
THE EFFICIENT MARKET HYPOTHESIS FOR BITCOIN IN THE CONTEXT OF NEURAL NETWORKS

A PREPRINT

Mike Kraehenbuehl
School of Engineering
Zurich University of Applied Sciences
Winterthur, Switzerland
kraehmik@students.zhaw.ch

Joerg Osterrieder *
School of Engineering
Zurich University of Applied Sciences
Winterthur, Switzerland
joerg.osterrieder@zhaw.ch

August 16, 2022

ABSTRACT

This study examines the weak form of the efficient market hypothesis for Bitcoin using a feedforward neural network. Due to the increasing popularity of cryptocurrencies in recent years, the question has arisen, as to whether market inefficiencies could be exploited in Bitcoin. Several studies we refer to here discuss this topic in the context of Bitcoin using either statistical tests or machine learning methods, mostly relying exclusively on data from Bitcoin itself. Results regarding market efficiency vary from study to study. In this study, however, the focus is on applying various asset-related input features in a neural network. The aim is to investigate whether the prediction accuracy improves when adding equity stock indices (S&P 500, Russell 2000), currencies (EURUSD), 10 Year US Treasury Note Yield as well as Gold&Silver producers index (XAU), in addition to using Bitcoin returns as input feature. As expected, the results show that more features lead to higher training performance from 54.6% prediction accuracy with one feature to 61% with six features. On the test set, we observe that with our neural network methodology, adding additional asset classes, no increase in prediction accuracy is achieved. One feature set is able to partially outperform a buy-and-hold strategy, but the performance drops again as soon as another feature is added. This leads us to the partial conclusion that weak market inefficiencies for Bitcoin cannot be detected using neural networks and the given asset classes as input. Therefore, based on this study, we find evidence that the Bitcoin market is efficient in the sense of the efficient market hypothesis during the sample period. We encourage further research in this area, as much depends on the sample period chosen, the input features, the model architecture, and the hyperparameters.

Keywords Neural Networks, Feedforward, Bitcoin, Efficient Market Hypothesis

*Financial support by the Swiss National Science Foundation within the project “Mathematics and Fintech - the next revolution in the digital transformation of the Finance industry” is gratefully acknowledged by the corresponding author. This research has also received funding from the European Union’s Horizon 2020 research and innovation program FIN-TECH: A Financial supervision and Technology compliance training programme under the grant agreement No 825215 (Topic: ICT-35-2018, Type of action: CSA) and from Innosuisse under the grant agreement “Strengthening Swiss Financial SMEs through Applicable Reinforcement Learning” (Innosuisse Innovation project 47959.1 IP-SBM). Furthermore, this article is based upon work from the COST Action 19130 Fintech and Artificial Intelligence in Finance, supported by COST (European Cooperation in Science and Technology), www.cost.eu (Action Chair: Joerg Osterrieder).

The authors are grateful to the management committee members of the COST (Cooperation in Science and Technology), Action Fintech, and Artificial Intelligence in Finance as well as speakers and participants of the 4th, 5th, and 6th European COST Conference on Artificial Intelligence in Finance and Industry.

Contents

1	Introduction	3
2	Literature Overview	3
2.1	Efficient Markets	3
2.2	Bitcoin and Efficiency	3
2.3	Machine Learning in Financial Markets	4
3	Data and Summary Statistics	5
3.1	Features	5
3.2	Missing Values	7
4	Experiment Approach	8
4.1	Feedforward Neural Networks	8
4.2	Data Preprocessing	9
4.3	Neural Network Architecture	10
5	Results	12
5.1	Accuracy & Loss	13
5.2	Performance	15
6	Discussion	15
7	Conclusion	16
	References	17
8	Appendix	19
8.1	Plots	19
8.2	Results of individual feature sets	21
8.2.1	ETH same day (to test model performance)	21
8.2.2	Feature set 0: random sample "rNorm"	22
8.2.3	Feature set 1: Bitcoin	23
8.2.4	Feature set 2: Bitcoin + S&P500	24
8.2.5	Feature set 3: Bitcoin + S&P500 + Russell 2000	25
8.2.6	Feature set 4: Bitcoin + S&P500 + Russell 2000 + EUR/USD	26
8.2.7	Feature set 5: Bitcoin + S&P500 + Russell 2000 + EUR/USD + 10Y Treasury Yield	27
8.2.8	Feature set 6: Bitcoin + S&P500 + Russell 2000 + EUR/USD + 10Y Treasury Yield + Gold/Silver	28

1 Introduction

Cryptocurrencies have become popular over the past few years, both among retail and institutional investors. Many studies, but not all, conclude that traditional markets such as stocks or commodities are highly efficient. Therefore it has become tricky to outperform these markets [1]. However, since crypto markets have emerged only in the past ten years, there is a possibility that the markets for this new asset class are not (yet) efficient.

In this study, we challenge the weak form of the efficient market hypothesis (EMH) for Bitcoin with the help of a neural network. This has already been tried in a large number of works with different machine learning techniques, various features, and different outcomes [2], [3], [4], [5], [6]. However, most of the attempts deal with features that are directly derived from the asset under investigation - namely technical features. Less focus is placed on analyzing the connection of other financial instruments with Bitcoin. By examining the following two research questions, we try to contribute to this field of research:

RQ 1: Do additional input features improve the prediction accuracy?

RQ 2: Is the Bitcoin market efficient as per the weak form of the efficient market hypothesis?

First, we give an overview of the relevant literature about the EMH. Second, we review some of the papers that examined Bitcoin and its efficiency. Third, we point out recent papers that used machine learning techniques to predict Bitcoin, as well as other asset classes. Finally, we investigate the research questions with a feedforward neural network with the main focus on analyzing if additional input features improve prediction accuracy.

2 Literature Overview

2.1 Efficient Markets

Fama introduces the efficient market hypothesis (EMH) in 1965 [7], [8]. Another contributor to the EMH is Samuelson [9]. According to Fama, markets are efficient, only if the prices always “fully reflect” all available information. Further, Fama separates the theory into three parts [10]: the weak form (prices reflect all information about past prices), semi-strong (prices reflect all public information: earnings, fundamentals, etc.), and strong form (all private information is reflected in the prices). The following statements can be derived from this:

- Market prices are accurate.
- A consistent outperformance by using already known information is not possible.

Fama describes the sufficient conditions for a capital market to be efficient as follows: (1) no transaction costs, (2) all information is available to all market participants without any costs, and (3) all market participants agree on the implications of the current information for the current price and the distributions of future prices of each security. Such a market would be efficient. These conditions are sufficient but not necessary for market efficiency.

It is also important to note that most economists agree that the EMH is not absolute and has its limitations [11]. Charlie Munger thinks that it is “roughly right” that the stock market is efficient, but not completely efficient [12]. The efficiency in the markets comes from its participants, the investors and traders, who try to make a profit out of every potential information advantage. This causes the prices to include all the information and fully reflect the intrinsic value of its underlying. As markets become more efficient, their returns start to become more and more unpredictable. Grossman and Stiglitz (1980) [13] conclude that markets can never be fully efficient. If a market was efficient, no trader would make the effort to gather exclusive information since no outperformance is possible anyway. On the other hand, for a market to be efficient, there have to be enough participants to drive prices to the intrinsic values and consequently make them efficient. According to Lo (2004) [14], the efficient market hypothesis is incomplete because, during times of market distortions, investors act irrationally which creates inefficiencies that can be exploited by others. To describe this circumstance, he proposes the adaptive market hypothesis (AMH). On the other hand, Hsieh (2010) [15] claims that market participants may be irrational at some points (e.g. by being systematically biased), but the market can still be efficient, as long as the rational market participants cannot predict the price movements.

2.2 Bitcoin and Efficiency

Nakamoto first describes the concept of Bitcoin in his paper in 2008 [16]. In the following years, Bitcoin becomes more and more popular and scientists start to study this new asset class. Garcia et al. (2014) [17] study social feedback related to Bitcoin. Hencic and Gourieroux (2014) [18] predict Bitcoin with a non-causal regressive model. Sapuric and Kokkinaki (2014) [19] measure the different volatilities of Bitcoin against several major exchange rates of currencies. Kristoufek (2015) [20] studies how the price of Bitcoin is formed and what the drivers of its price are. Osterrieder and Lorenz (2017) [21] conclude that the Bitcoin return distribution exhibits higher volatility, stronger non-normal characteristics, and heavier tails than those of the G10 currencies. Chan et al. (2017) [22] fit the exchange rates of the largest (by market capitalization) cryptocurrencies against the USD to various parametric distributions.

They find that none of the cryptocurrencies analyzed fit a normal distribution and that for Bitcoin the best fit is the hyperbolic distribution. Enoksen et al. (2020) [23] research the emergence of bubbles in the Bitcoin market and other cryptocurrencies. Most financial assets were already examined for their efficiency. However, these findings may not be valid for the new asset class of cryptocurrencies. One of the first to investigate whether the Bitcoin price behaves as it would in an efficient market is Urquhart (2016) [24]. The author conducts various statistical tests like Ljung-Box, AVR test, etc. over the years from 2010 until 2016. He concludes that the Bitcoin market is not weakly efficient during this period. However, it appears that the market becomes less inefficient during the second half of the test window and returns become random. A study that follows after Urquhart (2016) is "On the inefficiency of Bitcoin" by Nadarajah and Chu (2017) [25]. The authors use the same data as Urquhart but continue with a power transformation of the Bitcoin returns. With the transformed data the tests used by the authors show that the Bitcoin market is indeed efficient. Bariviera (2017) [26] uses the Hurst exponent to examine daily Bitcoin prices from August 2011 to February 2017. He concludes that the Bitcoin market becomes, at least weak form, efficient after 2014. Brauneis and Mestel (2018) [27] investigate the efficiency of a large number of cryptocurrencies. They find that efficiency is positively related to liquidity. Among all cryptocurrencies tested, Bitcoin is the one to pass the most tests of price randomness. Bitcoin is the least predictable and therefore the most efficient. Tiwari et al. (2018) [28] find that the Bitcoin market is efficient except for the periods of April–August 2013 and August–November 2016. Al-Yahyaee et al. (2018) [29] show that among stocks, gold, and currencies, Bitcoin is the least efficient. Sensoy (2018) [30] also analyzes the weak form of the EMH for Bitcoin and finds that on an intraday level Bitcoin becomes more informationally efficient. Additionally, he states that BTCUSD is slightly more efficient than BTCEUR in the sample period. Two important findings from the study are that the higher the frequency, the lower the pricing efficiency is, and liquidity (volatility) has a significantly positive (negative) effect on the efficiency of Bitcoin prices. Studies like [24] conclude that Bitcoin becomes more efficient over time and also the correlation with stock indices increases. This pattern is not surprising as young markets, like some emerging markets, tend to become more efficient over time as more investors and traders enter the market and price in the information that they believe to have. Vidal-Tomás and Ibañez (2018) [31] study the semi-strong form of efficiency. They find that Bitcoin does not respond to measures of central banks, but becomes increasingly efficient in regards to negative and positive news related to Bitcoin itself. Khuntia and Pattanayak (2018) [32] test the adaptive market hypothesis on the Bitcoin market where they verify the evolving efficiency in the Bitcoin market. Major events coincide with times of efficiency and inefficiency, therefore the AMH rather than the EMH seems to apply to the behavior of the Bitcoin price. Bitcoin is probably also less efficient/rational than the stock and bond markets because of the lack of intrinsic value. Investors can derive a fair price for a stock at least in part from the company's key figures, revenues, profits, etc. This is not possible for cryptocurrencies to the same extent, and therefore much more speculation is involved. However, this can be a sign of irrationality, but not necessarily of inefficiency.

All in all, there is mixed evidence for the predictability of Bitcoin. It seems as if the Bitcoin market is becoming more efficient as it matures.

2.3 Machine Learning in Financial Markets

Machine learning (ML) methods enjoy increasing popularity in applications used for empirical asset pricing due to their ability to learn complex relationships between features and target. Gu et al. (2020) [33] identify neural networks and regression trees are the best performing methods due to their advantage of accommodating nonlinear interactions in contrast to other ML methods. As the most powerful predictors for stock prices, the authors mention the ones showing price trends including return reversal and momentum, liquidity, volatility, and valuation ratios. Wildi and Bundi (2019) [2] analyze the inefficiency of Bitcoin using several methods including ARMA and neural networks. They use a feedforward neural net, with Bitcoin lags 1-6 as inputs. The architecture with only two hidden layers containing 6 and 3 nodes respectively is rather simple. However, they are still able to generate an outperformance when looking at the ensemble average of random nets. Overall they strongly reject the EMH for Bitcoin. In contrast to other before-mentioned studies, they observe increasing market inefficiency towards the end of the sample period which is 2019-01-10. Huang et al. (2019) [3] use Bitcoin price-based high dimensional technical indicators in their research. With a classification tree-based model and 124 technical indicators, they provide evidence for strong out-of-sample performance for narrow ranges of daily Bitcoin returns. They find that technical analysis can be useful when predicting returns for an asset like Bitcoin that is hard to value by fundamentals. Mudassir et al. (2020) [4] use the four models ANN, SANN, SVM, and LSTM to research different prediction horizons for Bitcoin. The classification models for next day forecasts score up to 65% accuracy. Jaquart et al. (2021) [5], however, criticize the authors for creating a potentially unbalanced training dataset. For financial time series, which are usually noisy, an unbalanced dataset can lead to a majority prediction of a class that is overrepresented in the training dataset. If the same class is also overrepresented in the test dataset, this can automatically lead to higher accuracy. Jaquart et al. express similar critics to Chen et al. (2020) [6] that find that simpler methods like logistic regression perform better in the 1-day forecast than more complex ones like recurrent neural networks. In their work, Jaquart et al. (2021) [5] analyze the intraday Bitcoin price on prediction horizons from 1 to 60 minutes. Their approach is to use different machine learning models like recurrent networks and gradient boosting among others. Moreover, they use four different feature sets; technical features derived from

historical Bitcoin data, blockchain-related features, sentiment-/interest-based features, and asset-based features other than Bitcoin (e.g. gold and MSCI World index returns). Their findings are that technical features are most important for most methods. An outperformance against a random classifier is achieved but this turns into negative returns after accounting for transaction costs. Therefore, this study does not violate the efficient market hypothesis. Feng and Giglio (2020) [34] introduce a methodology to systematically evaluate the marginal contribution of a new factor relative to the existing factors in explaining asset prices. They use techniques like the double-selection LASSO method [35] and two-pass regressions such as Fama-MacBeth for their procedure. The proposed method can be used to identify a good set of potentially market-predictive features.

Comparing performances of predictive models from different works is often difficult. The reason for this is that one can use different time horizons, target and explanatory variables, parameter specifications, and evaluation metrics.

3 Data and Summary Statistics

This section will discuss the input data that is used for the experiment in section 4. With an artificial neural network, we use different feature sets to predict the Bitcoin price movement for the next day. The predictive performances of the feature sets are subsequently compared.

3.1 Features

There are several studies that use technical features derived from its underlying to predict Bitcoin, and they show some evidence for their predictive power. [5], [6] use other asset-based features such as major stock indices, interest rates, or currency market data. With our study, we would like to contribute to research in this direction and therefore investigate the six features in Table 1 to validate the two research questions stated at the beginning.

The daily prices of selected features are imported from Yahoo Finance (YF) with a starting date of 2014-09-17. This period should be sufficient to answer the research questions because, as mentioned in the previous section, some studies conclude that Bitcoin was not efficient in the early years and has become more efficient over time.

Feature	Source	YF Symbol	Timestamp
BTC/USD	CoinMarketCap	BTC-USD	23:59 UTC [36]
S&P 500	Standard & Poors	^GSPC	20:00 UTC [37]
Russell 2000	FTSE Russell	^RUT	20:00 UTC [37]
EUR/USD	ICE Data Services	EURUSD=X	21:30 UTC [38]
10Y Treasury Note Yield Index	CBOE	^TNX	19:00 UTC [39]
PHLX Gold&Silver Index	NASDAQ	^XAU	21:30 UTC [40]

Table 1: Selected input features and source.

Bitcoin is the largest cryptocurrency by market cap and accounts for around 45% of the total USD 1.28 trillion in the crypto market [41]. On the one hand, we will use Bitcoin as a target for our predictive model, and on the other hand, we will also use past Bitcoin returns as input features. CoinMarketCap reports the Bitcoin price (BTC/USD) by calculating a volume-weighted average of all market pair prices reported for Bitcoin [42]. The prices are reported from crypto exchanges like Binance, Coinbase, and FTX among others.

The S&P 500 Index represents the market price of the 500 biggest (by market cap) companies in the United States. Its value is reported whenever the NYSE is open for trading. The closing price of the S&P 500 Index is calculated according to the closing prices and the weightings of the included stocks. Kim et al. (2020) [43] show in an empirical study that there exists a time-varying relationship between the Bitcoin market and the S&P 500. Therefore the S&P 500 price may also provide information to predict future Bitcoin price movements.

Fama and French (1993) [44] outline several risk factors in the returns of stocks and bonds. One of the factors is the size, or that small companies outperform large ones. Therefore, the Russell 2000 Index will be used in the experiment to account for this risk factor. The Russell 2000 Index is a subset of the 2000 smallest companies that are represented in the Russell 3000 Index [45]. The companies represented by the Russell 2000 Index only account for 10% of the total market capitalization of the Russell 3000 Index. Its closing price is calculated daily based on the closing prices at the NYSE of the stocks contained in the index.

Currencies may also add value to the prediction of cryptocurrencies. The EUR/USD exchange rate is among the biggest and most liquid currency pairs.

The CBOE 10 Year Treasury Note Yield Index (TNX) is based on the yield-to-maturity of the most recently auctioned 10 Year Treasury Note.

In the before-mentioned study by Kim et al. (2020) [43], they not only show a relationship between Bitcoin and the US stock market but also with gold. For reproducibility reasons, we use the PHLX Gold and Silver Index (XAU). We calculate a 0.95 Pearson correlation with the CME Globex Gold Future on 100 troy ounces [46] (YF-symbol: GC=F). If gold and silver prices share a dependency with Bitcoin, the XAU index should also serve as a good input feature for predicting Bitcoin. The Philadelphia Stock Exchange Gold and Silver Index is a capitalization-weighted index that includes the leading companies involved in the mining of gold and silver [40]. All securities in the index must be listed on The Nasdaq Stock Market, New York Stock Exchange, NYSE American, or Cboe BZX Exchange [47], [48].

Furthermore, we analyze the log returns of each input feature as seen in equation 1. The main reason for using log returns is that they are stationary and additive. The neural network is then trained with normalized log returns. To describe the distributions of the features, the skewness and kurtosis of each features log returns are calculated in Table 2.

$$\log(1 + r_t) = \log\left(\frac{P_t}{P_{t-1}}\right) = \log(P_t) - \log(P_{t-1}) \quad (1)$$

Feature	mean	min	max	skewness	kurtosis
BTC/USD	0.0015	-0.465	0.225	-0.767	11.112
S&P 500	0.0005	-0.099	0.089	-0.217	13.434
Russell 2000	0.0003	-0.118	0.089	-0.696	9.545
EUR/USD	0.0000	-0.028	0.028	0.069	3.236
10Y Treasury Yield	0.0003	-0.271	0.404	1.564	28.149
Gold/Silver	0.0002	-0.146	0.144	-0.059	3.314

Table 2: Stats of log returns.

Bitcoin is notorious for its high volatility compared to other assets. Among the selected features, Bitcoin had the highest daily loss of -0.465. The Fisher kurtosis of 11.112 also shows that the tails are thicker than in the normal distribution (normal distribution has a kurtosis of 0.00). Bitcoin, S&P 500, Russel 2000, and the Gold/Silver index have negative skewness. On the other hand, the 10Y Treasury Yield has a high positive skewness of 1.564 and the highest kurtosis of 28.149. Between 2020-03-06 and 2020-03-10, both the highest negative log return of -0.271 and the highest positive log return of 0.404 occurred. The S&P 500 and Russell 2000 distributions are very similar, which is not surprising as both indices represent US companies. Among the assets examined, the EUR/USD and Gold/Silver returns most closely resemble a normal distribution. Fig. 1 shows the distributions of the log returns.

The correlation heatmap in Fig. 17 (appendix) shows a strong correlation between the Russell 2000 and S&P 500 Index since both indices represent US companies. Other combinations show a moderate positive correlation with two exceptions. First, the 10Y Treasury Yield appears to not correlate with Bitcoin. Second, we can see a negative correlation of Gold/Silver with the 10Y Treasury Yield. This plot only shows a simple correlation and not a possible time-varying nonlinear dependence between the different features. In the appendix 8.1, we further discuss Bitcoin's autocorrelation and a reason why we do not consider it in this study.

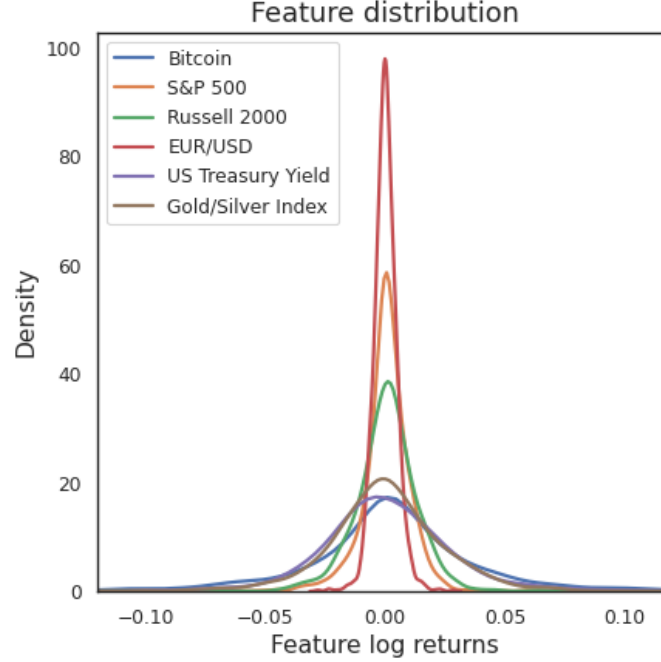


Figure 1: Density plot of the log returns of the features.

3.2 Missing Values

All six input features described above are combined into one dataset. The time series range from 2014-09-17 until 2022-04-30. At this point, the dataset has missing values for all features except Bitcoin, which is traded every day without exception. All other assets have no prices published on weekends and public holidays. Additionally, because we use log returns as inputs for the model in the experiment, we produce missing values on Mondays and all other dates with no value for the previous day because no one-day log returns can be calculated. There are multiple ways how to deal with missing values in machine learning. One method is to replace them with the last known value but this will distort the data. Another common approach is to interpolate to estimate the missing value. In practice, this is not very feasible because only the past values are known and this would also falsify the data. Since a neural network cannot be trained with missing values, the most reasonable consequence, therefore, is to delete all rows in the dataset, containing assets with missing values. After doing so, the initial size of the dataset shrinks from 2783 rows to 1495 rows. Furthermore, because we want to predict the next day's price movement of Bitcoin, the input dataset with all selected features contains data delayed by one day (lag 1). The problem of missing values intensifies when lag 2 is added to the input dataset. For illustration, if we take the target (Bitcoin signal) for 2022-03-23 (Wednesday) and choose the S&P 500 as the input feature, we have a value for 2022-03-22 (lag 1) but no value for 2022-03-21 (lag 2). To explain the current Bitcoin price movement, there has to be a lag 1 and lag 2 value for every input feature. If we create a feature dataset that besides lag 1 also contains lag 2, the dataset will shrink to 1092 rows. With each additional lag added for features other than Bitcoin, more rows contain missing values, therefore more rows have to be deleted which makes the already sparse daily data even smaller. Therefore we stay with lag 1 for this experiment. This is unfavorable for training but at least training and predictions are based upon real values. Estimating missing data would distort our analysis in particular the comparison between different feature sets. Using intraday data would be an interesting option as more data points are available and the problem becomes less extensive.

4 Experiment Approach

With an artificial neural network, the next day's Bitcoin price movements (signals) shall be predicted. **Neural networks have the advantage that they can deal with nonlinear dependencies.** With this experiment, we want to analyze the following research questions:

RC 1: Do additional input features lead to a better result? This will be measured by comparing the effects of different feature sets on prediction accuracy.

RC 2: Is the Bitcoin market efficient under the weak form of the EMH? If the correctly predicted price movements are statistically significant, this will be a sign that the Bitcoin market is not efficient.

The aim of this experiment is not to parameterize the model in such a way as to generate the best possible outperformance but to examine these two questions under the chosen scenario.

Experiment process

1. Import daily prices of selected features from Yahoo Finance
2. Calculate log returns
3. Create six different sets of input features
4. Shift (lag) features by one day
5. Drop all rows with missing values
6. Normalize lagged features to a range between 0-1
7. Create binary Bitcoin signals as the target variable
8. Split data into training-, validation-, and test-set
9. Train feedforward neural network with training set to predict next days Bitcoin signal
10. Evaluate model out-of-sample with testing set

4.1 Feedforward Neural Networks

A neural network is designed to mirror the functioning of the human brain. It contains many nodes that are interconnected in a specific way. The individual neurons receive information (input) from other neurons, transform the information and then pass it on (output) to other connected neurons. The basic idea is that the network attempts to construct a function $f(x)$ that maps a given input X to a given target value y . This is an example of supervised learning, but there are also neural networks that can be used for unsupervised learning [49]. An ever-increasing number of types of neural networks exist with each having its advantages; e.g. perceptron, multilayer perceptron (MLP), feedforward neural network (FNN), convolutional neural networks (CNN), recurrent neural networks (RNN), long short-term memory (LSTM) [50] and many more. For this experiment, we chose the rather basic feedforward neural network since the focus is not on optimizing to the best possible performance, but rather on comparing different feature sets. Fig. 3 shows the information flow in an FNN. In contrast to recurrent neural networks, feedforward networks only pass information in one direction (forward). The input layer represents the explanatory variables or features (x) that are entered into the system. After that comes any number of hidden layers with each containing any number of neurons. These layers are responsible for processing, transforming, and connecting information. Finally, there is an output layer that generates the predicted y . For the network to differ from a normal linear regression, at least one layer must have an activation function $g(x)$. This function must be either nonlinear or at least only partially linear function, but never completely linear. Equation 2 describes the transformation of information from layer $(n-1)$ to the next layer n [5].

$$a^{(n)} = g^{(n)}\left(W^{(n)}a^{(n-1)} + b^{(n)}\right) \quad (2)$$

where:

$a^{(n)}$ is the output of layer n

$g^{(n)}$ is the activation function

$W^{(n)}$ is the weight matrix that weight the values $a^{(n-1)}$ that have been forwarded from layer $n-1$ to layer n

$b^{(n)}$ is the bias for layer n

4.2 Data Preprocessing

Features

Data is imported with the Python package yfinance from Yahoo Finance so that the code for the neural network can easily be executed by any user without the need for external datasets. The used data starts on 2014-09-17 and ends on 2022-04-30. The imported dataset of Bitcoin and the selected features consist of prices. At this stage, rows that include N/A values are kept in the dataset. Note that Bitcoin is traded every day. Stocks and currencies on the other hand have missing values at weekends and public holidays. To answer the question of additional input features leading to higher prediction accuracy, the sets of input features outlined in Table 3 are created. Each feature set contains the log returns of the input features shifted by one day (lag 1) and all rows with missing values are deleted. For the experiment, each feature set will be trained separately with the neural network, and predictions are made individually. Then the results metrics from the different feature sets will be compared. We decide to only use lag 1 for two reasons. Firstly, if more lags are added, the dataset is reduced due to more missing values in the rows. For example, if lag 2 is added, the dataset is reduced by 27% from 1495 rows to 1092 rows. Since the return from the previous day will most likely be most important, we only choose lag 1 as the input feature. The reason for this decision is described in more detail in the appendix section 8.1.

Feature set	Included features (lag 1)
Feature set 1	Bitcoin
Feature set 2	Bitcoin, S&P 500
Feature set 3	Bitcoin, S&P 500, Russell 2000
Feature set 4	Bitcoin, S&P 500, Russell 2000, EUR/USD
Feature set 5	Bitcoin, S&P 500, Russell 2000, EUR/USD, US Treasury Yield
Feature set 6	Bitcoin, S&P 500, Russell 2000, EUR/USD, US Treasury Yield, Gold+Silver

Table 3: Feature sets containing log returns shifted by 1 day (lag 1).

The feature sets now consist of shifted log returns without missing values. It is common practice to normalize the input features for the training of neural networks. For this study, the log returns of the features are normalized to a range between 0-1 with the **MinMaxScaler** [51] equation 3. The distribution of the data remains the same.

$$X_{sc} = \frac{X - X_{min}}{X_{max} - X_{min}} \quad (3)$$

Target - binary Bitcoin signal

The goal of the neural network in this experiment is to predict the Bitcoin price movement of the next day. This will be defined as a binary classification problem in equation 4, where 0 represents a negative Bitcoin return and 1 represents a positive Bitcoin return. For further analysis or in practice applications, 0 can stand for sell and 1 for buy. Or short or long. It would have to be decided what it means when implementing a trading strategy. In this experiment, we use a long-only approach where the threshold between a positive and a negative signal is a probability of 0.5. One could also define that a long position is taken if the probability of a rising price is above 65%, and a short position is taken if the probability is less than 35%. In this study, this is ignored when analyzing the effects of the features on accuracy.

$$y_t = \begin{cases} 1 & , \text{ if return } \geq 0 \\ 0 & , \text{ if return } < 0 \end{cases} \quad (4)$$

where y_t is the Bitcoin signal

Train- validation- and test-set

Now that the different feature sets are defined, the data is split into separate datasets to meet the requirements of the neural network algorithm. In this case, a train level of 60% is applied, which means that the first 60% of data is used for training. For this purpose, the two datasets x_{train} and y_{train} are generated, which contain the normalized log returns of the selected input features shifted by one day and the binary target variable, respectively (Fig. 2). The remaining 40% are used for testing the out-of-sample performance. The test dataset is named x_{test} and y_{test} . Additionally, we define a validation split of 10%. Therefore, 10% of the training data is used as a validation set for hyperparameter tuning during the training process. The distribution of values may differ between train, validation, and test dataset. Since in practice this will most likely also be the case we do not see this as a major problem but it is discussed further with the results in section 5.1.

y_train		x_train					
Bitcoin	Signal	Bitcoin	S&P500	Russell2000	EURUSD	TrasuryYield	GoldSilver
	0	0.565	0.553	0.591	0.350	0.418	0.453
	1	0.569	0.525	0.518	0.597	0.377	0.416
	0	0.790	0.497	0.524	0.514	0.383	0.561
	0	0.631	0.568	0.610	0.503	0.421	0.466
	0	0.633	0.441	0.492	0.395	0.367	0.501

Figure 2: Example of the dataset with binary target and normalized explanatory variables used for training.

4.3 Neural Network Architecture

There exist multiple different types of neural networks, each having its advantages in certain application areas. For time series the most suitable types are feedforward, recurrent, and LSTM. Because the aim of this experiment is not to generate the best possible performance on the sample we use a simple feedforward neural network instead of more complicated models. Wildi and Bundi [2] use a feedforward net with two hidden layers with dimensions six and three. Even with such a simple architecture, the neural network is able to outperform a buy-and-hold strategy, an EqMa(6), and an MA(6) model. The following neural network, accessible on GitHub [52], is programmed in Python by using the Keras API [53]. The neural network is set up with the architecture in Table 4.

Layer	Neurons
Input Layer	with dimension according to the amount of input features
First hidden Layer	30 Neurons with ReLu activation function
Second hidden Layer	15 Neurons with ReLu activation function
Third hidden Layer	5 Neurons with ReLu activation function
Output Layer	1 Neuron with Sigmoid activation function

Table 4: Layers, neurons, and activation functions used for the experiment.

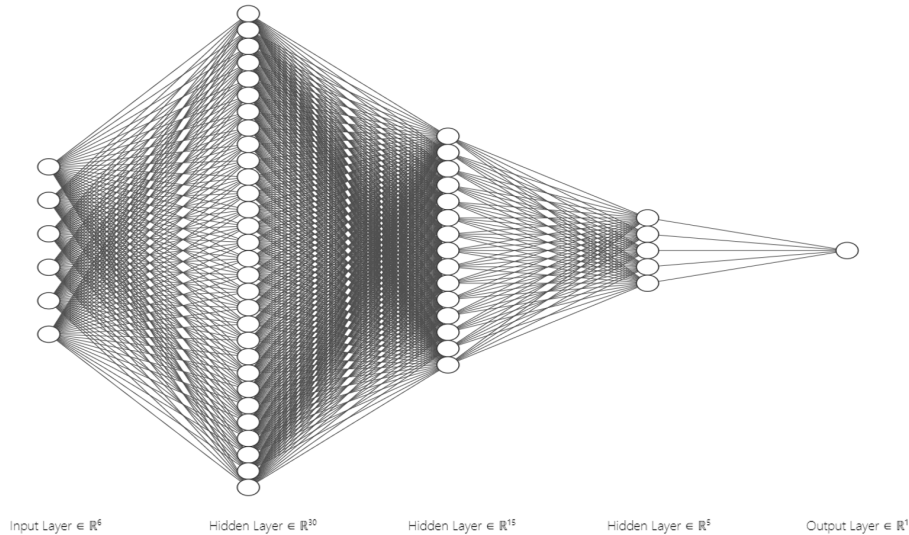


Figure 3: Architecture used for the model [54].

ReLU & Sigmoid activation functions

The Rectified Linear Units (ReLU) function 5 visualized in Fig. 4 is widely used in deep learning applications. Despite its simplicity, the function can account for interactions and nonlinearities. Traditional activation functions like tanh and sigmoid can lead to the problem of a vanishing gradient if they are used in many hidden layers. Because their derivatives are very small at most points of the curve, it makes it difficult to change the weights when applying gradient descent. ReLU helps reduce this problem since its derivative is either 0 or 1. This is why ReLU is commonly used as an activation for hidden layers [55].

$$\text{ReLU}(z) = \max(0, z) \quad (5)$$

The sigmoid function 6 visualized in Fig. 5 maps any real number to a value between 0 and 1. One could also say that it is a function that maps a weighted sum to probabilities. As the sigmoid function is used in the output layer, the output of the sigmoid function is between 0 and 1 and the threshold function 7 maps the value $\sigma(z)$ to 1 if the predicted probability for the Bitcoin signal being positive is greater than 0.5, and vice versa for a negative signal. In this case the prediction 0 means "not-in-the-market" and 1 means "in-the-market".

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (6)$$

$$t(p) = \begin{cases} 1 & p \geq 0.5 \\ 0 & p < 0.5 \end{cases} \quad (7)$$

where:

z is the weighted sum from equation 2

p is the probability output from the sigmoid function 6

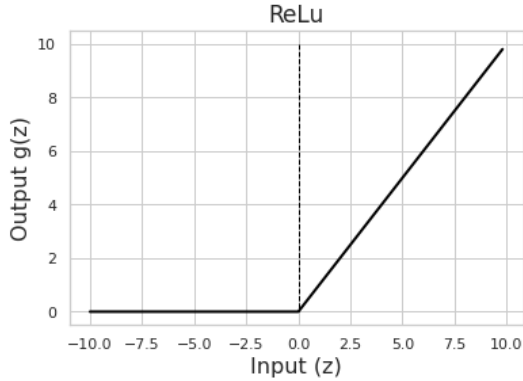


Figure 4: ReLU activation function used in hidden layers.

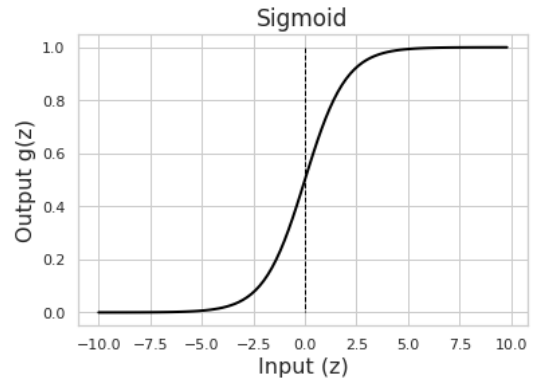


Figure 5: Sigmoid activation function used in the output layer.

Loss function - Binary cross-entropy

Different loss functions can be used for training a neural network. The purpose of the loss function is to calculate the error of a given state. A common loss function used for regression is the mean squared error (MSE). The binary cross-entropy (equation 8), also called **log-loss**, is another loss function, which is designed to measure the performance of a **binary classification problem**. This function is applied in the output layer after the sigmoid but before the threshold function. The input for the cross-entropy loss function is a value between 0 and 1, or a probability value for the signal being positive. Since signals (0 and 1) have to be predicted in this experiment, the cross-entropy loss function is a better choice than the MSE.

Example: If the target signal is 1 (positive log return of Bitcoin), the cross-entropy loss function illustrated in Fig. 6 calculates $-\log(a)$ as loss value. If the output of the sigmoid function is near 0, the loss will be large. This means that a wrong prediction that is near 0 instead of 1, will be punished strongly. On the other hand, it will drop drastically if it is

close to 1. This way, wrong predictions cause a high loss and will therefore update the weights strongly. A small loss will update the weights not much. Fig. 7 shows the curve when the target is 0.

$$\mathcal{L}(a, y) = -[y \log(a) + (1 - y) \log(1 - a)] \quad (8)$$

where:

y is the binary target class (0 or 1)

a is the predicted probability that z belongs to class 1 (equation 6)

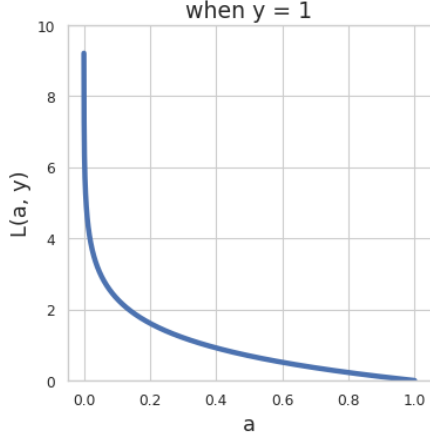


Figure 6: Cross-entropy when target signal is 1.

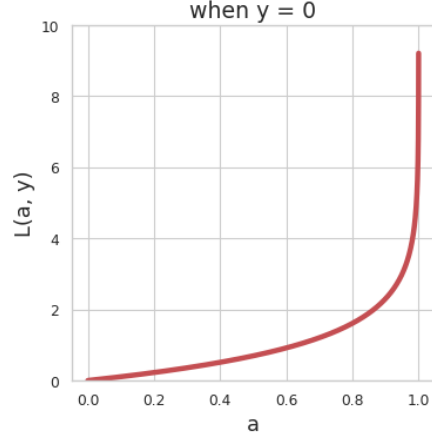


Figure 7: Cross-entropy when target signal is 0.

Optimizers and hyperparameters

As large as the selection of network types and architectures is, there is also an enormous choice of hyperparameters and optimizers. ADAM and BatchNormalisation is used in the model. Further, the model is trained with 300 epochs and a batch size of 20 at a learning rate of 0.00005.

5 Results

Influence of randomness

The neural network's algorithm starts with initial weights that are adjusted during training according to the training data and the defined network specifications. Depending on how these weights are initiated at the beginning, the outcome of the final weights and therefore the predictions and performance on the test dataset will be different. Therefore, one could just run the model many times until a good set of initial weights was initiated so that the performance on the test set will also be good. This would distort the meaning of out-of-sample testing. Therefore, for this experiment, each feature set is trained and tested 50 times. For later result comparison of the feature sets, we then take averages of the metrics of the 50 iterations. This is how we try to make a representative comparison. To show the importance of randomness, we take a look at the best (iteration 36/50) and worst (iteration 40/50) model outcomes of feature set 2 (Bitcoin lag 1 + S&P 500 lag 1). The data and model settings have been the same, except for the different initial weights. The best iteration achieves an out-of-sample accuracy of 54.5% and a logarithmic cumulative sum of 1.52 while the worst iteration only has an accuracy of 49% and a logarithmic cumulative sum of 0.13.

Test model quality on same day data

The model is tested on same-day Ethereum data to prove that the model would generate reasonable predictions with inputs that actually have a dependency on the target. Ethereum is the second-largest cryptocurrency and is highly correlated with Bitcoin. We create a set of training features that only consists of Ethereum log returns with the same timestamp as Bitcoin (same day). The data is processed and used in the model as described above. We refer to the appendix section 8.2.1 for the results of the same-day Ethereum feature set. Tested 10 times, the average prediction accuracy for the Bitcoin signal is 80.5%. Therefore, we assume that the model would recognize existing dependencies between input features and target variable (if there are any) and generate prediction accuracies higher than a buy-and-hold strategy.

Results of 50 iterations

The feature set 0 "rNorm" is created to serve as a benchmark. It consists of a sample of randomly drawn normally distributed values where the mean is 0 and the variance is equal to the variance of the Bitcoin log returns. Subsequently, the feature set rNorm is processed and evaluated in the same way as the other feature sets. Each feature set in Table 5 is trained and tested 50 times and the averages of the resulting metrics are saved in Table 6. The plots of the individual feature sets with their results can be found in the appendix section 8.2.

Feature set	Included features (lag 1)
Feature set 0	rNorm $\sim \mathcal{N}(\mu, \sigma^2)$, where $\mu = 0$ and $\sigma^2 = \text{var}(\text{Bitcoin log returns})$
Feature set 1	Bitcoin
Feature set 2	Bitcoin, S&P 500
Feature set 3	Bitcoin, S&P 500, Russell 2000
Feature set 4	Bitcoin, S&P 500, Russell 2000, EUR/USD
Feature set 5	Bitcoin, S&P 500, Russell 2000, EUR/USD, US Treasury Yield
Feature set 6	Bitcoin, S&P 500, Russell 2000, EUR/USD, US Treasury Yield, Gold+Silver

Table 5: Feature sets 1-6 containing log returns shifted by 1 day (lag 1).

Feature set	Accuracy (in-sample)	Accuracy (out-of-sample)	Loss (in-sample)	Loss (out-of-sample)	Log-cumsum	Sharpe ratio	Pos. / neg. signals forecasted
Feature set 0	0.549	0.525	0.687	0.698	0.922	0.603	534 / 64
Feature set 1	0.546	0.527	0.693	0.701	0.938	0.687	499 / 99
Feature set 2	0.574	0.519	0.678	0.712	0.870	0.696	472 / 126
Feature set 3	0.585	0.523	0.671	0.720	0.812	0.646	453 / 145
Feature set 4	0.595	0.530	0.663	0.719	1.109	0.866	424 / 174
Feature set 5	0.603	0.527	0.659	0.730	0.833	0.637	400 / 198
Feature set 6	0.610	0.525	0.654	0.733	0.853	0.656	392 / 206

Table 6: Results averaged over 50 iterations per feature set.

5.1 Accuracy & Loss

The in-sample-accuracy in the training set illustrated in Fig. 8 improves monotonically with additional inputs from 54.6% in feature set 1 up to 61% in feature set 6. This increase is expected because more input features enable the network to make more dependencies. However, this can also lead to over-sampling if the out-of-sample accuracy is not increasing with additional features. We can not observe a monotonically out-of-sample accuracy increase in Fig. 9 in contrast to the the training set. As mentioned above, it appears that the additional input features lead to overfitting rather than providing real added value. In feature set 4, where EUR/USD is added, an increase in accuracy is recorded, however, accuracy drops again when the 10Y Treasury Yield is added. All in all, we cannot identify any clear added value at this stage. Furthermore, the accuracy of the rNorm feature set is similar to that of the other feature sets. The 95% confidence intervals show that the true mean of the out-of-sample accuracy should be in these intervals with a probability of 95%. The lower limit of the intervals is above 50% for all feature sets. This might show that there is indeed a possible outperformance. However, it is important to note that in the test dataset (out-of-sample) there are 318 ones or positive signals (53.2%) out of the total 598 observations that are in the sample. Consequently, a buy-and-hold strategy would have an accuracy of 53.2%. As we can see in Fig. 11, the model predicts a positive signal more often than a negative one. Therefore, the accuracy of at least 51.9% that is achieved by the neural network is probably not due to the model, but rather to the data distribution itself. Barandela et al. [56] describe the problem of an imbalanced training dataset in supervised pattern recognition. The problem occurs when one class is overrepresented in the training set while the other class is less frequent. Especially with time series that are usually rather noisy, this can lead to the model generally predicting the overrepresented class more often. If the test dataset then also has a similar imbalance, the prediction accuracy will likely be higher without the features contributing anything. In the training set used in this study class 1 (positive signal) is overrepresented with 55.96%. Therefore it is likely that due to this small imbalance, the model will also predict class 1 more often than class 0. In the testing set, class 1 accounts for 53.18%. Therefore, it can be assumed that the out-of-sample prediction accuracy will most likely be above 50%, which applies to our case. Fig. 10 shows that the distributions of Bitcoin returns in the training and test datasets are not exactly the same. Fig. 11 shows that the number of predicted positive signals decreases with each added feature.

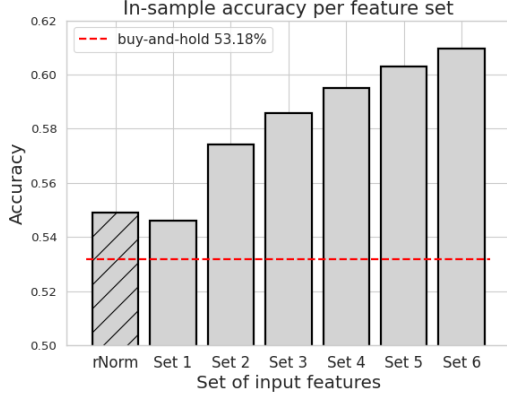


Figure 8: Average in-sample accuracy per feature set.

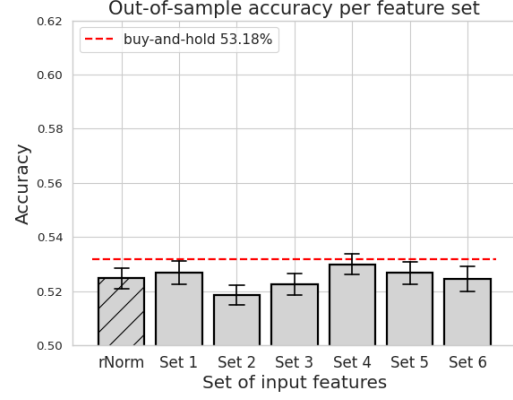


Figure 9: Average out-of-sample accuracy per feature set.

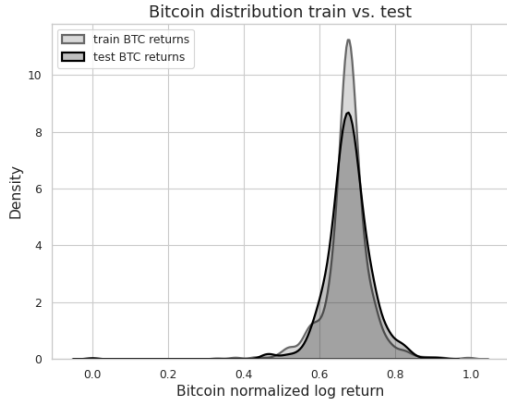


Figure 10: Distributions of returns in train and test set.

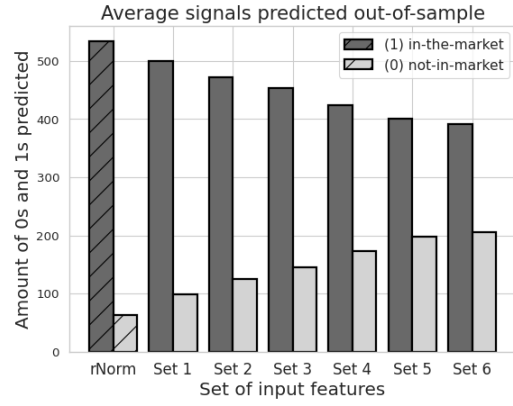


Figure 11: Average predicted signals per feature set.

The loss is described by the binary cross-entropy function 8. In order to compare the individual feature sets, the last loss value of each iteration is stored and then the average of the 50 iterations is taken. The in-sample loss decreases with more input features as Fig. 12 shows. The reason is the same as with the in-sample accuracy, that more features give more scope to recognize dependencies. This makes the "punishment" value for errors smaller. On the other hand, in Fig. 13 the loss in the test dataset increases with almost every feature added, which is a sign of overfitting.

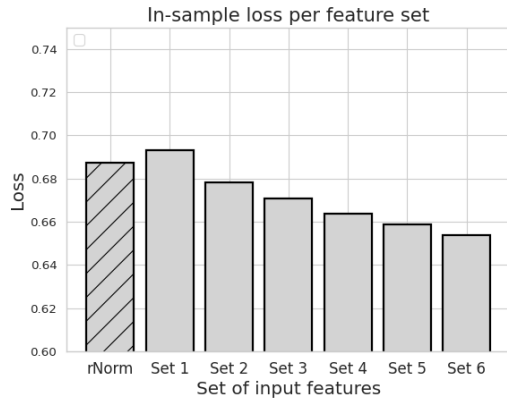


Figure 12: Average final in-sample loss.

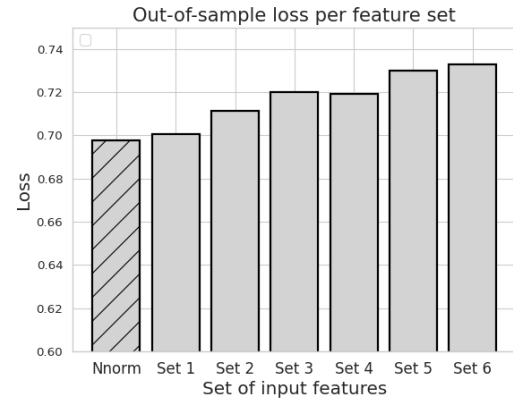


Figure 13: Average final out-of-sample loss.

5.2 Performance

The effective (buy-and-hold) cumulative sum of log returns of Bitcoin in the test dataset is 1.054. The cumulative sum of log returns achieved with the different feature sets is calculated with equation 9. Each feature set generates a vector of signals (0s and 1s) in the length of the test dataset. For simplicity reasons, we use a long-only strategy for the experiment. It is hereby defined that 0 stands for “not in the market” or “not invested” and 1 stands for “in the market” or “invested”. Therefore when calculating the performance, only log returns that occurred at a positive signal will be cumulated and all log returns at a negative signal are ignored. For comparison, we take the average of all cumulative sums of log returns in each feature set (e.g. 50 iterations of feature set 1 then take the average). Fig. 14 shows that only in feature set 4 is the achieved performance higher than it would be with a buy-and-hold approach. Furthermore, additional input features do not seem to linearly lead to better performance. As a second performance metric, the Sharpe ratio is calculated with equation 10. Fig. 15 shows, that the Sharpe ratio is sometimes better than with buy-and-hold, but this is probably also due to the fact that the strategy of the model is not always invested, and thus the volatility is lower. Both plots have a similar shape as the out-of-sample accuracy with the feature set 4 being significantly better than the others.

$$\text{Cumulative Log Returns} = \bar{P} \cdot \bar{R} \quad (9)$$

$$\text{Sharpe Ratio} = \frac{R_p}{\sigma_p} \cdot \sqrt{252} \quad (10)$$

where:

\bar{P} is the vector of signals (0s & 1s) predicted

\bar{R} is the vector of effective Bitcoin log returns during the test period

R_p is the average cumulative log returns over 50 iterations

σ_p is the average daily standard deviation over 50 iterations

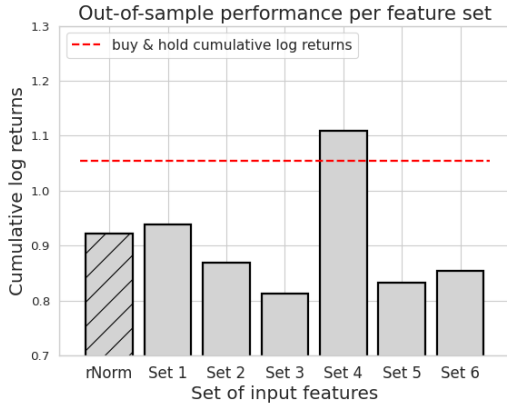


Figure 14: Average cumulative log returns.

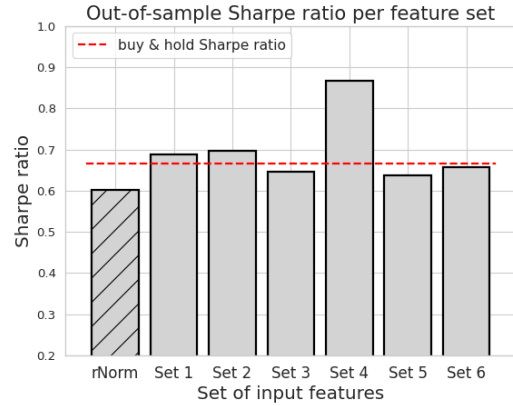


Figure 15: Average Sharpe ratios.

6 Discussion

General discussion

As expected, the in-sample accuracy increases with each input feature added, as the model has more scope to explain the signals. However, the accuracy that is generated out-of-sample on the test dataset does not improve consistently. The feature set 4 has the highest accuracy of 53% but in feature set 5 when the 10Y Treasury Yield was added, accuracy drops again, so the higher accuracy in feature set 4 appears to be random. This is also shown by the comparison with the benchmark rNorm, which achieves a higher out-of-sample predictive accuracy than the other feature sets in 3 out of 6 cases. Furthermore, a buy-and-hold strategy would be correct 53.2% of the time during the test sample period. The fact that the out-of-sample accuracy is not improving the same way as the in-sample accuracy is a sign that additional inputs lead to overfitting and, likely, they do not have a dependency on the target. We can also see the same behavior

with the loss on the training (decreasing) and test set (increasing). Additional input features seem to make the model more active. It predicts more negative Bitcoin signals with each input added. The random feature set rNorm predicts a positive signal 89% of the time, while the largest feature set 6 predicts a positive signal only 66% of the time. Exactly why this is the case would need to be investigated in a follow-up study.

The performance (cumulative log returns) is always worse than buy-and-hold, except again for the feature set 4, which is also better than the other feature sets in terms of accuracy. However, the performance of feature set 5 and 6 are again significantly worse than that of feature set 4. The Sharpe ratios are partly better or about the same as buy-and-hold, but with every feature set, the model is not always invested, which lowers the volatility compared to buy-and-hold. The fact that the feature set rNorm consisting of random values is not consistently outperformed also suggests that the features chosen in this study are not able to explain the Bitcoin price.

Research questions

RC 1: Do additional input features improve the prediction accuracy?

No consistent improvement in prediction accuracy is recorded in the test samples.

RC 2: Is the Bitcoin market efficient as per the weak form of the EMH?

Since buy-and-hold is not outperformed with any of the feature sets in terms of accuracy, the Bitcoin market appears to be weakly efficient under the chosen scenario.

These two assertions only refer to the selected test period, the network architecture, the input features, the number of feature lags, and hyperparameters for this experiment. If these parameters are changed, other results will be the consequence. Also, other combinations of input features may lead to other results.

7 Conclusion

Since much of this experiment is based on chance, we encourage further studies on this topic. Models based on and also tested on past data are unfortunately never 100% accurate or meaningful, as the author can adjust the constellation to produce results in the test sets that are favorable to him. We attempted to address this problem by specifying a model architecture and then averaging the values of the metrics over 50 iterations in Table 6 before evaluating the individual feature sets. The results obtained from this study lead us to the partial conclusion that weak market inefficiencies for Bitcoin cannot be detected using neural networks and the selected asset classes as input.

References

- [1] B. G. Malkiel, "Reflections on the efficient market hypothesis: 30 years later," *Financial Review*, vol. 40, no. 1, pp. 1–9, Feb. 2005.
- [2] M. Wildi and N. A. Bundi, "Bitcoin and market-(in)efficiency : a systematic time series approach," *Digital Finance*, vol. 1, pp. 47–65, Nov. 2019.
- [3] J.-Z. Huang, W. Huang, and J. Ni, "Predicting bitcoin returns using high-dimensional technical indicators," *The Journal of Finance and Data Science*, vol. 5, pp. 140–155, Sep. 2019.
- [4] M. Mudassir, S. Bennbaia, D. Unal, and M. Hammoudeh, "Time-series forecasting of bitcoin prices using high-dimensional features: a machine learning approach," *Neural Computing and Applications*, Jul. 2020.
- [5] P. Jaquart, D. Dann, and C. Weinhardt, "Short-term bitcoin market prediction via machine learning," *The Journal of Finance and Data Science*, vol. 7, pp. 45–66, Nov. 2021.
- [6] Z. Chen, C. Li, and W. Sun, "Bitcoin price prediction using machine learning: An approach to sample dimension engineering," *Journal of Computational and Applied Mathematics*, vol. 365, Feb. 2020.
- [7] E. F. Fama, "The behavior of stock-market prices," *The Journal of Business*, vol. 38, no. 1, pp. 34–105, Jan. 1965.
- [8] E. F. Fama, "Random walks in stock market prices," *Financial Analysts Journal*, vol. 21, no. 5, pp. 55–59, Sep. 1965.
- [9] P. A. Samuelson, "Rational theory of warrant pricing," *Industrial Management Review*, vol. 6, no. 2, pp. 13–39, 1965.
- [10] B. G. Malkiel and E. F. Fama, "Efficient capital markets: A review of theory and empirical work," *The Journal of Finance*, vol. 25, no. 2, pp. 383–417, May 1970.
- [11] E. Sinclair, "The Efficient Market Hypothesis and Its Limitations," in *Positional Option Trading*. John Wiley & Sons, Ltd, Sep. 2015, ch. 2, pp. 11–27.
- [12] C. Munger. Markets are fairly efficient; berkshire hathaway 1999 annual meeting. Accessed: 2022-05-31. [Online]. Available: <https://buffett.cnbc.com/video/1999/05/03/afternoon-session---1999-berkshire-hathaway-annual-meeting.html>
- [13] S. J. Grossman and J. E. Stiglitz, "On the impossibility of informationally efficient markets," *The American Economic Review*, vol. 70, no. 3, pp. 393–408, Jun. 1980.
- [14] A. W. Lo, "The adaptive markets hypothesis: Market efficiency from an evolutionary perspective," *Journal of Portfolio Management*, vol. 30, pp. 15–29, Oct. 2004.
- [15] J. R. Boatright and N.-H. Hsieh, "Efficiency and Rationality," in *Finance Ethics*. John Wiley & Sons, Ltd, 2010, ch. 4, pp. 63–83.
- [16] S. Nakamoto, "Bitcoin: A peer-to-peer electronic cash system," 2008. [Online]. Available: <https://bitcoin.org/bitcoin.pdf>
- [17] D. Garcia, C. Tessone, P. Mavrodiev, and N. Perony, "The digital traces of bubbles: Feedback cycles between socio-economic signals in the bitcoin economy," *J. R. Soc. Interface*, vol. 11, Aug. 2014.
- [18] A. Hencic and C. Gouriou, "Noncausal autoregressive model in application to bitcoin/usd exchange rates," *Studies in Computational Intelligence*, vol. 583, pp. 17–40, Jan. 2014.
- [19] S. Sapuric and A. Kokkinaki, "Bitcoin is volatile! isn't that right?" *Lecture Notes in Business Information Processing*, vol. 183, pp. 255–265, Oct. 2014.
- [20] L. Kristoufek, "What are the main drivers of the bitcoin price? evidence from wavelet coherence analysis," *PLOS ONE*, vol. 10, no. 4, pp. 1–15, Apr. 2015.
- [21] J. Osterrieder and J. Lorenz, "A statistical risk assessment of bitcoin and its extreme tail behavior," *Annals of Financial Economics*, vol. 12, no. 01, p. 1750003, Apr. 2017.
- [22] S. Chan, J. Chu, S. Nadarajah, and J. Osterrieder, "A statistical analysis of cryptocurrencies," *Journal of Risk and Financial Management*, vol. 10, no. 2, May 2017.
- [23] F. Enoksen, C. Landsnes, K. Lučivjanská, and P. Molnár, "Understanding risk of bubbles in cryptocurrencies," *Journal of Economic Behavior & Organization*, vol. 176, pp. 129–144, Aug. 2020.

- [24] A. Urquhart, “The inefficiency of bitcoin,” *Economics Letters*, vol. 148, pp. 80–82, Nov. 2016.
- [25] S. Nadarajah and J. Chu, “On the inefficiency of bitcoin,” *Economics Letters*, vol. 150, pp. 6–9, Jan. 2017.
- [26] A. F. Bariviera, “The inefficiency of bitcoin revisited: A dynamic approach,” *Economics Letters*, vol. 161, pp. 1–4, Dec. 2017.
- [27] A. Brauneis and R. Mestel, “Price discovery of cryptocurrencies: Bitcoin and beyond,” *Economics Letters*, vol. 165, pp. 58–61, Apr. 2018.
- [28] A. K. Tiwari, R. Jana, D. Das, and D. Roubaud, “Informational efficiency of bitcoin—an extension,” *Economics Letters*, vol. 163, pp. 106–109, Feb. 2018.
- [29] K. H. Al-Yahyaee, W. Mensi, and S.-M. Yoon, “Efficiency, multifractality, and the long-memory property of the bitcoin market: A comparative analysis with stock, currency, and gold markets,” *Finance Research Letters*, vol. 27, pp. 228–234, Dec. 2018.
- [30] A. Sensoy, “The inefficiency of bitcoin revisited: A high-frequency analysis with alternative currencies,” *Finance Research Letters*, vol. 28, pp. 68–73, Mar. 2019.
- [31] D. Vidal-Tomás and A. Ibañez, “Semi-strong efficiency of bitcoin,” *Finance Research Letters*, vol. 27, pp. 259–265, Dec. 2018.
- [32] S. Khuntia and J. Pattanayak, “Adaptive market hypothesis and evolving predictability of bitcoin,” *Economics Letters*, vol. 167, pp. 26–28, Jun. 2018.
- [33] S. Gu, B. Kelly, and D. Xiu, “Empirical asset pricing via machine learning,” *The Review of Financial Studies*, vol. 33, p. 2223–2273, Mai. 2020.
- [34] G. Feng and S. Giglio, “Taming the factor zoo: A test of new factors,” *The Journal of Finance*, vol. 75, pp. 1327–1370, Jun. 2020.
- [35] A. Belloni, V. Chernozhukov, and C. Hansen, “Inference on treatment effects after selection among high-dimensional controls,” *The Review of Economic Studies*, vol. 81, no. 2 (287), pp. 608–650, 2014.
- [36] CoinMarketCap. Market open and close times. Accessed: 2022-05-24. [Online]. Available: <https://support.coinmarketcap.com/hc/en-us/articles/360016193231-Market-Open-and-Close-Times#:~:text=The%20market%20opens%20at%2012,UTC%20time%20unless%20otherwise%20specified>
- [37] NYSE. Holidays & trading hours. Accessed: 2022-05-24. [Online]. Available: <https://www.nyse.com/markets/hours-calendars>
- [38] YahooFinance. Eur/usd (eurusd=x). Accessed: 2022-05-24. [Online]. Available: <https://finance.yahoo.com/quote/EURUSD=X/>
- [39] YahooFinance. Treasury yield 10 years (^tnx). Accessed: 2022-05-24. [Online]. Available: <https://finance.yahoo.com/quote/%5ETNX/>
- [40] Bloomberg. Philadelphia gold and silver index description. Accessed: 2022-04-28.
- [41] CoinMarketCap. Today’s cryptocurrency prices by market cap. Accessed: 2022-05-24. [Online]. Available: <https://coinmarketcap.com/de/>
- [42] CoinMarketCap. How are prices calculated on coinmarketcap? Accessed: 2022-05-24. [Online]. Available: <https://support.coinmarketcap.com/hc/en-us/articles/360015968632-How-are-prices-calculated-on-CoinMarketCap->
- [43] J.-M. Kim, S.-T. Kim, and S. Kim, “On the relationship of cryptocurrency price with us stock and gold price using copula models,” *Mathematics*, vol. 8, no. 11, Sep. 2020.
- [44] E. F. Fama and K. R. French, “Common risk factors in the returns on stocks and bonds,” *Journal of Financial Economics*, vol. 33, pp. 3–56, Feb. 1993.
- [45] FTSE. Index factsheet: Russell 2000 index. Accessed: 2022-05-24. [Online]. Available: <https://research.ftserussell.com/Analytics/FactSheets/temp/d5c141a6-bb80-4471-83f4-2bf18e473adb.pdf>
- [46] CMEGroup. Gold futures and options. Accessed: 2022-05-31. [Online]. Available: <https://www.cmegroup.com/markets/metals/precious/gold.html>
- [47] NASDAQ. Phlx gold/silver index methodology. Accessed: 2022-05-24. [Online]. Available: https://indexes.nasdaqomx.com/docs/methodology_XAU.pdf

- [48] NASDAQ. Nasdaq index methodology guide. Accessed: 2022-05-24. [Online]. Available: [https://indexes.nasdaqomx.com/docs/Nasdaq_Index_Methodology_Guide%20\(1\).pdf](https://indexes.nasdaqomx.com/docs/Nasdaq_Index_Methodology_Guide%20(1).pdf)
- [49] N. Japkowicz, “Supervised versus unsupervised binary-learning by feedforward neural networks,” *Machine Learning*, vol. 42, pp. 97–122, Jan. 2001.
- [50] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [51] Scikit-learn. Sklearn minmaxscaler. Accessed: 2022-05-24. [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html>
- [52] M. Kraehenbuehl. Code bachelorthesis neural network with bitcoin. Accessed: 2022-05-29. [Online]. Available: https://github.com/mikekraehenbuehl/BA_2022_Neural_Network_Bitcoin
- [53] Keras. Keras api reference. Accessed: 2022-05-24. [Online]. Available: <https://keras.io/api/>
- [54] NN-SVG. Nn-architecture schematics. Accessed: 2022-05-24. [Online]. Available: <https://alexlenail.me/NN-SVG/>
- [55] Kaggle. Rectified linear units ReLU in deep learning. Accessed: 2022-05-24. [Online]. Available: <https://www.kaggle.com/code/dansbecker/rectified-linear-units-relu-in-deep-learning/notebook>
- [56] R. Barandela, R. M. Valdovinos, J. S. Sánchez, and F. J. Ferri, “The imbalanced training sample problem: Under or over sampling?” in *Structural, Syntactic, and Statistical Pattern Recognition*. Springer Verlag, Aug. 2004, pp. 806–814.
- [57] A. Maravall, “An application of nonlinear time series forecasting,” *Journal of Business & Economic Statistics*, vol. 1, pp. 66–74, Jan. 1983.

8 Appendix

8.1 Plots

Bitcoin autocorrelation

In Fig. 16 only lag 6 and 10 appear to have a significant autocorrelation since the values are outside of the 95% confidence interval. In the Bitcoin paper of Wildi and Bundi [2] they also notice the autocorrelation at lag 6 and therefore use the first 6 lags as inputs for the neural network. However, despite this statistical significance, why should the Bitcoin price from 6 days or even 10 days ago have a dependency on today’s price, whereas the prices of the last two days are completely irrelevant. Autocorrelation shows the amount of linear dependency in a time series. However, a lack of autocorrelation does not imply independence since it can not identify a possible nonlinearity [57]. Hence, a neural network with its nonlinear activation functions is a suitable predictive model for this study and we only select lag 1 as feature because intuitively it should be most relevant.

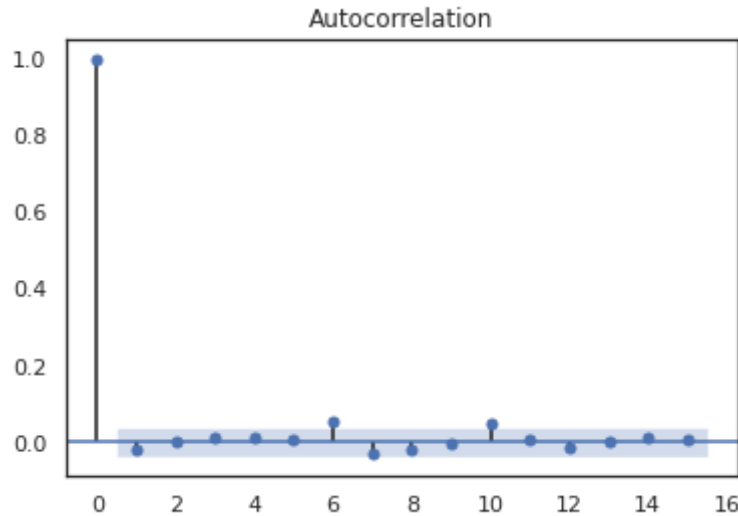


Figure 16: Plot shows autocorrelation for BTC at lag 6 & 10.

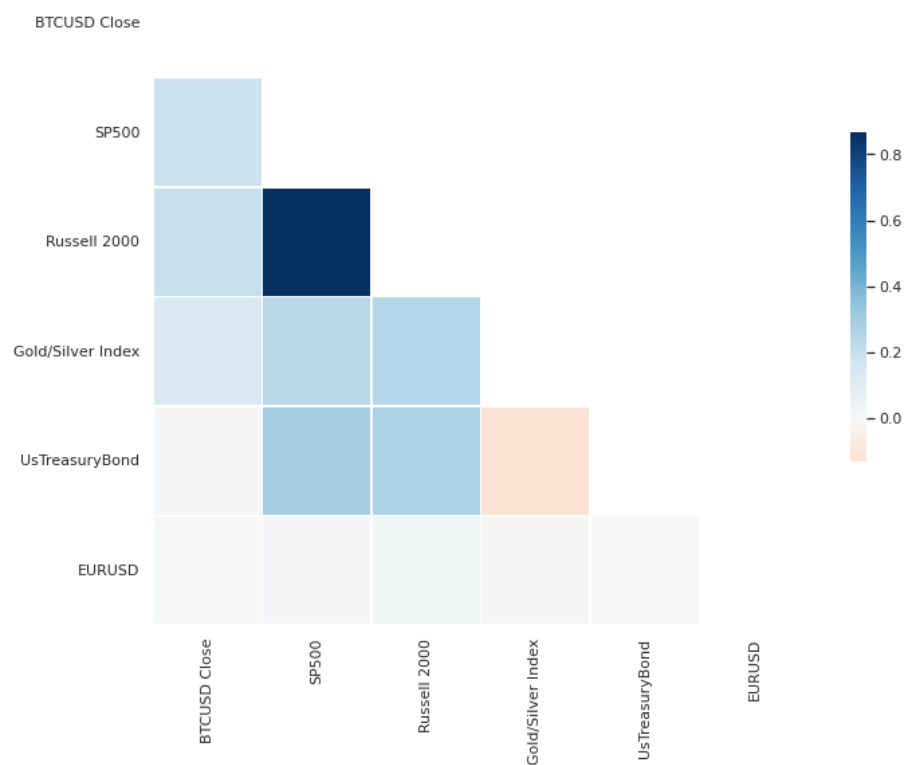


Figure 17: Correlation heatmap

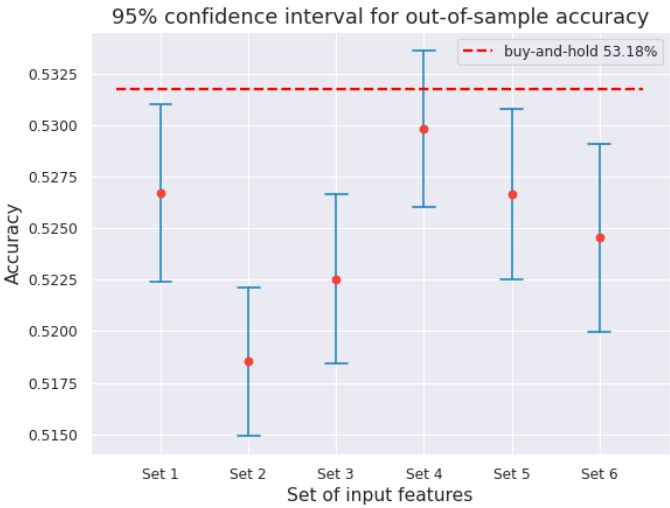


Figure 18: Confidence intervals for the true mean of the accuracies.

8.2 Results of individual feature sets

8.2.1 ETH same day (to test model performance)

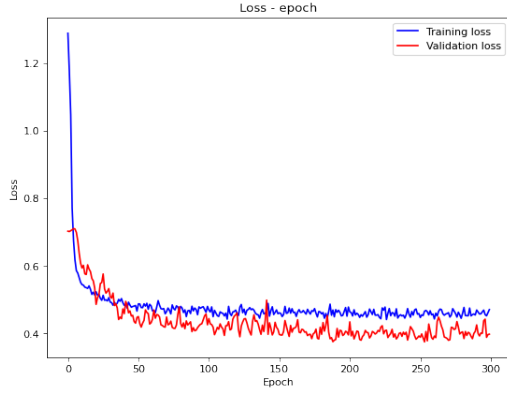


Figure 19: Train/Val. loss history of first iteration.

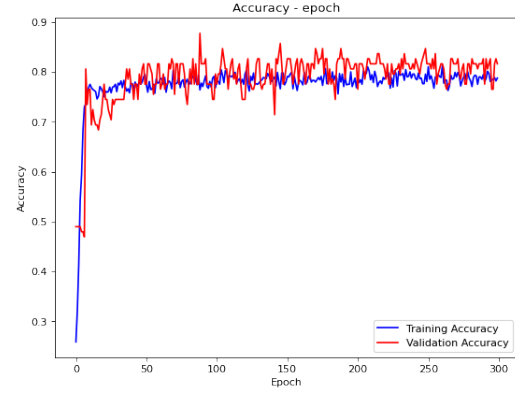


Figure 20: Train/Val. accuracy history of first iteration.

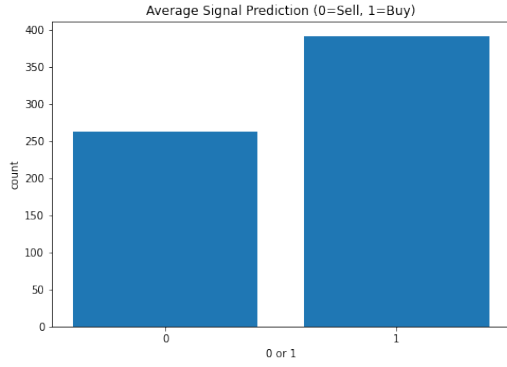


Figure 21: Average out-of-sample signals of 10 iterations.

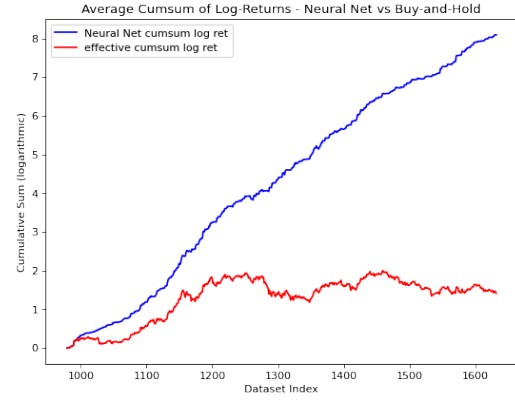


Figure 22: Average out-of-sample cumulative log returns of 10 iterations.

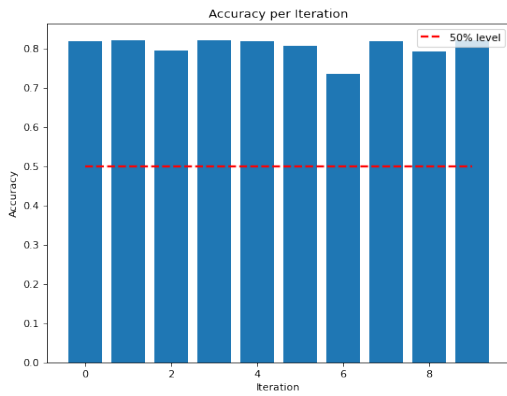


Figure 23: Out-of-sample accuracy per iteration.

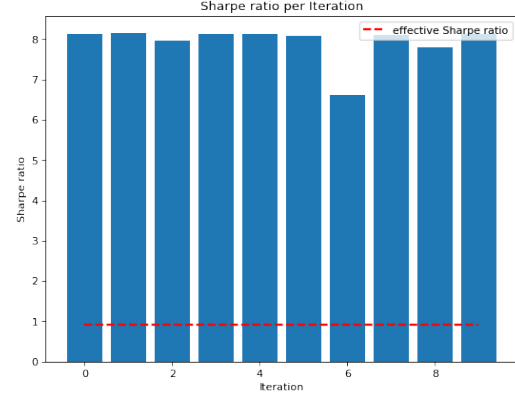


Figure 24: Out-of-sample Sharpe ratio per iteration.

8.2.2 Feature set 0: random sample "rNorm"

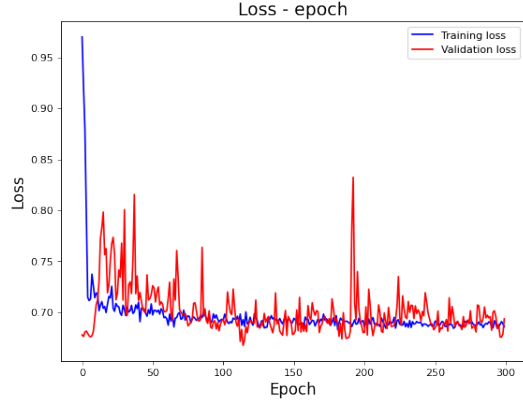


Figure 25: Train/Val. loss history of first iteration.

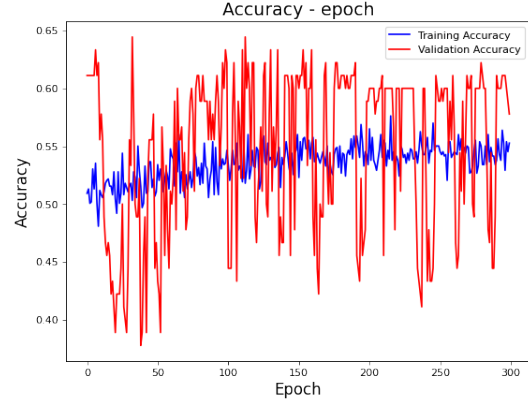


Figure 26: Train/Val. accuracy history of first iteration.

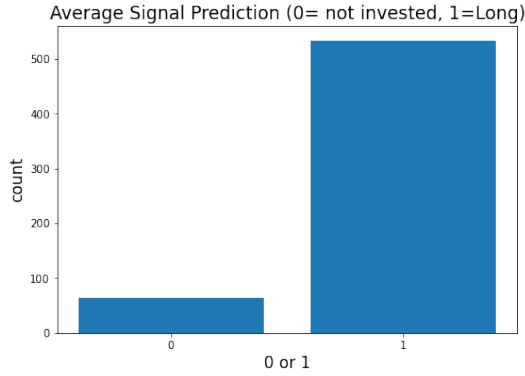


Figure 27: Average out-of-sample signals of 50 iterations.

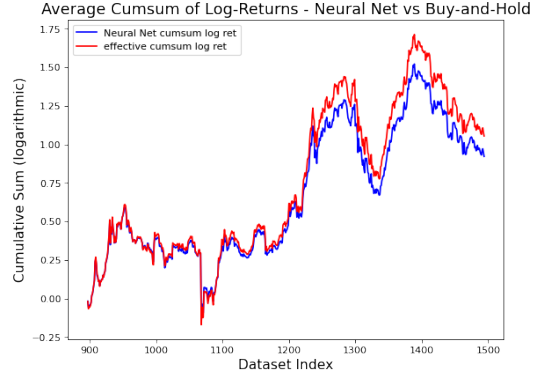


Figure 28: Average out-of-sample cumulative log returns of 50 iterations.

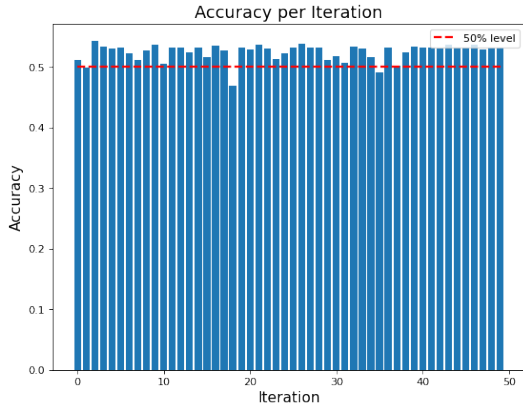


Figure 29: Out-of-sample accuracy per iteration.

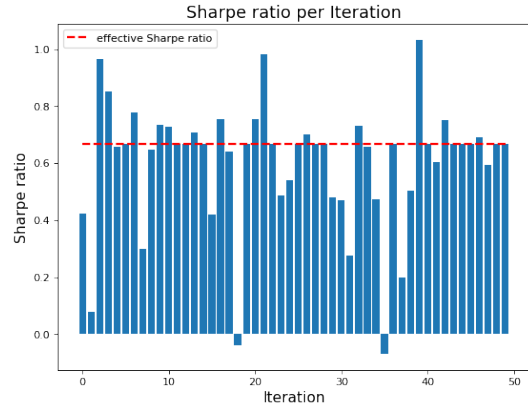


Figure 30: Out-of-sample Sharpe ratio per iteration.

8.2.3 Feature set 1: Bitcoin

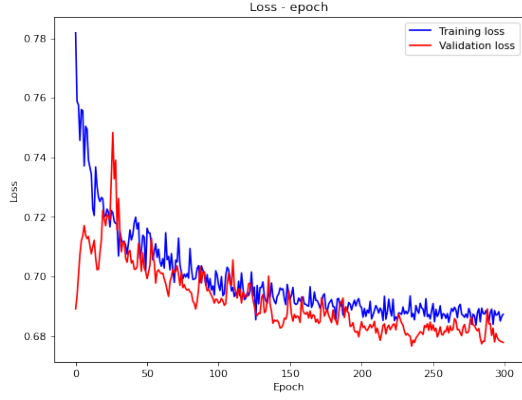


Figure 31: Train/Val. loss history of first iteration.

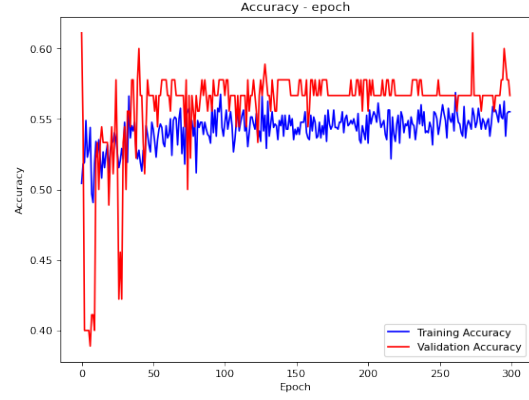


Figure 32: Train/Val. accuracy history of first iteration.

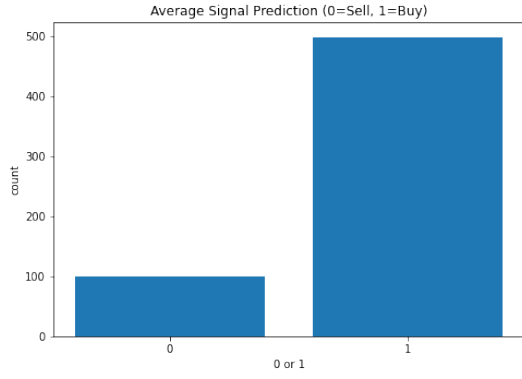


Figure 33: Average out-of-sample signals of 50 iterations.

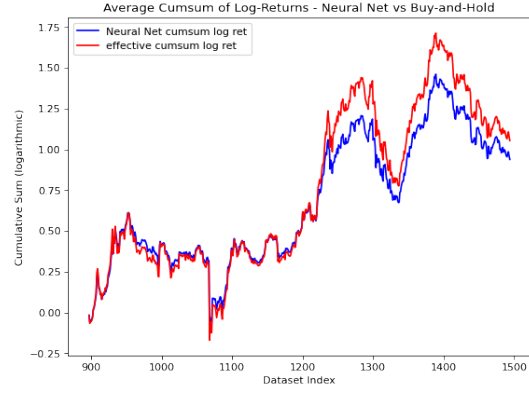


Figure 34: Average out-of-sample cumulative log returns of 50 iterations.

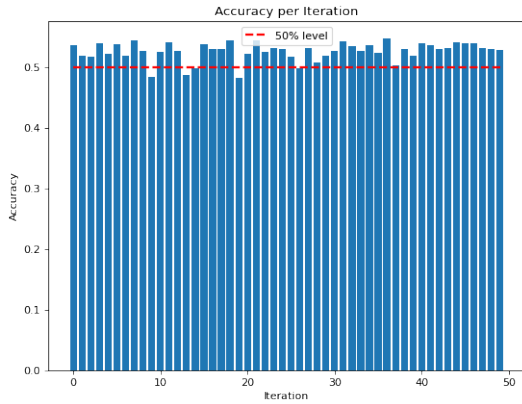


Figure 35: Out-of-sample accuracy per iteration.

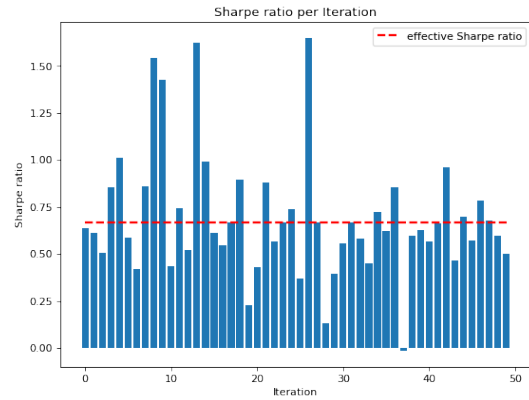


Figure 36: Out-of-sample Sharpe ratio per iteration.

8.2.4 Feature set 2: Bitcoin + S&P500

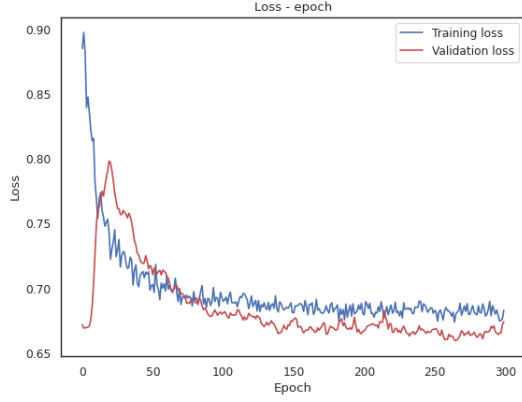


Figure 37: Train/Val. loss history of first iteration.

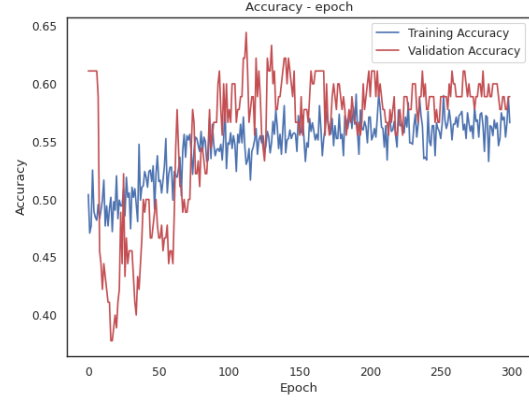


Figure 38: Train/Val. accuracy history of first iteration.

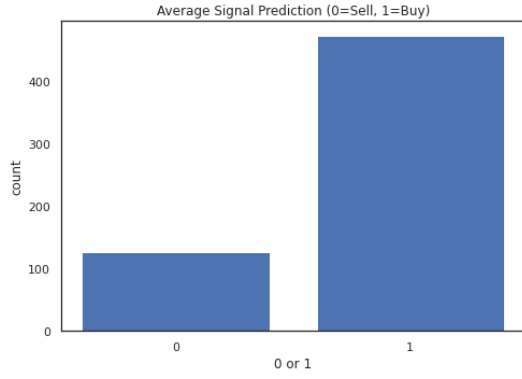


Figure 39: Average out-of-sample signals of 50 iterations.

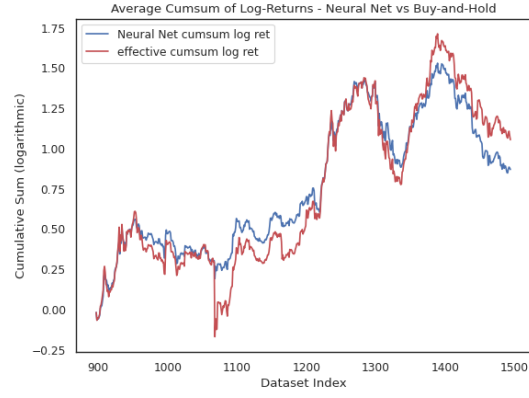


Figure 40: Average out-of-sample cumulative log returns of 50 iterations.

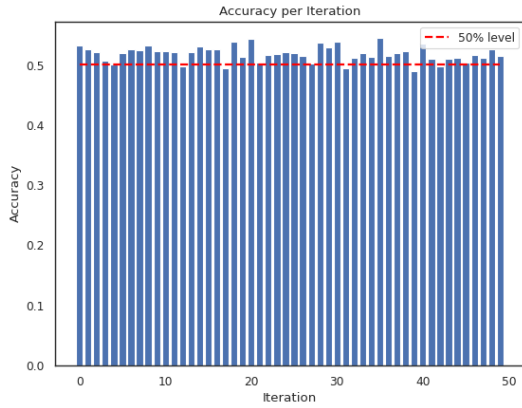


Figure 41: Out-of-sample accuracy per iteration.

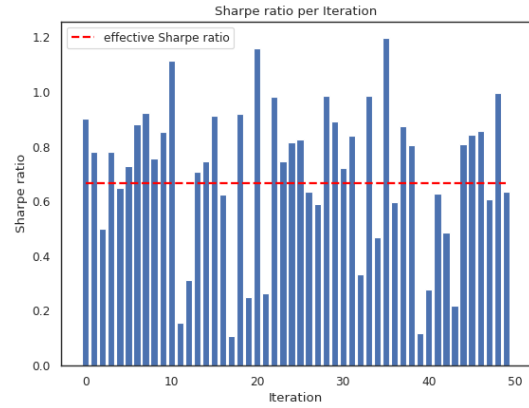


Figure 42: Out-of-sample Sharpe ratio per iteration.

8.2.5 Feature set 3: Bitcoin + S&P500 + Russell 2000

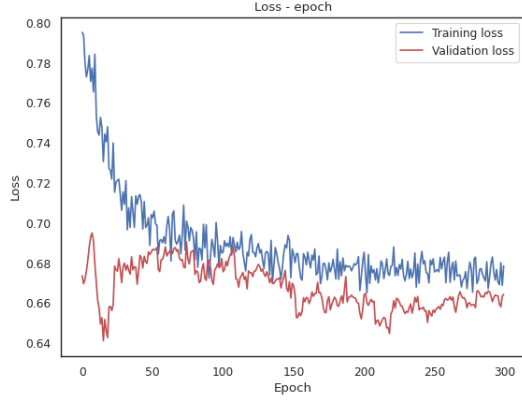


Figure 43: Train/Val. loss history of first iteration.

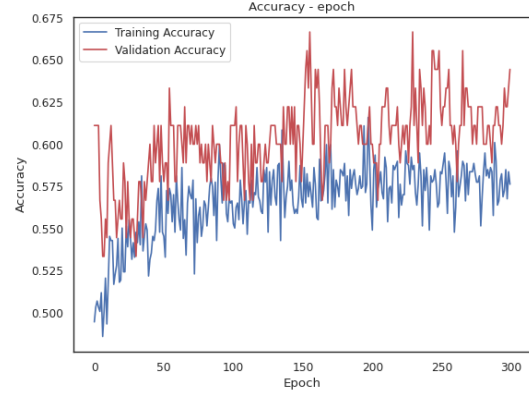


Figure 44: Train/Val. accuracy history of first iteration.

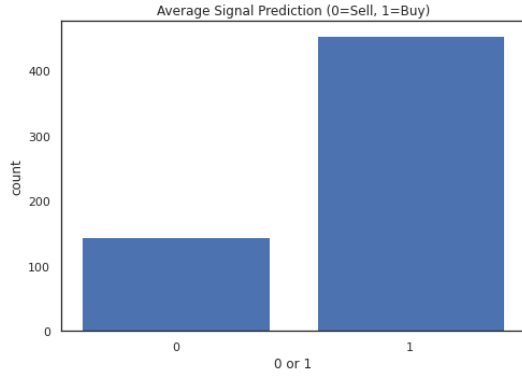


Figure 45: Average out-of-sample signals of 50 iterations.

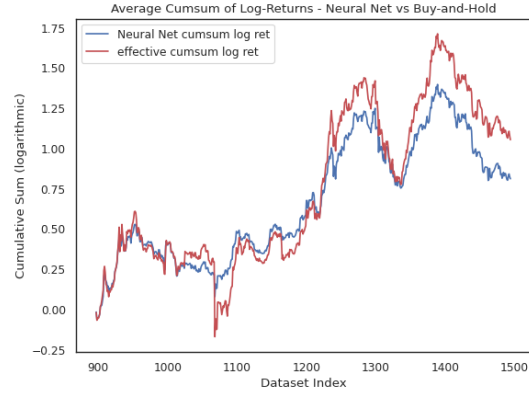


Figure 46: Average out-of-sample cumulative log returns of 50 iterations.

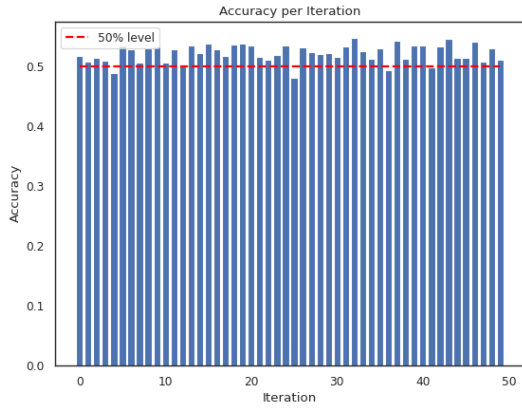


Figure 47: Out-of-sample accuracy per iteration.

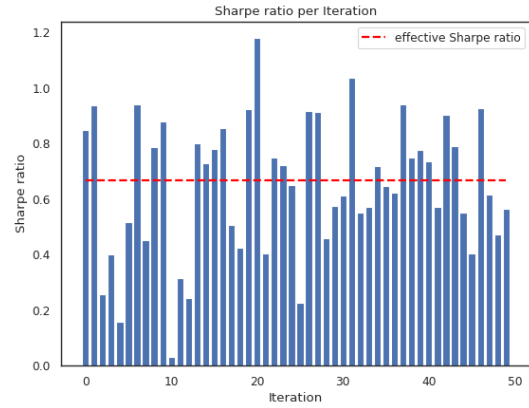


Figure 48: Out-of-sample Sharpe ratio per iteration.

8.2.6 Feature set 4: Bitcoin + S&P500 + Russell 2000 + EUR/USD

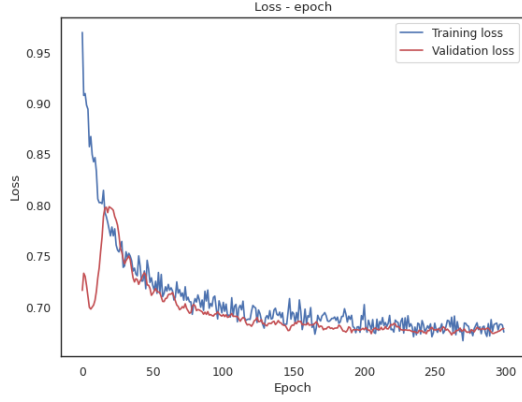


Figure 49: Train/Val. loss history of first iteration.

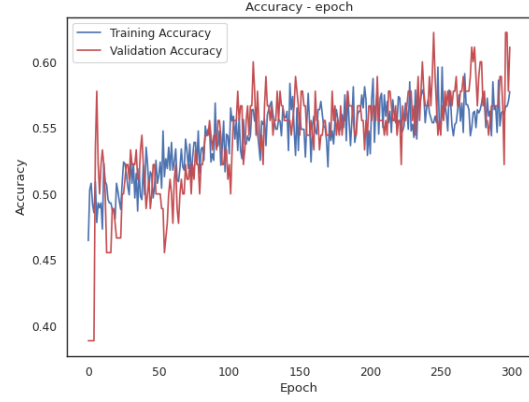


Figure 50: Train/Val. accuracy history of first iteration.

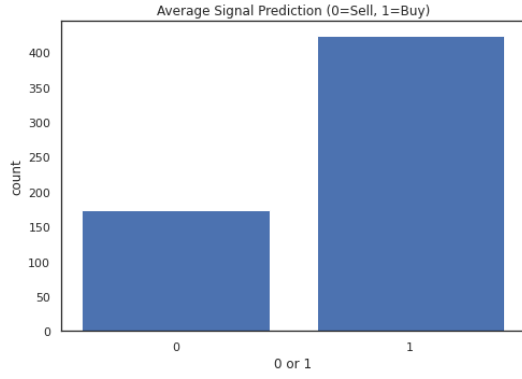


Figure 51: Average out-of-sample signals of 50 iterations.

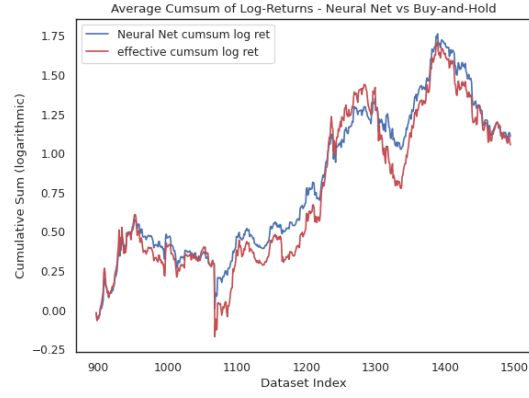


Figure 52: Average out-of-sample cumulative log returns of 50 iterations.

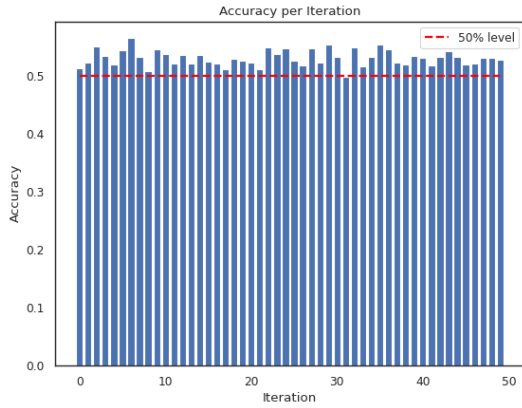


Figure 53: Out-of-sample accuracy per iteration.

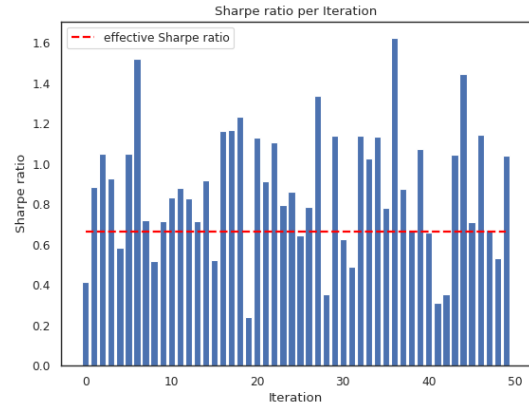


Figure 54: Out-of-sample Sharpe ratio per iteration.

8.2.7 Feature set 5: Bitcoin + S&P500 + Russell 2000 + EUR/USD + 10Y Treasury Yield

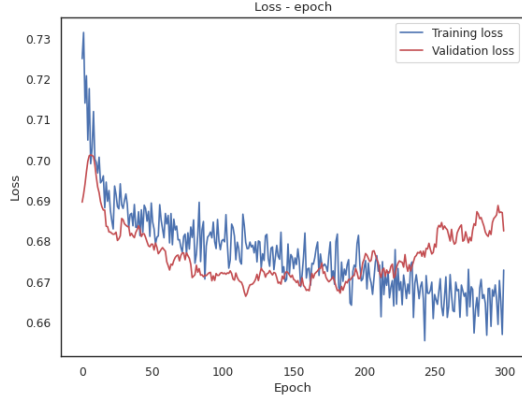


Figure 55: Train/Val. loss history of first iteration.

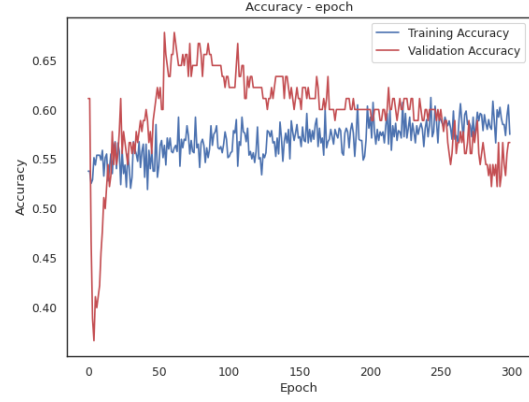


Figure 56: Train/Val. accuracy history of first iteration.

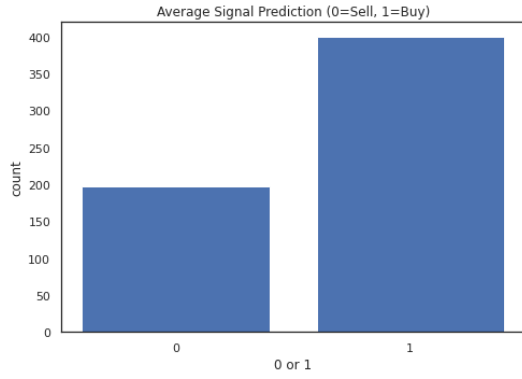


Figure 57: Average out-of-sample signals of 50 iterations.

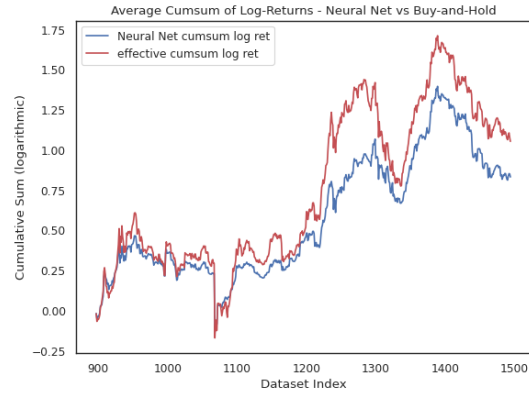


Figure 58: Average out-of-sample cumulative log returns of 50 iterations.

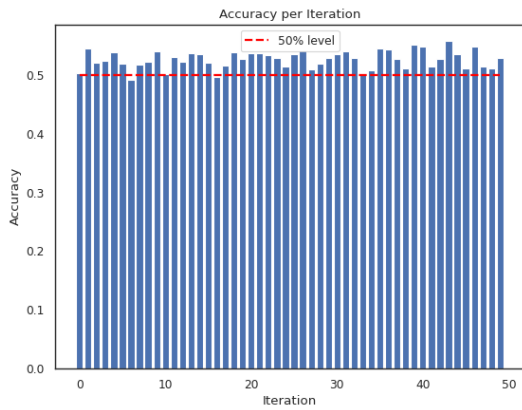


Figure 59: Out-of-sample accuracy per iteration.

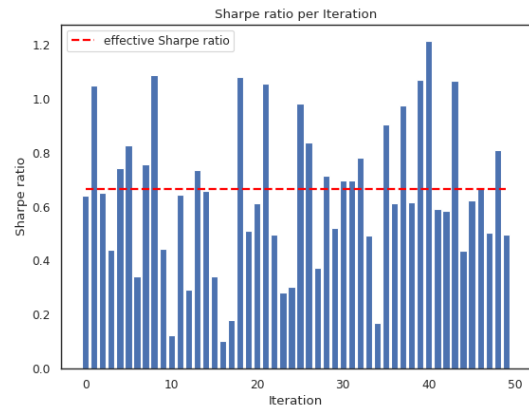


Figure 60: Out-of-sample Sharpe ratio per iteration.

8.2.8 Feature set 6: Bitcoin + S&P500 + Russell 2000 + EUR/USD + 10Y Treasury Yield + Gold/Silver

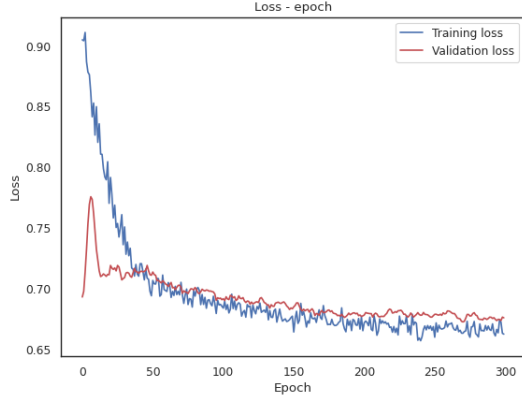


Figure 61: Train/Val. loss history of first iteration.

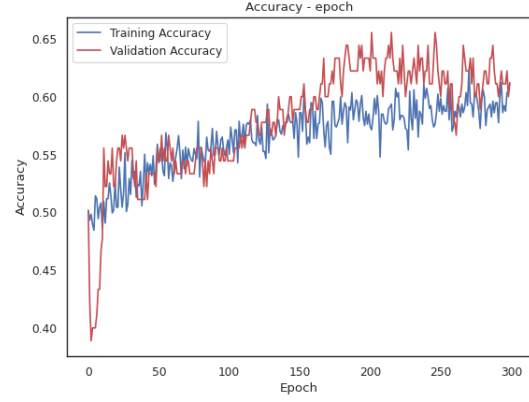


Figure 62: Train/Val. accuracy history of first iteration.

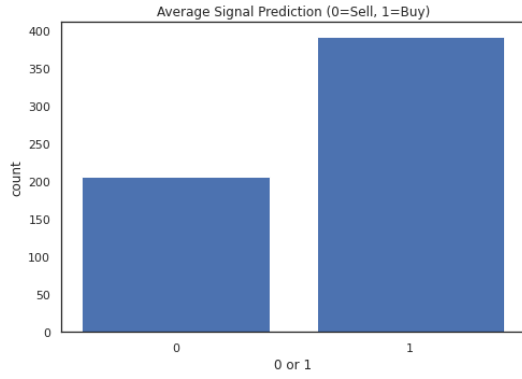


Figure 63: Average out-of-sample signals of 50 iterations.

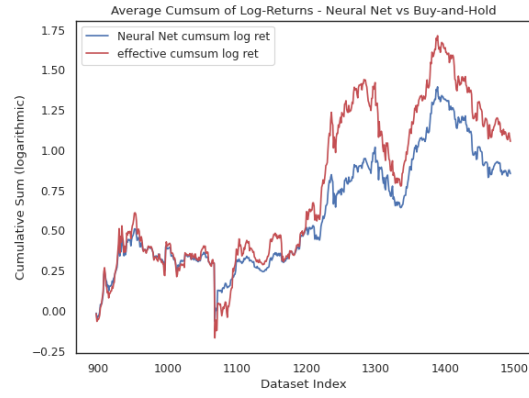


Figure 64: Average out-of-sample cumulative log returns of 50 iterations.

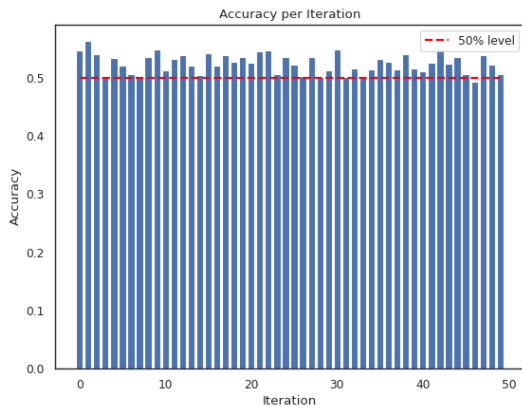


Figure 65: Out-of-sample accuracy per iteration.

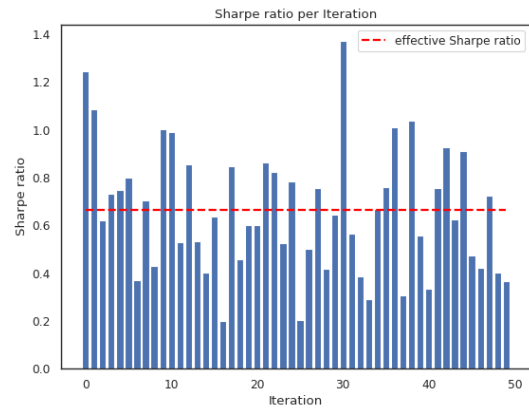


Figure 66: Out-of-sample Sharpe ratio per iteration.