

# МОДЕЛИ ПРЕДИКТИВНОЙ ЗАЩИТЫ ИНФОРМАЦИИ АВТОМАТИЗИРОВАННОЙ СИСТЕМЫ УПРАВЛЕНИЯ ВОДОСНАБЖЕНИЕМ НА ОСНОВЕ ВРЕМЕННЫХ РЯДОВ С ИСПОЛЬЗОВАНИЕМ ТЕХНОЛОГИЙ МАШИННОГО ОБУЧЕНИЯ<sup>1</sup>

Применительно к задаче прогнозирования кибератак рассмотрены модели, основанные на методах предиктивного обслуживания, а также сформированы гипотезы о границах применимости метода предиктивной защиты информации. Для проверки гипотез проанализирован набор данных «Water\_rump» (сайт Kaggle), состоящий из количественных и качественных характеристик автоматизированной системы управления водоснабжением, записанных в течение шести месяцев. Проведена предобработка, очистка и группировка экспериментальных данных для обучения нейронной сети, а также проведен анализ основных свойств данных и поиск в них общих закономерностей, распределений и аномалий. Метод предиктивной защиты информации реализован с применением технологий машинного обучения. Для каждой модели выполнена настройка гиперпараметров и проведена оценка по метрикам качества «precision», «recall» и «accuracy». Полученные результаты позволяют сделать вывод о применимости реализованного метода предиктивной защиты информации на практике для анализа данных автоматизированных систем управления технологическими процессами.

**Ключевые слова:** автоматизированная система управления технологическим процессом (АСУ ТП), временной ряд, задача прогнозирования, информационная безопасность, кибератака, машинное обучение, предиктивная защита информации.

<sup>1</sup> Исследование выполнено при финансовой поддержке Минобрнауки России (грант ИБ МТУСИ) в рамках научно-го проекта № 40469-29/2021-К.

# PREDICTIVE INFORMATION PROTECTION MODELS OF AUTOMATED WATER MANAGEMENT SYSTEM BASED ON TIME SERIES USING MACHINE LEARNING TECHNOLOGIES

*Models based on predictive maintenance methods are considered in relation to the problem of predicting cyberattacks, and also formed hypotheses about the limits of applicability of the predictive information protection method. To test the hypotheses, the "Water\_pump" dataset (Kaggle website) was analyzed, consisting of quantitative and qualitative characteristics of an automated water supply management system recorded over six months. Preprocessing, cleaning and grouping of experimental data for training a neural network, and also analyzed the main properties of the data and searched for general patterns, distributions and anomalies in them. The preprocessing, cleaning and grouping of experimental data for training the neural network was carried out, as well as the analysis of the main properties of the data and the search for general patterns, distributions and anomalies in them. The method of predictive information protection is implemented using machine learning technologies. For each model, the hyperparameters were adjusted and the quality metrics "precision", "recall" and "accuracy" were assessed. The results obtained allow us to draw a conclusion about the applicability of the implemented predictive information protection method in practice for analyzing data from Industrial Control Systems.*

**Keywords:** Industrial Control Systems (IDS), time series, forecasting problem, Information Security, cyberattack, machine learning, predictive information protection.

Проблематике обнаружения и предотвращения кибератак на объекты автоматизированных систем управления технологическими процессами (АСУ ТП) посвящено большое количество публикаций. Однако описанные подходы, как правило, позволяют выявить аномальное поведение системы [1], когда злоумышленник уже проник в систему, либо совершает попытки несанкционированного доступа к ней. Однако более интересной и полезной с точки зрения практического применения задачей является задача обнаруживать кибератаку до того, как она началась, и прогнозировать время, через которое система даст сбой при её реализации. Подобный подход имеет много общего с новыми

технологиями предиктивного обслуживания, основной целью которых является обеспечение надежности критичных для деятельности предприятия производственных и технологических процессов [2]. В отличие от традиционных подходов, направленных на поддержание каждой единицы оборудования в безупречном состоянии, технологии предиктивного обслуживания не требуют неоправданно высоких затрат. Под предиктивной защитой информации при этом понимается деятельность, позволяющая по косвенным признакам (параметрам) системы определить возможность наступления кибератаки, спрогнозировать время, через которое она наступит, а также выбрать адекватные превентив-

ные меры защиты. В рамках проведенного исследования разработана модель предиктивной защиты объекта автоматизированной системы управления водоснабжением, направленной на предотвращение нарушения одного из свойств информационной безопасности объекта – доступности информации. Под доступностью при этом понимают свойство информации быть готовой к использованию по запросу авторизованного субъекта, имеющего на это право [6].

Для исследования поставленной задачи был взят датасет «Water\_pump» с сайта Kaggle. Данные представляют собой количественные значения, принимаемые с 51 датчика (сенсора). Размер датасета 220320x55. Столбец «timestamp» представляет собой временной

интервал, показывающий, что измерения проводились каждую минуту (рис. 1). Анализ показал, что датасет имеет достаточно большое количество признаков с пропусками данных (более 50%). Вследствие появления зашумленных значений при обучении нейронной сети исключены столбцы, у которых процент пропущенных значений составлял более 60. Чтобы отразить максимально реалистичную ситуацию по работе с датчиками, выпадающие значения восполнены медианными значениями последних 20 минут. Кроме того, вследствие разной размерности датчиков/сенсоров (например, sensor\_00 изменяет свои значения от 1 до 5, в то время как sensor\_01 – от 47 до 48) данные были нормированы с использованием функции `StandardScaler`.

timestamp	sensor_00	sensor_01	sensor_02	sensor_03	sensor_04	sensor_05	sensor_06	sensor_07	...	sensor_43	sensor_44	sensor_45
2018-04-01 00:00:00	2.465394	47.09201	53.2118	46.310760	634.3750	76.45975	13.41146	16.13136	...	41.92708	39.641200	65.68287
2018-04-01 00:01:00	2.465394	47.09201	53.2118	46.310760	634.3750	76.45975	13.41146	16.13136	...	41.92708	39.641200	65.68287
2018-04-01 00:02:00	2.444734	47.35243	53.2118	46.397570	638.8889	73.54598	13.32465	16.03733	...	41.66666	39.351852	65.39352
2018-04-01 00:03:00	2.460474	47.09201	53.1684	46.397568	628.1250	76.98898	13.31742	16.24711	...	40.88541	39.062500	64.81481
2018-04-01 00:04:00	2.445718	47.13541	53.2118	46.397568	636.4583	76.58897	13.35359	16.21094	...	41.40625	38.773150	65.10416
2018-04-01 00:05:00	2.453588	47.09201	53.1684	46.397568	637.6157	78.18568	13.41146	16.16753	...	42.70833	38.773150	63.65741
2018-04-01 00:06:00	2.455556	47.04861	53.1684	46.397568	633.3333	75.81614	13.43316	16.13136	...	43.22916	38.194440	61.92130
2018-04-01 00:07:00	2.449653	47.13541	53.1684	46.397568	630.6713	75.77331	13.25231	16.12413	...	42.96875	38.194443	59.60648
2018-04-01	2.463426	47.09201	53.1684	46.397568	631.9444	74.58916	13.28848	16.13136	...	42.18750	38.194440	57.87037

Рис. 1. Представление данных в датафрейме

Наличие в датасете информации с большого количества сенсоров приводит к существенному проценту зашумленности данных, что может сказаться на качестве распознавания будущей модели. Было решено проанализировать временное распределение каждого датчика и выявить взаимосвязь в их работе. На рис. 2 видно, что sensor\_01, sensor\_02, sensor\_03 имеют похожее распределение. Аналогичный вывод можно сделать и для датчиков sensor\_4 – sensor\_12.

Для уменьшения размерности данных датчики были сгруппированы по степени схожести их временного распределения по 7 блокам (табл. 1).

Категориальным признаком в датасете является состояние датчика в текущий момент времени:

- «NORMAL» (нормальная работа система): 205836 значений;
- «RECOVERING» (восстановление работы датчика): 14477 значений;
- «BROKEN» (датчик не работает): 7 значений.

Как показано на гистограмме распределения категориальной переменной (рис. 3) выборка достаточно несбалансированная, и преобладают записи с нормальным поведением системы. Однако приоритетом формируемой модели является умение предсказывать имен-

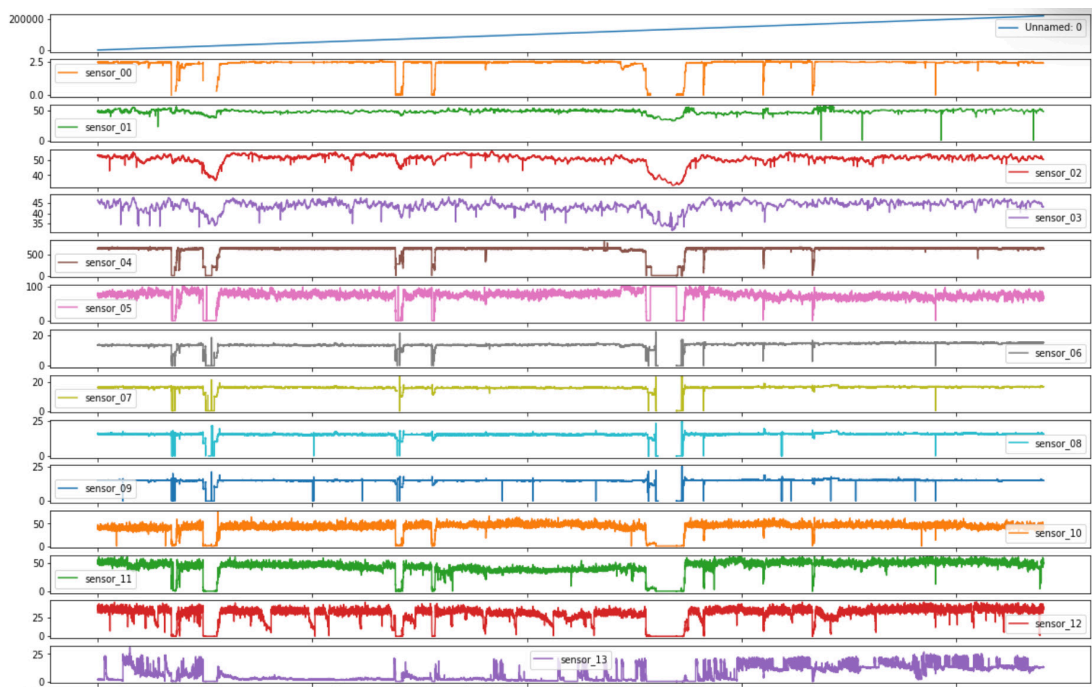


Рис. 2. Временное распределение различных сенсоров

Таблица 1

### Группировка сенсоров по блокам

№ блока	Номер сенсора
I блок	sensor_01, sensor_02, sensor_03
II блок	sensor_04, sensor_05, sensor_06, sensor_07, sensor_08, sensor_09
III блок	sensor_10, sensor_11, sensor_12
IV блок	sensor_14, sensor_16, sensor_17, sensor_18
V блок	sensor_19, sensor_20, sensor_21, sensor_22, sensor_23, sensor_24
VI блок	sensor_25, sensor_26, sensor_28, sensor_29, sensor_30, sensor_31, sensor_32, sensor_33
VII блок	sensor_34, sensor_35

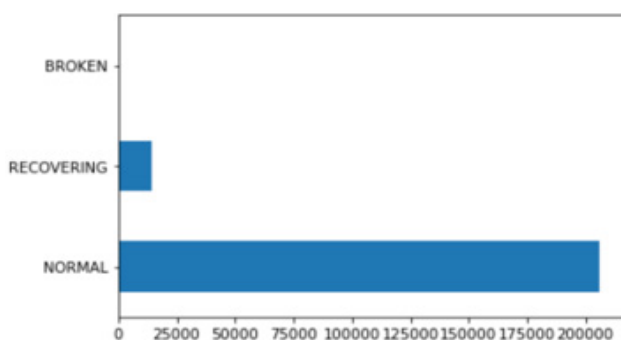


Рис. 3. Гистограмма распределения категориального признака

но выход датчика из строя (нарушение доступности информации). Данные были разбиты на обучающую и тестовую выборку. В обучающей выборке помимо нормальных состояний системы, состояний восстановления были использованы 2 записи состояний поломки датчика.

Так как модели машинного обучения не

умеют работать с категориальными признаками (за исключением модели Catboost от компании Яндекс), то было решено заменить состояния системы на числа: «1» - нормальная работа датчика, «0.5» - восстановление работы датчика, «0» - датчик не работает (рис. 4).

Стоит отметить, что такие модели машин-

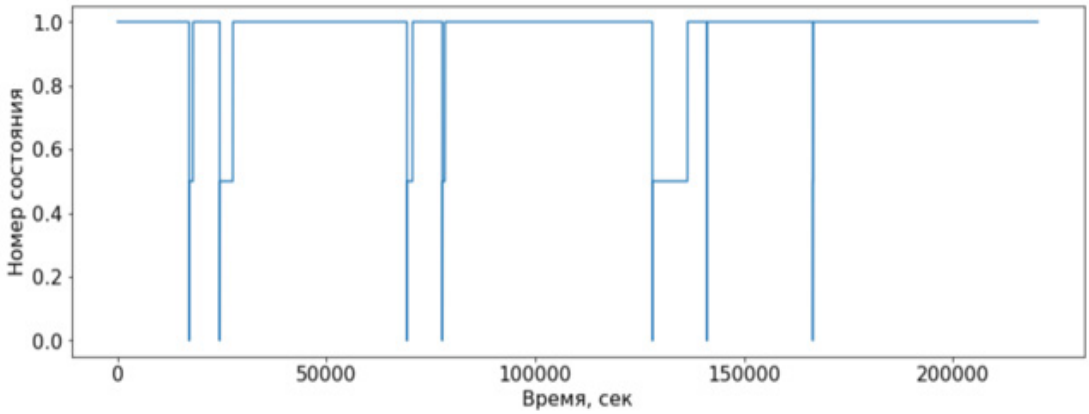


Рис. 4. Временное распределение сенсора с заменой категориальной переменной

ного обучения для работы с временными данными, как модель скользящего среднего (ARIMA) и разновидность данной модели, учитывающая сезонность (SARIMA), не смогли предсказать, через какое время датчик выйдет из строя. Поэтому была построена рекуррентная сеть с долгой краткосрочной памятью (LSTM) (рис. 5) [7]. Для борьбы с переобучением использовалась Lasso регуляризация – Dropout(0.3). Количество эпох для обучения составило 50. В качестве алгоритма оптимизации [9] был выбран Adaptive moment estimation (Adam). А в качестве функции ошибки – среднеквадратичная ошибка (RMSE) [10]:

$$RMSE = \sqrt{\frac{1}{n} \sum (y_i - \gamma_i)^2}, \quad (1)$$

где  $n$  – количество записей,  $y_i$  – предсказанный ответ моделью,  $\gamma_i$  – истинный ответ. На тестовом датасете среднеквадратичная

ошибка составила 0.016. На рис. 6 показан процесс обучения модели на тренировочном наборе данных, где после 30-й эпохи модель вышла на плато. Так как на тестовом наборе данных с увеличением эпохи ошибка не увеличивается, сделан вывод о том, что модель не была подвержена переобучению.

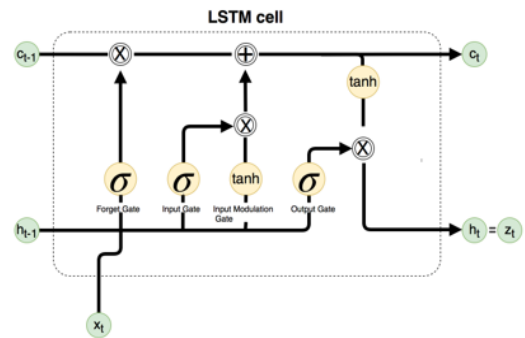


Рис. 5. Архитектура LSTM

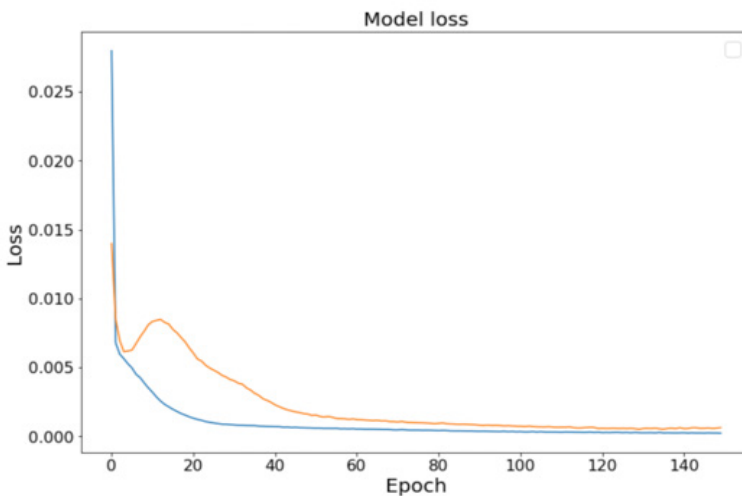


Рис. 6. График обучения модели на тренировочном и тестовом датасетах

На рис. 7 представлены графики значений, предсказанных разработанной моде-

лью, и истинных значений. При сравнении графиков видно, что модель предсказала все

5 отказов в системе и спрогнозировала время ремонта (замены) датчиков, близкое к реальным значениям. Кроме этого, на тестовых данных отсутствуют ошибки второго рода, что является преимуществом данной модели.

Стоит также отметить, что разработанная модель позволяет спрогнозировать состояние системы на достаточно длительный период времени.

Полученные результаты позволяют сде-

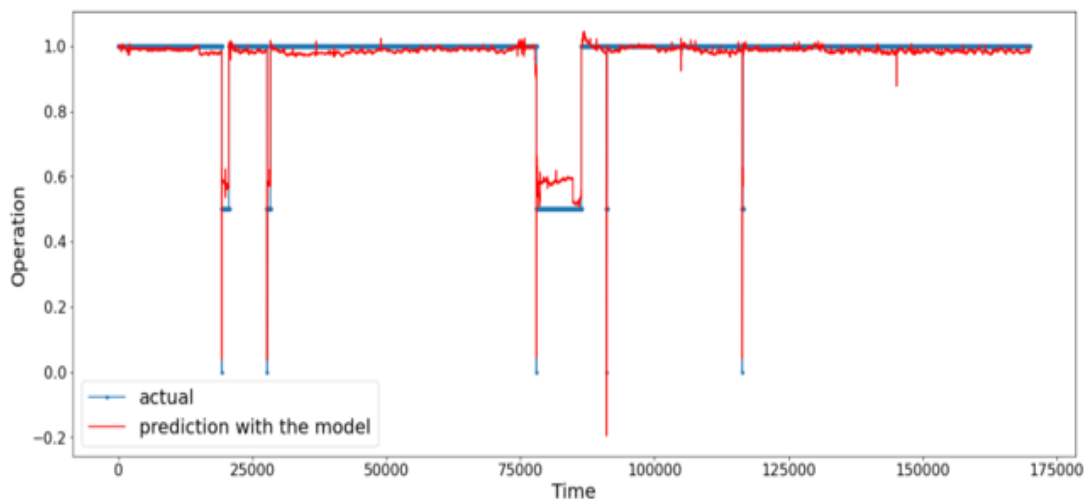


Рис. 7. Графики значений, предсказанных моделью (prediction with the model), и истинных значений (actual)

лать вывод о применимости реализованного метода предиктивной защиты информации на практике для анализа данных АСУ ТП. Раз-

работанные модели могут использоваться для предотвращения нарушений других свойств безопасности.

## Литература

1. Гарбук С.В., Правиков Д.И., Полянский А.В., Самарин И.В. Обеспечение информационной безопасности АСУ ТП с использованием метода предиктивной защиты // Вопросы кибербезопасности, 2019. №3(31). С. 30-36.
2. Боровков А.И. «Умные» цифровые двойники – основа новой парадигмы цифрового проектирования и моделирования глобально конкурентоспособной продукции нового поколения. Трамплин к успеху// Журнал АО «ОДК». 2018. № 13. С. 12-18.
3. Правиков Д.И. Об одном подходе к обеспечению информационной безопасности автоматизированных систем // Вопросы защиты информации. 2007. № 3. С. 17-19.
4. Гарбук С.В., Бурцев А.Г. Методические основы исследования уязвимостей компонентов АСУ ТП // Защита информации. Inside. 2012. № 3. С. 34-38.
5. Гарбук С.В. Перспективы применения интеллектуальных технологий для решения задач безопасности // Национальная безопасность / 2016. № 4. С. 451-457.
6. Pump sensor data (2021). Доступ к ресурсу по ссылке <https://www.kaggle.com/nphantawee/pump-sensor-data> (23 мая 2021).
7. Асеев Г.Д. Обнаружение вторжений на основе анализа аномального поведения локальной сети с использованием алгоритмов машинного обучения с учителем / Г.Д. Асеев, А.Н. Соколов // Вестник УрФО. Безопасность в информационной сфере. – 2020 № 1(35). – С.77-83
8. Luzhnov V.S Simulation of Protected Industrial Control Systems Based on Reference Security Model using Weighted Oriented Graphs / V.S. Luzhnov, A.N. Sokolov, A.E. Barinov // Proceedings - 2019 International Russian Automation Conference, RusAutoCon 2019. – 2019
9. Sokolov A.N. Applying Methods of Machine Learning in the Task of Intrusion Detection Based on the Analysis of Industrial Process State and ICS Networking / A.N. Sokolov, I.A. Pyatnitsky, S.K. Alabugin // FME Transactions. –2019. –Vol. 47 No. 4. – P.782-789
10. Соколов А.Н. Применение методов одноклассовой классификации для обнаружения вторжений / А.Н. Соколов, С.К. Алабугин, И.А. Пятницкий // Вестник УрФО. Безопасность в информационной сфере. – 2018. – Том - № 2(28). – С.43-48

## References

1. Garbuk S.V., Pravikov D.I., Polyanskiy A.V., Samarin I.V. Obespecheniye informatsionnoy bezopasnosti ASU TP s ispol'zovaniyem metoda prediktivnoy zashchity // Voprosy kiberbezopasnosti, 2019. No3(31). S. 30-36.
2. Borovkov A.I. «Umnyye» tsifrovyye dvoyniki – osnova novoy paradigmy tsifrovogo proyektirovaniya i modelirovaniya global'no konkurentosposobnoy produktsii novogo pokoleniya. Trampolin k uspekhу// Zhurnal AO «ODK». 2018. No 13. S. 12-18.
3. Pravikov D.I. Ob odnom podkhode k obespecheniyu informatsionnoy bezopasnosti avtomatizirovannykh sistem // Voprosy zashchity informatsii. 2007. No 3. S. 17-19.
4. Garbuk S.V., Burtsev A.G. Metodicheskiye osnovy issledovaniya uyazvimostey komponentov ASU TP // Zashchita informatsii. Inside. 2012. No 3. S. 34-38.
5. Garbuk S.V. Perspektivy primeneniya intellektual'nykh tekhnologiy dlya resheniya zadach bezopasnosti // Natsional'naya bezopasnost' / 2016. No 4. S. 451-457.
6. Pump sensor data (2021). Accessed at <https://www.kaggle.com/nphantawee/pump-sensor-data> (May 23, 2021).
7. Asyayev G.D. Obnaruzheniye vtorzheniy na osnove analiza anomal'nogo povedeniya lokal'noy seti s ispol'zovaniyem algoritmov mashinnogo obucheniya s uchitelem / G.D. Asyayev, A.N. Sokolov //Vestnik UrFO. Bezopasnost' v informatsionnoy sfere. – 2020 № 1(35). – С.77-83
8. Luzhnov V.S. Simulation of Protected Industrial Control Systems Based on Reference Security Model using Weighted Oriented Graphs / V.S. Luzhnov, A.N. Sokolov, A.E. Barinov //Proceedings - 2019 International Russian Automation Conference, RusAutoCon 2019.-2019
9. Sokolov A.N. Applying Methods of Machine Learning in the Task of Intrusion Detection Based on the Analysis of Industrial Process State and ICS Networking / A.N. Sokolov, I.A. Pyatnitskiy, S.K. Alabugin //FME Transactions. -2019.-Vol. 47 No. 4.- P.782-789
10. Sokolov A.N. Primeneniye metodov odnoklassovoy klassifikatsii dlya obnaruzheniya vtorzheniy / A.N. Sokolov, S.K. Alabugin, I.A. Pyatnitskiy //Vestnik UrFO. Bezopasnost' v informatsionnoy sfere. – 2018. – Tom - № 2(28). – С.43-48

---

**АСЯЕВ Григорий Дмитриевич**, аспирант кафедры защиты информации высшей школы электроники и компьютерных наук ФГАОУ ВО «Южно-Уральский государственный университет (национальный исследовательский университет)». Россия, 454080, г. Челябинск, проспект Ленина, д. 76. E-mail: [asiaevgd@susu.ru](mailto:asiaevgd@susu.ru).

**СОКОЛОВ Александр Николаевич**, кандидат технических наук, доцент, заведующий кафедрой защиты информации высшей школы электроники и компьютерных наук ФГАОУ ВО «Южно-Уральский государственный университет (национальный исследовательский университет)». Россия, 454080, г. Челябинск, проспект Ленина, д. 76. E-mail: [sokolovan@susu.ru](mailto:sokolovan@susu.ru).

**ASYAEV Grigorii Dmitrievich**, postgraduate student of the department of information security of the school of electrical engineering and computer science in FSAEI HE «South Ural State University (national research university)». 76, Lenin prospect, Chelyabinsk, Russia, 454080. E-mail: [asiaevgd@susu.ru](mailto:asiaevgd@susu.ru).

**SOKOLOV Alexander Nikolaevich**, Ph.D., Associate professor, Head of the department of information security of the school of electrical engineering and computer science in FSAEI HE «South Ural State University (national research university)». 76, Lenin prospect, Chelyabinsk, Russia, 454080. E-mail: [sokolovan@susu.ru](mailto:sokolovan@susu.ru).