockchain.com. 2021. *Blockchain.com Wallets*. December 12. Accessed December 13, 2021. <https://www.blockchain.com/charts/my-wallet-n-users>.
- Bohte, Rossini. 2019. "Comparing the Forecasting of Cryptocurrencies by Bayesian Time-Varying Volatility Models." *Journal of Risk and Financial*.
- Boldt, Martin, and Anton Borg. 2020. "Using VADER sentiment and SVM for predicting customer response sentiment." *Expert Systems with Applications* 162 11.



- Borges, Tomé Almeida, and Rui Ferreira Neves. 2020. "Ensemble of machine learning algorithms for cryptocurrency investment with different data resampling methods." *Applied Soft Computing Journal* 90.
- Bouri, Elie, Peter Molnár, Georges Azzi, David Roubaud, and Lars Ivar Hagfors. 2017. "On the hedge and safe haven properties of Bitcoin: Is it really more than a diversifier?" *Finance Research Letters* 20: 192-198.
- Breiman, Leo. 2001. "Random Forests." *Machine Learning* 5–32.
- Breiman, Leo. 1996. "Bagging Predictors." *Machine Learning* 24: 123-140.
- Brownlee, Jason. 2021. *Machine Learning Mastery*. Februar 17. Accessed December 4, 2021. <https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning/>.
- Bugár, Gyöngyi, and Márta Somogyvári. 2020. "Bitcoin: Digital Illusion or a Currency of the Future?" *Financial and Economic Review, Vol. 19 Issue 1* 132–153.
- Carvalho, Diogo V, Eduardo M. Pereira, and Jaime S. Cardoso. 2019. "Machine Learning Interpretability: A Survey on Methods and Metrics." *Electronics* 2019, 8, 832–844.
- Chan, Ernest P. 2021. *Quantitative Trading, How to Build Your Own Algorithmic Trading Business*. New Jersey: John Wiley & Sons, Inc., Hoboken.
- Cheah, Eng-Tuck, and John Fry. 2015. "Speculative bubbles in Bitcoin markets? An empirical investigation." *Economics Letters* 130: 32-36.
- Chen, Zheshi, Chunhong Li, and Wenjun Sun. 2020. "Bitcoin price prediction using machine learning: An approach to sample dimension engineering." *Journal of Computational and Applied Mathematics* 365.
- Chevallier, Julien, Bangzhu Zhu, and Lyuyuan Zhang. 2021. "Forecasting Inflection Points: Hybrid Methods with Multiscale Machine Learning Algorithms." *Computational Economics* 537-575.

- Chong, Eunsuk, Chulwoo Han, and Frank Chongwoo Park. 2017. "Deep learning networks for stock market analysis and prediction: Methodology, data representations, and case studies." *Expert Systems with Applications* 83: 187-205.
- Chu, Jeffrey, Stephen Chan, and Yuanyuan Zhang. 2020. "High frequency momentum trading with cryptocurrencies." *Research in International Business and Finance* 52 101176.
- Chu, Xiaojun, Chongfeng Wu, and Jianying Qiu. 2015. "A nonlinear Granger causality test between stock returns and investor sentiment for Chinese stock market: a wavelet-based approach." *Applied Economics*, 48:21 1915-1924.
- Chung, Hyejung, and Kyung-shik Shin. 2018. "Genetic Algorithm-Optimized Long Short-Term Memory Network for Stock Market Prediction." *Sustainability* 10(10): 3765-3783.
- Chung, Junyoung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling." *NIPS 2014*.
- Claesen, Marc, Jaak Simm, and Dusan Popovic. 2014. *Optunity*. Accessed December 14, 2021. <https://optunity.readthedocs.io/en/latest/user/solvers/TPE.html>.
- Clower, Eric. 2021. *APTECH*. October 4. Accessed October 20, 2021. <https://www.aptech.com/blog/introduction-to-granger-causality/>.
- Cocco, Luisanna, Roberto Tonelli, and Michele Marchesi. 2021. "Predictions of bitcoin prices through machine learning based frameworks." *PeerJ Computer Science* 7:e413.
- Cohen, Jerome Bernard, Edward D Zinbarg, and Arthur Zeikel. 2014. *Investment analysis and portfolio management*<sup>^</sup>. Homewood: R.D. Irwin.
- Coinbound. 2021. *Coinbound*. Accessed October 25, 2021. <https://coinbound.io/best-crypto-influencers-on-twitter/>.
- CoinMarketCap. 2021. *CoinMarketCap*. December 07. Accessed December 07, 2021. <https://coinmarketcap.com/de/>.

- Cointaxlist. 2021. *Cointaxlist*. 12 14. Accessed 12 14, 2021. <https://cointaxlist.com/blog/is-portugal-really-a-tax-haven-for-crypto>.
- Comment\_Picker. 2021. *Twitter ID & Follower Count*. Accessed October 10, 2021. <https://commentpicker.com/twitter-id.php>.
- Cybenko, George. 1989. "Approximation by superpositions of a sigmoidal function." *Mathematics of Control, Signals and Systems* 2: 303–314.
- Dastgir, S., E. Demir, G. Downing, G. Gozgor, and C. K. Lau. 2019. "The causal relationship between Bitcoin attention and Bitcoin returns: Evidence from the Copula-based Granger causality test." *Finance Research Letters Volume 28* 160-164.
- de Souza, Matheus José Silva, et al. 2019. "Can artificial intelligence enhance the Bitcoin bonanza." *The Journal of Finance and Data Science* 5, 83-98.
- Dickey, David A., and Wayne A. Fuller. 1981. "Likelihood Ratio Statistics for Autoregressive Time Series with a Unit Root." *Econometrica Vol. 49, No. 4* 1057-1072.
- Diks, Cees, and Valenty Pancheko. 2006. "A new statistic and practical guidelines for nonparametric Granger causality testing." *Journal of Economic Dynamics & Control* 30 1647–1669.
- Dorsey, Jack. 2006. *Twitter*. March 21. Accessed October 20, 2021. <https://twitter.com/jack/status/20>.
- Dr. Rajab, Muhanad, and Feda Hassan Jahjah. 2020. "Impact of Twitter Sentiment Related to Bitcoin on Stock Price Returns." *University of Baghdad Engineering Journal* 26 (6) 60-71.
- Dutta, Aniruddha, Saket Kumar, and Meheli Basu. 2020. "A Gated Recurrent Unit Approach to Bitcoin Price Prediction." *Journal of Risk and Financial Management* 13(2):23.
- Eisner, Ben, Tim Rocktäschel, Isabelle Augenstein, Matko Bosnjak, and Sebastian Riedel. 2016. "emoji2vec: Learning Emoji Representations from their Description." *4th*

*International Workshop on Natural Language Processing for Social Media at EMNLP 2016* 7.

- Elbagir, Shihab, and Jing Yang. 2019. "Twitter Sentiment Analysis Using Natural Language Toolkit and VADER Sentiment." *Proceedings of the International MultiConference of Engineers and Computer Scientists 2019* 5.
- Elsaraiti, Meftah, and Adel Merabet. 2021. "A Comparative Analysis of the ARIMA and LSTM Predictive." *Energies* 14: 6782-6798.
- Fang, Fan, Carmine Ventrea, Michail Basiosb, Hoilong Kong, Leslie Kanthan, David Martinez-Rego, Fan Wu, and Li Lingbo. 2020. "Cryptocurrency Trading: A Comprehensive Survey." *arXiv preprint arXiv:2003.11352*.
- Farrar, Donald E., and Robert R. Glauber. 1967. "Multicollinearity in Regression Analysis: The Problem Revisited." *The Review of Economics and Statistics* 49(1): 92-107.
- FIA. 2021. *Commodities*. September 24. Accessed October 21, 2021. <https://www.fia.org/commodities>.
- Finance, Yahoo. 2021. *Top Cryptos by Volume*. November 11. Accessed November 11, 2021. <https://ca.finance.yahoo.com/u/yahoo-finance/watchlists/crypto-top-volume-24hr>.
- Fischer, Thomas Günter, Christopher Krauss, and Alexander Deinert. 2019. "Statistical Arbitrage in Cryptocurrency Markets." *Journal of Risk and Financial* 12(1), 31.
- Garcia, David, and Frank Schweitzer. 2015. "Social signals and algorithmic trading of Bitcoin." *Royal Society open science* 2.9.
- Genesis, Cloud. 2021. *Pricing*. December 12. Accessed December 12, 2021. <https://www.genesiscloud.com/pricing>.
- Gil-Alana, Luis Alberiko, Emmanuel Joel Aikins Abakah, and María Fátima Romero Rojo. 2020. "Cryptocurrencies and stock market indices. Are they related?" *Research in*

doi:<https://doi.org/10.1016/j.ribaf.2019.101063>.

- Granger, C. W. J. 1969. "Investigating Causal Relations by Econometric Models and Cross-spectral Methods." *Econometrica* Vol. 37, No. 3 424-438.
- Greaves, Alex, and Benjamin Au. 2015. "Using the Bitcoin Transaction Graph to Predict the Price of Bitcoin." *Computer Science*.
- Gupta, Nitin, Shashank Mujumdar, Hima Patel, Satoshi Masuda, Naveen Panwar, Sambaran Bandyopadhyay, Sameep Mehta, et al. 2021. "Data Quality for Machine Learning Tasks." *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'21)*. Singapore: ACM. 2.
- Hannun, Awni, Guo Chuan, and van der Maaten Laurens. 2021. "Measuring Data Leakage in Machine-Learning Models with Fisher Information." *Conference on Uncertainty in Artificial Intelligence*.
- Heimerl, F., Lohmann S., S. Lange, and Ertl T. 2014. "Word Cloud Explorer: Text Analytics Based on Word Clouds." *47th Hawaii International Conference on System Sciences*. 1833-1842.
- Henckaerts, Roel, Marie-Pier Côté, Katrien Antonio, and Roel Verbelen. 2021. "Boosting Insights in Insurance Tariff Plans with Tree-Based Machine Learning Methods." *North American Actuarial Journal* (25:2): 255-285.
- Hochreiter, Sepp. 1997. "Long Short-term Memory." *Neural Computation* 9(8): 1735-1780.
- Hochreiter, Sepp. 1998. "The Vanishing Gradient Problem During Learning Recurrent Neural Nets and Problem Solutions." *International Journal of Uncertainty Fuzziness and Knowledge-Based Systems* 6(2): 107-116.
- Hoi, Steven C.H., Doyen Sahoo, Jing Lu, and Peilin Zhao. 2021. "Online learning: A comprehensive survey." *Neurocomputing* 459: 249-289.

- Huang, Jing-Zhi, William Huang, and Jun Ni. 2019. "Predicting bitcoin returns using high-dimensional technical indicators." *The Journal of Finance and Data Science* 5 140-155.
- Hutto, C.J., and Eric Gilbert. 2014. "VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text." *Eighth International AAAI Conference on Weblogs and Social Media* 216-225.
- Ibrahim, Ahmed, Rasha Kashef, Menglu Li, Esteban Valencia, and Eric Huang. 2020. "Bitcoin Network Mechanics: Forecasting the BTC Closing Price Using Vector Auto-Regression Models Based on Endogenous and Exogenous." *Journal of Risk and Financial Management* 13(9), 189.
- Jacobs, AM., B. Herrmann, G. Lauer, J. Lüdtkke, and S. Schroeder. 2020. "Sentiment Analysis of Children and Youth Literature: Is There a Pollyanna Effect?" *Frontiers in Psychology. Frontiers Media SA.* 8.
- Jacobs, Arthur M. 2019. "Sentiment Analysis for Words and Fiction Characters From the Perspective of Computational (Neuro-)Poetics." *Frontiers in Robotics and AI. Frontiers Media SA.* 13.
- Jain, Achin, and Vanita Jain. 2019. "Sentiment classification of twitter data belonging to renewable energy using machine learning." *Journal of Information and Optimization Sciences (Vol. 40, Issue 2)* 521-533.
- Ji, Suhwan, Jongmin Kim, and Hyeonseung Im. 2019. "A Comparative Study of Bitcoin Price Prediction Using Deep Learning." *Mathematics* 7(10): 898-918.
- John, Alun, Samuel Shen, and Tom Wilson. 2021. *Future of Money*. September 24. Accessed December 07, 2021. <https://www.reuters.com/world/china/china-central-bank-vows-crackdown-cryptocurrency-trading-2021-09-24/>.

- Khoo, Christopher, and Sathik Basha Johnkhan. 2018. "Lexicon-based sentiment analysis: Comparative evaluation of six sentiment lexicons." *Journal of Information Science*, Volume 44 (4): 22 - Aug 1, 2018.
- Kim, J., H. Wimmer, Liu Kim, and S. Kim. 2021. "A Streaming Data Collection and Analysis for Cryptocurrency Price Prediction using LSTM." *IEEE/ACIS 6th International Conference on Big Data, Cloud Computing, and Data Science (BCD)* 45-52.
- Kraaijeveld, Oliver, and Johannes De Smedt. 2020. "The predictive power of public Twitter sentiment for forecasting." *Journal of International Financial Markets, Institutions and Money* Volume 65 22.
- Kristjanpoller, Werner, and Marcel C. Minutolo. 2018. "A hybrid volatility forecasting framework integrating GARCH, artificial neural network, technical analysis and principal components analysis." *Expert Systems With Applications* 109 1-11.
- Kuhn, Max, and Kjell Johnson. 2013. *Applied Predictive Modeling*. New York: Springer-Verlag.
- Kunesh, Andrew. 2020. *Traject*. February 25. Accessed November 25, 2021. <https://fanbooster.com/blog/social-media-post-lengths/>.
- Lahmiri, Salim, and Stelios Bekiros. 2021. "Deep Learning Forecasting in Cryptocurrency High-Frequency Trading." *Cognitive Computation* 13: 485–487.
- Li, Susan. 2020. *towardsdatascience*. December 23. Accessed November 26, 2021. <https://towardsdatascience.com/a-quick-introduction-on-granger-causality-testing-for-time-series-analysis-7113dc9420d2>.
- Lim, Bryan, and Stefan Zohren. 2020. "Time Series Forecasting With Deep Learning: A Survey." *Philosophical Transactions A*.

- Liu, Yang, Jian-Wu Bi, and Zhi-Ping Fan. 2017. "Multi-class sentiment classification: The experimental comparisons of feature selection and machine learning algorithms." *Expert Systems with Applications (Vol. 80)* 323–339.
- Liu, Zhen, Ruoyu Wang, Ming Tao, and Xianfa Cai. 2015. "A class-oriented feature selection approach for multi-class imbalanced network traffic datasets based on local and global metrics fusion." *Neurocomputing* 168: 365–381.
- Lu, Haitian, Bingzhong Wang, Qing Wu, and Jing Ye. 2020. "Fintech and the Future of Financial Service: A Literature Review and Research Agenda." *China Accounting and Finance Review* 107-136.
- Mallqui, Dennys C.A., and Ricardo A.S. Fernandes. 2019. "Predicting the direction, maximum, minimum and closing prices of daily Bitcoin exchange rate using machine learning techniques." *Applied Soft Computing Journal* 75: 596–606.
- Matsuo, Yutaka. 2003. "Prediction, Forecasting, and Chance Discovery." In *Advanced Information Processing*, by Y. Ohsawa and P. McBurney, 30-43. Berlin, Heidelberg: Springer.
- McCabe, Michael. 2019. *intelligencefusion*. June. Accessed November 15, 2021. <https://www.intelligencefusion.co.uk/insights/resources/article/top-30-most-followed-news-accounts-on-twitter/>.
- Medhat, Walaa, Ahmed Hassan, and Hoda Korashy. 2014. "Sentiment analysis algorithms and applications: A survey." *Ain Shmas Engineering Journal Volume 5, Issue 4* 1093-1113.
- Mitra, A. 2020. "Sentiment Analysis Using Machine Learning Approaches (Lexicon based on movie review dataset)." *Journal of Ubiquitous Computing and Communication Technologies (Vol. 2, Issue 3)* 145–152.
- Mordor Intelligence. 2020. *Algorithmic Trading Market - Growth, Trends, COVID-19 Impact, and Forecasts (2021 - 2026)*. Online: Mordor Intelligence.



- Mudassir, Mohammed, Shada Bennbaia, Devrim Unal, and Mohammad Hammoudeh. 2020. "Time-series forecasting of Bitcoin prices using high-dimensional features: a machine learning approach." *Neural Computing and Applications*.
- Munim, Ziaul Haque, Mohammad Hassan Shakil, and Ilan Alon. 2019. "Next-Day Bitcoin Price Forecast." *Journal of Risk and Financial* 12(2).
- Murphy, J. J. 1999. "Technical analysis of the financial markets: A comprehensive guide to trading methods and applications." *Penguin*.
- Nakano, Masafumi, Akihiko Takahashi, and Soichiro Takahashi. 2018. "Bitcoin technical trading with artificial neural network." *Physica A* 510: 587–609.
- Neufeld, Dorothy. 2021. *visualcapitalist*. January 27. Accessed November 25, 2021. <https://www.visualcapitalist.com/the-50-most-visited-websites-in-the-world/>.
- Niu, Tong, Jianzhou Wang, Haiyan Lu, Wendong Yang, and Pei Du. 2020. "Developing a deep learning framework with two-stage feature selection for multivariate financial time series forecasting." *Expert Systems With Applications* 148.
- Ossinger, Joanna, and Jim Silver. 2021. *Bloomberg*. June 14. Accessed Oktober 17, 2021. <https://www.bloomberg.com/news/articles/2021-06-13/musk-says-tesla-sold-about-10-of-bitcoin-holdings#:~:text=Elon%20Musk%20said%20Tesla%20Inc,use%20about%2050%25%20clean%20energy.&text=Musk%20has%20whipsawed%20Bitcoin%20and,in%20the%20past%20few%20months>.
- Pano, Toni, and Rasha Kashef. 2020. "A Complete VADER-Based Sentiment Analysis of Bitcoin (BTC) Tweets during the Era of COVID-19." *In Big Data and Cognitive Computing (Vol. 4, Issue 4, p. 33)* 17.
- Pant, Dibakar Raj, Prasanga Neupane, Anuj Poudel, Anup Kumar Pokhrel, and Bishnu Kumar Lama. 2018. "Recurrent Neural Network Based Bitcoin Price Prediction by Twitter

- Sentiment Analysis." *IEEE 3rd International Conference on Computing, Communication and Security (ICCCS)* 6.
- Patel, Mohil Maheshkumar, Sudeep Tanwar, Rajesh Gupta, and Neeraj Kumar. 2020. "A Deep Learning-based Cryptocurrency Price Prediction Scheme for Financial Institutions." *Journal of Information Security and Applications* 55 (2020) 102583 12.
- Petukhina, Alla A., Raphael C. G. Reule, and Wolfgang Karl Härdle. 2021. "Rise of the machines? Intraday high-frequency trading patterns of cryptocurrencies." *The European Journal of Finance* 8-30.
- Phaladisailoed, Therasak, and Thanisa Numnonda. 2018. "Machine Learning Models Comparison for Bitcoin Price Prediction." *2018 10th International Conference on Information Technology and Electrical Engineering (ICITEE)*. Bali, Indonesia.
- Raschka, Sebastian. 2015. *Python Machine Learning*. Birmingham: Packt Publishing Ltd.
- Saxena, Anshul, and T.R. Sukumar. 2018. "Predicting bitcoin price using lstm And Compare its predictability with arima model." *International Journal of Pure and Applied Mathematics* 119(17): 2591-2600.
- Schafer, Cullen. 1993. "Technical Note: Selecting a Classification Method by Cross-Validation." *Machine Learning* 13: 135-143.
- Sebastião, Helder, and Pedro Godinho. 2021. "Forecasting and trading cryptocurrencies with machine learning under changing market." *Financial Innovation* 7:3.
- Shintate, Takuya, and Lukáš Pichl. 2019. "Trend Prediction Classification for High Frequency Bitcoin Time Series with Deep Learning." *Journal of Financial Risk and Financial Management* 1-15.
- Silver, Caleb. 2020. *Investopedia*. December 24. Accessed October 21, 2021. <https://www.investopedia.com/insights/worlds-top-economies/>.

- Sovbetov, Yhlas. 2018. "Factors Influencing Cryptocurrency Prices: Evidence from Bitcoin, Ethereum, Dash, Litecoin, and Monero." *Journal of Economics and Financial Analysis* 2(2): 1-27.
- Srishthi, Vashishtha, and Susan Seba. 2020. "Fuzzy Interpretation of Word Polarity Scores for Unsupervised Sentiment Analysis." *11th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*. Kharagpur, India: IEEE.
- Statista. 2021. "Biggest cryptocurrency exchanges based on 24h volume in the world on November 22, 2021." *Statista*. November 22. Accessed December 7, 2021. <https://www.statista.com/statistics/864738/leading-cryptocurrency-exchanges-traders/>.
- Staudemeyer, Ralf C., and Eric Rohstein Morris. 2019. "Understanding LSTM -- a tutorial into Long Short-Term Memory Recurrent Neural Networks."
- Stenqvist, Evita, and Jacob Lönnö. 2017. *Predicting Bitcoin price fluctuation with Twitter sentiment analysis*. Stockholm, Sweden.
- Swarnkar, Naman. 2020. *quantinsti*. May 01. Accessed October 12, 2021. <https://blog.quantinsti.com/vader-sentiment/>.
- Taboada, Maite, Julian Brooke, Milan Tofiloski, Kimberly Voll, and Manfred Stede. 2011. "Lexicon-Based Methods for Sentiment Analysis." *Computational Linguistics. MIT Press - Journals* 267–307.
- Tao, Ran, Chi-Wei Su, Yidong Xiao, Ke Dai, and Fahad Khalid. 2021. "Robo advisors, algorithmic trading and investment management: Wonders of fourth industrial revolution in financial markets." *Technological Forecasting & Social Change* 163 6.
- Timmins, Fiona, and Catherine McCabe. 2005. "How to conduct an effective." *art & science* 20(11): 41-47.

- Tsai, Chih-Fong, and Yu-Chieh Hsiao. 2010. "Combining multiple feature selection methods for stock prediction: Union, intersection, and multi-intersection approaches." *Decision Support Systems* 50: 258–269.
- Turner, Zane, Kevin Labille, and Susan Gauch. 2021. "Lexicon-based sentiment analysis for stock movement prediction." *Journal of Construction Materials. Institute of Construction Materials* 12.
- Twitter. 2021. *Fiscal Year 2020 Annual Report*. Annual Report, Twitter, Inc. Accessed November 20, 2021. [https://s22.q4cdn.com/826641620/files/doc\\_financials/2020/ar/FiscalYR2020\\_Twitter\\_Annual\\_Report.pdf](https://s22.q4cdn.com/826641620/files/doc_financials/2020/ar/FiscalYR2020_Twitter_Annual_Report.pdf).
- Twitter. 2021. *Twitter Developer*. Accessed December 14, 2021. <https://developer.twitter.com/en/docs/twitter-api/tweets/lookup/migrate/standard-to-twitter-api-v2>.
- Twitter. 2021a. *Developer Platform*. Accessed November 29, 2021. <https://developer.twitter.com/en/docs/twitter-api>.
- Twitter. 2021b. *Developer Platform*. Accessed December 2, 2021. <https://developer.twitter.com/en/docs/twitter-api/tweets/search/integrate/build-a-query>.
- Twitter. 2021c. *Developer Platform*. December 12. Accessed December 2021, 2021. <https://developer.twitter.com/en/pricing/search-30day>.
- Twitter. 2021d. *Developer Platform*. Accessed November 25, 2021. <https://developer.twitter.com/en/docs/counting-characters>.
- Twitter. 2021e. *Developer Platform*. Accessed November 25, 2021. <https://developer.twitter.com/en/docs/twitter-api/rate-limits>.

- Valencia, Franco, Alfonso Gómez-Espinosa, and Benjamín Valdés-Aguirre. 2019. "Price Movement Prediction of Cryptocurrencies Using Sentiment Analysis and Machine Learning." *Entropy* 12.
- Vo, Au, and Christopher Yost-Bremm. 2018. "A High-Frequency Algorithmic Trading Strategy for." *Journal of Computer Information Systems* 60(6): 555-568.
- Wall Street Journal. 2021. *Wall Street Journal*. Accessed 12 14, 2021. <https://www.wsj.com/market-data/quotes/TSLA/historical-prices>.
- Walther, Thomas, Klein, Tony, & Bouri, Elie. 2019. "Exogenous drivers of Bitcoin and Cryptocurrency volatility—A mixed data sampling approach to forecasting." *Journal of International Financial Markets*.
- Wang, Tai-Yue, and Huei-Min Chiang. 2007. "Fuzzy support vector machine for multi-class text categorization." *Information Processing and Management* 43: 914–929.
- Wei, Wang. 2016. *Achieving Inclusive Growth in China Through Vertical Specialization*. Elsevier Science & Technology.
- Williams, Sarah C., and Anita K. Foster. 2011. "Promise Fulfilled? An EBSCO Discovery Service." *Journal of Web Librarianship* 5(3): 179-198.
- Woebeking, Fabian. 2021. "Cryptocurrency volatility markets." *Digital Finance* 273–298.
- Wollams, Ben. 2021. *The Drum*. 09 06. Accessed November 25, 2021. <https://www.thedrum.com/opinion/2021/09/06/why-twitter-making-comeback-one-the-most-popular-social-channels>.
- Xiaolei, Sun, Liu Mingxi, and Sima Zeqian. 2020. "A novel cryptocurrency price trend forecasting model based on LightGBM." *Finance Research Letters* 32.
- Zhou, Liyi, Kaihua Qin, Christof Ferreira Torres, Duc V Le, and Arthur Gervais. 2021. "High-Frequency Trading on Decentralized On-Chain Exchanges." *IEEE Symposium on Security and Privacy (SP)* 18.