

University of Bonn  
Master of Science in Economics

# **Identifying Arbitrage in Cryptomarkets with Algorithmic Trading: A Machine Learning Approach**

Submitted by  
Raphael Redmer

Supervisor: TBA

June 1, 2020

## Abstract

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Data and Software</b>	<b>1</b>
2.1	Data . . . . .	1
2.2	Software . . . . .	1
<b>3</b>	<b>Methodology</b>	<b>2</b>
3.1	Training and Trading Set . . . . .	2
3.2	Feature and Target Generation . . . . .	2
3.3	Model Training . . . . .	3
3.3.1	Logistic Regression . . . . .	3
3.3.2	Random Forest . . . . .	3
3.3.3	Support Vector Machine . . . . .	3
3.3.4	Artificial Neural Network . . . . .	3
3.4	Trading . . . . .	3
<b>4</b>	<b>Results</b>	<b>3</b>
4.1	General Results . . . . .	3
4.2	Strategy Performance . . . . .	4
4.3	Further Analyses . . . . .	4
<b>5</b>	<b>Discussion</b>	<b>4</b>

# 1 Introduction

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

## 2 Data and Software

### 2.1 Data

In this setup, we use minute-binned OHLC data of crypto/USD-pairs obtained from the cryptocurrency exchange [Bitfinex](#) via its API ranging from 01.01.2019 to 31.12.2019. For each minute-bin, we collect *Open*, *High*, *Low*, *Close*, *Volume* and *Timestamp* data. *High* and *Low* denote the highest and lowest price respectively that was traded within this timeframe. *Open* and *Close* denote the first and last traded price. *Volume* denotes the total volume traded within the respective minute-bin. *Timestamp* denotes the point in time for each minute-bin as a UNIX-Timestamp, i.e., is the number of seconds that have passed since 01.01.1970 (Cite: [IEEE](#) and [The Open Group 2018](#)).

Even though Bitfinex is the largest exchange for cryptocurrency with a [market capitalization](#) and [different tradable coins](#), for most coins, the trading frequency is so low such that many crypto/USD-trading pairs have a considerable amount of minute-bins in which no volume was traded. In case of a crypto-pair having no volume for a particular minute, the API leaves out this bin when requesting its data resulting in missing bins. We resolved this issue by propagating price values from the last active minute-bin and setting the volume to zero. Further, we restricted the number of crypto-USD-pairs to the top ten pairs by market capitalization (see [plot](#)). In addition, we decided to only take data from 01.01.2019 to 31.12.2019, since for most coins 2019 was the most active year in terms of trading frequency. Thus, the resulting data set contains roughly  $10 \times 365 \times 24 \times 60 = 5.256.000$  rows.

### 2.2 Software

The programming language used for conducting this study is Python 3.7 ([cite](#)). For data preparation and feature engineering, we used Pandas and numpy ([cite](#)). Data Visualization was done via Matplotlib ([cite](#)). For the training of the models Logistic Regression, Random Forest and Support Vector Machine, we used the respective Scikit-learn implementation ([cite](#)). The Artificial Neural Network was trained using the Keras framework with Tensorflow backend to enable GPU calculation.

### 3 Methodology

Similarly to [Krauss et al. (2017)] and [cite arbitrage], the methodology of this paper consists of the following steps:

1. The entire data set is split into a training, a validation and a trading set.
2. The respective features (explanatory variables) and targets (dependent variables) are created
3. Each model is trained on the training set
4. Conduct out-of-sample predictions on the trading set for each model
5. Evaluate its accuracy and trading-performance on the trading set respectively
6. Go to Step 2, and repeat the same steps for a different feature- and target-specification

#### 3.1 Training and Trading Set

In our application to minute-binned data, the test set, i.e. trading set, contains all observations from 01.11.2019 to 31.12.2019. The training set ranges from 01.01.2019 to 14.09.2019, and the remaining 15.09.2019 to 31.10.2019 is reserved for the validation set. We decided against the usual k-fold cross-validation approach in order to emphasize the importance of future observation for the model, since its performance only gets evaluated on the future trading set.

#### 3.2 Feature and Target Generation

Broadly following [cite], we generate the feature space as follows:

**Input:** Let  $P^c = (P_t^c)_{t \in T}$  denote the price process of coin-USD-pair  $c$ , with  $c \in \{1, \dots, n\}$ . The price itself is the average between *Open* and *Close*.

**Features:** From the data set we obtain the following features:

**Returns:** Let  $R_{t,t-m}^c$  be the simple return for coin  $c$  over  $m$  periods defined as

$$R_{t,t-m}^c = \frac{P_t^c}{P_{t-m}^c} - 1 \quad (1)$$

**Volumes:** Let  $V_t^c$  be the traded volume for coin  $c$  in minute-bin  $t$  scaled by Quantile-Transformer fitted separately for each coin

**Target:** Let  $Y_{t+1,t}^c$  be a binary response variable for each coin  $c$ . It assumes value 1 (class *up*) if its future 120 min return  $R_{t+120,t+1}^c$  is greater than its cross-sectional median across all pairs  $(R_{t+120,t+1}^c)_{c=1}^n$ , else -1 (class *down*).

We decided for the inclusion of volume such that the model has a measure for taking trading activity into account without breaking vital assumptions needed for testing the 1.

Efficient Market Hypothesis (jcite<sub>l</sub>). In addition, the volume got scaled for each coin in order to make the measure more comparable across coins, since we are training a single universal model for each of the selected coins. Further, the Quantile-Transformation handles outliers (jcite<sub>l</sub>) and restricts the feature to an intervall ranging from 0 to 1.

### **3.3 Model Training**

#### **3.3.1 Logistic Regression**

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

#### **3.3.2 Random Forest**

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

#### **3.3.3 Support Vector Machine**

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

#### **3.3.4 Artificial Neural Network**

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

### **3.4 Trading**

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

## **4 Results**

### **4.1 General Results**

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

## 4.2 Strategy Performance

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

## 4.3 Further Analyses

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

## 5 Discussion

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.