

Bitcoin Price Forecasting via Ensemble-based LSTM Deep Learning Networks

[†]MyungJae Shin, [§]David Mohaisen, and ^{○,‡}Joongheon Kim

[†]Seoul National University Hospital, Seoul, Republic of Korea

[§]Department of Computer Science, University of Central Florida, Orlando, FL, USA

[○]School of Electrical Engineering, Korea University, Seoul, Republic of Korea

[‡]Artificial Intelligence Engineering Research Center, College of Engineering, Korea University, Seoul, Republic of Korea

E-mails: mjshin.cau@gmail.com, mohaisen@ucf.edu, joongheon@korea.ac.kr

Abstract—Time series prediction plays a significant role in the Bitcoin market because of volatile characteristics. Recently, deep neural networks with advanced techniques such as ensembles have led to studies that show successful performance in various fields. In this paper, an ensemble-enabled Long Short-Term Memory (LSTM) with various time interval models is proposed for predicting Bitcoin price. Although hour and minute data set are shown to provide moderate shifts, daily data has relatively a deterministic shift. As such, the ensemble-enabled LSTM network architecture learned the individual characteristics and impact on price predictions from each data set. Experimental results with real-world measurement data show that this learning architecture effectively forecasts prices, especially in risky time such as sudden price fall.

I. INTRODUCTION

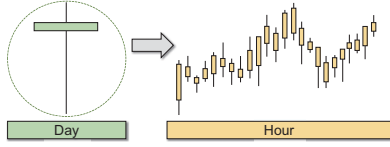
Time series forecasting has been considered as a challenging problem that predicts future values based on the past. This problem is widely seen in many real-world applications addressed by research, including finance, climate forecasting, and power generation. Especially, financial time series forecasting is one of the major difficult problems; thus leading to many scientific interests in the subject for decades. However, traditional statistical models have been limited in predicting financial time series data due to inherently unstable financial time series characteristics. Financial time series, especially stocks, are very difficult to predict because of the heavy-tailed distributions. Meanwhile, the predictability of stock market returns is still an open problem and the results are controversial [1]–[5]. As the price of cryptocurrencies, such as Bitcoin, has soared, interest in their price forecasting has increased. Because Bitcoin price forecasting is an important issue for investors, since the accuracy of price prediction will result in the acquisition of greater income, as in the stock market. However, the nature of time series price data calls for the use of more sophisticated methods to capture complex nonlinear relationships between data variables; such sophisticated models include deep learning [6]. As a result, in the machine learning research community, the prediction of investment markets has been one of the major topics. As investors began to focus on Bitcoin as an investment asset, Bitcoin price forecasting becomes of interest [7]. Many classical statistical methods, such as the Moving Average (MA), Auto-Regression (AR), and Auto-Regressive Integrated

Moving Average (ARIMA), have been applied to solve this problem.

In 2017, the Bitcoin price increases from less than 1,000 at the beginning of 2017 to almost 20,000 in December. However, the Bitcoin price has fallen back to around 6,500. Since the investment market is one of the most important investment channels, the ability to predict Bitcoin market means being able to generate reasonable decisions by avoiding critical financial losses and making financial gains. The price data of Bitcoin market is characterized as dynamic “time series data” [8], and its inherent noisy, non-stationary and complex attributes make much more difficult to construct forecasting models.

In this paper, we propose various time scales and the associated price prediction technique for Bitcoin price forecasting. Furthermore, we propose a new technique for learning weights between ensemble models. The Long Short-Term Memory (LSTM) aggregated with an ensemble technique are used to make scalable and accurate forecasts in Bitcoin trading. In our proposed learning architecture, three LSTM-based neural network models are designed for differentiated time scales such as short-term (i.e., minute), mid-term (i.e., hour), and long-term (i.e., day) data sets in order to reflect Bitcoin price’s considerable fluctuation. Based on the three LSTM-based learning results, an ensemble technique is used for aggregating the results. To the best of our knowledge, the major research efforts in Bitcoin are focused on security and privacy issues, because protecting financial data, securing user privacy, and reducing peer-to-peer system operation risks are obviously important [9]–[13]. In addition, there exists research contributions on stock market forecasting with ensemble-based algorithms [14]–[16]. However, none of this work consider how to differentiate price forecasting with granularity whereas our proposed algorithm conducts Bitcoin price forecasting based on various time scales. Most of the previous work have used information such as price, volume, text contents for price forecasting, and focused on data sets that have moderate shifts. However, in the Bitcoin trading market, there are no moderate shifts during the trading time. This is mainly because there is no mediator for the transaction, as in traditional financial trade system, and the global market is tied together. As such, the price fluctuations are intertwined with diverse interests, and

Fig. 1: Daily data and Hourly data. Note that the daily data can be converted into hourly data through dispersing to shed light on the intrinsic dynamics in price at lower granularity.



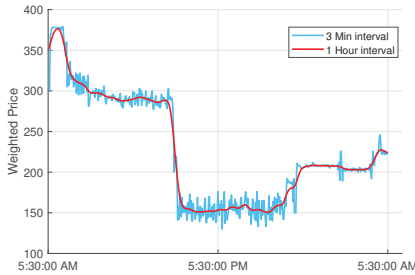
it is almost impossible to find relationships between variables and factors that affect the price change.

Our proposed LSTM based price forecasting algorithm embodies various notable contributions, as follows:

- First, we propose a new multi-cycle time based price prediction technique for Bitcoin price forecasting. We further propose a new technique for learning weights between ensemble models to facilitate such a prediction technique.
- Second, the ensemble model shows stable forecast results in the period of rapid price change, as shown through various experiments and systematic evaluations.
- Finally, using the pricing data of different time intervals, the trained model shows that the same input information also shows different prediction patterns, also validated through various experiments and evaluations.

II. MODELING AND APPROACH

Fig. 2: 2014-12-01 - 2014-12-02 : Difference in price fluctuation according to different time interval



The historical data consists of a series of features, such as the opening and closing prices. The traditional approach to trade market prediction is to analyze these features. As shown in Fig. 1, the market data is represented in the form of candlesticks, which are constructed based on high level features (i.e., highest price, lowest price, open price, close price, etc.). The highest price and lowest price are the highest and lowest price at which a Bitcoin has been traded during a time interval. The closing price generally refers to the last price at which a Bitcoin is traded during a time interval. The open price is the price when the new time interval is started.

The features of candlesticks consist of the historical statistic data within a certain time interval. This means that the

candlestick is reorganized when the time interval change. As a result, the reorganized candlesticks sometimes have different tendencies and patterns, as shown in Fig. 2. In Fig. 2, we notice that there is no clear trend of price at 5:30:00 PM.

As such, we consider multiple time intervals at the same time to adjust the bias of predictions. Our Bitcoin price data organization, data measurement plans, and related descriptions are explained as follows:

- We define the unit time interval for the price forecasting computation. After that, we extract the behaviors and classify them into five price dynamics, namely *High*, *Low*, *Open*, *Close*, and *Volume*.
- To compare predicted data and original data (ground truth), close price is used for original data. As shown in Fig 1, the number of data samples changes when the time stamp changes, and in turn the prices significantly change.
- As illustrated in Fig. 2, we aggregate all hour data values in one day data. However, this aggregation cannot reflect the instantaneous price fluctuation within a day. Furthermore, this is obviously harmful in terms of price forecasting, because the Bitcoin price fluctuation shows deterministic shifts. Therefore, we differentiate the times as minute, hour, and day, in order to reflect short-term, mid-term, and long-term characteristics for more precise Bitcoin price forecasting.
- For gathering data experimentally, Bitcoin trading information, from December 1, 2014 (2014-12-01) to November 11, 2018 (2018-11-11), has been extracted and pre-processed; and it was used for learning and testing.
- Based on the extracted measurement data, customized structuring in order to produce input streams for artificial neural networks is required. For this purpose, the values of input and output of the artificial neural networks should be between 0 and 1 with

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}, \quad (1)$$

and the results are used for learning and forecasting.

- The data we extract from the trading market (e.g., Coinbase) shows the different tendencies of training data set in Fig. 2. As the time interval becomes large, the tendency of the price's change becomes stable. That is, the price change does not fluctuate over a long time. In this paper, and using these characteristics, the LSTM models that perform price predictions are constructed based on those different tendencies.

III. BITCOIN PRICE FORECASTING USING LSTM-BASED ENSEMBLE LEARNING

A. Model Overview: LSTM-based Ensemble Learning

In this work, LSTM, which is one of recurrent neural network (RNN) variants, is used for time series forecasting [17]. The RNN-based algorithms are able to remember time series information patterns, and this is one of its main differences with other feedforward neural networks. The LSTM network

- *Memory Block*: The memory block in LSTM, called the *cell state*, has a self-connection that stores (remembers) the temporal state of the network as well as a special multiplicative unit, called gate, in order to control the flow of information [17]. This cell state is used to repeatedly inference output as shown in (10).
- *Input Gate*: There exist two inputs in this gate, i.e., *current input data* (input in Fig. 3) and the *feedback from previous cell* (recurrent in Fig. 3), respectively. This is calculated as (10). After getting the two inputs, our proposed LSTM cell determines how much of the input will be considered.
- *Output Gate*: This determines the final output of the cell values where the cells are located in the LSTM unit.
- *Forget Gate*: The forget gate scales the internal state of the cell before adding it back to the cell as input through self-recurrent connection, therefore adaptively forgetting or resetting the cells memory [17].

In this section, the equations for the activation (forward pass) and gradients calculation (backward pass) of an LSTM layer are shown. The error gradients are used to train the LSTM network in backpropagation. The LSTM equations are given for a single unit. For multiple units and layers, the calculations are repeated for each block, in any order with following variables.

- $f(x), g(x), h(x)$: the activation function of the gates (sigmoid), the input activation of functions, and the output activation function of gates, respectively, and
- I : the number of inputs, K the number of outputs, H the number of cells in hidden layer.

$$\delta_j^t = \frac{\partial}{\partial a_k^t} O(a_k^t), \quad (2)$$
$$a_{\nu}^t = \sum_{i=1}^I w_{i,\nu} x_i^t + \sum_{h=1}^H w_{h,\nu} b_h^{t-1} + \sum_{c=1}^C w_{c,\nu} s_c^{t-1} \quad (3)$$

$$a_{\phi}^t = \sum_{i=1}^I w_{i,\phi} x_i^t + \sum_{h=1}^H w_{h,\phi} b_h^{t-1} + \sum_{c=1}^C w_{c,\phi} s_c^{t-1} \quad (5)$$

$$a_w^t = \sum_{i=1}^I w_{i,w} x_i^t + \sum_{h=1}^H w_{h,w} b_h^{t-1} + \sum_{c=1}^C w_{c,w} s_c^{t-1} \quad (7)$$

$$a_c^t = \sum_{i=1}^I w_{i,c} x_i^t + \sum_{h=1}^H w_{h,\phi} b_h^{t-1}, \quad (9)$$

$$b_c^t = b_w^t h(s_c^t), \quad (11)$$

$$g(x) = \frac{4}{1+e^{-x}} - 2, \quad (12)$$

$$h(x) = \frac{2}{1 + e^{-x}} - 1, \quad (13)$$

Authorized licensed use limited to: Univ of Calif Santa Barbara. Downloaded on June 16, 2021 at 02:26:06 UTC from IEEE Xplore. Restrictions apply.

Fig. 4: Proposed Ensemble-based Learning Architecture.

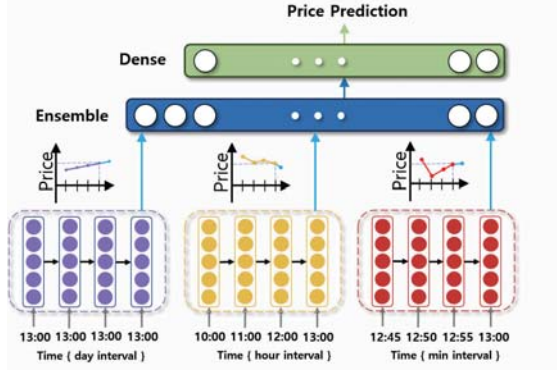
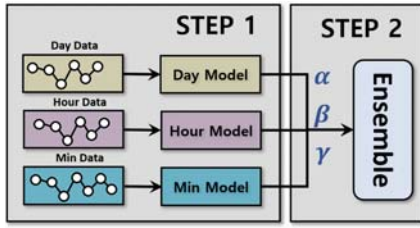


Fig. 5: Training Procedure



C. Ensemble Model

In this work, we propose an ensemble model, which consists of three different LSTM layers with individually and independent time-interval data sets. The benefit of our approach is that the characteristic of the data, which is splitted by different time intervals, can be considered concurrently. That is, the proposed model does not take into account only the short-term trends, represented by one fine granularity, but long-term trends, captured by coarse granularity of time, and using trends in terms of various time intervals. As such, in this work, we use independent LSTM models that use trends in minutes, hours, and days intervals. In Fig. 4, \mathcal{D} is the LSTM model which is trained through the 24-hours interval data, \mathcal{H} is the LSTM model which is trained through 1-hours interval data, and \mathcal{M} is the LSTM model that is trained based on 1-minutes interval data. Throughout this paper, \mathcal{T} is used to refer to the “combined” model representing the total data set, consisting of 24-hours, 1-hours, 1-minutes interval data. The models trained independently for various time granularities are used in tandem (aggregated) as a single larger model. In Fig. 4, the aggregation is illustrated and the presented variables are defined as follows:

- α : the trainable variable for adjusting LSTM trained with minute data set,
- β : the trainable variable for adjusting LSTM trained with hour data set,
- γ : the trainable variable for adjusting LSTM trained with day data set, and
- δ : the variable used to solve the error that arises even after calculation through the above three variables. In addition, trainable bias variable for alleviating error between

forecasting value and real value.

Fig. 5 shows the training procedure of the proposed model. For training, the aggregation parameters and the parameters of the LSTM models are fixed; this training process is called an aggregation training procedure. In the aggregation training procedure, the data is uniformly selected from \mathcal{T} , and the aggregation variables range from 0 to 1. The parameter α determines the effect of the day model when ensemble-model forecasts the price, while the parameter β determines the effect of hour model, and γ determines the effect of minute model on the final ensemble prediction. δ is a bias to adjust for the shift of the predicted values.

Based on this model and corresponding variables, our final output can be calculated as follows:

$$o = \alpha M + \beta H + \gamma D + \delta, \quad (14)$$

where o denotes the final output of the proposed model.

As illustrated in Fig. 4, multiple independent LSTM networks are used for learning at all given three data sets (i.e., day, hour, and minute). After the individual LSTM-based learning procedure, the results are aggregated for more precise forecasting based on ensemble techniques. In Fig. 4, LSTM-based learning with day data is beneficial for the prediction of deterministic shifts in the Bitcoin value over time. On the other hand, the LSTM-based learning with minute data is beneficial for the prediction of moderate shifts. Therefore, conducting ensemble-based learning for the given three LSTM networks with day, hour, and minute data sets is able to utilize the benefits for the given LSTM networks, i.e., high forecasting performance for long-term and short-term price estimation.

D. Prediction

In this paper, the proposed model is trained to predict the 3-minute time interval price, and thus the results of the experiment represent the predicted value after 3 minutes of the input data. In the proposed model, three different time intervals are used, where the corresponding data set for each model is fed into that model for prediction. Each data set uses the most recent candlestick data for prediction. For example, if the 5 : 30 PM data of 3 minutes-interval is fed into the proposed model to predict the price at 5 : 33 PM, the data of 1 hour-interval is fed into the model up to 5 : 00PM.

IV. EXPERIMENTS

A. Data from the Wild

In this section, we describe the data used and then present the results. Furthermore, we provide our analysis of the results. In this paper, the market data has samples with three time intervals, ranging from December 1, 2014 (2014-12-01) to November 11, 2018 (2018-11-11). The market data is adjusted to the volume weighted average price of the three individual time intervals. The volume weighted average price is the ratio of the price traded to the total volume traded over a particular time interval. Such a ratio is used as a measure of the average price at which a stock is traded over the trading time. In our evaluation, the whole data set is separated into a training set

(80%) and a test set (20%). The volume weighted average price is calculated as follows:

$$WP_{VWAP} = \frac{\sum_{t=1}^{T_{interval}} P_t Q_t}{\sum_{t=1}^{T_{interval}} Q_t} \quad (15)$$

where the P_t is the price traded at time t and the Q_t is the quantity traded at time t . The $T_{interval}$ means that the time interval which we set (i.e., 3 minutes, 1 hours and 1 day).

Historic price data for Bitcoin are measured and collected in the format of time series (open, close, high, low and volume) in the granularity intervals of 3 minutes, 1 hour and 1 day. In addition, close series data set is fed into our model for forecasting. In this work, data from 2014 to 2018 was gathered for training from the Bittrex cryptocurrency exchange.

B. Validation Setting

Experiments with real-world historical data (from May of 2015 to November of 2018) have been conducted for the verification of the proposed ensemble-enabled LSTM learning networks. For testing the three LSTM layers in the proposed ensemble-enabled LSTM model, 450000 points of data are used for the ground truth verification and the rest of the data is used for training. This means that the forecasting results by the proposed learning architecture are compared with real-world 450000 points of data. Based on this comparison, the prediction accuracy and errors are observed. The number of nodes in the LSTM was 256. The dense network has 128 as the number of nodes with the *ReLU* function. We used 128 as batch size as well as the Adam optimizer.

C. Evaluation

The performance of the proposed ensemble-enabled LSTM network architecture is evaluated with its *Accuracy* and *Error* where they are calculated using following (16) and (17):

$$\text{Accuracy} \triangleq \frac{1}{N} \sum_{t=1}^N \left(1 - 100 * \frac{|y[t] - y'[t]|}{y[t]} \right), \quad (16)$$

$$\text{Error} \triangleq |y[t] - y'[t]|, \quad (17)$$

where $y[t]$ and $y'[t]$ are the real measured data and the forecasting data by the proposed model.

In Table I, the results of accuracy and root mean square error (RMSE) computation are presented. Table I shows the comparison between the proposed algorithm and ARIMA. Those results are based on real-world measurement data when using four neural network learning architectures, described as follows: (i) Minute LSTM NN: LSTM neural network which is trained only with minute data, (ii) Hour LSTM NN: LSTM neural network which is trained only with hour data, (iii) Day LSTM NN: LSTM neural network which is trained only with day data, and (iv) Ensemble M+H+D LSTM NN: the proposed ensemble-enabled LSTM neural networks with day, hour, and minute data.

This experiment is using the models trained with the different intervals as an input data set. For example, the Hour

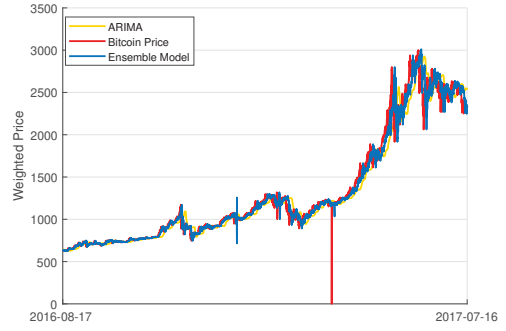


Fig. 6: Forecasting Graph for the Term (2016/08/17 – 2017/07/16), during 3 minute.

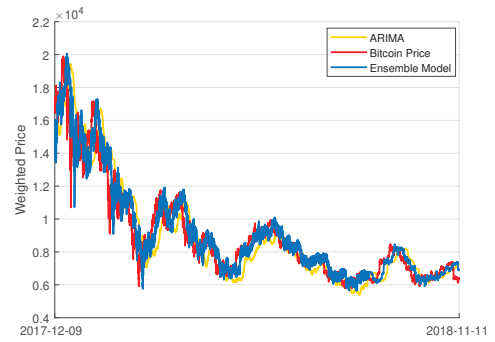


Fig. 7: Forecasting Graph for the Term (2017/12/09 – 2018/11/11), during 3 minute.

LSTM NN is trained using hour-interval data set; and the 3 minute-interval data is fed into Hour LSTM NN to predict the next weighted price. With the accuracy and RSME results in Table I, the following results are observed and reported:

- Day LSTM NN presents the lowest accuracy performance as shown in Tab. I. This is due to the fact that the day-based data cannot reflect the Bitcoin market trends which dramatically changes the prices in a short term. As a result, the results of the trained model showed smooth changes. However, this model shows similar accuracy in Fig. 7 cases. This is because the data set used in Fig. 7 didn't have dramatic changes.
- Minute LSTM NN shows better prediction performance than Day LSTM NN because it is suitable for fast changing price forecasting.
- Hour LSTM NN shows a reasonable performance in both Fig. 6 and Fig. 7 cases. However, in Fig. 6, the Hour LSTM NN shows lowest accuracy. This model shows the most large differences between two data set, and highlights why considering the time interval is very important for price prediction.
- Ensemble M+H+D LSTM NN shows a reasonable performance compared to the other networks. In addition to this performance gain, we can observe the actual benefits

TABLE I: Model evaluation using the accuracy and error (measured by the RMSE).

	Term (1) in Fig.6		Term (2) in Fig.7	
	Accuracy	RMSE	Accuracy	RMSE
Minute LSTM NN	94.91	36.64	95.86	31.75
Hour LSTM NN	89.46	39.26	95.81	33.48
Day LSTM NN	94.95	37.70	94.86	37.58
Ensemble M+H+D LSTM NN	95.56	37.24	96.86	31.60
ARIMA	88.32	42.32	93.21	41.75

(1) 2015/08/17–2017/07/17 [3 min interval] (2) 2017/12/09–1028/11/11 [3 min interval].

of this ensemble-based LSTM networks in terms of risk management. In real Bitcoin trading market, the moments when the financial losses actually happen when the short-term price down-falling (e.g., in Term (1) and Term (2)). In Table. 1, it is shown that the prediction of the rise in price has high error, while the prediction of the down-falling point shows a reasonable performance. In a highly volatile times, our proposed prediction model shows the minimum prediction forecasting errors. This verifies that our proposed learning architecture shows the best performance in terms of reliable Bitcoin market price dynamics estimation especially in risky moments.

V. CONCLUSION AND DISCUSSION

In this work, we address the limitations of the existing Bitcoin price prediction models and propose a more realistic model that acknowledges instantaneous price fluctuations. Particularly, we observe that popular prediction models [1] generalize Bitcoin as a stablecoin and apply machine learning techniques for predictions. However, as shown in Fig1:twograph, Bitcoin is not a stablecoin and Bitcoin price can significantly change in a short duration. Moreover, we also observed that the existing Long Short-Term Memory (LSTM) models, in the same problem space, partially characterize the price change by considering trends in a single time interval.

To address all these limitations, we propose ensemble-enabled Long Short-Term Memory (LSTM) neural network architectures that interpret and predict instantaneously changing Bitcoin price. We use three individual LSTM models for long-term (i.e., day data set), mid-term (i.e., hour data set), and short-term (i.e., minute data set) price prediction. Our experiments show that individual models that consist of LSTM layers, provide different results by training different time-interval data even though the dataset arrives at the same price in different timestamps. The proposed united model aggregates the LSTM networks with ensemble techniques and predicts the Bitcoin price with high accuracy while demonstrating best performance especially in risky periods, i.e., the moments when the Bitcoin price shows deterministic shifts. Therefore, our model accurately captures the real world price fluctuations while providing a high prediction accuracy.

ACKNOWLEDGMENT

This research was supported by the National Research Foundation of Korea (2019R1A2C4070663); J. Kim is the

corresponding author of this paper (joongheon@korea.ac.kr)

REFERENCES

- [1] M. Saad, J. Choi, D. Nyang, J. Kim, and A. Mohaisen, "Towards Characterizing Blockchain-based Cryptocurrencies for Highly-Accurate Predictions," *IEEE Systems Journal*, 14(1):321–332, March 2020.
- [2] J. Kim and J. Kim, "Demo: Light-Weight Programming Language for Blockchain," in *ACM International Conference on Mobile Systems, Applications, and Services (MobiSys)*, Seoul, Korea, June 2019.
- [3] M. Saad, L. Njilla, C.A. Kamhoua, J. Kim, D. Nyang, and A. Mohaisen, "Mempool Optimization for Defending Against DDoS Attacks in PoW-based Blockchain Systems," in *Proc. IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*, Seoul, Korea, May 2019.
- [4] H. Lee, M. Shin, K.S. Kim, Y. Kang, and J. Kim, "Recipient-Oriented Transaction for Preventing Double Spending Attacks in Private Blockchain," in *Proc. IEEE International Conference on Sensing, Communication and Networking (SECON)*, Hong Kong, China, June 2018.
- [5] S. Kim and J. Kim, "POSTER: Mining with Proof-of-Probability in Blockchain," in *Proc. ACM Asia Conference on Computer and Communications Security (AsiaCCS)*, Incheon, Korea, June 2018.
- [6] A. Einstein, B. Podolsky, and N. Rosen, "Can Quantum-Mechanical Description of Physical Reality be Considered Complete?," *Phys. Rev.*, 47:777–780, 1918.
- [7] M. Amjad and D. Shah, "Trading Bitcoin and Online Time Series Prediction," in *Proc. Conference on Neural Information Processing Systems (NIPS) Time Series Workshop*, Barcelona, Spain, December 2016.
- [8] N.I. Sapankevych and R. Sankar, "Time Series Prediction using Support Vector Machines: A Survey," *IEEE Computational Intelligence Magazine*, 4(2):23–38, 2009.
- [9] S. Bag, S. Ruj, and K. Sakurai, "Bitcoin Block Withholding Attack: Analysis and Mitigation," *IEEE Transactions on Information Forensics and Security*, 12(8):1967–1978, 2017.
- [10] M. Apostolaki, A. Zohar, and L. Vanbever, "Hijacking Bitcoin: Routing Attacks on Cryptocurrencies," in *Proc. ACM/IEEE Symposium on Security and Privacy (S&P)*, San Jose, CA, USA, May 2017.
- [11] E.B. Sasson, A. Chiesa, C. Garman, M. Green, I. Miers, E. Tromer, and M. Virza, "Zerocash: Decentralized Anonymous Payments from Bitcoin," in *Proc. ACM/IEEE Symposium on Security and Privacy (S&P)*, San Jose, CA, USA, May 2014.
- [12] M. Andrychowicz, S. Dziembowski, D. Malinowski, and L. Mazurek, "Secure Multiparty Computations on Bitcoin," in *Proc. ACM/IEEE Symposium on Security and Privacy (S&P)*, San Jose, CA, USA, May 2014.
- [13] C. Decker and R. Wattenhofer, "Information Propagation in the Bitcoin Network," in *Proc. IEEE International Conference on Peer-to-Peer Computing (P2P)*, Trento, Italy, September 2013.
- [14] P.K. Mahato and V. Attar, "Prediction of Gold and Silver Stock Price using Ensemble Models," in *Proc. IEEE International Conference on Advances in Engineering and Technology Research*, August 2014.
- [15] M. Asad, "Optimized Stock Market Prediction using Ensemble Learning," in *Proc. International Conference on Application of Information and Communication Technologies*, Rostov-on-Don, Russia, October 2015.
- [16] J. Yang, R. Rao, P. Hong, and P. Ding, "Ensemble Model for Stock Price Movement Trend Prediction on Different Investing Periods," in *Proc. International Conference on Computational Intelligence and Security*, Wuxi, China, December 2016.
- [17] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, 9(8):1735–1780, 1997.