МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ

НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ
УКРАЇНИ "КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ
ІНСТИТУТ ІМЕНІ ІГОРЯ СІКОРСЬКОГО"

НАВЧАЛЬНО-НАУКОВИЙ
ФІЗИКО-ТЕХНІЧНИЙ ІНСТИТУТ

# THEORETICAL AND APPLIED CYBERSECURITY

Перша Всеукраїнська
науково-практична конференція,
присвячена 100-річному ювілею
академіка В.М. Глушкова

Матеріали конференції

Київ – 2023

# OSINT TIME SERIES FORECASTING METHODS ANALYSIS

Feher A., Lande D.

National Technical University of Ukraine "Igor Sikorsky Kyiv Polytechnic Institute" Kyiv, Ukraine

Time series forecasting is an important niche in the modern decision-making and tactics selection process, and in the context of OSINT technology, this approach can help predict events and allow for an effective response to them.

For this purpose, LSTM, ARIMA, LPPL~(JLS), N-gram were selected as time series forecasting methods, and their simple forms were implemented based on the time series of quantitative mentions of starlink systems obtained and generated using OSINT technology. Based on this, their overall effectiveness and the possibility of using them in combination with OSINT technology to form a forecast of the future were investigated.

**Keywords:** time-series, prediction, forecasting, OSINT

## Introduction

Business, finance, logistics, medicine, biology, and chemistry, use forecasting as one of the most applied methods of science that help to effectively solve typical problems and contribute to overall developments. At the same time, in the modern world, the latest neural network developments find their fits in various cybersecurity fields, such as threat intelligence, malware detection, and endpoint protection, which use probabilistic forecasting concepts for training, as well as show overall needs in the chosen topic.

Time-series forecasting methods as a scientific attitude use historical and current data to predict future values over a period of time or at a certain point in the future. By analysing the available data stored in the past, forecasting helps to understand future trends and allows you to respond to them in the most effective way.

In today's world, a well-designed forecasting system frees up hands and gives freedom in the field of the targeted application, even within the framework of national and cyber security. From the point of view of military and civilian security, such a system allows for the correct construction and adjustment of tactics and strategy at different time intervals in accordance with the forecasted events.

The task of the study is, first of all, to create a basis of the most effective forecasting methods for effective further research, and make qualitative comparisons between methods of its different nature. The methods themselves are analysed and used in conjunction with Open Source Intelligence (OSINT) technology to prove the application probability concept. The time series considered for the forecasting study represent quantitative collected information obtained using OSINT technologies.

## Methodology

The selected time series for the study represents a complex dependence of the number of selected events obtained using OSINT technology on the time interval of one year. The selected event for analysis was presented as a dependence on the quantitative characteristics of mentions of Starlink systems in news, blogs, and articles over the Internet on the corresponding time period of the 2022 year.

To create a comparative base, only 333 days out of 364 were used to train the selected models, where the last month in a count of 31 days of the selected year was used as the predicted outcome values for further analysis. The main modern approaches to forecasting are considered to be the following: neural network, statistical, econometric, and linguistic. Each of them is actively used in their respective industries, and in some cases, a combination of several approaches or tuning modifications is used to obtain the most relevant values of needs.

As typical modern representatives of the described approaches, the following methods have been chosen to study time series and build forecasts series corresponding to them:

- Long Short-Term Memory (LSTM) as the most common neural network method;
- Autoregressive Integrated Moving Average (ARIMA) as the most widely used statistical method;
- Log Periodic Power Law (LPPL) or Johansen-Ledoit-Sornette (JLS) as an econometric method, which is subject to criticism and is not popular, but is used in some cases;
- N-gram as a linguistic one, which is already quite strongly implemented in modern technologies and life aspects.

Depending on the selected dataset, chosen forecasting models can predict the quantitative characteristics of selected events based on corresponding values there is possible to calculate the average error between the real and predicted data.

To determine the accuracy of the forecasting models, we used the mean square error (MSE) using the formula \eqref{eq:mse} and the root mean square error (RMSE) using the formula \eqref{eq:rmse} as the most accurate methods for determining such kind of errors.

$$MSE = \frac{1}{n}\Sigma_{i=1}^{n}\left(x_i - \overline{x_i}\right)^2$$

$$RMSE = \sqrt{\frac{1}{n}\Sigma_{i=1}^{n}\left(x_i - \overline{x_i}\right)^2}$$

## LSTM

LSTM is a kind of recurrent neural network model, with the difference that LSTM can handle long time series of data. In addition, the conventional recurrent model has a vanishing gradient problem for long data sequences, while LSTM can prevent this problem during training and perform qualitative results.

The model can recall previous long-time series of data [1] and has automatic controls to keep relevant features or discard irrelevant features. It is because of these factors that LSTM was

chosen among other recurrent methods as a method for the study.

For LSTM, we used its single-layer configuration with 32 units, and the Adam optimizer as an extension to stochastic gradient descent which gave more relevant results, with batch size and epoch values of 512 which defines a number of predictions at the time, the output data was (inverse) transformed by normalization to obtain the predicted series.

## ARIMA

ARIMA is an autoregressive integrated moving average model, where the AR part shows that the time series is regressed on its own past data. The MA part shows that the forecast error is a linear combination of past corresponding errors. The I part shows that the data values have been replaced by different values of order $d$ to obtain stationary data, which is a requirement of the ARIMA approach.

It is because of this complexity that the ARIMA model is effective in re-examining past data using this combined learning approach and helps to effectively predict future points in the time series [2]. This attitude creates a base of popularity for the method and its practical value.

For ARIMA, we used its one-layer configuration with $p = 33$ which defines the number of lag observations included in the model, $d = 2$ which defines the number of times that the raw observations are differenced as a degree of differencing, $q = 0$ which defines the size of the moving average window, the values of which were determined empirically according to the more relevant output forecast values.

## LPPL

The LPPL – or Johansen-Ledoyt-Sornett (JLS) model -- attempts to diagnose, time, and predict the end of financial bubbles, a common term in the financial industry for crisis points when the majority of participants lose confidence during speculative growth.

Despite the widespread criticism [3], the creators of the model provide a motivation based on some natural assumptions,

including risk-neutral assets, rational expectations, local self-reinforcing imitation, and probabilistic critical moments for the algorithm to calculate the stages of bubble development directly [4] with a simple equation.

This way, we can see how the chosen forecasting algorithm work with an atypical for it time series.

For the LPPL (JLS), were used its modification using the Covariance Matrix Adaptation Evolution Strategy (CMA-ES), which gives a more varied and relevant forecasted series.

## N-gram

N-grams represent a continuous sequence of $N$ elements from a given set of texts. The N-grams technique has found its main application in the field of probabilistic language models. They estimate the probability of the next element in a sequence of words, and this is the basis of the theoretical approach to assume study time series forecasting.

This approach to language modeling estimates a close relationship between the position of each element in the string, calculating the occurrence of the next word in relation to the previous one and the frequency of their occurrence.

In a broad sense, these elements do not necessarily mean strings of words, they can also be phonemes, syllables, or letters [5], depending on what exactly is required, and it is thanks to this flexibility that the work was able to be based on numeric time series as well.

There is an additional variation in modeling by creating semantically connected elements in turn, in this paper, the unigram was studied as N-gram with one connection inside, to provide a complete forecast of 31 days, with other values of $N$ the model could not produce a chain of values with a length of 31 values, and a simple general type of tokenization of all elements was used.

## Result and Discussion

The software was developed for each method, and the time series was adjusted to obtain the predicted results. The graphs shown in Figure 1 of actual (real) and predicted values were

modeled according to the dataset of chosen time series, and the processes were repeated to obtain the most relevant predicted series.

Neural network and statistical approaches proved to be the most effective for forecasts, while econometric and linguistic methods proved to be rather limited in their use in forecasting such time series.
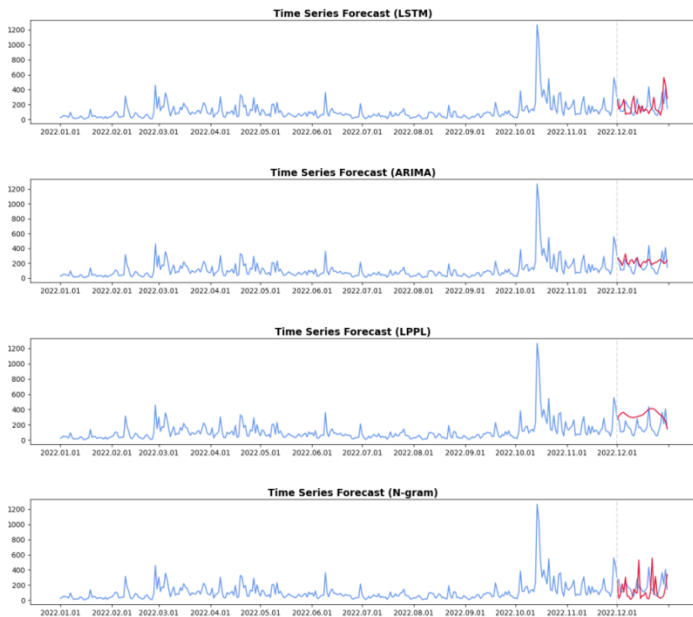


Figure 1. Time series and predictions

## LSTM

The LSTM method is quite flexible and can be easily adjusted to the specifics of the time series, due to its complexity, the method works stably, without fail, and the predicted results are quite close to the real ones. It is also possible to adjust additional parameters [6], so it is possible to create a multilayer model with stronger rejection, which can give more accurate predicted results.

## ARIMA

ARIMA was a good choice, it is less flexible in use, but with the correct selection of parameters *p, d, q* it makes its forecasts quite accurately according to different kinds of time series, among the selected options it showed itself to be the best.

## LPPL

In another way, LPPL performs rather poorly as a method for forecasting time series, which is not surprising due to its narrow focus on solving other mentioned problems. The model is still evolving over time, partly in response to valid criticism, and in the course of the study was found that the strategy of evolutionary adaptation of the covariance matrix CMA-ES is a good improvement that allows for more accurate results, but despite this, when using generative algorithms such as CMA-ES to improve the forecast, the complexity of the calculation itself increases proportionally. It turned out that the calculation of individual large numerical values is also problematic, which requires taking their logarithmic representation, which can also affect the distortion of the forecast.

## N-gram

The N-gram model presented a rather limited version of time series forecasting due to the limited number of previous possible values according to which the forecast can take them. That is, considering this method within the framework of non-stationary series, the forecast is limited by the threshold values of the time series and cannot go beyond it, which reduces its accuracy. Therefore, in studying time series, there is wide room for improvement when using the model with the N-1 algorithm, in which the forecast distortion at short intervals will be much smaller and gradually graduated with respect to time, and the addition of a recurrent component that can increase the accuracy at longer time intervals.

It is also worth noting the creation of joint or separate dictionaries for different series, which will increase the accuracy of joint series if there is an appropriate semantic

correlation, but vice versa in the absence of such correlations. To determine the accuracy of the models were calculated their mean square error (MSE) and root mean square error (RMSE). The results can be found in Table. where a lower number reflects a higher accuracy of the forecast.

|  | LSTM | ARIMA | LPPL | N-gram |
|---|---|---|---|---|
| MSE | 21009.85 | 10242.77 | 36911.94 | 33618.03 |
| RMSE | 144.9477 | 101.2065 | 192.1248 | 183.3522 |

## Conclusion

Based on the practical part, it can be noted that each of the considered methods satisfies the task, despite the low accuracy of effective time series forecasting of such models as LPPL and N-gram, they provide much more creative space for further study and optimisation. In turn, LSTM and ARIMA models have proved to be quite effective, so it is not surprising that these models and their approaches are dominant in terms of time series forecasting.

Thanks to the study carried out, there is a basis for further study of the topic, forecasting various types of events obtained from open sources, and in particular the models themselves. In the context of this study, having the means of automated OSINT data collection, it is possible to confirm the effectiveness of their use for building predictive options for the future.

## References

1. *Sudriani Y.*, *Ridwansyah I.*, *Rustini H. A.* Long short term memory (LSTM) recurrent neural network (RNN) for discharge level prediction and forecast in Cimandiri river, Indonesia. — 2019. — DOI: 10.10 88/1755-1315/299/1/012037.

2. *Brownlee J.* Introduction to Time Series Forecasting with Python. — 1st ed. — 2020. — 365 p.

3. *Fantazzini D.*, *Geraskin P.* Everything You Always Wanted to Know about Log Periodic Power Laws for Bubble Modelling but Were Afraid to Ask // European Journal of Finance. — 2011. — Jan. — Vol. 19. — P. 11–13. — DOI: 10.1080/1351847X.2 011.601657.

4. *Shu M.*, *Zhu W.* Diagnosis and Prediction of the 2015 Chinese Stock Market Bubble. — 2019. — arXiv: 1905.09633 [q-fin.ST].

5. *Jurafsky D.*, *Martin J. H.* Speech and Language Processing. — 3rd ed. — 2023. — 636 p.

6. *Staudemeyer R. C.*, *Morris E. R.* Understanding LSTM – a tutorial into Long Short-Term Memory Recurrent Neural Networks. — 2019. — arXiv: 190 9.09586 [cs.NE].

# РОЛЬОВА МОДЕЛЬ: ВПЛИВ НА БЕЗПЕКУ ТА ДЕЦЕНТРАЛІЗАЦІЮ БЛОКЧЕЙНУ RONIN

Гузенко Г. С., Гальчинський Л. Ю.
Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського», НН Фізико-технічний інститут, Київ, Україна

Наводиться аналіз структури протоколу Ronin, пояснюється сутність моделі консенсусу Proof of Authority (PoA). Фокусується на рольовій політиці доступу до функцій смартконтрактів у протоколі Ronin та вивчається вплив такого доступу на децентралізацію та загальну безпеку блокчейн системи.
**Ключові слова:** Ronin, PoA, безпека, рольова політика

## Вступ

Протокол Ronin – це інноваційна сайд-блокчейнплатформа, призначена для гри Axie Infinity, яка здійснюється на базі неперервної блокчейн технології. Завдяки своїм унікальним функціям та особливостям, Ronin забезпечує широкий спектр можливостей для гравців, а також стабільну та ефективну криптовалютну економіку, що допомагає забезпечити безпеку, швидкість та масштабованість гри. [1]

Ronin походить від Ethereum і спочатку використовував консенсусний алгоритм Proof-of-Authority з низькими комісіями та високою швидкістю транзакцій. Зараз протокол вдосконалили, використовуючи гібридний