# METHODS OF INTELLIGENT DATA ANALYSIS USING NEURAL NETWORKS IN DIAGNOSIS

## Volodymyr Lyfar[1], Olena Lyfar[1], Volodymyr Zynchenko[2]

[1]Volodymyr Dahl East Ukrainian National University, Kyiv, Ukraine, [2]Institute of Telecommunications and Global Information Space, Kyiv, Ukraine

*Abstract. The considered methods make it possible to develop the structure of diagnostic systems based on neural networks and implement decision support systems in classification diagnostic problems. The study uses general special methods of data mining and the principles of constructing an artificial intelligence system based on neural networks. The problems that arise when filling knowledge bases and training neural networks are highlighted. Methods for developing models of intelligent data processing for diagnostic purposes based on neural networks are proposed. The authors developed and verified an activation function for intermediate neural levels, which allows the use of weighting coefficients as probabilities of diagnostic processes and avoids the problem of local minima when using gradient descent methods. The authors identified special problems that may arise during the practical implementation of a decision support system and the development of knowledge bases. An original activation function for intermediate layers is proposed, obtained based on the modernization of the Gaussian error function. The experience of using the considered methods and models allows us to implement artificial intelligence diagnostic systems in various classification problems.*

Keywords: artificial intelligence, neural networks, diagnostics, data analysis, knowledge bases, machine learning

## METODY INTELIGENTNEJ ANALIZY DANYCH Z WYKORZYSTANIEM SIECI NEURONOWYCH W DIAGNOZIE

*Streszczenie. Rozważane metody pozwalają na opracowywanie struktury systemów diagnostycznych opartych na sieciach neuronowych oraz wdrażanie systemów wspomagania decyzji w klasyfikacji problemów diagnostycznych. W pracy zastosowano ogólnie specjalistyczne metody eksploracji danych oraz zasady budowy systemu sztucznej inteligencji opartego na sieciach neuronowych.. Zwrócono uwagę na problemy pojawiające się przy wypełnianiu baz wiedzy i szkoleniu sieci neuronowych. Zaproponowano metody opracowywania modeli inteligentnego przetwarzania danych do celów diagnostycznych w oparciu o sieci neuronowe. Autorzy opracowali i zweryfikowali funkcję aktywacji dla pośrednich poziomów neuronowych, która pozwala na wykorzystanie współczynników ważących jako prawdopodobieństw procesów diagnostycznych i pozwala uniknąć problemu minimów lokalnych przy stosowaniu metod gradientowego opadania. Autorzy zidentyfikowali szczególne problemy, które mogą pojawić się podczas praktycznego wdrażania systemu wspomagania decyzji i rozwoju baz wiedzy. Zaproponowano oryginalną funkcję aktywacji warstw pośrednich, otrzymaną w oparciu o modernizację funkcji błędu Gaussa. Doświadczenie w stosowaniu rozważanych metod i modeli pozwala na wdrażanie systemów diagnostycznych sztucznej inteligencji w różnych problemach klasyfikacyjnych.*

Słowa kluczowe: sztuczna inteligencja, sieci neuronowe, diagnostyka, analiza danych, bazy wiedzy, uczenie maszynowe

## Introduction

When modeling complex systems, the possibilities of using deterministic models and analytical mathematical description of their behavior are significantly limited. In this case, the way out is the application of fuzzy logic models with the help of sets of system behavior rules obtained taking into account practical experience.

The main problem of using such experience is the uncertainty of the mathematical apparatus of modeling and possible conclusions arising from this alogism. This work considers the possibility of applying methods of intelligent data analysis and artificial intelligence based on neural networks in the most common practice of diagnosing the state of complex systems. Such systems primarily include the human body [1-3].

The highest complexity of human structure and functioning does not allow creating a closed class of models based on the understanding of biophysical, chemical, psychological and other functional interactions both within the body and with external sources. That is why, since the time of Avicenna, medicine has been considered more of an art than a science. Attempts to differentiate various interconnected communication functions of the human organism in the case allow to locally describe and quantify some parameters of human states.

No fewer difficulties arise when trying to model or diagnose complex information systems, control processes, electro-mechanical systems, and other complex mechanisms that are widely used in technology. In this work, the authors tried to propose methods of constructing intelligent diagnostic systems based on examples of medical practice, comparing the experience of traditional Chinese medicine and modern medical science [4].

As you know, the fundamental difference between these two approaches in medicine lies in the ideology of the principles of functioning of the human body. If modern medicine considers a person as an extremely complex set of functional organs, their connections (including reverse ones), cause-and-effect processes of incoherence of internal interactions of organs and systems, then traditional Chinese medicine considers the human body as a single system of material and energy flows.

No single approach provides a complete and unambiguous description of the cause-and-effect relationships of the occurrence and development of the disease. It is as if the same processes and phenomena were described not only in different languages, but also in different concepts. In this regard, a comparative analysis of the clusters of symptoms in relation to the disease syndrome by two methods is impossible at the same time due to differences in concepts. However, the mathematical apparatus of artificial intelligence using neural networks makes it possible to apply universal methods of diagnostic processes, pushing aside differences in the semantic perception of medical experience [5].

## 1. Models, methods and research problems of diagnostic systems

Models, methods and tasks of researching diagnostic methods. Diagnostic tasks are mostly related to classification (sometimes clustering) tasks. This affects the architecture of the neural model, the choice of activation function for different neural layers, the depth and complexity of the input information representation, and the methods of interpreting the results of neural models. The architecture of such a neural network can be built on the basis of a combination of Rosenbath perceptrons (Fig. 1). The XOR problem for perceptrons is easily solved by using the signature of Cantor's algebra for three logical operations: AND, OR, NOT. The input layer of the set of symptoms S provides the output combinations of independent, non-antagonistic manifestations of syndromes as symptoms understandable to specialists [6-8].
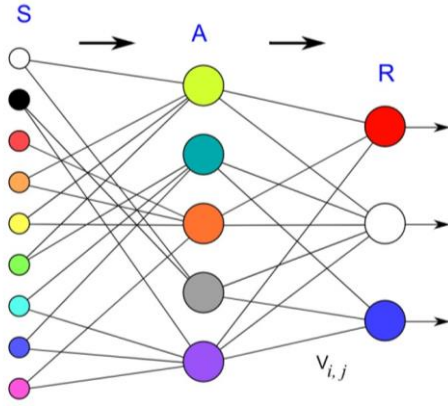
*Fig. 1. The structure of the Rosenbat Perceptron with an intermediate association layer*

The intermediate layers of neurons A provide the definition of the set of states of the key "entities" of the syndrome. The final syndromes are formed in the output layer R, and are a concatenation of "entities". In fact, this is a syndrome formula composed by conjunctions of individual entities, the set of which is reflected in the set of output neurons. The set of entities contains subsets of the key determinants of the general syndrome formula.

For example: in traditional Chinese medicine, the hybrid syndrome "Liver fever with rising YIN and lung wind" is a combination of entities: (fever ∧ liver) ∪ (rising YIN) ∪ (wind ∧ lungs). In this case, there are elements of the essence of "substance" – fever, wind; essence of organs – liver, lungs; essence of "energy" – yin, yang, etc. Depending on the complexity of the disease and the manifestation of symptom sets, the syndrome may be hybrid and combine a number of separate syndromes. In modern medicine, "organs" (appendix, glands...), "functional systems" (immunity, gastrointestinal tract...), condition determinants (fracture, inflammation, ischaemia...) and others can be used as intermediate layers. For example, the syndrome "acute appendicitis" is the intersection of the state "acute (inflammation)" and the organ "appendix" [9].

A set of k symptom vectors is used as input for medical diagnosis $S=\cup_{J=1}^{k}(\cap_{m=1}^{n}s_m)$; where is the $s_m \in S$ of the input layer. Such a set contains subsets of hierarchical attributes built in a tree structure using the xml structured markup language. The nested hierarchy structure allows you to establish the relationship of the corresponding symptom (tree leaf) to the localization of the symptom, its belonging to the type of symptom and weight characteristics that allow you to take into account the intensity and reliability of its manifestation.

- body – nested subsets of signs localizing the symptom to a certain part of the body (for example, "left hand") → "forearm" → "phalanx of the index finger");
- Stype – a set of symptom types (for example, pain, swelling, inflammation). Nested subsets of symptom type attributes (pain – acute, aching, throbbing).

The set of mappings of the elements of the input layer $s_i$ to the sets of the intermediate layer of neurons $a_k$ is weighted by the value $w_i$. The sum of the weighted input elements determines the state of the k-th neuron of the intermediate layer.

$$Ns_k = \sum_{i=0}^{n}\left(s_i \cdot w_i\right)_k \qquad (1)$$

At the same time, the activation function that determines the output value of the neuron at the output of the intermediate layer may differ from the activation function of the original layer.

$$Na_r = \sum_{j=0}^{m}\left(a_j \cdot v_j\right)_r \qquad (2)$$

The neuron's state is used as an argument to the activation function. Since the output layer completes the determination of the results of the classifiation task, the well-proven Leaky ReLU function can be used as the activation function FR($Na_r$). This piecewise linear function is differentiated over the range of probability of the final syndrome from zero to one and is not problematic for classification tasks. The activation function for intermediate layers FA($Ns_k$) should be different in priority.

The authors have studied a number of activation functions and can offer an original variant, a normalised range of weights [0;1]:

$$P(x) = \left(\frac{1}{\sqrt{2\pi}} \cdot \int_{0}^{6x-3} \exp\left(\frac{-t^2}{2}\right)dt\right) + 0.5 \qquad (3)$$

Fig. 2 shows the graph of the proposed activation function Fa~P(x) (x∈[0:1]) the normalising parameter for the intermediate layer of neurons).

This function can be obtained by normalising the Gaussian error function (erf), is smooth and differentiable over the entire range of values from 0 to 1. The presence of an integral and the complexity of the function are not an obstacle to its fast calculation with an accuracy of $10^{-3}$, which is much less than the possible minimum calculation error. The proposed function is smooth, monotonic, and differentiable over the entire range from 0 to 1. The argument of the function can be easily normalised for any non-negative values and can be related to the probability value of the weight of connections between neuronal layers.
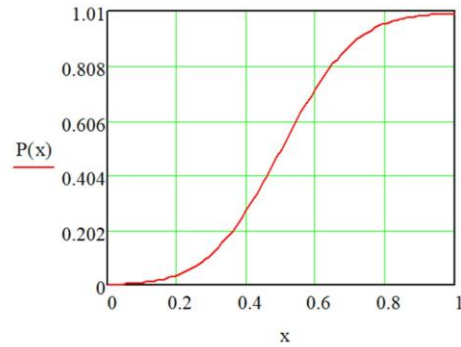


*Fig. 2. Activation function for intermediate layers of neurons*

## 2. Results and discussion

**Methods of working with the diagnostic system.** The proposed architecture of the neural network can be used in any classification tasks of diagnostics. Before using such a structure, it is necessary to train the network initially. With an effective interface, such training can be performed by a non-programmer and even a highly qualified diagnostic specialist. To create and fill in the knowledge base of syndromes, it is necessary to have a minimal initial set of many final syndromes, each of which is a reflection of the intersection of the entities of the syndrome formulas belonging to the collection of intermediate layers. The conjunctions of entity values fragmented from the formula of the corresponding syndrome, which have a default linkage weight of 1, are correlated as the knowledge base is filled. The values of the collection of intermediate entities are determined, in turn, by the conjunctions of symptom sets. The set of connections of symptom sets $s_i$ of entity $P_i$ and further with syndromes is weighted by relations, hereinafter represented as events $e_{ij}$.

$$
\begin{array}{cccc}
& s_1 & s_2 & \cdots & s_n \\
P_1 & \begin{pmatrix} e_{11} & e_{12} & \cdots & e_{1n} \\ P_2 & e_{21} & e_{22} & \cdots & e_{2n} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ P_m & e_{m1} & e_{m2} & \cdots & e_{mn} \end{pmatrix}
\end{array}
$$

$P_i$ – represents the probability of occurrence of event $e_{ij}$. The connections of elementary events are reflected in their probability weights $e_{ij} \rightarrow p_{ij}$; $p_{ij} \in P_i$. The neural network is trained directly by the expert during the creation and filling of the knowledge base. The neurons are tuned in the reverse order starting with the construction of sets of syndromes R, vectors of subsets of "entities" A, vectors of the set of symptoms S. In fact, the resulting knowledge base is a consistent surjective mapping of functions: (fA: S→A; F→R: A→R). During training, the well-known method of backpropagation is used. At the same time, the well-differentiated activation function of layer A allows the use of a fast gradient descent method, bypassing the problem of local minima. The correction of weights is determined by the formula [10]:

$$
w_i(m+1) = w_i(m) - u\alpha \tag{4}
$$

$u$ – significance coefficient of the algorithm; $\alpha$ – error rejection.

In the process of training a neural network, an expert generates a limited number of symptom vectors, identifying symptoms as antagonists. Each created vector is matched with the resulting vector of entities and corresponds to the correct syndrome. Since no categories of significance are assigned to individual symptoms when filling in data from different sources, the problem of unreliability may arise. We will consider this issue separately in the following. Training and correction of weights are carried out in parallel when creating each new link and each new syndrome. This can lead to the problem of dynamic changes in the knowledge base.

**Problems of constructing neural models in diagnostics and suggestions for overcoming them.**

The diagnostic methods discussed in this article rely on non-deterministic models and data on the objects of diagnosis, which leads to a number of problems. We will briefly discuss the main ones below.

**The problem of knowledge blurring.** The diagnostic system uses data obtained from different sources. This theoretically allows for contradictory diagnoses to be added to the knowledge base. Alogisms may arise that cannot be found by any means other than expert intervention. When creating knowledge bases, there may be problems with the reliability of the knowledge entered. Since any elementary symptoms are equilibrium and mutually independent, different symptom vectors for the same syndromes obtained from different sources (from different authors) may cause deviations from the true values of these syndromes (the so-called knowledge dilution). To effectively search for such situations, it is proposed to use a quantitative indicator defined as the cumulative share of the symptom vector in many diagnostic formulas. The final probability of the validity of the symptom vector is determined by logical OR and shows the frequency of repetition for syndromes of the same type. If the total proportion of the total symptom vector is less than, for example, 0.90, it can be assumed that this scenario may have a 10 per cent probability of error. However, it should be borne in mind that this is only an average and may not always correspond to the level of reliability of the syndrome in question.

It is possible that different end syndromes will correspond to the same symptom vector. This is the problem of exclusive OR. It is important that this illogicality cannot be eliminated by any means other than an expert's decision. Software tools that implement the knowledge bases under consideration should be equipped with XOR problem search functions and prevent the simultaneous use of contradictory statements.

**Dynamic changes in the knowledge base.** As the knowledge base is filled in, the weighting coefficients and mapping vectors change, changing the knowledge logic as well, which leads to the problem of dynamic data changes. In simple terms, the conclusions of the network can change significantly as the knowledge base is filled. Change beyond recognition. This can affect the perception of the reliability of the diagnosis. In addition, only relational databases can be used to build holistic data, as the chronology of the data warehouse does not allow for changes over time.

**The problem of weight correction.** The weighting coefficients of neural connections in a network are correlated during network training by the back-propagation method. The authors propose to set the initial weights of the coefficients to only 0.5. At the first step of machine learning, when establishing the relationship between the symptom vector and the final syndrome, a preliminary weight correction operation is performed. The weights $w_i$ can take on values in the range of ]0:1].

In this case, the initial correlation of the weight corresponds to the frequency of its manifestation in the resulting syndrome for different sources of information. For a syndrome obtained from a single source of information, the frequency of the relationship is one. The reliability of the information in this case cannot be determined. To eliminate such uncertainty in the formulation of such a syndrome, it is necessary to manually adjust the values of $u\alpha$.

**Reliability of neural network training. Antagonism of syndromes.**

If we consider a syndrome as a mapping of symptom vectors to a set of syndromes, considering such a mapping to be a statement, then there is no prohibition on logical contradiction between statements about the same syndrome by different authors. The problem of antagonism of statements arises. Since any statement, in addition to semantic perception, also carries a deep logical meaning and is essentially knowledge about the diagnostic model, there are no mathematical ways to detect such illogisms. In such a situation, different (possibly contradictory) syndrome formulas are assigned to the same symptom vector. This creates an OR or XOR problem. The task of the decision support system is to search for such situations and suggest corrections to the expert [11].

## 3. Conclusion

1. The paper considers the task of building diagnostic models based on neural networks. The proposed methods and models have undergone verification procedures and some of the problems described in the article have been identified. The proposed methods of building problem-solving networks can be easily implemented by effective means of mathematical modelling and software development of decision support systems in both medicine and technical diagnostic systems.

2. The authors have established that the problem of knowledge blurring that arises when creating knowledge bases by processing information from many independent sources can be partially solved and regulated by methods of verifying the reliability of classification statements. Along with the problems that exclude, or, through the use of alternative syndromes, the positive side of this approach is manifested in the form of objectification of the resulting classification task. At the same time, it is important to use the methods of antagonistic statements and the correction of such statements by ranked majoritarian methods. At the same time, it is impossible to eliminate the role of experts in the final decision-making on the reliability of the data entered into the knowledge base. In this regard, it is assumed that it is impossible to obtain new reliable knowledge after intensive mass training of the network.

The developers' attempts to obtain new knowledge (creation of new syndromes not previously mentioned in the knowledge bases) based on the processing of critical material from the knowledge bases were unsuccessful. However, the neural networks based on the proposed models performed the main task of classifying symptom states and identifying reliable syndromes perfectly.

3. The activation function for the intermediate layer proposed by the authors has proven to be positive and is recommended for use in diagnostic systems where quantitative indicators of the probability of realisation of the final syndromes of diagnostic tasks are important.

## References

[1] Balogh E. P. et al. (eds.): Improving Diagnosis in Health Care. National Academies Press (US), Washington 2015 [https://doi.org/10.17226/21794].

[2] Caliskan A., Yuksel M. E.: Classification of Coronary Artery Disease Data Sets by Using a Deep Neural Network. Euro Biotech J 1(4), 2017, 271–277.

[3] Checcucci E.: Applications of neural networks in urology: a systematic review. Current Opinion in Urology 30(6), 2020, 788–807.

[4] Glover E.: Artificial Intelligence [https://builtin.com/artificial-intelligence] (available: 28.01.2022).

[5] Kharlamova N. V. et al.: The use of artificial intelligence to diagnose diseases and predict their outcomes in newborns. Russian Bulletin of Perinatology and Pediatrics, 2023, 108–114.

[6] Lins A.J.C.C. et al.: Using Artificial Neural Networks to Select the Parameters for the Prognostic of Mild Cognitive Impairment and Dementia in Elderly Individuals. Computer methods and programs in biomedicine 152, 2017, 93–104 [https://doi.org/10.1016/j.cmpb.2017.09.013].

[7] Mantzaris D. et al.: Artificial Neural Networks for Estimation of Dementias Types. Artif Intell Appl 1(1), 2014, 74–82.

[8] Mirbabaie M., Stieglitz M.: Artificial intelligence in disease diagnostics: A critical review and classification on the current state of research guiding future direction. Health Technol. 11, 2021, 693–731.

[9] Rasmy L. et al.: Recurrent neural network models (CovRNN) for predicting outcomes of patients with COVID-19 on admission to hospital: model development and validation using electronic health record data. Lancet Digit Health 4(6), 2022, e415–e425 [https://doi.org/10.1016/S2589-7500(22)00049-8].

[10] Sanoob M. U. et al.: Artificial Neural Network for Diagnosis of Pancreatic Cancer. IJCI 5(2), 2016, 41–49.

[11] Sutton R. T., Pincock D., Baumgart D. C.: An overview of clinical decision support systems: benefits, risks, and strategies for success. NPJ Digit. Med. 3(17), 2020.

**Prof. Lyfar Volodymyr**
e-mail: lifar@snu.edu.ua

Lyfar Volodymyr graduated Physics Faculty of the university in 1984. He received his Ph.D. in Information Technology in 2007 and his Doctor of Engineering in 2017. He is a professor of the Department of Information Technologies of the East Ukrainian National University. His areas of interest include information technologies of artificial intelligence and decision support systems.



https://orcid.org/0000-0002-7860-9663

**M.Sc. Lyfar Olena**
e-mail: eklyfar@gmail.com

Lyfar Olena graduated Physics Faculty of the university in 1984. She is a senior lecturer of the Department of Information Technologies of the East Ukrainian National University. Her areas of interest include information technologies of artificial intelligence and decision support systems.



https://orcid.org/0000-0002-3014-5521

**M.Sc. Zinchenko Volodymyr**
e-mail: zinchenko@outlook.com

Volodymyr Zinchenko graduated from Donbass State Technical University in 2015. He received a full higher education in the specialty "Electronic Systems" and qualified as an engineer in the field of electronics and telecommunications. In 2021, he entered the postgraduate program at the Institute of Telecommunications and Global Information Space, specializing in Computer Science. His research interests include artificial intelligence information technologies and sensor networks.



https://orcid.org/0000-0001-6081-4848