

Abstracts of Dissertations

Institute of Information and
Communication Technologies

BULGARIAN ACADEMY OF
SCIENCES



1 / 2021



METHODS AND MEANS OF
DATA ANALYSIS IN
INFORMATION SYSTEMS
USING
TIME SERIES

Ivan Blagoev

МЕТОДИ И СРЕДСТВА ЗА
АНАЛИЗ НА ДАННИ В
ИНФОРМАЦИОННИ
СИСТЕМИ С ИЗПОЛЗВАНЕ
НА ВРЕМЕВИ РЕДОВЕ

Иван Благоев

Автореферати на дисертации

Институт по информационни и
комуникационни технологии

БЪЛГАРСКА АКАДЕМИЯ НА НАУКИТЕ

ISSN: 1314-6351

Поредицата „Автореферати на дисертации на Института по информационни и комуникационни технологии при Българската академия на науките“ представя в електронен формат автореферати на дисертации за получаване на научната степен „Доктор на науките“ или на образователната и научната степен „Доктор“, защитени в Института по информационни и комуникационни технологии при Българската академия на науките. Представените трудове отразяват нови научни и научно-приложни приноси в редица области на информационните и комуникационните технологии като Компютърни мрежи и архитектури, Паралелни алгоритми, Научни пресмятания, Лингвистично моделиране, Математически методи за обработка на сензорна информация, Информационни технологии в сигурността, Технологии за управление и обработка на знания, Грид-технологии и приложения, Оптимизация и вземане на решения, Обработка на сигнали и разпознаване на образи, Интелигентни системи, Информационни процеси и системи, Вградени интелигентни технологии, Йерархични системи, Комуникационни системи и услуги и др.

Редактори

Геннадий Агре

Институт по информационни и комуникационни технологии, Българска академия на науките
E-mail: agre@iinf.bas.bg

Райна Георгиева

Институт по информационни и комуникационни технологии, Българска академия на науките
E-mail: rayna@parallel.bas.bg

Даниела Борисова

Институт по информационни и комуникационни технологии, Българска академия на науките
E-mail: dborissova@iit.bas.bg

Настоящото издание е обект на авторско право. Всички права са запазени при превод, разпечатване, използване на илюстрации, цитирания, разпространение, възпроизвеждане на микрофилми или по други начини, както и съхранение в бази от данни на всички или част от материалите в настоящето издание. Копирането на изданието или на част от съдържанието му е разрешено само със съгласието на авторите и/или редакторите

The series Abstracts of Dissertations of the Institute of Information and Communication Technologies at the Bulgarian Academy of Sciences presents in an electronic format the abstracts of Doctor of Sciences and PhD dissertations defended in the Institute of Information and Communication Technologies at the Bulgarian Academy of Sciences. The studies provide new original results in such areas of Information and Communication Technologies as Computer Networks and Architectures, Parallel Algorithms, Scientific Computations, Linguistic Modelling, Mathematical Methods for Sensor Data Processing, Information Technologies for Security, Technologies for Knowledge management and processing, Grid Technologies and Applications, Optimization and Decision Making, Signal Processing and Pattern Recognition, Information Processing and Systems, Intelligent Systems, Embedded Intelligent Technologies, Hierarchical Systems, Communication Systems and Services, etc.

Editors

Gennady Agre

Institute of Information and Communication Technologies, Bulgarian Academy of Sciences
E-mail: agre@iinf.bas.bg

Rayna Georgieva

Institute of Information and Communication Technologies, Bulgarian Academy of Sciences
E-mail: rayna@parallel.bas.bg

Daniela Borissova

Institute of Information and Communication Technologies, Bulgarian Academy of Sciences
E-mail: dborissova@iit.bas.bg

This work is subjected to copyright. All rights are reserved, whether the whole or part of the materials is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in other ways, and storage in data banks. Duplication of this work or part thereof is only permitted under the provisions of the authors and/or editor.

e-ISSN: 1314-6351

© ICT-BAS 2012

www.iict.bas.bg/dissertations



BULGARIAN ACADEMY OF SCIENCES

Abstract of PhD Thesis

METHODS AND MEANS OF DATA ANALYSIS IN INFORMATION SYSTEMS USING TIME SERIES

Ivan Ivanov Blagoev

Supervisor: Assoc. Prof. Tatiana Atanasova

Approved by Supervising Committee:

Prof. Ivan Garvanov

Prof. Radoslav Yoshinov

Prof. Vladimir Monov

Assoc. Prof. Desislava Ivanova

Assoc. Prof. Velizar Shalamanov



**INSTITUTE OF INFORMATION AND
COMMUNICATION TECHNOLOGIES**

Department of Modelling and Optimization

Introduction

Advances in technology are so obvious that it can only be mentioned without the need for a factual description. In this respect, a significant difference from recent times is the highly expansive digital transformation. Due to COVID-19, the threat to human health, the speed of technology entering our lives has accelerated greatly, leading to a total change in many activities, and in the coming years will be even more noticeable as humanity transforms and adapts to this new way of life.

All that has been mentioned so far leads with it and very completely new to science and processes not previously researched. The collection and processing of time series and big data will be expanded by infiltrating new processes. The need for research and new discoveries will be crucial for the development of science and technology in the coming years. For this development of new methods and tools for time series research and big data processing, it is extremely important and will be a major tool for research and development of science and technology in the future.

This dissertation, through time-series research, contributes to achieving better results in methods of forecasting financial instruments, processing big data and improving cryptography and cyber security.

Purpose and tasks of the dissertation

The purpose of this dissertation is to develop new methods and means of data analysis in information systems using timelines.

For this purpose, the following tasks are defined:

- 1 develop a method for analyzing and predicting price movements in the financial field using time series;
- 2 to propose an algorithm for the training of artificial neural networks in forecasting financial time lines;
- 3 propose solutions to increase cryptographic protection in information systems by applying methods of analysis of time lines;
- 4 to conduct experimental studies to verify the proposed methods of enhancing cryptographic protection in solving cyber security tasks.

- 5 develop programmatic methods to overcome problems when working with big data in time series.

Structure of the dissertation

Thesis work is structured in four chapters.

Chapter 1 provides an overview of the current themes in the field of data science, especially when these data are presented as time series.

Chapter 2 presents the developed methods for researching and forecasting financial time series using different mathematical apparatus.

Chapter 3 also describes the developed solutions to provide cryptographic protection in the provision of information services by examining random number generators representing sequences of time series. The practical application of the proposed approaches to cyber security is presented. The actual results of the tests carried out, demonstrating the successful resolution of the tasks assigned, are shown.

Chapter 4 overcomes loading and processing a big data problems in limited computer resources when examining time series with the means of programming language R.

A summary of the results of the development has been presented in the Conclusion.

The dissertation work contains 125 pages, 33 figures, 1 table and 122 bibliography sources.

Chapter 1. Analysis of the state of the tests.

If you look at the world through the eye of technology, the first thing that would impress any specialist is how much data that is. This is a side effect of mass digital transformation and automation (Wang, 2020), leaving a digital trace of the real process. Time series in communications, technology, business come as a result of measuring characteristics from technical, natural, social, economic and other systems (Mikalef,2020), (Ciampi,2020).

1.1 Time series

The time series is a representing row of data collected at equal or uneven time intervals. A key feature of the time series is that each subsequent value is depending on

the previous values. This dependency can be both very complex and relatively simple. Currently, many forecasting methods that act as effective tools are widely accepted for evaluating and analyzing data from time line models. Of these, the most commonly used model is an integrated method of seasonal component auto-aggression (SARIMA - Seasonal ARIMA), which essentially belongs to a linear model.

1.2. Applying time series to financial instruments

Market price movements are described through time lines and are subject to analysis by finance strategists, economists and market strategists. Types of financial analyses currently used to analyze financial instruments:

- Fundamental - based on analyzing events happening around the world and concerning financial and commodity markets (Wafi, 2015);
- The technical analysis - based primarily on statistical methods with time series calculations - is based on statistical methods. Allows of forecasting to be described by statistical means and mathematical algorithms (Plummer, 1991), (Scott, 2016).

1.2.1 Neural networks

Approaches for examining time series can be divided into two categories: statistical methods and computational intelligence. Statistical methods investigate dependencies between baseline and relevant factors after studying past data, while the other group of methods mimics the human way of thinking and logical conclusion in order to gain knowledge from past experience (such as artificial neural networks) and to predict future values (Atanasova, 2017). Artificial neural networks (ANN) are used in various scientific and daily tasks. In the simplest case, it is a multilayered perceptron. Time series are attractive for research with artificial neural networks (Tomov, 2016).

1.3. Application of time lines in cryptography and cyber security

Meeting cyber security requirements is a prerequisite for the safety and security of IT infrastructures, digital resources and the protection of personal data. In its foundation is cryptography, which provides a number of processes, such as authentication, identification, encryption, data approving processes and etc. The main root of

cryptology is random numbers, and in the most frequent case for modern cryptography needs, two types of random number generators are used:

- True Random Number Generator (TRNG);
- Pseudo Random Number Generator (PRNG).

Traditional RNG measures are mainly aggregated statistics relating to deviations from mathematical chance. In order to help check the quality of a random number generator, its output may be stored in a time series and the data may be subjected to specialized mathematical analyses.

1.4 Conclusions

As a result of the conclusions drawn, it should be summarized that time series studies in different fields and applications need to develop specific methods and means to achieve the specific objectives.

Chapter 2. Methods for examining and forecasting financial time series

This chapter examines the widespread Momentum Indicator that belongs to the oscillators group. Its calculation is based on a mathematical time series processing apparatus. The dissertation aims to improve its effectiveness.

2.1.1 Momentum Oscillator

Momentum is a major oscillator that shows whether the price trend is accelerating, slowing down or moving at the same rate. The function of this oscillator is to take into account the acceleration of the price trend. As confirmation of the signals for torque differentiation, figure formations from technical analysis are included in the study to combine and confirm the current price reversal. In the specific example of Fig. 2.2. Momentum and the signal for divergence 1C confirm an upcoming price turnaround, through a multiple peak of 1D.



Fig. 2.2. EUR/USD Forex historical chart between September 2006 to August 2008 period on a daily time frame

2.1.2 Weaknesses in market trend analysis through Momentum

However, it is clear in the study that the rear incidents in which Momentum may make an exception and not take into account a divergence at the completion of an ongoing market trend, where it is depicted with 1C in Fig. 2.3.



Fig. 2.3. USD/CAD Forex historical chart between September 2006 and August 2008 period on a daily time frame

According to the example presented by the real forex market, the value of Momentum reaches a higher peak even than the previous one. But the price then makes a significant adjustment of about 60% without the oscillator's divergence. The question that excites the study is therefore, can the accuracy of the Momentum oscillator be improved?

2.1.3 Method for increasing the accuracy of Momentum

This dissertation provides a non-traditional method for signaling a market turnaround, namely the **MA Volatility Indicator** developed. And the MA is divided into two under the form:

- Simple moving average (SMA):

To calculate the SMA, a time series is used to sum the data of the last periods(t) where, for example, $t=10$ for 10 days, according to the time frame (may be a different value, optional). Then divide by the number of t periods. This calculation is made for the period of each bar of the chart. The SMA formula is as follows:

$$SMA_t = \sum_{n=1}^t price_n / t$$

- Exponential moving average (EMA):

To reduce the lagging effect of SMA, technical analysis users often prefer Exponential Moving Average (EMA). They reduce the backlog by adding new values to

the latest prices depending on the length of the MA. The shortest EMA will be of greater value than will be applied to most MA.

$$X = K * (C - P) + P,$$

where X – current EMA, C – current price, P – EMA from the previous period (a value of SMA is used for the calculation of the first period), K – smoothing coefficient.

The adjustment coefficient applies an appropriate coefficient to the newer prices, which are related to the previous EMA prices. Formula for smoothing coefficient:

$$K = 2 / (1 + N),$$

where N – number of previous EMA prices.

A conventional approach to MA trading is a higher time frame for the price not to cross MA, as in case of market adjustment, reaching MA from the market price is considered a strong support for the current trend. In case of a break in the MA price, it is taken as a signal for a reversal and, in the case of a bounce, as a signal to confirm the current market trend. The other method is an analysis with more than one MA, all of which are at different speeds.

The MA Volatility Indicator method developed in the dissertation relies on determining extreme values for moving the price away from MA, On the basis of which to determine the sentiment of market participants at the present time.



Fig.2.8 Combining Momentum with the proposed method - MA Volatility Indicator
(historical Forex USD/CAD chart, on a daily time frame, period is between 1997 and 1999)

On the Figure 2.8, the MA Volatility Indicator method is applied and combined with the values of Momentum, and the simulation data is historically real from the Forex market. Figure 2.8 makes it clear that Momentum, depicted with a green line, does not take into account lower values at the last market peak. This is highlighted by the straight red line of its trend.

In conclusion, it can be said that Momentum is an effective oscillator, which has become part of many automated systems and trading strategies in financial markets. But the MA Volatility Indicator method manages to improve forecasting accuracy. Therefore, it could be applied both to automated systems and to analysis of market trends by man.

2.2 Forecasting financial time series via neural networks

Multilayer perceptron is the most commonly used type of artificial neural network that can be presented as an oriented weighted graph. In this study, the basic idea is that instead of the number of hidden layers, the number of neurons at the entrance increases and the hidden layers expand during neural network training. The extension of the input layer is related to the fact that each time line grows with the appearance of a new measurement. The purpose of the training is to make the size of the input layer as large as the size of the full time line.

2.2.1 Modeling prerequisites

The proposed model uses a set of artificial neural subnets and these subnets are joined into a common artificial neural network. The smallest artificial neural subnet has a 1-1-1 topology (Fig. 2.2.6 - left). The network is trained with examples that have only one value. The target in the model is a forecast of only one value ahead of time. Therefore, all sub-modes have only one output. All input values are provided as examples of elastic learning to reverse the error. The training stops at a certain *level of epsilon* for complete change of neural network errors.

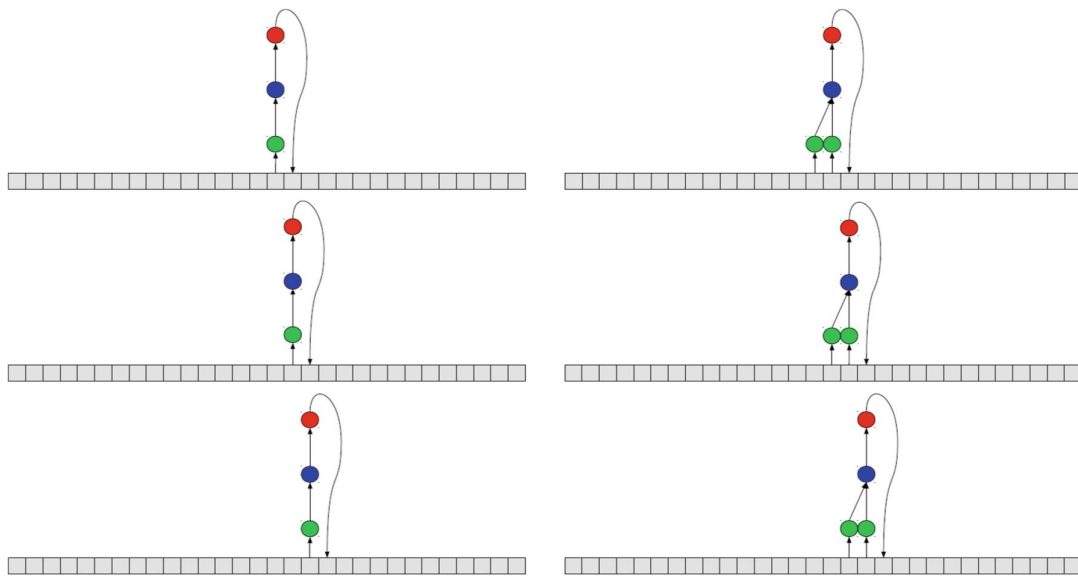


Fig. 2.2.6. Training of artificial neural subnets with 1-1-1 topology (left) and 2-1-1 topology (right).

After training the 1-1-1 topology, the weight values of the first subnet are loaded in the second subnet with a 2-1-1 topology (Fig. 2.2.6-right). A third subnet has a 3-2-1 topology. The size of the hidden layer is selected automatically by a gradual pruning algorithm implemented in the Encog Machine Learning Framework (<http://www.heatonresearch.com/encog/>). The topologies of the subnets are formed by adding one neuron to the input layer and adjusting the size of the hidden layer with an incremental trimming algorithm. The ultimate goal is to reach n-m-1 topology (Fig. 2.8), which covers all known values of time lines.

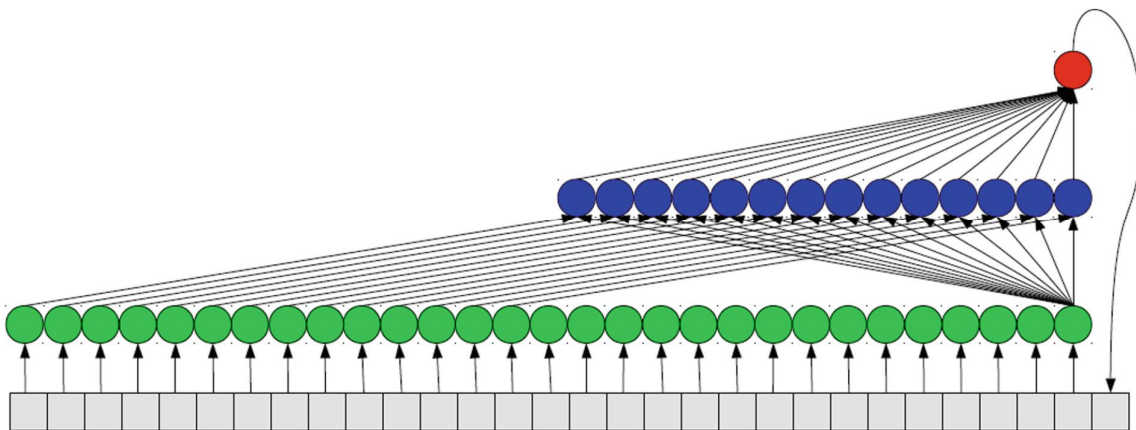


Fig. 2.2.8. Training of artificial neural subnet with n-m-1 topology. Some of the links between the input layer and the hidden layer are not displayed for a better look.

The general idea behind the proposed model is the gradual training of racing in size artificial neural networks. The common problem in the training of artificial neural networks is the size of the network. The proposed model has a higher degree of self-adaptation, since when a new value appears in the time line, the size of the artificial neural network increases, which means that the training phase and the work phase are simultaneous.

2.2.2 Experiments on the study

Experiments are done through a JAVA program where artificial neural networks are implemented through the API provided by the Encog Machine Learning Framework. As input data for experiments, financial time lines are used on the FOREX market. Data are taken from daily two-month trading for EUR/USD and USD/JPY currency pairs. Time line values are scaled in the range from -0.99 to +0.99 with the Min Max zoom rule. The results of the experiments are still within the range of statistical error that comes from the complexity of financial processes and the high-frequency noise inside the data.

The proposed model for self-build three-layer MLP for predicting time series is a promising approach to speeding up the training of artificial neural networks.

2.3 Conclusions

This chapter offers new methods for analyzing and forecasting market price movements through time series and neural networks.

As a result, the following conclusions have been drawn:

- 1 The study so far covers the main aspects of the analysis process – from defining the problem and placing the tasks, to presenting methods for solving them. In each of the stages, real evidence is presented to identify weaknesses or the need to find a more rational approach in the area under question.
- 2 The methods allow them to be integrated into automated processing and decision-making systems. The developed method (MA Volatility Indicator) improves oscillator precision (Momentum) and works in the combination of two EMA or

SMA instruments, offering a new methodology for interpreting results in market analyses and helping to reduce the risk of losses and increase success in automated trading.

- 3 The proposed self-build training algorithm in three-layer MLP accelerates ANN's training in forecasting financial time lines.

The methods presented so far can be applied by specialists in different fields in systems of an estimated nature, for decision-making, analyzing events and processes based on time lines.

Chapter 3. Solutions for providing cryptographic protection by applying time series in cryptography and cyber security

The dissertation proposes time series approach be applied to the quality analysis of a random number generation system (RNG) to ensure cryptographic protection in information systems. For the current RNG survey, a numeric array is retrieved to analyze the values from random numbers in time rows. The results are displayed graphically, where the vulnerable random numbers produced by the generator become more prominent.

3.1 Application of time line techniques for analyzing a random number generator in the field of cyber security

RSA is an asymmetric encryption algorithm that allows anyone to send encrypted messages that only the private key holder can decode. The principle of operation can be explained briefly by generating a very large random number p , then generating another such number q and calculating their work $x=p*q$, in fact x is known as a public key.

3.1.1 The researchers of the (almost) secret algorithm – weaknesses due to insufficient RNG entropy

On the surface, RSA encryption seems invincible. But according to the study presented, the problem lies in the random number generators that provide the algorithm. The vulnerability is fundamental and comes from the fact that RSA needs very large numbers to create encryption keys, and generators in mass computer systems have significantly less capacity. By producing from the number generator, a starting value called seed, entered into the pseudo generator and after computing-intensive computing, cryptographic RSA keys are generated. The problem is that of devices such as phones, IoT, small routers, etc. small systems are even more pronounced because they often do not have sufficient resources for this laborious work. This greatly speeds up the process of generating RSA keys as needed, but opens up a major vulnerability in cryptography security. And considering these disturbing observations, they are reason enough to do research on the subject.

The modern criteria for a reliable RSA key is a minimum of 2048 bits, the recommended length being even 4096 bits. Other research has also found that between 4096, 8192 and 16384 bits of RSA key, greater security of larger keys is minimal. The

reason also comes from the limitations of random number generators. For larger RSA keys, extremely large real random numbers are required. Which in a computer system is extremely difficult to obtain.

And if weaknesses in cryptographic functions are not illuminated, we run the risk of being discovered and exploited by malicious individuals without it being known to others. In conclusion, it can be said that the weaknesses do not proceed from an error in RSA arithmetic. They come from the technological weakness with which RSA is applied.

3.2 Method for assessing the vulnerability of random number generators for cryptographic protection in information systems

The subject of the study includes the technology of the widespread PHP programming language. For the needs of systems developed with this technology to ensure the need for random numbers, PHP has the following means:

1. Lyne congressional generator (LCG), e.g. `lcg_value()`
2. The Marsenne-Twister algorithm, e.g. `mt_rand()`
3. Locally supported function C, i.e. `rand()`

They are also reused for functions such `asarray_rand()` and `uniqid()`, and the downside of entropy and random number generators of the above functions consists in easily predicting future PRNG values. This is because the initial internal states or PRNG SEED are limited and the output of values is in an insufficient range, and this is predictable from readily available modern computational resources. Often, to get a SEED value in PHP, developers use `mt_rand()` or the following script to use automatically:

```
<? php
mt_srand(3231153718);
for ($i=1; $i < 15; $i++) {
    echo mt_rand(), PHP_EOL;
}
```

Which, due to the weak entropy of the tools on offer, risks the recovery of SEED by an attacker. For this purpose, a simulation of a real information system is created in the study, which uses the following source code to generate a token for the different purposes of the application:


```
$newtoken = hash('sha512', mt_rand());
```

Generating a token in the presented way is a nice example, such as a single conversion to `mt_rand()`, which is the hash with SHA512. The fact is that, in fact, if a programmer assumes that the functions of the random PHP values are "random enough", he will be much more inclined to embed a simple usage model. But the method used above to generate mark erssuffers from one flaw - random values are limited to numbers (i.e. its uncertainty or entropy is close to negligible). If you check the output of `mt_getrandmax()`, it will be found that the maximum random number `mt_rand()` can generate only 2.147 billion. This limited number of options makes it vulnerable to a violent attack. In the presence of a modern good video card (GPU) and with the help of specialized brute force attack software such as hash cat, such a calculation can be completed in just a few minutes. Therefore, the use of hash to hide the output of `mt_rand()` is useless.

To protect this type of system, random values of higher quality must be generated. For use in non-trivial tasks, PHP requires sources of high-end entropy that can be provided by the operating system. In Linux is usually used with `/dev/urandom`, unless devices with even high erentropy are installed. In Linux, with the correct setting, a regular random number generator that is of the PRNG type (which is a pseudo random number generator), is often loaded from a source of high entropy `/dev/random`, which here it makes it resistant to attacks. Therefore, any software system developed with PHP inorder to be well protected should be redirected to the reuse functions of the `mcrypt_create_iv openssl_pseudo_random_bytes` OpenSSL external library. They are optimized to use a cryptographically protected pseudo-random generator. Which is tailored and integrated with the operating system.

3.2.2. Understanding RNG Entropy in Linux

In the Linux operating system, the random numbering architecture has the following type:

1. `/dev/random` is a real random number generator if the entropy ends.
2. `/dev/urandom` is a pseudo random number generator (PRNG) and it is not blocked due to entropy depletion.

3. `/dev/hwrng` is an additional hardware for true random numbers that is specialized and not installed in the default computer systems. It provides entropy noise to maintain random numbers;

The accumulated entropy in Linux system can be verified by the following command:

```
$ cat /proc/sys/kernel/random/poolsize
```

```
4096
```

```
$ cat /proc/sys/kernel/random/entropy_avail
```

```
3868
```

Where:

`/proc/sys/kernel/random/poolsize` is used to declare the size (in bits) of the Entropy Pool buffer, for example: How many random numbers should we store before we stop "pumping" for more.

`/proc/sys/kernel/random/entropy_avail` shows the quantity (in bits) of currently stored random numbers in the pool.

Through the user activity and operation of the computer system, such as network, disks, memory state, central processor, peripherals, etc. With the penal functions in the Linux kernel, they have functions for continuously procuring random numbers. Which is designed to compensate for the constant need for them when the computer system is working. For the purposes of the study, such a situation can be easily triggered in order for this process to be observed. By next command, just discard everything that is in `/dev/random` random generator of random numbers and displayed on the screen:

```
$ hexdump /dev/random
```

```
00000000 d5c4 ff0a b8ef 9bdc ad95 480b e853 f0ef
```

```
00000010 e0cb 7c08 4bc4 daef 2b21 ea62 0eac 2c6c
```

```
00000020 d6bd 70e6 5d6f a7e3 0874 d52f 77df 6a2b
```

```
00000030 1909 efe8 9964 acee 2aad 2522 4ddb 1d0b
```

At the same time, the entropy buffer status may be displayed in a parallel open command terminal, with the content refreshed every second. To do this, it is necessary to run the following combination of commands:

```
$ watch -n 1 cat /proc/sys/kernel/random/entropy_avail
```

As a result, the presence of entropy will begin to decline, and its condition will reach critical values, even to zero. By pressing Ctrl-C, this pointless waste is stopped. Perhaps this should never be done in practice - especially on a real server system- except for research purposes of course. But often the systems have problems with the accumulation of entropy in the buffer, and the result seems disturbing:

```
$ cat /proc/sys/kernel/random/entropy_avail
```

```
96
```

From the example presented, the machine produced an entropy result of 96 bits and the increase in this value is too slow and insufficient. The reasons for this can be heterogeneous. For example, from a lack of specific hardware, incorrect settings, virtualization, too much activity with random system numbers, and inability to compensate for the contagion of random values, etc. One possible solution is to launch specialized software to help collect random numbers. This is a daemon that is designed to use any events that can be considered relatively random when the machine is working to produce more and better random numbers. For example, the cpu "flicker", the change in memory status, input output operations, network traffic can add more entropy to the system buffer. Installing this solution and the basic setting in the system are as follows:

```
# apt install haveged
# systemctl start rngd
# update-rc.d haveged defaults
# rngd -r /dev/urandom
```

On a system with relatively moderate traffic:

```
# pv /dev/random > /dev/null
```

```
40 B 0:00:15 [ 0 B/s] [ <=> ]
```

```
52 B 0:00:23 [ 0 B/s] [ <=> ]
```

```
58 B 0:00:25 [5.81 B/s] [ <=> ]
```

```
64 B 0:00:30 [6.05 B/s] [ <=> ]
```

```
^C
```

```
# systemctl start haveged
```

```
# pv /dev/random > /dev/null
```

```
7.12MiB 0:00:05 [1.43MiB/s] [ <=> ]
```

```
15.7MiB 0:00:11 [1.44MiB/s] [ <=> ]
```

```
27.2MiB 0:00:19 [1.46MiB/s] [ <=> ]
```

```
43MiB 0:00:30 [1.47MiB/s] [ <=> ]
```

```
^C
```

Using the `pv` command, you can see how much data is transmitted for this purpose. It is clear from the data flow shown that 2.1 bits per second (B/s) were obtained before, while then ~ 1.5 MB / sec was obtained.

3.2.3. Time series for random number generators

The specifics of RNG and PRNG allow them to be analyzed using time series analysis and forecasting techniques, such as capturing the flow of output numerical values as a sequence and in itself arranged sequentially over time. Such a flow of numerical values can be described as follows:

$$N = T * V$$

where: N - the length of the number row, T - time (duration) of the generation of numbers, V - *number* of generated numbers per unit of time.

So through the time lines it is possible to determine the quality of entropy over time. If a random number generator is not very reliable, then its weaknesses could be found for a shorter time series of data, for which fewer resources will be needed for processing and analysis. For the needs of the current study, a number array will be used,

which will not be created by a high-quality random number generator, but by a mediocre one. The idea is to apply the approach and analyze the time series of random values of a mid-range computer system, which everyone usually has.

3.2.4. Study of random number generators with time series

Time series as stochastic process can also be used for analysis of RNG/PRNG. For this purpose, an algorithm has been developed to detect repetitive patterns (patterns) of data in the time lines generated by RNG. A specially written program is used to collect data from random number generators in a time series for the needs of the survey. Using a re-written program, the random number data collected is presented graphically, helping to make it easier to spot important elements of the time series (Fig. 3.2 and Fig. 3.3). At first glance, with the results of the data from System.Random in Fig. 3.2, everything is fine and it is possible to think that they have a good quality of entropy. But let me suggest another way in another graphical view to make sure of our judgment.

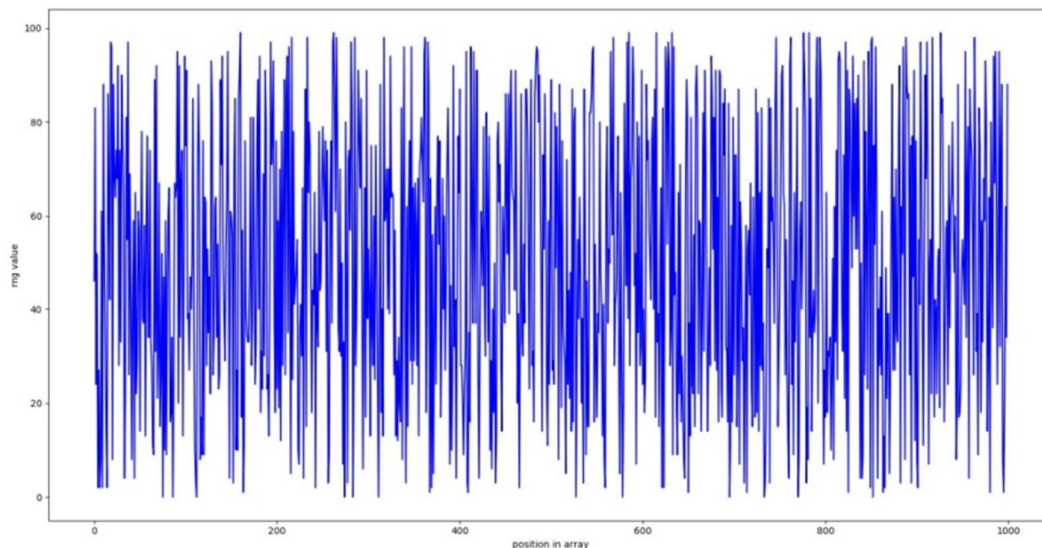


Fig. 3.2. Visualization of data obtained by System.Random as a noise diagram.

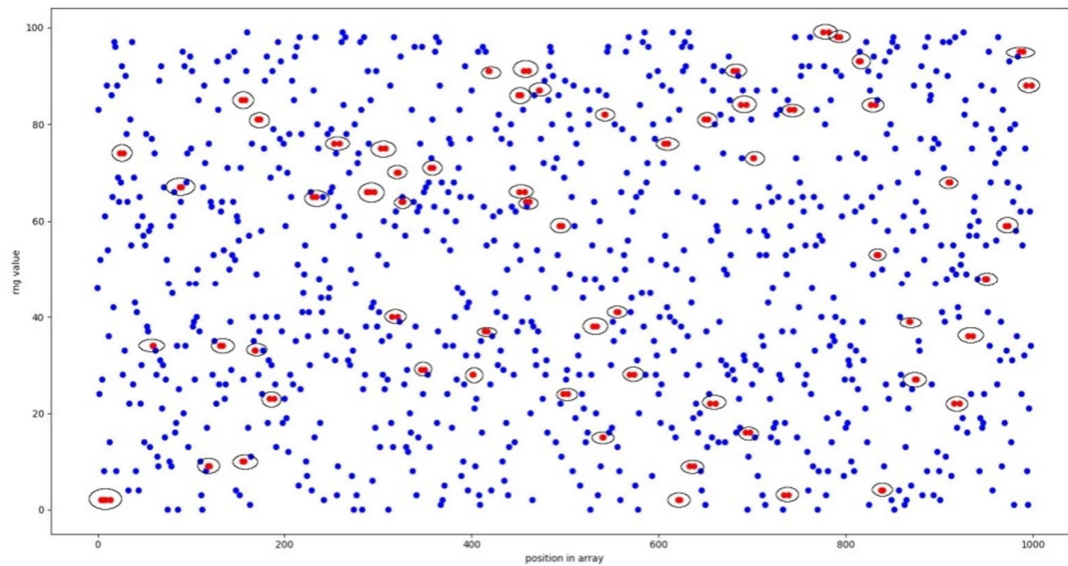


Fig. 3.3. Presentation of System.Random data in a field of dots preview

Presenting the same data with a different graphical interpretation can help to reveal some quality problems with the values studied. In Fig. 3.2 and Fig. 3.3. They are recognized and stained red and circled for better visibility.

The visualization of Fig. 3.3 shows the weaknesses of the processed results. Patterns of recurrence occur periodically over time. These cases are colored red by a program developed using predefined forecasting models, as mentioned earlier. If such a random number generator is used in cryptography, the SEED values produced by it can be successfully attacked by predicting the next SEED value or by monitoring encrypted data transmitted, the values underlying the encryption system can be adopted at a certain point.

3.3 Neglected cyber security risks in public internet hosting service providers

So far, RNG quality analytics research and problem areas have been able to affect cryptographic algorithms, programming languages and operating systems. Now the focus is shifting to massively offered public hosting services. This research uses a web hosting provider that is one of the most popular in the industry. The web application service is installed on a mass-marketed shared hosting. Web certificate added and SSL access is enabled, all running on standard communication ports. On the first line of client-server security, the critical ciphers supported by the hosting server come out. If they are up to

date and there are no vulnerable and already outdated and time-compromised ones, it can be considered that the communication protocol is sufficiently secured.

A test was carried out by scanning the cryptographic protocols that provide the connection between client and server (hosting service). Another protocol that the server supports is TLS v1.2, which is still up to date and approved for use, It contains cryptographic ciphers that need to be removed, but the server offers them for communication, which is also a significant security vulnerability of the service provided. The analysis of protocols and ciphers also identified another significant flaw. The TLS v1.3 protocol is not supported at all, this is currently the most up-to-date and secure protocol of the TLS family for tunnel connectivity.

After checking the cryptographic protocols and ciphers supported for communication, and the follow-up moved on to the more sensitive theme – random number generators. To perform this analysis, a computer program was created that establishes server connectivity in the available cryptographic security protocols between client and server. In this case, TLS v1.2 was used, and in the connection phase, the program takes the generated random numbers from the server and saves them to a file as a time line. The program in question runs in a loop until it collects a sufficient amount of data for analysis.

The data collected from random numbers is analyzed using the specialized open source software to analyze random numbers used in Robert G. Brown's cryptography Dieharder (Brown,2021). A simulation of 114 tests, as well as a check of the quality of numbers and the cyber security standard of FIPS-140 random number generators have been performed. In summary, the data from the random number simulation test are:

- Only 25 tests have passed successfully;
- Failed, which have compromised /predictable/ value and therefore detectable cryptography are 76;
- Vulnerable where cryptography can be revealed with relatively good computer hardware are 13;

From the results presented, it can be concluded that the that due to weaknesses in random numbers and the identified violation of cryptographic protection, the risk of success in cyber-attacks for compromising cryptography is critically high.

The solutions that are allowed in this case are to use private hosting on its own infrastructure, which will not allow excessive load of the type described. However, if it is not possible to provide continuity for a hardware configuration and a suitable location, such as a server room, it is better to rent a VPS server that will only be under the control of one client and also avoid the problem. However, action can also be taken on the part of the hosting provider to increase the capacity of cyber security. The techniques for configuring proper functioning and enhancing the capacity of Linux entropy described in section 3.2.2 "Understanding RNG Entropy in Linux" of this thesis should be applied.

After the correct system setup, you can resort to another unconventional approach, noting the principle of working on collecting entropy in its buffers from the Linux operating system. A program can be written that generates a series of events that will not particularly harass the system, but will create numerous processes supporting the collection of entropy:

```
#!/bin/sh

## list of sites using round-robin DNS
ROUND_ROBINS="www.yahoo.com google.com twitter.com outlook.com"

## Entropy start and end value limits
STOP_LIMIT="3800"
START_LIMIT="3000"

until [ "$(cat /proc/sys/kernel/random/entropy_avail)" -gt "$STOP_LIMIT" ]
do while [ "$(cat /proc/sys/kernel/random/entropy_avail)" -lt "$START_LIMIT" ]
do for thing in "/tmp/loyeyoung" "/tmp/sueellen" "/tmp/rootdev" "/tmp/files"
do echo $thing =====
        touch /tmp/toss

for robins in $ROUND_ROBINS
do nslookup "$robins" 8.8.8.8 > /tmp/toss
```



```

nslookup "$robins" 9.9.9.9 >> /tmp/toss

nslookup "$robins" 192.168.2.3 >> /tmp/toss

nslookup "$robins" >> /tmp/toss

cat /tmp/toss

mkdir $thing -p

cp /tmp/toss $thing/toss

cat $thing/toss

rm -f /tmp/toss

rm -f $thing/toss

done

done

done

done

```

The presented program script is a quite basic and could be upgraded and compiled in other programming or scripting languages. The rate of entropy build-up has improved. Which contributes to the system in question to bear greater loads on the generation of RNG values. The mode of action is as currently specified is that the additional operations in memory, processor, disk and network will be activated when a value in the entropy buffer reaches less than 3000. Also, the solution provided could be used in combination with hardware solutions supporting cryptographic algorithms and random number entropy, which Intel also provides in its processors.

The name of the random number generation module is Intel Secure Key, its previous code name is Bull Mountain Technology. Therefore, it must be verified whether the current system has such processors and its configuration could be upgraded. In the presence of a computer system with a Linux operating system, verification may be done except through the technical documentation of the chips from the manufacturer and through the following combination of commands:

```
$ cat /proc/cpuinfo | grep -i rdrand | echo $?
```

0

As a result, 0 means that an RDRAND flag is available and the processor can be turned on to improve the cryptographic functions of the system as follows:

```
# apt install rng-tools-debian
```

```
# /etc/init.d/rng-tools-debian start
```

```
# /etc/init.d/rng-tools-debian status
```

```
* rng-tools-debian.service - LSB: rng-tools (Debian variant)
```

```
Loaded: loaded (/etc/init.d/rng-tools-debian; generated)
```

```
Active: active (running) since Fri 2020-11-28 17:30:54 EET; 3min 10s ago
```

```
Docs: man:systemd-sysv-generator(8)
```

```
Tasks: 4 (limit: 4915)
```

```
Memory: 1.3M
```

```
CGroup: /system.slice/rng-tools-debian.service
```

```
'-3597 /usr/sbin/rngd -r /dev/hwrng
```

```
$ cat /proc/sys/kernel/random/entropy_aval
```

```
4096
```

The results show that the rate of entropy collection for our case exceeds the rate of its consumption.

3.4 Results in real technological infrastructure

The proposed approach to improving cyber security in cryptography and random number generators in busy server systems with public services has been applied in the technological infrastructure of the IICT-BAS. The hardware configuration used is mid-range, taking into account the complexity of the task performed. The server is equipped with one six-core Xeon(R) E-2236 V6, 32GB RAM and two hard drives in Raid1 configuration. The public service operating server is operating on Linux and all services are entirely open source software. The services as running from the virtual machine are:

- server, currently with 242 user accounts. Available through SMTP, POP3, IMAP, all of which are protected by cryptographic communication protocol TLS v1.2 and TLS v1.3. Certified with a server certificate to establish TLS sessions with an asymmetric algorithm of the type elliptical curve secp384r1. Connecting to the service cannot take place without encryption of communication;
- Web mail that allows all 242 users to operate their mail and through a web browser. The communication is protected by the cryptographic communication protocol TLS v1.2 and TLS v1.3. Certified by a server certificate for establishing TLS sessions with an asymmetric algorithm of the type elliptical curve secp384r1. Connecting to the service cannot be done without encryption of communication;
- Web portal of the Institute of Information and Communication Technologies at the Bulgarian Academy of Sciences, which is the main web space of the institute. It contains activity information, two scientific journals, and structural information. Communication is protected by cryptographic communication protocol TLS v1.2 and TLS v1.3 and server certificate for establishing sessions with asymmetric algorithm with elliptical curve secp384r1.
- SSH remote administration service with the highest degree of cryptographic protection currently offered by the protocol.
- FTP Web Content Remote Management Service. Communication is protected by cryptographic communication protocol TLS v1.2 and TLS v1.3 and server certificate for establishing sessions with asymmetric algorithm with elliptical curve secp384r1.

For all services, the priority protocol for encrypted connectivity is the latest and most secure protocol TLS v1.3, but if it turns out that the client does not support it passes the protocol TLS v1.2. The latter is left only for compatibility, removing all cryptographic algorithms in which vulnerabilities are detected.

It has been found that at present a server is used significantly intensively by IICT users and external Internet users. The level of cryptographic protection is at the highest at the moment according to established standards and no compromises have been made with cryptographic protocols or ciphers. As evidence of the quality of cyber security by cryptographic means, the result of an SSL Labs scanner test was applied to the level of encryption in the TLS protocols offered on the available services. The results of this test

are derived on the basis of current cryptographic protection requirements at present, which have been validated by the international cryptographic protection laboratories FIPS and NIST for the US and Common criteria for Europe. The results show that the protocols and encryption tools are of the current level. The assessment of all tests is the highest possible A+. The level of protection of the HTTP protocol that communicates with the browser via the TLS tunnel is also the maximum level of protection A+.

Tests have also been carried out on the quality of entropy, using the methods presented in the dissertation. Two well-established methods have been applied by the command shell server, the first checking the quality of FIPS entropy with RNG test, and the second with the Dieharder analysis tool. All tests of the entropy of random numbers shall pass with the highest possible result according to the criteria of the programs in question.

Despite the good results, another study has been carried out which shows whether, at times of high user activity and intensive cryptography load, the capacity of random numbers is exhausted. This statistic is collected within half a month. To make it clear whether there are moments that lead to the depletion of the entropy buffer faster than its charging by the system. As a result, after half a month of data collection in a time series, 2137 values were obtained. These values are displayed in graphical form and are depicted in Fig. 3.10:

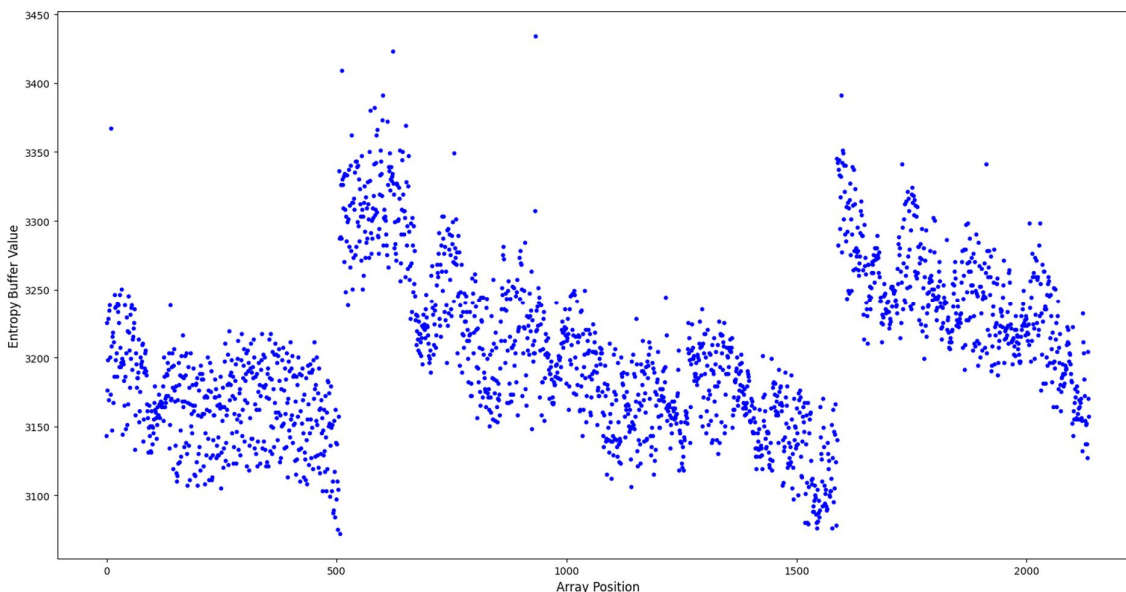


Fig. 3.10 Level of entropy at different moments in time

Fig. 3.10 shows that the system had peaks of more intense activity, which resulted in a strong drawing from the entropy accumulated in the buffer. However, the values at times fell sharply. However, the approach proposed in the study was able to compensate for the high consumption of entropy. While the graph amplitudes between maximum and minimum values appear wide - the values are in a narrow range, with an entropy level between 3000 and just under 3450. The lack of values below 3000 indicates that the system is in very good health and is able to absorb greater loads because the values are far from critical. Taking into account all these results from the actual working environment of the server, the proof of effectiveness of the proposed approach is available. It can therefore be considered that the proposed approach can be beneficial and support the different Internet systems and solutions.

3.4 Conclusions

The research of the services presented and their level of cyber security is key to a more secure transition to modern digital transformation. The rapid transfer of all social and economic activities to digital platforms proves that but in terms of cyber security, many of the current IT services are still lagging behind. Increasing the success rate of cybercrime can lead to a loss of trust in technology and obstruction of these processes, which will also affect scientific and technical progress.

Applying mathematical and statistical analyses with time series to solve cyber security problems is effective. The approaches proposed here can also be combined with other cyber security analysis techniques and methods to be more complex and effective. In this dissertation work, a method for examining the quality of RNG and PRNG in an information system has been developed by applying time series. The method allows to increase the quality of entropy in the use of cryptography, providing various Internet services. The algorithm for detecting repetitive patterns of data generated by RNG has been developed. A study of cryptographic tests and the quality of entropy on real-world busy server systems with public Internet services has been conducted.

Chapter 4. Software approaches to working with big data sets and limited computer resources with programming language R

4.1 The programming language R

Programming language R is a product with powerful tools for statistical calculation and analysis. R is both a programming language and a software environment (Borcard, 2011),(The R, 2017). Language R offers a wide variety of statistical techniques such as linear and non-linear modeling, classical statistical tests, analysis of time lines, classification, grouping, etc., as well as graphic techniques and is extremely expandable (Long, 2015).

Despite the many advantages of R, the richness of its statistical models and data processing tools, as well as powerful visualization capabilities, problems arise when working with large amounts of data. R's limitations stem from the fact that it is designed to operate in single process calculation mode only (single CPU core) and data loaded at once in RAM.

4.2 Overcoming problems with big data using a multi-core microprocessor

Parallel programming calculation of more than one CPU core is possible by recompiling and adding some program components to R. This is possible due to the fact that R is an open source system and this is one of the advantages that this concept brings.

4.3 Methods for optimizing data volumes

One of the well-known features of the R language is that it loads all the data it operates in the RAM of the computer system, which would be critical even on powerful systems with a large resource. As a way to solve this problem, the dissertation looks at ways to load data into memory by excluding data with incorrect content at the time of loading it.

In some statistical surveys, it is not necessary to load all the data, but only certain time frames in order to make an approximate statistical analysis in a time frame. In this case, positioned reading methods suggested in the dissertation work can be applied. This makes it possible to process only a certain snippet of the data located in a big data file. Another common problem is when loading large data is that after loading in memory and processing, some of the data is no longer needed, but it continues to occupy significant

amounts of memory. The dissertation presents a way to reduce memory data by removing excess data and releasing memory.

4.4 Conclusions

The author's contribution is that this material, with the means of programming language R, helps solve problems when working with large data sets with limited computer resources.

In conclusion, it can be said that with the examples presented so far, the topic of optimized data loading cannot be exhausted when working with programming language R. Working with real data is always a challenge (Baumer, 2017). But the techniques presented so far are between good practices and are often used, they could also be combined with other approaches to problem solving in this area.

Conclusion - summary of the results obtained

The dissertation examined in detail methods and means of using time lines in solving various tasks arising in modern applications of information technology and systems.

A method titled MA Volatility Indicator has been proposed to improve precision in oscillator (Momentum). MA Volatility Indicator works in the combination of two EMA or SMA tools and offers a new methodology for interpreting results, which contributes to the detection of levels of over-purchase and over-selling in the market trend. All EMA, SMA and Momentum tools used in the study, as well as MA Volatility Indicator use time series.

The applicability of the neural networks apparatus for forecasting time lines in the financial field has been examined. It has been shown that a new model of presentation of input data characteristic of financial indicators results in a higher degree of self-adaptation in neural network training. Experiments conducted confirm the complexity of the financial processes and the presence of high-frequency noise in the data.

A method for examining the quality of RNG and PRNG in an information system has been developed by applying time lines to increase entropy in the use of cryptography providing various Internet services. This contributes to better cyber security of it infrastructure for digital resources and data protection. In the dissertation, the topic of cryptography received special attention due to its critical importance.

The practical results of the actual experiment showed that the golden ratio between mass services and actual cyber security requirements was found.

In view of the work carried out in this dissertation and the results obtained in the course of the studies and set out above, the following scientific and applied results may be formulated:

1. A method entitled MA Volatility Indicator has been developed to combine indicators for detecting price movements with new approaches when using time lines of financial data.

2. The apparatus of artificial neural networks shall be applied for the purpose of examining financial time lines. An algorithm has been developed to train the neural network by increasing the size of the neural network input and creating a hybrid structure, and a model for self-build three-layer MLP has been proposed.

3. A method has been developed to increase cryptographic protection in information systems based on studies on the quality of random number generators.

4. Experimental research has been carried out to solve cyber security problems in public widespread hosting services. The results obtained confirm the validity of the proposed method of enhancing cyber security.

5. Programming methods have been developed for efficient operation with large data with means in the R language.

6. The developed methods for increasing cryptographic protection are implemented in the technological infrastructure of IICT-BAS. A study of cryptographic tests and the quality of entropy on real-world busy server systems with public Internet services was conducted.

Guidelines for future research

The guidelines for future research on the subject of the dissertation include:

- Implementation of the MA Volatility Indicator method and its application in combination with other methods of analysis and forecasting of market price trends;
- Application of the MA Volatility Indicator method to automated systems for analyzing market trends and extracting decision-making signals;
- Conducting more research in the field of training algorithms and neural network systems for analyzing and forecasting time lines;

- The development of new methods to increase cryptographic protection in information systems;
- Research of a combination of developed method with other RNG methods and systems for analysis in cryptography and other technological fields to help create and improve RNG, as well as to more accurately determine the range of tasks that the generator can perform well;
- Find more approaches to loading and filtering big data to make it more efficient.

Publications on the subject of the dissertation work

- 1 **Ivan Blagoev**, Nikolay Dokev, Combining Momentum with one method for predicting market price movements for more accurate results (Combination of Momentum with One Method for Forecasting Market Trends to Improve the Results), International Scientific Conference "UNITH'17" – Gabrovo, 2017 Selected papers, ISSN 2603-378X, pp. II-265-II-270
- 2 **Ivan Blagoev**, Methods for Optimized Use of Computer Memory during Data Loads with R Programming Languages, International Conference "Automatics and Informatics'2017", 4-6 October 2017, Sofia, Bulgaria, ISSN:1313-1850, pp.213-215.
- 3 **Blagoev I.**, Improving the Momentum Oscillator Accuracy by a Method for Forecasting of Market Price Movements, Collection of Reports from An International Conference, Vasil Levski National University, 14-15 June 2018, vol. 9, p. 1. 177-185. (ceeol.com)
- 4 **Blagoev I.**, Method for More Reliable Users' Authentication on the Internet, Collection of Reports from International Conference, Vasil Levski National University, 14-15 June 2018, Vol. 9, p. 1. 167-176. (ceeol.com)
- 5 **Blagoev, I.**, Using R Programming Language for Processing of Large Data Sets, Proc. Int. Conf. Big Data, Knowledge and Control Systems Engineering – BdKCSE'2018, 21-22 November 2018, Sofia, Bulgaria ISSN 2367-6450, pp. 91-98.
- 6 **Ivan Blagoev**, Application of Time Series Techniques for Random Number Generator Analysis, Proceedings of XXII Int. Conference DCCN 2019, September 23-27, 2019, Moscow, Russia, pp.437-446. ISBN 978-5-209-09683-2, 2019 (RINZ).

- 7 **Blagoev I.**, Neglected Cybersecurity Risks in the Public Internet Hosting Service Providers. *Information&Security International Journal* - ISIJ, 47, no. 1, pp. 62-76 (2020)
- 8 Balabanov T.D., **Blagoev I.I.**, Dineva K.I. (2018) Self Rising Tri Layers MLP for Time Series Forecasting. In: Vishnevskiy V., Kozyrev D. (eds) Distributed Computer and Communication Networks. DCCN 2018. Communications in Computer and Information Science, vol 919. Springer, Cham. https://doi.org/10.1007/978-3-319-99447-5_50, pp. 577-584, **SJR:0.188**
- 9 **Blagoev I.** (2020) Method for Evaluating the Vulnerability of Random Number Generators for Cryptographic Protection in Information Systems. In: Dimov I., Fidanova S. (eds) Advances in High Performance Computing. HPC 2019. Studies in Computational Intelligence, vol 902. Springer, Cham. https://doi.org/10.1007/978-3-030-55347-0_33 **SJR:0.215**

Citations noted

- I **Blagoev, I.**, 2018. Using R Programming Language for Processing of Large Data Sets, Proc. Int. Conf. Big Data, Knowledge and Control Systems Engineering – BdkCSE'2018, pp. 91-98.

It is quoted in:

- 1 Dineva, K., Atanasova, T.: Regression Analysis on Data Received from Modular IoT System. ESM'2019, EUROSIS-ETI, ISBN: 978-9492859-09-9, EAN: 9789492859099, pp.114-118, 2019
 - 2 Ivaylo Blagoev, G. Vassileva and V. Monov, "Methodology for content preparation of online courses," 2020 International Conference Automatics and Informatics (ICAI), Varna, Bulgaria, 2020, pp. 1-4, doi: 10.1109/ICAI50593.2020.9311364.
- II **Blagoev I.**, Neglected Cybersecurity Risks in the Public Internet Hosting Service Providers. *Information&Security International Journal* - ISIJ, 47, no. 1, pp. 62-76 (2020)

It is quoted in:

- 3 M Terzieva, D Karastoyanov, ICT for Innovation in Advanced Banking, PROBLEMS OF ENGINEERING CYBERNETICS AND ROBOTICS • 2020 • Vol. 73, pp. 47-54 p-ISSN: 2738-7356; e-ISSN: 2738-7364, doi: 10.7546/PECR.73.20.05
- III **Blagoev I.**, Method for more reliable users' authentication on the Internet, Collection of reports from international conference, Vasil Levski National University, 14-15 June 2018, vol. 9, p. 1. 167-176.

It is quoted in:

- 4 Ivaylo Blagoev, Gergana Vassileva and Vladimir Monov, "Methodology for content preparation of online courses," 2020 International Conference Automatics and Informatics (ICAI), IEEE, Varna, Bulgaria, 2020, pp. 1-4, doi: 10.1109/ICAI50593.2020.9311364.
- 5 Dineva, K., Atanasova, T.: Security in IoT Systems. Proceedings 19th International Multidisciplinary Scientific Geoconference SGEM 2019, 19, 2.1, ISBN:978-619-7408-79-9, ISSN:1314-2704, DOI:10.5593/sgem2019/2.1, 576-577. SJR (Scopus):0.232 Q4

Participation in projects

- 1 National Scientific Program "Information and Communication Technologies for a Digital Single Market in Science, Education and Security" (ICT in the IA) -2018-2021.
- 2 Project Zora on Order No 147/14.06.2019 "Digital and Cyber Sustainable IICT"

Awards

1. Award of IICT-BAS for excellence in 2019 in the category "PhD students".

Bibliography

- 1 Atanasova, T., Barova, M.: Exploratory analysis of Time Series for hypothesized feature values. In: International Scientific Conference UniTech 2017, vol. II, pp. 399-403, University publishing house V. Aprilov, Gabrovo (2017)
- 2 Balabanov T.D., Blagoev I.I., Dineva K.I. Self Rising Tri Layers MLP for Time Series Forecasting. In: Vishnevskiy V., Kozyrev D. (eds) Distributed Computer and

- Communication Networks. DCCN 2018. Communications in Computer and Information Science, vol 919. Springer, Cham. https://doi.org/10.1007/978-3-319-99447-5_50 (2018)
- 3 Balabanov, T., Atanasova, T., Blagoev, I., Activation Function Permutation for Multilayer Perceptron Training, International Conference on Big Data, Knowledge and Control Systems Engineering BdKCSE'2018, Sofia, Bulgaria, ISSN 2367-6450, pp. 9-14 (2018)
 - 4 Blagoev I., Dokev N.: A Method for Investigating the Alterations in the Price Trends of the Currency Markets and Forecasting of Probable Future Alterations, *Problems of Engineering Cybernetics and Robotics*, vol.65, pp.39-48 (2012)
 - 5 Blagoev I., Neglected Cybersecurity Risks in the Public Internet Hosting Service Providers. *Information&Security International Journal - ISIJ*, 47, no. 1, pp. 62-76 <https://doi.org/10.11610/isij> (2020)
 - 6 Blagoev I.: Method for Evaluating the Vulnerability of Random Number Generators for Cryptographic Protection in Information Systems. In: Dimov I., Fidanova S. (eds) *Advances in High Performance Computing. HPC 2019. Studies in Computational Intelligence*, vol 902. Springer, Cham. https://doi.org/10.1007/978-3-030-55347-0_33. (2021)
 - 7 Blagoev, I., Using R Programming Language for Processing of Large Data Sets, Proc. Int. Conf. Big Data, Knowledge and Control Systems Engineering – BdKCSE'2018, 21-22 November 2018, Sofia, Bulgaria ISSN 2367-6450, pp. 91-98.
 - 8 Camara C., Martín H., Peris-Lopez P., Aldalaien M., Design and Analysis of a True Random Number Generator Based on GSR Signals for Body Sensor Networks, *Sensors* 19, 2033; doi:10.3390/s19092033 (2019)
 - 9 Plummer T., *Forecasting Financial Markets: The Psychology of Successful Investing*, January (2010)
 - 10 Pseudo-Random Number Generators, <https://crypto.stanford.edu/psc/notes/crypto/prng.html>
 - 11 Zhao P., R with Parallel Computing from User Perspectives, <https://www.r-bloggers.com/r-with-parallel-computing-from-user-perspectives/> (2016)
 - 12 Brown R. G.: Dieharder: A Random Number Test Suite, <https://webhome.phy.duke.edu/~rgb/General/dieharder.php> (2021)
 - 13 Ciampi F., G. Marzi, S. Demi, M. Faraoni, The big data-business strategy interconnection: a grand challenge for knowledge management. A review and future perspectives, *Journal of Knowledge Management*, Vol. 24, Issue 5 (2020).

- 14 Koeune F. Pseudo-random number generator. In: van Tilborg H.C.A. (eds) Encyclopedia of Cryptography and Security. Springer, Boston, MA . https://doi.org/10.1007/0-387-23483-7_330 (2005)
- 15 Borcard, D., Gillet, F., Legendre, P. Numerical Ecology with R, Springer, pp. 9 – 30 (2011)
- 16 Baumer B. S., Kaplan D. T., Nicholas J., Modern Data Science with R, Horton Chapman & Hall/CRC, Boca Raton, (2017)
- 17 Long C. (Ed.) Data Science & Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data, John Wiley & Sons, Inc., (2015)
- 18 Martínez-Acosta L., Medrano-Barboza J.-P., López-Ramos Á., López J., López-Lambrano Á., SARIMA Approach to Generating Synthetic Monthly Rainfall in the Sinú River Watershed in Colombia, *Atmosphere*, 11, 602; doi:10.3390/atmos11060602 (2020)
- 19 Mikalef P., Krogstie J., Examining the interplay between big data analytics and contextual factors in driving process innovation capabilities, *European Journal of Information Systems*, Volume 29, - Issue 3: Business Process Management and Digital Innovation <https://doi.org/10.1080/0960085X.2020.1740618> (2020)
- 20 Scott G., Carr M., Cremonie M. , Technical Analysis: Modern Perspectives, e CFA Institute Research Foundation (2016)
- 21 The R Journal, ISSN: 2073-4859, <https://journal.r-project.org/> (2017)
- 22 Tomov, P., Monov, V., Artificial Neural Networks and Differential Evolution Used for Time Series Forecasting in Distributed Environment, Proc. of Int. conference Automatics and Informatics, ISSN 1313-1850, pp.129-132, Sofia, Bulgaria, (2016)
- 23 Wafi A.S., Hassan H., Mabrouk A., Fundamental Analysis Models in Financial Markets – Review Study, *Procedia Economics and Finance*, Vol. 30, 939 – 947. Elsevier (2015)
- 24 Wang, W., Y. Wang, Analytics in the era of big data: The digital transformations and value creation in industrial marketing, *Industrial Marketing Management*, Vol. 86, pp. 12-15, ISSN 0019-8501, <https://doi.org/10.1016/j.indmarman.2020.01.005> (2020)
- 25 Li C., Zhang J., Sang L., Gong L., Wang L., Wang A., Wang Y., Deep Learning-Based Security Verification for a Random Number Generator Using White Chaos, *Entropy*, 22, 1134; doi:10.3390/e22101134 (2020)



БЪЛГАРСКА АКАДЕМИЯ НА НАУКИТЕ

АВТОРЕФЕРАТ НА ДИСЕРТАЦИЯ

за присъждане на образователна и научна степен “доктор” по научна специалност “Информатика“

МЕТОДИ И СРЕДСТВА ЗА АНАЛИЗ НА ДАННИ В ИНФОРМАЦИОННИ СИСТЕМИ С ИЗПОЛЗВАНЕ НА ВРЕМЕВИ РЕДОВЕ

Иван Иванов Благоев

Ръководител: Доц. Татяна Атанасова

Научно жури:

**Проф. Иван Гарванов
Проф. Радослав Йошинов
Проф. Владимир Монов
Доц. Десислава Иванова
Доц. Велизар Шаламанов**



**Институт по информационни и
комуникационни технологии
Секция „Моделиране и оптимизация“**

Увод

Напредъкът в технологиите е толкова очевиден, че може само да се спомене, без нужда от фактологично описание. В това отношение, значителна разлика от последно време е силно експанзиращата дигитална трансформация. Поради COVID-19 заплахата за човешкото здраве, скоростта на навлизане на технологиите в нашият живот силно се ускори, което води до тотална промяна в множество дейности, а в следващите години ще се забелязва още по-силно, когато човечеството се трансформира и адаптира към този нов начин на живот.

Всичкото споменато до тук, води със себе си и до много напълно нови за науката и неизследвани до сега процеси. Събирането и обработката на времеви поредици и големи данни ще се разшири с проникване и в новите процеси. Нуждата от изследване и нови открития ще е решаваща за развитието на науката и технологиите в следващите години. За това разработката на нови методи и средства за изследванията с времеви редове и обработка на големи данни и е изключително важна и ще бъде основен инструмент за изследванията и развитието на науката и технологиите в бъдеще.

Настоящият дисертационен труд, чрез изследвания с времеви редове допринася за постигане на по-добри резултати при методи за прогнозиране на финансови инструменти, обработката на големи данни и подобряване на криптографията и киберсигурността.

Цел и задачи на дисертацията

Целта на настоящата дисертация е да се разработят нови методи и средства за анализ на данни в информационни системи с използване на времеви редове.

За тази цел се дефинират следните задачи:

- 1 да се разработи метод за анализ и предсказване на ценови движения във финансовата област с използване на времеви редове;
- 2 да се предложи алгоритъм за обучение на изкуствени невронни мрежи при прогнозиране на финансови времеви редове;
- 3 да се предложат решения за повишаване на криптографската защита в информационните системи чрез прилагане на методи за анализ на времеви редове;
- 4 да се проведат експериментални изследвания за верификация на предложените методи за повишаване на криптографска защита при решаване на задачите за осигуряване на киберсигурността.

- 5 Да се разработят програмни методи за преодоляване на проблеми при работа с големи обеми от данни във времеви редове.

Структура на дисертацията

Дисертационният труд е структуриран в четири глави.

В **първа** глава е направен преглед на актуалните теми в областта на науката на данните, особено когато тези данни се представят като времеви редове. Мотивирана е необходимостта от разработване на нови методи и средства за анализ на данни в информационни системи с използване на времеви редове.

Във **втора** глава са представени разработените методи за изследване и прогнозиране на финансовите времеви редове с използване на различни математически апарати.

В **трета** глава са описани разработените решения за осигуряване на криптографска защита при предоставяне на информационни услуги чрез изследване на генератори на случайни числа, представляващи поредици от времеви редове. Представено е практическото приложение на предложените подходи за обезпечение на киберсигурността. Показани са реалните резултати от проведените тестове, доказващи успешното решаване на поставените задачи.

В **четвърта** глава преодоляването на проблеми при работа с големи масиви от данни и ограничени компютърни ресурси при изследване на времеви редове е направено с разработените софтуерни подходи и със средствата на език за програмиране R.

В **Заклучението** е представено резюме на получените резултати от разработката. Определени са насоки за бъдещи изследвания и развитие. Представен е списък с научни публикации по темата и забелязани цитирания.

Дисертационният труд съдържа 125 страници, 33 фигури, 1 таблица и 122 литературни източника.

Глава 1. Анализ на състоянието на изследванията.

Ако се погледне в околния свят през окото на технологиите, първото което би впечатлило всеки специалист е колко много данни са това. Това е страничен ефект от масовата дигитална трансформация и автоматизацията (Wang, 2020), оставяйки цифрова следа от изпълнението на реалния процес. Тези цифрови следи отразяват случващото се в реалния свят и позволяват задълбочен анализ на основните процеси. Динамичните времеви редове в комуникациите, технологии, бизнеса идват в резултат

на измерване на характеристики от технически, природни, социални, икономически и други системи (Mikalef, 2020), (Ciampi, 2020).

1.1 Времеви редове

Времевите редове представляват редици от данни, събрани на равни или неравни интервали от време. Основна характеристика на времевия ред е, че всяка следваща стойност е в зависимост от предходните стойности. Тази зависимост може да бъде, както много сложна, така и относително проста. Понастоящем много методи за прогнозиране, които действат като ефективни инструменти, са широко приети за оценка и анализ на данни от модели на времеви редове. От тях, най-често използваният модел е интегриран метод на авторегресия със сезонен компонент (SARIMA - Seasonal ARIMA), който по същество принадлежи към линеен модел. Но на практика при решаване на различни задачи в информационни системи най-често срещаното е, че процесът на генериране на данни е силно нелинеен и прогнозите получени с тези модели, не позволяват да се достигне до точните резултати (Martínez-Acosta, 2020).

1.2. Приложение на времеви редове върху финансови инструменти

Пазарните ценови движения се описват чрез времеви редове и са предмет на анализ от финансисти, икономисти и пазарни стратегии. Видове финансови анализи, които към настоящия момент се използват за анализ на финансови инструменти:

- Фундаменталният - основава на анализиране на събитията, случващи се по света и касаещи финансовите и стокови пазари (Wafi, 2015);
- Техническият анализ – основава се предимно на статистически методи с изчисления върху времевите редове. Позволява методите за прогнозиране да бъдат описани чрез статистически средства и математически алгоритми (Plummer, 1991), (Scott, 2016).

1.2.1 Невронни мрежи

Подходите за изследване на времеви редове могат да бъдат разделени на две категории: статистически методи и изчислителна интелигентност. Статистическите методи изследват зависимости между изходните и съответните фактори след изучаване на минали данни, докато другата група методи имитира човешкия начин на мислене и логическо заключение, за да придобие знания от миналия опит (като изкуствени невронни мрежи) и да предвиди бъдещи стойности (Atanasova, 2017). Изкуствените невронни мрежи (ANN) се използват в различни научни и ежедневни задачи. Обикновено ANN се представят като претеглен насочен граф и има много различни

конфигурации на тази схема. В най-простия случай това е многослоен персептрон. Времеви редове са атрактивни за изследвания с изкуствени невронни мрежи (Tomov, 2016).

1.3. Приложение на времеви редове при криптографията и кибер сигурността

Изпълнението на изискванията за киберсигурност е предпоставка за безопасността и сигурността на ИТ инфраструктурите, цифровите ресурси и защитата на личните данни. В нейният фундамент е криптографията, която осигурява редица процеси, като идентификация, удостоверяване, кодиране, автентикация, потвърждение за състояние на процеси и данни и др. Основният корен на криптографията са случайните числа, като в най-честия случай за съвременните нужди на криптографията се използват два вида генератори на произволни числа:

- Генератор на случайни числа (RNG);
- Псевдо генератор на случайни числа (PRNG).

Традиционните мерки за RNG са предимно обобщена статистика, отнасяща се до отклонения от математическата случайност. За да се подпомогне проверката на качеството на генератор на случайни числа, може неговият изход да се запише във времеви ред и данните да бъдат подложени на специализирани математически анализи.

1.4 Изводи

В резултат на направените изводи следва да се обобщи, че изследванията на времеви редове в различни области и приложения се нуждаят от разработка на специфични методи и средства за постигане на конкретните цели.

Глава 2. Методи за изследване и прогнозиране на финансовите времеви редове

В тази глава е разгледано изследване върху широко разпространения индикатор Моментум, който принадлежи към групата на осцилаторите. Изчисляването му се базира на математически апарат за обработка на времеви редове. В дисертацията се цели подобряването на неговата ефективност.

2.1.1 Осцилатор Моментум (Momentum Oscillator)

Моментум е основен осцилатор, който показва дали ценовата тенденция се ускорява, забавя или се движи със същата скорост. Той обикновено достига максималната си стойност преди върха на цените и минимума преди дъното на спада.

Функцията на този осцилатор е да отчита ускорението на ценовата тенденция. При изтощаване на текущата тенденция и наличие на вероятност от промяна на същата, Моментум дава сигнал за дивергенция. Това е момент при който цената продължава да се движи в посока на тенденцията, но стойностите на Моментум намаляват при възходяща ценова тенденция или се повишават при низходяща ценова пазарна тенденция. Като потвърждение на сигналите за дивергенция на Моментум в изследването са включени фигурни формации от техническия анализ за комбиниране и потвърждение на текущия ценови обрат. В конкретния пример на фиг.2.2. Моментум и сигналът за дивергенция 1С потвърждава предстоящ ценови пазарен обрат, чрез множествен връх с 1D.



Фиг. 2.2. Дивергенция на Моментум при пазарен тренд EUR/USD на дневна база
(исторически данни от Forex пазарът)

2.1.2 Слабости при анализ на пазарния тренд чрез Моментум

В изследването обаче се вижда, че има случаи при които Моментум може да направи изключение и да не отчете дивергенция при завършване на текуща пазарна тенденция, което е изобразено с 1С на фиг.2.3.



Фиг. 2.3. Валутна двойка USD/CAD на дневна база (исторически данни от Forex пазарът)

Според изнесения пример от реалния форекс пазар, стойността на Моментум начертава по-висок връх даже от предходния, но цената след това прави значителна корекция от около 60% без да е на лице дивергенция при осцилатора. Дори напротив, според сигнала, текущата ценова тенденция се потвърждава с ускорението на осцилатора. Въпросът, който вълнува изследването следователно е, дали може да се подобри точността на осцилатора Моментум?

2.1.3 Метод за повишаване точността на Моментум

В този дисертационен труд се предоставя нетрадиционен метод за получаване на сигнал за пазарен обрат, а именно разработеният **MA Volatility Indicator**. Базира се на нетрадиционен начин на използване на индикатор от тип пълзяща средна линия Moving Average (MA). Индикаторът MA се разделя на два под вида:

- Simple moving average (SMA):

За пресмятане на SMA, се използва времеви ред, при който се сумират данните на последните периоди (t), където например $t=10$ за 10 дена, според времевата рамка (може да бъде различна стойност, по избор). След това се дели на броя t периодите. Такова пресмятане се прави за всеки един бар за период от графиката. Формулата за SMA е, както следва:

$$SMA_t = \sum_{n=1}^t price_n / t$$

- Exponential moving average (EMA):

С цел да се намали изоставащия ефект на SMA, ползващите технически анализ често предпочитат Exponential Moving Average (ЕМА). Те намаляват изоставането чрез добавяне на нови стойности върху най-новите цени, зависещи от дължината на МА. Най-кратката ЕМА ще е с по-голяма стойност, отколкото ще бъде приложена за повечето МА. Формула за пресмятане на ЕМА:

$$X = K * (C - P) + P,$$

където X – настояща ЕМА, C – настояща цена, P – ЕМА от предния период (за пресмятането на първия период се използва стойност от SMA), K – изглаждащ коефициент.

Изглаждащият коефициент прилага подходящ коефициент към по-новите цени, които са свързани с предходните цени на ЕМА. Формула за изглаждащия коефициент:

$$K = 2 / (1 + N),$$

където N – брой на предходните ЕМА цени.

Конвенционален подход за търговия чрез МА е на по-висока времева рамка цената да не пресича МА, като при пазарна корекция достигането на МА от пазарната цена се счита за силна подкрепа за текущата тенденция. При пробив на цената на МА се приема, като сигнал за обрат, а при отскачане, като сигнал за потвърждение на текущата пазарна тенденция. Другият метод е анализ с повече от една МА, като всичките МА са с различна скорост. При пробив или отскачане на по-бързата МА към по-бавната МА, се тълкува аналогично за сигнал за потвърждение или обрат в текущата тенденция.

Разработеният в дисертация метод MA Volatility Indicator разчита на определяне на екстремни стойности за отдалечаване на цената от МА, на база на което да се определи сантимента на участниците на пазара към текущия момент. При екстремно високи стойности на отдалечаване на цената от МА по посока на текущата тенденция, трябва да се счита, че това е сигурен сигнал за предстоящ обрат или дълбока корекция на текущата тенденция. Интересното е, че това явление се наблюдава добре в моменти, когато Моментум не дава сигнали за дивергенция и край на ускорението на текущата ценова тенденция.



Фиг.2.8 Комбиниране на Моментум с предлаганото решение (исторически данни от Forex пазарът USD/CAD)

На фигура 2.8 методът MA Volatility Indicator е приложен и комбиниран със стойностите на Моментум, като данните за симулацията са реални исторически от форекс пазара. На фигура 2.8 ясно личи, че Моментум, изобразен със зелена линия, не отчита по-ниски стойности на последния пазарен връх. Това е подчертано с правата червена линия на неговата тенденция. В същото време личи и най-високата на отдалечаване на цената от пурпурната линия на MA със стойност 554 пипса.

В заключение може да се каже, че Моментум е един ефективен осцилатор, който е станал част от множество автоматизирани системи и стратегии за търговия на финансовите пазари. Но разработения в дисертацията методът MA Volatility Indicator, успява да подобри точността на прогнозиране. Следователно, той би могъл да се прилага, както при автоматизирани системи, така и при анализа на пазарните тенденции и от човек.

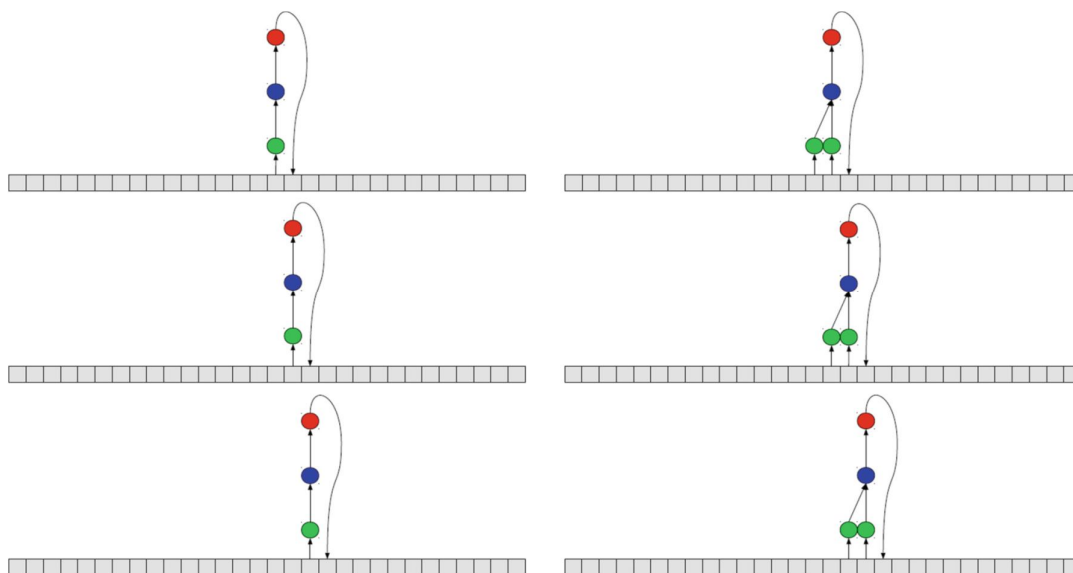
2.2 Прогнозиране на финансови времеви редове чрез невронни мрежи

Многослойният персептрон е най-често използваният вид на изкуствени невронни мрежи, който може да се представи като ориентиран претеглен граф. В това проучване основната идея е, че вместо броя на скрити слоеве, увеличава се броят на невроните на входа и скритите слоеве се разширяват по време на обучението на невронна мрежа. Удължаването на входния слой е свързано с факта, че всеки времеви

ред расте с поява на ново измерване. Целта на обучението е размерът на входния слой да бъде толкова голям, колкото размерът на пълната времева редица.

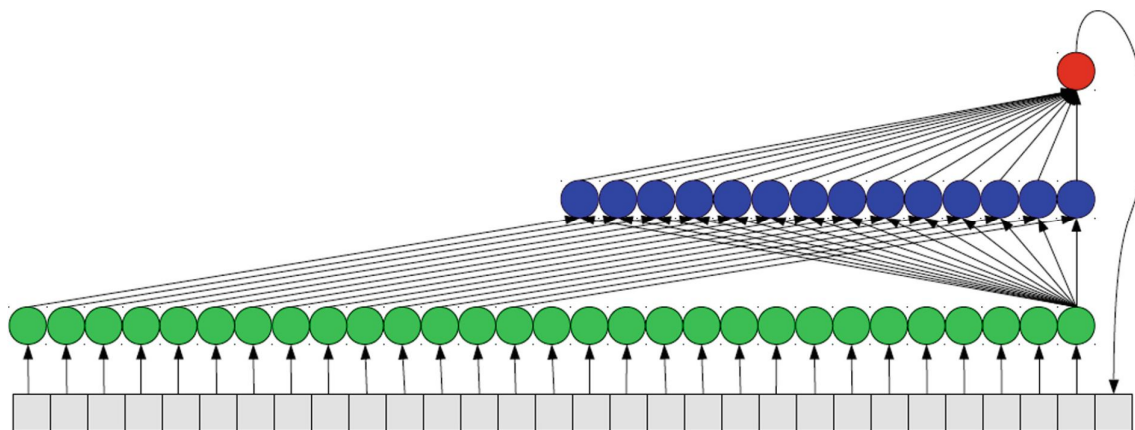
2.2.1 Предпоставки при моделирането

В предложения модел се използва набор от изкуствени невронни подмрежи и тези подмрежи се обединяват в обща изкуствена невронна мрежа. Най-малката изкуствена невронна подмрежа има 1-1-1 топология (фиг. 2.2.6 - вляво). Мрежата е обучена с примери, чийто вход има само една стойност. Целта в модела е прогноза за само една стойност напред във времето. Ето защо всички подмрежи имат само един изход. Всичките входни стойности се предоставят като примери за еластично обучение за обратно разпространение на грешката. Обучението спира на определено ниво на *epsilon* за пълна промяна на грешките на невронната мрежа.



Фиг. 2.2.6. Обучение на изкуствени невронни подмрежи с 1-1-1 топология (вляво) и 2-1-1 топология (вдясно).

След обучение на 1-1-1 топология стойностите на теглата на първата подмрежа се зареждат във втората подмрежа с 2-1-1 топология (Фиг. 2.2.6-вдясно). Трета подмрежа има 3-2-1 топология. Размерът на скрития слой се избира автоматично чрез алгоритъм за постепенно подрязване, внедрен в Encog Machine Learning Framework (<http://www.heatonresearch.com/encog/>). Топологиите на подмрежите се формират чрез добавяне на един неврон във входния слой и коригиране на размера на скрития слой с алгоритъм за инкрементално подрязване. Крайната цел е да се достигне n-m-1 топология (фиг. 2.8), която обхваща всички известни стойности на времеви редове.



Фиг. 2.2.8. Обучение на изкуствена невронна подмрежа с $n-m-1$ топология.

Някои от връзките между входния и скрития слой не се визуализират за по-добър вид.

Общата идея зад предложения модел е постепенното обучение на състезателни по размер изкуствени невронни мрежи. Често срещаният проблем при обучението на изкуствени невронни мрежи е размерът на мрежата. Чрез разделянето на най-голямата мрежа в много по-малки мрежи се постига ускоряване на процеса на обучението. Предложеният модел има по-висока степен на самоадаптация, тъй като когато се появи нова стойност във времеви ред, размерът на изкуствената невронна мрежа нараства, което означава, че фазата на обучение и фазата на работа са едновременни.

2.2.2 Експерименти върху изследването

Експериментите се правят чрез JAVA програма, където изкуствените невронни мрежи се изпълняват чрез API, предоставен от Encog Machine Learning Framework. Като входни данни за експериментите се използват финансови времеви редове на FOREX пазар. Данните се вземат от ежедневна двумесечна търговия за валутни двойки EUR / USD и USD / JPY. Стойностите на времевите редове се мащабират в диапазона от -0,99 до +0,99 с правилото за мащабирание MinMax. Изходът на изкуствената невронна мрежа се пренастройва до първоначалния диапазон със същото правило, което се използва в обратна посока. Резултатите от експериментите все още са в диапазона на статистическата грешка, която идва от сложността на финансовите процеси и високочестотния шум вътре в данните.

Предложеният модел за самонадграждащи се трислойни MLP за прогнозиране на времеви редове е обещаващ подход за ускоряване на обучението на изкуствени невронни мрежи. Нарастващият размер на входния слой включва максимална информация, налична във времевите редове, но предложената процедура за обучение

на изкуствена невронна мрежа отчита, че по-старите стойности трябва да бъдат по-малко информативни.

2.3 Изводи

В тази глава се предлагат нови методи за анализ и прогнозиране на пазарни ценови движения чрез времеви редове и невронни мрежи.

В резултат на това са направени следните заключения:

- 1 В изследването до тук се обхващат основните аспекти от процеса по анализ – от дефиниране на проблема и поставяне на задачите, до представянето на методи за решаването им. Във всеки един от етапите се извършва представяне на реални доказателства, чрез които може да се идентифицира наличието на слабости или необходимост от намиране на по-рационален подход в разглежданата област.
- 2 Методите позволяват да се интегрират в системи за автоматизирана обработка и вземане на решения. Разработеният метод (MA Volatility Indicator) подобрява прецизността в осцилатор (Моментум) и работи в комбинацията от два инструмента ЕМА или SMA, като предлага нова методика за интерпретиране на резултатите при пазарни анализи и спомага за намаляване на риска от загуби и увеличаване на успех при автоматизираната търговия.
- 3 Предложеният алгоритъм за обучение чрез самонадграждане в трислойни MLP ускорява обучението на ANN при прогнозиране на финансови времеви редове.

Методите представени до тук, могат да бъдат прилагани от специалисти в различни области в системи с прогнозен характер, за вземане на решения, анализиращи събития и процеси базирани на времеви редове.

Глава 3. Решения за осигуряване на криптографска защита чрез приложение на времеви редове при криптографията и киберсигурността

В дисертацията се предлага подходът на времеви редове да се приложи към анализа на качеството на система за генериране на произволни числа (RNG) за осигуряване на криптографската защита в информационните системи. За текущото изследване от RNG се извлича числов масив, за да може да се анализират стойностите от случайните числа във времеви редове. Резултатите се изобразяват графично, където по-ясно стават видни произведените от генератора уязвими случайни числа.

3.1 Приложение на техники от времеви редове за анализ на генератор на произволни числа в областта на киберсигурността

RSA е асиметричен алгоритъм за криптиране, който позволява на всеки да изпраща криптирани съобщения, които само притежателя на частния ключ може да декодира. Принципът на работа може да се обясни накратко, като се генерира едно много голямо произволно число p , след това се генерира още едно такова число q и се изчислява тяхното произведение $x=p*q$, всъщност x е известен като публичен ключ.

3.1.1 Изследователите на (почти) секретния алгоритъм – слабости поради недостатъчна ентропия на RNG

На повърхността, RSA криптирането изглежда неуязвимо. Но според представеното изследване проблемът се крие в генераторите на случайни числа, които обезпечават алгоритъма. Уязвимостта е фундаментална и идва от там, че на RSA са необходими много големи числа, за да се създадат ключовете за криптиране, а генераторите в масовите компютърни системи са със значително по-малък капацитет. За това се налага да се използват псевдо генератор, който да се комбинира в качествените източници на ентропия, за да се изпълнят нуждите на алгоритъма. Чрез произлязла от генератора на числа началната стойност наречена семе (Seed), вложена в псевдо генератора и след трудоемки за компютърната система изчисления криптографските RSA ключове се генерират. Проблемът при устройства, като телефони, IoT, малки рутери и др. малки системи е още по-силно изразен, защото често те нямат достатъчно ресурс за тази трудоемка работа. За това в тях се залагат на готово предварително изчислени основни фактори необходими за съставянето на ключовете. Това значително ускорява процеса по генериране на RSA ключове при необходимост, но отваря голяма уязвимост в сигурността на криптографията. И отчитайки тези обезпокоителни наблюдения, те са достатъчно основание да се направи изследване по темата.

Съвременните критерии за надежден RSA ключ е минимум 2048 бита, като препоръчителната дължина е даже 4096 бита. При други изследвания също е установено, че между 4096, 8192 и 16384 бита RSA ключ, по-голямата сигурност на по-големите ключове е минимална. Причината също идва от ограниченията при генераторите на случайните числа. При по-големи RSA ключове са необходими

изключително големи истински случайни числа. Които в една компютърна система е крайно трудно да се получат.

Ако слабостите в криптографски функции не се осветяват, рискуваме да бъдат открити и да се използват от злонамерени лица без това да е известно на останалите. В заключение може да се каже, че слабостите не изхождат от грешка в аритметиката на RSA. Те идват от технологичната слабост, с която се прилага RSA.

3.2 Метод за оценка на уязвимостта на генераторите на случайни числа за криптографска защита в информационните системи

В предмета на изследването попада технологията на широко разпространения език за програмиране PHP. За нуждите на системи, разработвани с тази технология, за да се обезпечават необходимостта от случайни числа, PHP разполага със следните средства:

1. Линеен конгресен генератор (LCG), напр. `lcg_value()`
2. Алгоритъмът Marsenne-Twister, напр. `mt_rand()`
3. Локално поддържана функция `C`, т.е. `rand()`

Те се използват повторно и за функции като `array_rand()` и `uniqid()`, като недостатъкът на ентропията и генераторите на произволни числа на гореописаните функции се състои в лесното прогнозиране на бъдещите стойности на PRNG. Причината е, че първоначалните вътрешни състояния или SEED на PRNG са ограничени и изходът на стойности е в недостатъчен диапазон и това е предвидимо от лесно достъпните съвременни изчислителни ресурси. Често, за да получат стойност за SEED в PHP, разработчиците използват `mt_rand()` или следния скрипт, за да се използва автоматично:

```
<?php
mt_srand(3231153718);
for ($i=1; $i < 15; $i++) {
    echo mt_rand(), PHP_EOL;
}
```

Което поради слабата ентропия на предлаганите инструменти, води до риск възстановяването на SEED от нападател. Въпреки пасивния си характер, това всъщност е истинска уязвимост. За целта в изследването се създава симулация на истинска информационна система, което използва следния изходен код за генериране на токен за различните цели на приложението:

```
$newtoken = hash('sha512', mt_rand());
```

Генериране на токен по представения начин е хубав пример, като единично обръщане към `mt_rand()`, което се хешира с SHA512. Факт е, че в действителност, ако програмист приеме, че функциите на случайните стойности на PHP са "достатъчно случайни", той ще бъде много по-склонен да вгради прост модел на използване. Което е срещано многократно в практиката. Но използваният по-горе метод за генериране на маркери страда от един недостатък - случайните стойности са ограничени до цифри (т.е. неговата несигурност или ентропия е близка до незначителна). Ако се провери продукцията на `mt_getrandmax()`, ще се открие, че максималният произволен брой `mt_rand()` може да генерира само 2,147 милиарда. Този ограничен брой опции го прави уязвим за груба атака. При наличие на съвременна добра видео карта (GPU) и с помощта на специализиран софтуер за атака с груба сила като `hashcat`, такова изчисление може да се завърши само в рамките на няколко минути. Следователно използването на хеш за скриване на изхода на `mt_rand()` е безполезно.

За да се защити този тип система, трябва да се генерират случайни стойности с по-високо качество. За използване в нетривиални задачи, PHP изисква източници на ентропия от висок клас, които могат да бъдат осигурени от операционната система. В Linux обикновено се използва с `/dev/urandom`, освен ако не са инсталирани устройства с още по-висока ентропия. В Linux, с правилната настройка, редовен генератор на произволни числа, който е от типа PRNG (който е псевдо генератор на произволни числа), често се зарежда от източник на висока ентропия `/dev/random`, което го прави устойчив на атаки. Следователно всяка една софтуерна система разработвана с PHP, за да бъде добре защитена следва да се пренасочи към функциите за повторно използване на външната библиотека на OpenSSL. Като се извикват функциите `openssl_pseudo_random_bytes()` и `mcrypt_create_iv()`. Те са оптимизирани да използват криптографски защитен псевдослучайни генератор. Който е съобразен и интегриран с операционната система.

3.2.2. Разбиране на RNG Entropy в Linux

В операционната система Linux архитектурата за получаване на случайните числа има следния вид:

1. `/dev/random` е истински генератор на случайни числа, ако свърши ентропията блокира. Този инструмент си осигурява ентропията, събрана от

системните параметри по време на работата му като достъп до диск, мрежов трафик, състояние на паметта, преместване на мишката и други прекъсвания на системата;

2. `/dev/urandom` е генератор на псевдо произволни числа (PRNG) и той не се блокира поради изчерпването на ентропията. Може да се използва за рандомизиране на неограничен поток. Случайният поток се осигурява от PRNG структури и необходимите начални SEED стойности ще се презареждат периодично от `/dev/random`;
3. `/dev/hwrng` е допълнителен хардуер за истински случайни числа, който е специализиран и не е инсталиран в компютърните системи по подразбиране. Той осигурява шум от ентропия за поддържане на случайни числа;

Натрупаната ентропия в Linux система може да бъде проверена чрез следната команда:

```
$ cat /proc/sys/kernel/random/poolsize
4096
$ cat /proc/sys/kernel/random/entropy_avail
3868
```

където:

`/proc/sys/kernel/random/poolsize` се използва за деклариране на размера (в битове) на буфера Entropy Pool, например: Колко произволни числа трябва да съхраним, преди да спрем да „помпаме“ за повече.

`/proc/sys/kernel/random/entropy_avail` показва количеството (в битове) на текущо съхранени случайни числа в пула.

Чрез потребителската активност и работата на компютърната система, като мрежа, дискове, състояние на паметта, централен процесор, периферия и др. специалните функции в ядрото на Linux имат функции за непрекъснато набавяне на случайни числа. Коего има за цел да компенсира непрестанната нужда от такива, при работа на компютърната система. Факт е, че колкото и да се опитва ядрото на операционната система да компенсира със случайни числа буферите, в определени моменти е възможно те да бъдат източвани много по-бързо. За нуждите на изследването лесно може да се предизвика такава ситуация, за да може да бъде наблюдаван този процес. Чрез следващата команда, просто се изхвърли всичко, което е в `/dev/random` генератора на произволни числа и се извежда на екрана:

```
$ hexdump /dev/random
00000000 d5c4 ff0a b8ef 9bdc ad95 480b e853 f0ef
00000100 e0cb 7c08 4bc4 daef 2b21 ea62 0eac 2c6c
00000200 d6bd 70e6 5d6f a7e3 0874 d52f 77df 6a2b
00000300 1909 efe8 9964 acee 2aad 2522 4ddb 1d0b
```

В същия момент може в паралелно отворен команден терминал да се изведе състоянието на буфера на ентропия, като съдържанието се обновява всяка секунда. За целта е необходимо да се стартира следната комбинация от команди:

```
$ watch -n 1 cat /proc/sys/kernel/random/entropy_avail
```

Като резултат наличието на ентропия ще започне да спада, като неговото състояние ще стигне критични стойности, дори и до нула. С натисне на Ctrl-C се спира това безсмислено разхищение. Може би никога не трябва да се прави това на практика, особено на реална сървърна система, освен с изследователска цел разбира се. Но често системите имат проблеми с натрупването на ентропия в буфера и резултатът изглежда смущаващ:

```
$ cat /proc/sys/kernel/random/entropy_avail
96
```

От представения пример машината произведе ентропиен резултат от 96 бита и увеличаването на тази стойност е твърде бавно и недостатъчно. Причините за това могат да са разнородни. Например от липса на специфичен хардуер, неправилни настройки, виртуализация, твърде голяма активност със случайните числа на системата и невъзможност да се компенсира консумацията на случайни стойности и др. Едно възможно решение е да се стартира специализиран софтуер подпомагащ събирането на случайни числа. Това е демон, който е проектиран да използва всякакви събития, които могат да се считат за сравнително случайни при работата на машината, за да се произведат повече и по-качествени случайни числа. Например процесорното „трептене“, промяната в състоянието на паметта, входно изходни операции, мрежов трафик могат да добавят още ентропия към буфера на системата. Инсталирането на това решение и основната настройка в системата са следните:

```
# apt install haveged
# systemctl start rngd
# update-rc.d haveged defaults
```

```
# rngd -r /dev/urandom
```

На система със сравнително умерен трафик:

```
# pv /dev/random > /dev/null
 40 B 0:00:15 [ 0 B/s] [          <=>          ]
 52 B 0:00:23 [ 0 B/s] [          <=>          ]
 58 B 0:00:25 [5.81 B/s] [          <=>          ]
 64 B 0:00:30 [6.05 B/s] [          <=>          ]
```

^C

```
# systemctl start haveged
```

```
# pv /dev/random > /dev/null
7.12MiB 0:00:05 [1.43MiB/s] [          <=>          ]
15.7MiB 0:00:11 [1.44MiB/s] [          <=>          ]
27.2MiB 0:00:19 [1.46MiB/s] [          <=>          ]
 43MiB 0:00:30 [1.47MiB/s] [          <=>          ]
```

^C

С помощта на командата `pv` може да се види колко данни се предават за целта. От показания поток на данните се вижда, че преди `haveged` се получаваха 2.1 бита в секунда (B / s), докато след това се получават ~ 1.5 MB / sec.

3.2.3. Времеви редове за генератори на случайни числа

Спецификата на RNG и PRNG позволява те да бъдат анализирани чрез техники за анализ и прогнозиране на времеви редове, тъй като улавянето на потока на изходните числови стойности е последователност и само по себе си е подредена последователно във времето. Такъв поток от числови стойности може да бъде описан, както следва:

$$N = T * V$$

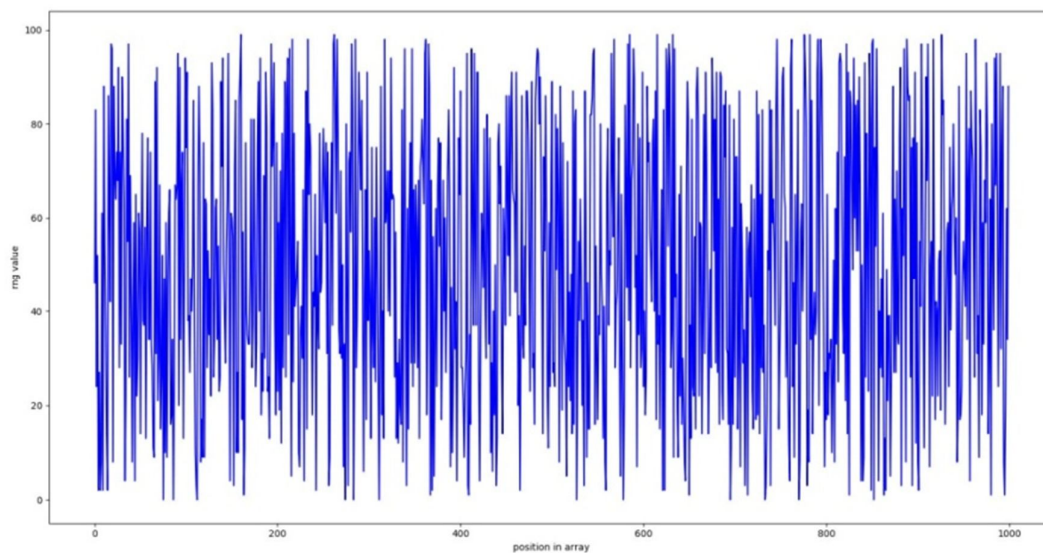
където: N - дължината на числовия ред, T - време (продължителност) на генерирането на числа, V - брой генерирани числа за единица време.

Така, че чрез времевите редове е възможно да се определи качеството на ентропията във времето. Ако даден генератор на случайни числа не е много надежден, то неговите слабости биха могли да се намерят за по-кратък времеви ред с данни, за които ще са необходими по-малко ресурси за обработка и анализ. За нуждите на текущото изследване ще се използва числов масив, който няма да бъде създаден от висококачествен генератор на случайни числа, а от посредствен такъв. Идеята е да се

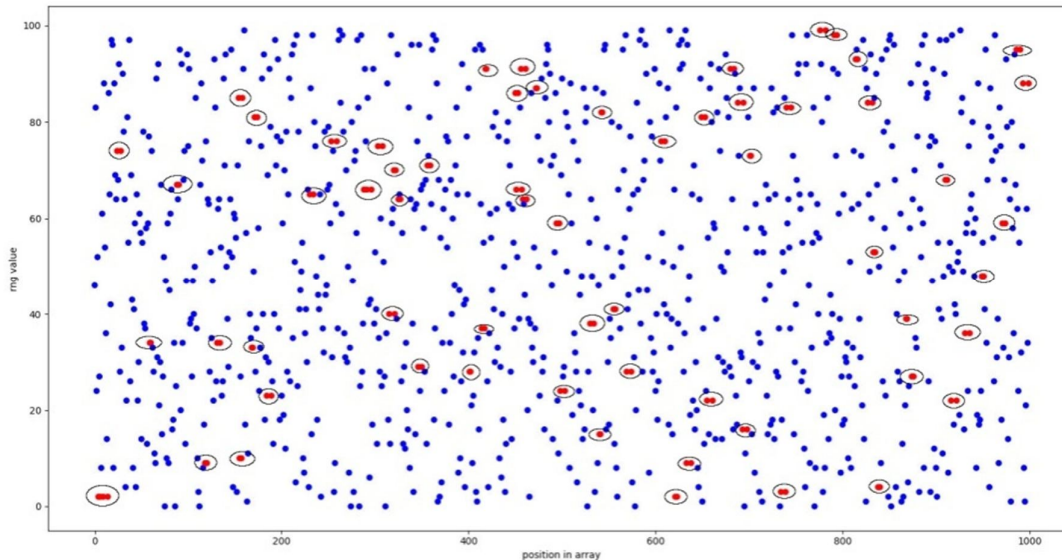
приложи подходът и да се анализира времевата линия от случайни стойности от среден клас компютърна система, с каквато най-често всеки разполага.

3.2.4. Проучване на генератори на случайни числа с времеви редове

Времевите редове като стохастичен процес могат да бъдат използвани за анализ и на RNG/PRNG. За тази цел е разработен алгоритъм за откриване на повтарящи се модели (patterns) от данни в генерираните от RNG времеви редове. За да се съберат данни от генераторите на случайни числа във времеви ред за нуждите на изследването се използва специално написана програма. Посредством отново написана за целта програма, събраните данни със случайни числа се представят графично, което помага за по-лесно забелязване на важните елементи от времевите редове (фиг.3.2 и фиг.3.3). На пръв поглед с резултатите от данните от System.Random на фиг. 3.2 всичко е наред и е възможно да се мисли, че имат добро качество на ентропия. Но нека предложи още един начин в друг графичен изглед, за да се уверим в преценката си.



Фиг. 3.2. Визуализиране на данните получени от System.Random, като линия на шум.



Фиг. 3.3. Представяне на данните от System.Random във визуализация тип поле с точки

Представянето на едни и същи данни с различна графична интерпретация, може да помогне за разкриването на някои проблеми с качеството на изследваните стойности. На фиг.3.2 и фиг.3.3 са изобразени графики на едни и същи данни. Като при фиг.3.2. качеството на случайните числа според графиката може да се приема за добро. Но след преглед на графиката от фиг.3.3 става ясно, че има чести случаи на повторения, които програмата прихваща, като модел на повторяемост. Съответните стойности се разпознават и оцветяват с червено и са оградени в кръг за по-добра видимост.

Визуализацията на фиг. 3.3 показва слабостите на обработените резултати. Моделите на повторното появяване се срещат периодично във времето. Тези случаи са оцветени в червено от разработената програма, като се използват предварително определени модели за прогнозиране, както беше споменато по-рано. В конкретния случай на разглежданите времеви редове, тези случаи са 65 случая от масив с 1000 стойности. Може да се каже, че 6,5% прогнозни числа от генерирания RNG масив са значителен резултат. Обикновено последователност от псевдослучайни числа се иницира от SEED вътре в PRNG (Koeune, 2005). Ако такъв генератор на произволни числа се използва в криптографията, произведените от него SEED стойности могат да бъдат атакувани успешно. Чрез прогнозиране на следваща стойност SEED или чрез наблюдение на предадени криптирани данни, стойностите в основата на системата за криптиране могат да бъдат възприети в определен момент.

3.3 Пренебрегнати рискове от киберсигурност в доставчиците на услуги за публичен Интернет хостинг

До тук изследването за анализи на качеството на RNG и областите, в които се срещат проблеми, успя да засегне криптографските алгоритми, езици за програмиране и операционни системи. Сега фокусът се измества върху масово предлагани публични хостинг услуги. В тази дисертация за настоящото изследване е използван уеб хостинг доставчик, който е един от популярните в бранша. Услугата за уеб приложения е инсталирана на масово предлаган нает споделен хостинг. Добавен е уеб сертификат и SSL достъпът е активиран, като всичко работи на стандартните портове за комуникация. На първа линия на защита между клиент и сървър излизат криптографските шифри, поддържани от хостинг сървъра. Ако те са актуални и между тях не се срещат уязвими и вече остарели и компрометирани във времето такива, може да се счита, че протоколът за комуникация е достатъчно добре подсигурен.

Извършен бе тест, като са сканирани криптографските протоколи, обезпечаващи връзката между клиент и сървър (хостинг услугата). Установено е, че от списъка с криптографски протоколи, които сървърът предлага участват: TLSv1.0 и TLSv1.1, които изобщо не трябва да се поддържат и предлагат, тъй като имат отдавна установени слабости и не трябва да се използват. Друг протокол, който сървърът поддържа е TLSv1.2, който е все още актуален и е одобрен за използване, но не в пълния му вид. В него се съдържат криптографски шифри, които трябва да бъдат извадени, но сървърът ги предлага за комуникация, което също е съществена уязвимост в сигурността на предоставяната услуга. Анализът на протоколите и шифрите също така установи още един съществен недостатък. Протоколът TLSv1.3 не се поддържа изобщо, това за момента е най-актуалния и сигурен протокол от семейството на TLS за тунелна свързаност.

След проверката на криптографските протоколи и шифри, поддържани за комуникация, изследването се премести върху по-чувствителната тема – генераторите на случайни числа. Понеже дори и при най-актуалните протоколи и шифри за защита, ако случайните числа не са достатъчно качествени и случайни, рискът да падне цялата криптираща защита е много голям. За да се извърши този анализ е създадена компютърна програма, която установява свързаност до сървъра по наличните криптографски протоколи за защита между клиент сървър. В конкретния случай е използван TLSv1.2 и във фазата на установяване на свързаност, програмата взима

генерираните случайни числа от сървъра и ги записва във файл, като времеви ред. Въпросната програма се изпълнява в цикъл, докато събере достатъчно количество данни за анализ.

Събраните данни от случайни числа са подложени на анализ чрез специализираният софтуер с отворен код за анализ на случайни числа използвани в криптографията Dieharder на Робърт Г. Браун (Brown, 2021). Изпълнени са симулация на 114 теста, както и проверка на качеството на числата и по стандарта за киберсигурност на генератори на случайни числа FIPS-140. Обобщено данните от теста за симулация на случайни числа са:

- Само 25 теста са преминали успешно;
- Неуспешни, които имат компрометирана /предсказуема/ стойност и следователно откриваема криптография са 76;
- Уязвими, където криптографията може да бъде разкрита с относително добър компютърен хардуер са 13;

От представените резултати може да се направи заключение, че заради слабостите в случайните числа и при установеното нарушаване на криптографската защита, рискът за успех при кибератаки за компрометиране на криптографията е критично висок. Причините за това може да са разнородни, но най-често срещаната от тях е, че хостинг доставчиците често поемат повече клиенти с техните приложения, от колкото капацитета на киберзащитата на сървърите им може да поеме. Работата на много приложения и клиенти едновременно, непрестанно източват криптографията и случайните числа на сървъра.

Решенията от страна на клиента, които се допускат в случая, е да се използва частен хостинг върху собствена инфраструктура, където няма да се допусне прекомерното натоварване от описания вид. При невъзможност да се осигури обаче непрекъсваемост за хардуерна конфигурация и подходящо място, като сървърно помещение. По-добре е да се наеме VPS сървър, който ще е само под контрола на един клиент и също проблемът ще се избегне. От страна на хостинг доставчика обаче, също може да се предприемат действия за повишаване капацитета на киберзащитата. Следва да се приложат техниките за конфигуриране на правилното функциониране и повишаване капацитета на ентропия в Linux, описани в раздел 3.2.2 „Разбиране на RNG Entropy в Linux“ на тази дисертация.

След правилната настройка на системата, може да се прибегне до друг нетрадиционен подход, познавайки принципа на работа на събиране на ентропия в

буферите си от операционната система Linux. Може да се напише програма, която да генерира редици от събития, които няма да затормозят особено системата, но ще създадат множество процеси подпомагащи събирането на ентропия:

```
#!/bin/sh

## list of sites using round-robin DNS
ROUND_ROBINS="www.yahoo.com google.com twitter.com outlook.com"

## Entropy start and end value limits
STOP_LIMIT="3800"
START_LIMIT="3000"

until [ "$(cat /proc/sys/kernel/random/entropy_avail)" -gt
"$STOP_LIMIT" ]

    do while [ "$(cat /proc/sys/kernel/random/entropy_avail)" -lt
"$START_LIMIT" ]

        do for thing in "/tmp/loyeyoung" "/tmp/sueellen"
"/tmp/rootdev" "/tmp/files"

            do echo $thing =====
                touch /tmp/toss
                for robins in $ROUND_ROBINS
                    do nslookup "$robins" 8.8.8.8 > /tmp/toss
                        nslookup "$robins" 9.9.9.9 >> /tmp/toss
                        nslookup "$robins" 192.168.2.3 >> /tmp/toss
                        nslookup "$robins" >> /tmp/toss
                        cat /tmp/toss
                        mkdir $thing -p
                        cp /tmp/toss $thing/toss
                        cat $thing/toss
                        rm -f /tmp/toss
                        rm -f $thing/toss
                    done
                done
            done
        done
    done
```

Представения програмен скрипт е съвсем базов и би могъл да бъде надграден и съставян и на други програмни или скриптов езици. Въпреки семплия вид, успява да даде очакваните резултати и покрие нуждите на текущото изследване. Скоростта на натрупване на ентропия се подобри. Което допринася въпросната система да понася големи натоварвания върху генерирането RNG стойности. Начинът на действие е както е зададен в момента е, че изпълнението на допълнителните операции в памет, процесор, диск и мрежа, ще се активират при достигане на стойност в буфера за ентропия под 3000. Също така, би могло предоставеното решение да се използва в комбинация с хардуерни решения, подпомагащи криптографските алгоритми и ентропията на случайните числа, което и компанията Intel предлага при своите процесори.

Наименованието на модула за подпомагане генерирането на случайни числа е Intel Secure Key, предишното му кодово име е Bull Mountain Technology. С това името Intel определя в процесорите си разширението за архитектура Intel64 и IA-32 RDRAND и свързаната с него хардуерна реализация на Digital Random Number Generator (DRNG). Освен всичко друго, DRNG, използвайки инструкцията RDRAND може да е изключително полезен при генериране на висококачествени ключове за криптографски протоколи. Следователно трябва да се провери, дали текущата система разполага с такива процесори и би могло нейната конфигурация да бъде обновена. При наличие на компютърна система с Linux операционна система, проверката може да стане освен чрез техническата документация на чиповете от производителя и чрез следната комбинация от команди:

```
$ cat /proc/cpuinfo | grep -i rdrand | echo $?
0
```

Като резултат 0 означава, че е наличен флаг RDRAND и процесорът може да бъде включен за подобряване на криптографските функции на системата по следния начин:

```
# apt install rng-tools-debian
# /etc/init.d/rng-tools-debian start
# /etc/init.d/rng-tools-debian status
* rng-tools-debian.service - LSB: rng-tools (Debian variant)
   Loaded: loaded (/etc/init.d/rng-tools-debian; generated)
```

```

Active: active (running) since Fri 2020-11-28 17:30:54 EET; 3min
10s ago
   Docs: man:systemd-sysv-generator(8)
  Tasks: 4 (limit: 4915)
  Memory: 1.3M
  CGroup: /system.slice/rng-tools-debian.service
          └─3597 /usr/sbin/rngd -r /dev/hwrng
$ cat /proc/sys/kernel/random/entropy_avail
4096

```

Резултатите показват, че скоростта на събиране на ентропия за нашият случай надхвърля скоростта на нейното консумиране.

3.4 Резултати в реална технологична инфраструктура

Предложеният подход за подобряване на киберсигурността в криптографията и генераторите на случайни числа при натоварени сървърни системи с публични услуги е приложен в технологичната инфраструктурата на института ИИКТ-БАН. Използваната хардуерна конфигурация е от среден клас, като е съобразена със сложността на изпълняваната задача. Сървърът е оборудван с един шест ядрен процесор Xeon(R) E-2236 от второ поколение и версия 6, 32GB RAM и два твърди диска в конфигурация с RAID1. Оперативния сървър с публичните услуги функционират върху Linux и всичките услуги са изцяло и от софтуер с отворен код. Функционират върху виртуална машина, като физическата машина е само виртуален хост, което е еквивалентно със ситуацията с разглежданите масови услуги, които са в предмета на текущото изследване за киберустойчивостта на криптографската защита. Сървърните услуги, изпълнявани от виртуалната машина са:

- мейл сървър, към момента с 242 потребителски акаунта. Достъпен чрез SMTP, POP3, IMAP, като всички те са защитени с криптографски комуникационен протокол TLSv1.2 и TLSv1.3. Удостоверяват се със сървърен сертификат за установяване на TLS сесии с асиметричен алгоритъм от типа елиптична крива `secp384r1`. Свързването до услугата не може да се осъществи без криптиране на комуникацията;
- Уеб мейл, който позволява на всичките 242 потребителя да оперират с пощата си и през уеб браузър. Комуникацията е защитена чрез криптографския комуникационен протокол TLSv1.2 и TLSv1.3. Удостоверяват се със сървърен

сертификат за установяване на TLS сесии с асиметричен алгоритъм от типа елиптична крива $secp384r1$. Свързването до услугата не може да се осъществи без криптиране на комуникацията;

- Уеб портал на Институтът по информационни и комуникационни технологии към Българската академия на науките, което е основното уеб пространство на института. Съдържа информация за дейността, два научни журнала, както и структурна информация. Комуникацията е защитена чрез криптографски комуникационен протокол TLSv1.2 и TLSv1.3 и съвършен сертификат за установяване на сесии със асиметричен алгоритъм с елиптична крива $secp384r1$. Удостоверяват се със съвършен сертификат за установяване на TLS сесии с асиметричен алгоритъм от типа елиптична крива $secp384r1$. Не се позволява свързване до услугата по не криптиран канал;
- Услуга за отдалечена администрация SSH с най-високата степен на криптографска защита, предлагана от протокола към момента. Идентификацията на потребител по SSH е само чрез криптографски ключове, не се допускат пароли;
- Услуга за отдалечено управление на Уеб съдържанието FTP. Комуникацията е защитена чрез криптографски комуникационен протокол TLSv1.2 и TLSv1.3 и съвършен сертификат за установяване на сесии със асиметричен алгоритъм с елиптична крива $secp384r1$. Удостоверяват се със съвършен сертификат за установяване на TLS сесии с асиметричен алгоритъм от типа елиптична крива $secp384r1$. Не се позволява свързване до услугата по не криптиран канал;

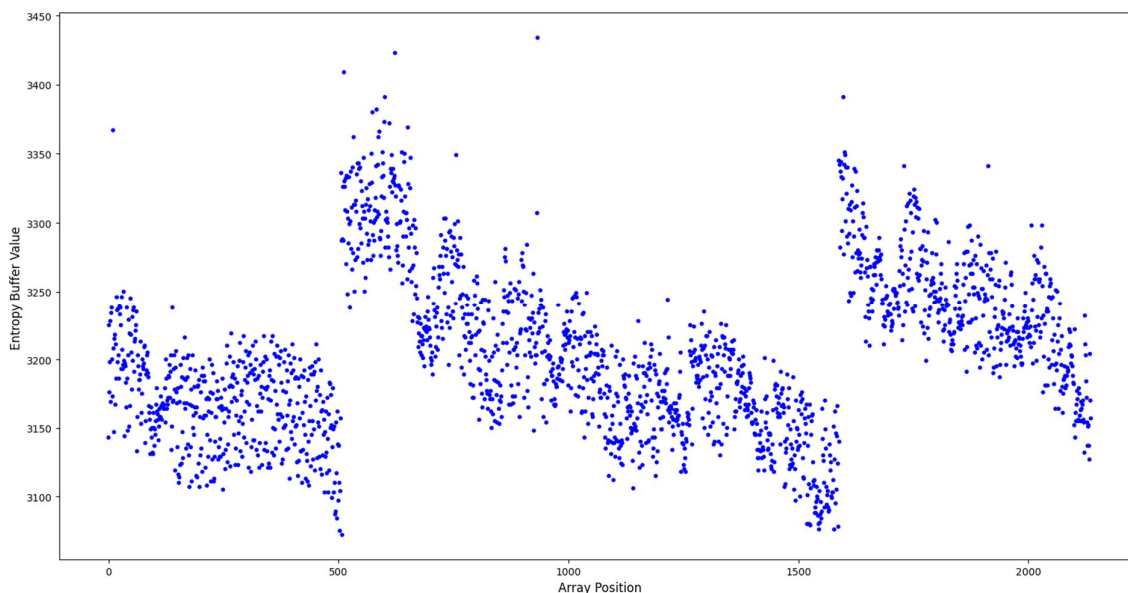
За всички услуги приоритетен протокол за криптирана свързаност е най-новият и сигурен протокол TLSv1.3, но ако се окаже, че клиента не го поддържа се минава на протокол TLSv1.2. Последният е оставен само за съвместимост, като от него са премахнати всички криптографски алгоритми в които са открити уязвимости.

Установено е, че към настоящия момент съвърхът се ползва значително интензивно от потребителите на ИИКТ и външни потребители на Интернет. Нивото на криптографска защита е на най-високото към момента според установените стандарти и не са правени никакви компромиси с криптографските протоколи или шифри. Като доказателство за качеството на киберзащитата с криптографски средства е приложен резултат от тест чрез скенер на SSL Labs за нивото на криптиране при предлаганите TLS протоколи върху наличните услуги. Резултатите от този тест са извлечени на базата на актуалните изискванията за криптографска защита към момента, които са

утвърдени от международните лаборатории по криптографска защита FIPS и NIST за САЩ, и Common criteria за Европа. От резултатите е видно, че протоколите и средствата за криптиране са от най-актуалните към сегашния момент. Оценката на всички тестове е най-високата възможна A+. Нивото на защита на HTTP протокола, който комуникира с браузъра, чрез TLS тунелът също е с максималното ниво на защита A+.

Извършени са тестове и на качеството на ентропия, чрез методите представени в дисертацията. От командния shell на сървъра са приложени два утвърдени метода, първият проверява качеството на ентропия по FIPS с `rngtest`, а вторият с инструмента за анализ `dieharder`. Всички тестове на ентропията на случайните числа издържат с най-високия възможен резултат според критериите на въпросните програми.

Въпреки добрите резултати е направено още едно изследване, според което става ясно, дали в моменти на висока потребителската активност и интензивното натоварване на криптографията ще доведат до изчерпване капацитета на случайните числа. За целта е съставен програмен скрипт, който на всеки 10 минути записва размерът на натрупаната в буфера ентропия. Тази статистика се събира в рамките на половин месец. За да стане ясно дали има моменти, които водят до изчерпване на буфера с ентропия по-бързо от колкото е неговото зареждане от системата. Като резултат след половин месец събиране на данни във времеви ред, се получиха 2137 стойности. Тези стойности са изведени в графичен вид и са изобразени на фиг.3.10:



Фиг.3.10 Ниво на ентропия в различни моменти от времето

От фиг.3.10 личи, че системата е имала пикове на по-интензивна дейност, които са водели до силно черпене от натрупаната в буфера ентропия. За това стойностите на моменти рязко падат. Предложеният в изследването подход върху системата обаче, успява да компенсира високата консумация на ентропия. Макар на графиката амплитудите между максимални и минимални стойности да изглеждат широки, то стойностите са в тесен диапазон, с ниво на ентропия между 3000 и малко под 3450. Липсата на стойности под 3000 показва, че системата се намира в много добро здраве и дори е способна да поеме по-големи натоварвания, защото стойностите са далеч от критичните. Като се вземат предвид всички тези резултати от реалната работна среда на сървъра, е на лице доказателството за ефективността на предлагания подход. Следователно може да се счита, че предлаганият подход може да бъде от полза и да подпомага различните Интернет системи и решения.

3.4 Изводи

Изследването на представените услуги и тяхното ниво на киберсигурност е от ключово значение за по-сигурен преход към съвременната дигитална трансформация. Бързото прехвърляне на всички социални и икономически дейности към дигитални платформи доказва, че съществуващата технологична инфраструктура може да отговори на днешните предизвикателства за дигитална трансформация. Ползите от това в икономическо и екологично отношение са неоспорими. Но по отношение на киберсигурността, много от настоящите ИТ услуги все още изостават. Увеличаването на степента на успех на киберпрестъпленията може да доведе до загуба на доверието в технологиите и възпрепятстването на тези процеси, което ще засегне и научно-техническия прогрес. Също така ще повлияе на забавянето в развитието и на много други свързани области в икономиката, сигурността, технологиите.

Прилагането на математически и статистически анализи с времеви редове за решаване на проблеми в киберсигурността е ефективно. Предлаганите тук подходи, може да се комбинират и с други техники и методи за анализ на киберсигурността, за да са по-комплексни и ефективни. В тази дисертационна работа е разработен метод за изследване на качество на RNG и PRNG в информационна система чрез прилагане на времеви редове. Методът позволява да се повиши качеството на ентропия при използването на криптография, осигуряваща различни Интернет услуги. Разработен е алгоритъмът за откриване на повтарящи се модели (patterns) от данни в генерираните от RNG времеви редове. Проведено е изследване на криптографски тестове и качеството

на ентропия върху работещи в реални условия натоварени сървърни системи с публични Интернет услуги.

Глава 4. Софтуерни подходи при работа с големи масиви от данни и ограничени компютърни ресурси с език за програмиране R

4.1 Програмният език R

Програмният език R е продукт, разполагащ с мощни инструменти за статистически изчисления и анализи. R едновременно е програмен език и софтуерна среда (Borcard, 2011), (The R, 2017). Компилира се и работи на различни операционни системи, като UNIX платформи, Linux, Windows и MacOS. Езикът R предлага широко разнообразие от статистически техники като, например линейно и нелинейно моделиране, класически статистически тестове, анализ на времеви редове, класификация, групиране и др., както и графични техники и е изключително разширяем (Long, 2015).

Въпреки многото предимства на R, богатството на неговите статистически модели и инструменти за обработка на данни, както и мощните способности за визуализация, изникват проблеми при работа с големи обеми от данни. Ограниченията на R произлизат от това, че той е проектиран да оперира в режим на изчисления само в единичен процес (на единично ядро на процесор) и при данните, заредени наведнъж в оперативната памет.

4.2 Преодоляване на проблеми на работа с големите данни чрез използване на микропроцесор с много ядра

Паралелното програмно изчисляване на повече от едно ядро на процесор е възможно чрез прекомпилиране и добавяне на някои програмни компоненти в R. Това е възможно, поради факта, че R е система с отворен код и това е едно от предимствата, което носи тази концепция.

4.3 Методи за оптимизиране на обеми от данни

Една от всеизвестните особености на езика R е, че зарежда всички данни с които оперира в RAM паметта на компютърната система, което при работа с големи данни би било критично, дори и на мощни системи с голям ресурс. Като начин за решаване на този проблем, в дисертацията се разглеждат начини за зареждане на данните в паметта,

като се изключват данните с некоректно съдържание още в момента на тяхното зареждане.

В някои статистически изследвания не е необходимо да се зареждат всичките данни, а само определени времеви рамки, за да се направи приблизителен статистически анализ в отрязък от време. В такъв случай, може да се приложат методи за позиционирано прочитане, предложени в дисертационния труд. Така става възможно да се обработи само определен отрязък от данните, разположени във файл с големите данни. Друг често възникващ проблем е, при зареждане на големи данни е, че след зареждане в паметта и обработка, някои от данните вече не са необходими, но продължават да заемат значителни обеми памет. В дисертацията е представен начин за редуциране на данните в паметта, като се премахват излишните от тях и се освобождава памет.

4.4 Изводи

Приносът на автора е, че чрез този материал със средствата на език за програмиране R се подпомага решаването на проблеми при работа с големи масиви от данни при ограничени компютърни ресурси. В тази дисертация са разработени софтуерни техники за оптимизиране на компютърната памет при работата с големи данни.

В заключение може да се каже, че с представените до тук примери не може да се изчерпа темата за оптимизираното зареждане на данни при работа с език за програмиране R. Работа с реални данни винаги е предизвикателство (Baumer, 2017). Но представените до тук техники са между добрите практики и са често използвани, те биха могли да се комбинират и с други подходи за решаването на проблеми в тази област.

Заклучение - резюме на получените резултати

В дисертационния труд подробно са изследвани методи и средства за използване на времеви редове при решаване на различни задачи, възникващи в съвременните приложения на информационни технологии и системи.

Предложен е метод озаглавен MA Volatility Indicator за подобряване прецизността в осцилатор (Моментум). MA Volatility Indicator работи в комбинацията от два инструмента EMA или SMA и предлага нова методика за интерпретиране на резултатите, което допринася за откриване на нива за свръх покупка и свръх продажба

при пазарната тенденция. Всички използвани в изследването инструменти ЕМА, SMA и Моментум, както и MA Volatility Indicator използват времеви редове.

Разгледана е приложимостта на апарата на невронните мрежи за прогнозиране на времеви редове във финансовата област. Показано е, че с нов модел на представяне на входните данни, характерни за финансови показатели, се получава по-висока степен на самоадаптация при обучение на невронната мрежа. Проведените експерименти потвърждават сложността на финансовите процеси и наличието на високочестотен шум в данните.

Разработен е метод за изследване на качество на RNG и PRNG в информационна система чрез прилагане на времеви редове за да се повиши ентропия на при използването на криптография, осигуряваща различни Интернет услуги. По този начин се допринася за по-добрата киберсигурност на ИТ инфраструктура за цифрови ресурси и защитата на данни. В дисертационната работа темата за криптографията получи специално внимание, поради нейното критичното значение. При пропускането само на един риск в киберсигурността е възможно да бъдат компрометирани всички ИТ услуги.

Практическите резултати от реалния експеримент показаха, че е намерено златното съотношение между масови услуги и действителните изисквания за киберсигурност.

С оглед на работата, извършена в този дисертационен труд и резултатите, получени в хода на изследванията и изложени по-горе, могат да бъдат формулирани следните научно-приложни резултати:

1. Разработен е метод, озаглавен MA Volatility Indicator, за комбиниране на индикатори за откриване на ценови движения с нови подходи при използване на времеви редове от финансовите данни.

2. Приложен е апаратът на изкуствени невронни мрежи с цел изследване на финансови времеви редове. Разработен е алгоритъм за обучение на невронната мрежа чрез увеличаване на размера на входа на невронна мрежа и създаване на хибридна структура, като е предложен модел за самонадграждащи се трислойни MLP.

3. Разработен е метод за повишаване на криптографската защита в информационните системи на базата на изследвания на качеството на генераторите на произволни числа.

4. Проведени са експериментални изследвания за решаване на проблемите с киберсигурността в публични широко разпространени хостинг услуги. Получените

результати потвърждават валидността на предложения метод за повишаване на киберсигурността.

5. Разработени са програмни методи за ефективна работа с големи данни със средства на езика R.

6. Разработените методи за повишаване на криптографска защита са имплементирани в технологичната инфраструктурата на ИИКТ-БАН. Проведено е изследване на криптографски тестове и качеството на ентропия върху работещи в реални условия натоварени сървърни системи с публични Интернет услуги.

Насоки за бъдещи изследвания

Насоките за бъдещи изследвания по тематиката на дисертацията включват:

- Имплементиране на методът MA Volatility Indicator и прилагането му в комбинация и с други методи за анализ и прогнозиране на пазарни ценови тенденции;
- Прилагане на методът MA Volatility Indicator към автоматизирани системи за анализ на пазарни тенденции и извличане на сигнали за взимане на решения;
- Провеждане на още изследвания в областта на обучаващи алгоритми и системи с невронни мрежи за анализ и прогнозиране на времеви редове;
- Развиването на нови методи за увеличаване на криптографската защита в информационните системи;
- Изследване на комбинация на разработен метод с други методи и системи за анализ на RNG в криптография и други технологични области, което да съдейства за създаване и усъвършенстване на RNG, както и за по-точно определяне на спектъра от задачи, които генераторът може да изпълнява добре;
- Намиране на още подходи за зареждане и филтриране на големите данни с цел по-ефективната им обработка.

Публикации по темата на дисертационния труд

- 1 **Иван Благоев**, Николай Докев, Комбиниране на Моментум с един метод за прогнозиране на пазарни ценови движения за по-точни резултати (Combination of Momentum with One Method for Forecasting of Market Trends to Improve the Results), Международна научна конференция “УНИТЕХ’17” – Габрово, 2017 Selected papers, ISSN 2603-378X, pp. II-265-II-270
- 2 **Ivan Blagoev**, Методи за оптимизирано използване на компютърна памет при зареждане на данни със средствата на език за програмиране R (Methods for

- Optimized Use of Computer Memory during Data Loads with R Programming Languages), International Conference “Automatics and Informatics’2017”, 4-6 October 2017, Sofia, Bulgaria, ISSN:1313-1850, pp.213-215.
- 3 **Blagoev I.**, Improving the Momentum Oscillator Accuracy by a Method for Forecasting of Market Price Movements, Сборник доклади от международна конференция, НВУ "Васил Левски", 14-15 юни 2018, Том 9, стр. 177-185. (ceeol.com)
 - 4 **Blagoev I.**, Method for more reliable users’ authentication in internet, Сборник доклади от международна конференция, НВУ "Васил Левски", 14-15 юни 2018, Том 9, стр. 167-176. (ceeol.com)
 - 5 **Blagoev, I.**, Using R Programming Language for Processing of Large Data Sets, Proc. Int. Conf. Big Data, Knowledge and Control Systems Engineering – BdkCSE’2018, 21-22 November 2018, Sofia, Bulgaria ISSN 2367-6450, pp. 91-98.
 - 6 **Ivan Blagoev**, Application of Time Series Techniques for Random Number Generator Analysis, Proceedings of XXII Int. Conference DCCN 2019, September 23-27, 2019, Moscow, Russia, pp.437-446. ISBN 978-5-209-09683-2, 2019 (РИИЦ).
 - 7 **Blagoev I.**, Neglected Cybersecurity Risks in the Public Internet Hosting Service Providers. *Information&Security International Journal* - ISIJ, 47, no. 1, pp. 62-76 (2020)
 - 8 Balabanov T.D., **Blagoev I.I.**, Dineva K.I. (2018) Self Rising Tri Layers MLP for Time Series Forecasting. In: Vishnevskiy V., Kozyrev D. (eds) Distributed Computer and Communication Networks. DCCN 2018. Communications in Computer and Information Science, vol 919. Springer, Cham. https://doi.org/10.1007/978-3-319-99447-5_50, pp. 577-584, **SJR:0.188**
 - 9 **Blagoev I.** (2020) Method for Evaluating the Vulnerability of Random Number Generators for Cryptographic Protection in Information Systems. In: Dimov I., Fidanova S. (eds) Advances in High Performance Computing. HPC 2019. Studies in Computational Intelligence, vol 902. Springer, Cham. https://doi.org/10.1007/978-3-030-55347-0_33 **SJR:0.215**

Забелязани цитирания

- I **Blagoev, I.**, 2018. Using R Programming Language for Processing of Large Data Sets, Proc. Int. Conf. Big Data, Knowledge and Control Systems Engineering – BdkCSE’2018, pp. 91-98.

Цитира се в:

- 1 Dineva, K., Atanasova, T.: Regression Analysis on Data Received from Modular IoT System. ESM'2019, EUROSIS-ETI, ISBN: 978-9492859-09-9, EAN: 9789492859099, pp.114-118, 2019
 - 2 Ivaylo Blagoev, G. Vassileva and V. Monov, "Methodology for content preparation of online courses," 2020 International Conference Automatics and Informatics (ICAI), Varna, Bulgaria, 2020, pp. 1-4, doi: 10.1109/ICAI50593.2020.9311364.
- II **Blagoev I.**, Neglected Cybersecurity Risks in the Public Internet Hosting Service Providers. *Information&Security International Journal* - ISIJ, 47, no. 1, pp. 62-76 (2020)
- Цитира се в:
- 3 M Terzieva, D Karastoyanov, ICT for Innovation in Advanced Banking, PROBLEMS OF ENGINEERING CYBERNETICS AND ROBOTICS • 2020 • Vol. 73, pp. 47-54 p-ISSN: 2738-7356; e-ISSN: 2738-7364, doi: 10.7546/PECR.73.20.05
- III **Blagoev I.**, Method for more reliable users' authentication in internet, Сборник доклади от международна конференция, НБУ "Васил Левски", 14-15 юни 2018, Том 9, стр. 167-176.

Цитира се в:

- 4 Ivaylo Blagoev, Gergana Vassileva and Vladimir Monov, "Methodology for content preparation of online courses," 2020 International Conference Automatics and Informatics (ICAI), IEEE, Varna, Bulgaria, 2020, pp. 1-4, doi: 10.1109/ICAI50593.2020.9311364.
- 5 Dineva, K., Atanasova, T.: Security in IoT Systems. Proceedings 19th International Multidisciplinary Scientific Geoconference SGEM 2019, 19, 2.1, ISBN:978-619-7408-79-9, ISSN:1314-2704, DOI:10.5593/sgem2019/2.1, 576-577. SJR (Scopus):0.232 Q4

Участие в проекти

- 1 Национална научна програма „Информационни и комуникационни технологии за единен цифров пазар в науката, образованието и сигурността“ (ИКТ в НОС) - 2018-2021.
- 2 Проект Зора по Заповед Но 147/14.06.2019 "Цифров и кибер устойчив ИИКТ"

Награди

1. Награда на ИИКТ-БАН за отлични научни постижения през 2019 г. в категория „Докторанти“.

Библиография

- 1 Atanasova, T., Barova, M.: Exploratory analysis of Time Series for hypothesized feature values. In: International Scientific Conference UniTech 2017, vol. II, pp. 399-403, University publishing house V. Aprilov, Gabrovo (2017)
- 2 Balabanov T.D., Blagoev I.I., Dineva K.I. Self Rising Tri Layers MLP for Time Series Forecasting. In: Vishnevskiy V., Kozyrev D. (eds) Distributed Computer and Communication Networks. DCCN 2018. Communications in Computer and Information Science, vol 919. Springer, Cham. https://doi.org/10.1007/978-3-319-99447-5_50 (2018)
- 3 Balabanov, T., Atanasova, T., Blagoev, I., Activation Function Permutation for Multilayer Perceptron Training, International Conference on Big Data, Knowledge and Control Systems Engineering BdKCSE'2018, Sofia, Bulgaria, ISSN 2367-6450, pp. 9-14 (2018)
- 4 Blagoev I., Dokev N.: A Method for Investigating the Alterations in the Price Trends of the Currency Markets and Forecasting of Probable Future Alterations, *Problems of Engineering Cybernetics and Robotics*, vol.65, pp.39-48 (2012)
- 5 Blagoev I., Neglected Cybersecurity Risks in the Public Internet Hosting Service Providers. *Information&Security International Journal - ISIJ*, 47, no. 1, pp. 62-76 <https://doi.org/10.11610/isij> (2020)
- 6 Blagoev I.: Method for Evaluating the Vulnerability of Random Number Generators for Cryptographic Protection in Information Systems. In: Dimov I., Fidanova S. (eds) Advances in High Performance Computing. HPC 2019. Studies in Computational Intelligence, vol 902. Springer, Cham. https://doi.org/10.1007/978-3-030-55347-0_33. (2021)
- 7 Blagoev, I., Using R Programming Language for Processing of Large Data Sets, Proc. Int. Conf. Big Data, Knowledge and Control Systems Engineering – BdKCSE'2018, 21-22 November 2018, Sofia, Bulgaria ISSN 2367-6450, pp. 91-98.
- 8 Camara C., Martín H., Peris-Lopez P., Aldalaien M., Design and Analysis of a True Random Number Generator Based on GSR Signals for Body Sensor Networks, *Sensors* 19, 2033; doi:10.3390/s19092033 (2019)
- 9 Plummer T., Forecasting Financial Markets: The Psychology of Successful Investing, January (2010)
- 10 Pseudo-Random Number Generators, <https://crypto.stanford.edu/psc/notes/crypto/prng.html>
- 11 Zhao P., R with Parallel Computing from User Perspectives, <https://www.r-bloggers.com/r-with-parallel-computing-from-user-perspectives/> (2016)
- 12 Brown R. G.: Dieharder: A Random Number Test Suite, <https://webhome.phy.duke.edu/~rgb/General/dieharder.php> (2021)

- 13 Ciampi F., G. Marzi, S. Demi, M. Faraoni, The big data-business strategy interconnection: a grand challenge for knowledge management. A review and future perspectives, *Journal of Knowledge Management*, Vol. 24, Issue 5 (2020).
- 14 Koeune F. Pseudo-random number generator. In: van Tilborg H.C.A. (eds) *Encyclopedia of Cryptography and Security*. Springer, Boston, MA . https://doi.org/10.1007/0-387-23483-7_330 (2005)
- 15 Borcard, D., Gillet, F., Legendre, P. *Numerical Ecology with R*, Springer, pp. 9 – 30 (2011)
- 16 Baumer B. S., Kaplan D. T., Nicholas J., *Modern Data Science with R*, Horton Chapman & Hall/CRC, Boca Raton, (2017)
- 17 Long C. (Ед.) *Data Science & Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data*, John Wiley & Sons, Inc., (2015)
- 18 Martínez-Acosta L., Medrano-Barboza J.-P., López-Ramos Á., López J., López-Lambraño Á., SARIMA Approach to Generating Synthetic Monthly Rainfall in the Sinú River Watershed in Colombia, *Atmosphere*, 11, 602; doi:10.3390/atmos11060602 (2020)
- 19 Mikalef P., Krogstie J., Examining the interplay between big data analytics and contextual factors in driving process innovation capabilities, *European Journal of Information Systems*, Volume 29, - Issue 3: Business Process Management and Digital Innovation <https://doi.org/10.1080/0960085X.2020.1740618> (2020)
- 20 Scott G., Carr M., Cremonie M., *Technical Analysis: Modern Perspectives*, e CFA Institute Research Foundation (2016)
- 21 The R Journal, ISSN: 2073-4859, <https://journal.r-project.org/> (2017)
- 22 Tomov, P., Monov, V., Artificial Neural Networks and Differential Evolution Used for Time Series Forecasting in Distributed Environment, Proc. of Int.conference Automatics and Informatics, ISSN 1313-1850, pp.129-132, Sofia, Bulgaria, (2016)
- 23 Wafi A.S., Hassan H., Mabrouk A., Fundamental Analysis Models in Financial Markets – Review Study, *Procedia Economics and Finance*, Vol. 30, 939 – 947. Elsevier (2015)
- 24 Wang, W., Y. Wang, Analytics in the era of big data: The digital transformations and value creation in industrial marketing, *Industrial Marketing Management*, Vol. 86, pp. 12-15, ISSN 0019-8501, <https://doi.org/10.1016/j.indmarman.2020.01.005> (2020)
- 25 Li C., Zhang J., Sang L., Gong L., Wang L., Wang A., Wang Y., Deep Learning-Based Security Verification for a Random Number Generator Using White Chaos, *Entropy*, 22, 1134; doi:10.3390/e22101134 (2020)

Abstracts of Dissertations

Number 1, 2021

INSTITUTE OF INFORMATION AND COMMUNICATION TECHNOLOGIES
BULGARIAN ACADEMY OF SCIENCES

БЪЛГАРСКА АКАДЕМИЯ НА НАУКИТЕ

ИНСТИТУТ ПО ИНФОРМАЦИОННИ И КОМУНИКАЦИОННИ ТЕХНОЛОГИИ

Брой 1, 2021

Автореферати на дисертации