

Mathematical Model of a Social Network User Profile Based on Interval Data Analysis

MYKOLA DYVAK, TYANDE PAN, OLEKSANDR KINDZERSKYI

West Ukrainian National University, Ternopil, Ukraine

Corresponding author: Mykola Dyvak (e-mail: mdy@wunu.edu.ua).

ABSTRACT It is proposed and substantiated to use a mathematical model for making decisions regarding the credibility of content posted on social networks, based on establishing a relationship between the outcome on which the decision about the credibility or unreliability of the content is made and the factors influencing it. Quantitative factors for assessing the user profile in the network have been justified: the number of posts, shares, or likes made by users within a short time after the content appears; the number of comments or reactions at certain time intervals; the number of participants interacting with the content within half a day after publication; the viral spread coefficient of the content. The resulting indicator of such a model is proposed to be the degree of credibility of specific content, ranging from 0 to 1. It is proposed and justified that methods of interval data analysis should be used to represent and analyze this indicator based on expert assessment of the content. An optimization problem is formulated for the two-stage identification of the model based on interval data analysis: forming the current structure based on candidate models (model structure synthesis), estimating its parameters and verifying the adequacy of the model. A hybrid method for identifying interval models of user profiles in a social network is proposed and substantiated. This method combines a metaheuristic algorithm for synthesizing the model structure based on the behavioral model of a bee colony with gradient methods for identifying the parameters of candidate models. Examples of applying the proposed method and mathematical model for making decisions about the credibility of content are presented.

KEYWORDS structural identification, parametric identification, artificial bee colony, social network user profile.

I. LITERATURE REVIEW ON THE GIVEN SUBJECT AREA

Social networks are one of the most common means of obtaining information [1], [2]. However, the information disseminated on social networks is not always reliable [3]. Modern social networks often host and spread fake or false information [4], [5]. Moreover, this is often done deliberately. Therefore, detecting false information in social networks is a relevant and important task. Existing methods can be classified into the following groups:

- Content analysis methods using artificial neural networks.
- Methods for analyzing information sources.
- Methods for verifying content credibility by comparing it with known facts in databases or with a priori verified content.
- Community-based verification, where users mark false information.
- Propagation network analysis methods.

- Community behavior analysis methods.
- In addition, combined methods are also possible.

The most widely used are content analysis methods employing artificial neural networks [6], [7], [8]. However, they have significant drawbacks, namely high sensitivity to the quality and size of the training dataset [9], [10] and high computational complexity [11], [12].

The disadvantages of source analysis methods include outdated data, lack of access to metadata, and the emergence of new sources or their dynamic changes, which reduce the effectiveness of these methods.

Methods that verify content credibility by comparing it with existing databases or a priori verified content also have limitations [13]. For instance, verified facts may not exist or may be difficult to find in databases, or there may be tight time constraints for retrieving such data.

Community-based verification [14], where users label content as false, can suffer from bias, lack of expertise within

the community, or susceptibility to manipulative statements. The disadvantages of propagation network analysis methods include complex structures and the possibility of manipulating participant communities.

In [15], the modeling of community participants' behavior in response to different types of content is examined. This approach can be used to identify content credibility [16], [17]. These methods are based on studying and modeling the typical behavior of a community when encountering certain information [2], [3]. Of course, these methods also have several drawbacks [18]. However, if the community's profile can be quantified — for example, in the form of the dynamics of reactions to specific content — then the initial reaction of the community can be used to model its further behavior and, on this basis, determine the credibility of the content [19], [20], [21]. Such an approach is considered in [10], [22]. However, it focuses only on capturing the dynamics of user reactions without analyzing the underlying reasons for such reactions. It is also possible to model user behavior quantitatively. For instance, one can estimate the number of posts in a short period following the appearance of certain content [23], [24].

One challenge in modeling user behavior is the lack of large datasets, which makes it impossible to use neural networks [7], [25]. In this case, it is advisable to use causal models that represent the relationship between the outcomes — on which the decision about content credibility is based — and the influencing factors in the form of algebraic equations. Moreover, the quantitative data used for decision-making often belongs to the category of imprecise data. In such cases, it is reasonable to use an interval representation of the result and, on this basis, build interval models [13], [26]. The advantage of this approach is that it does not require large datasets and guarantees the mathematical model's accuracy.

Thus, the aim of this paper is to develop an interval model for decision-making regarding the credibility of content posted on social networks by establishing the relationship between the outcome — on which the credibility decision is based — and the influencing factors. At the same time, to identify such a model, it is necessary to develop a hybrid method for identifying interval models of user portraits in social networks.

Quantitative methods make it possible to develop predictive models based on the analysis of user behavior observed in previous cases. For example, one can account for which topics generate the most interest or how quickly content can become viral.

To build such a quantitative model, let us introduce the indicator y , which characterizes the degree of credibility of specific content on a scale from 0 to 1. Here, 0 means completely false content, and 1 means fully credible content. It should be noted that this indicator is not interpreted as a probability but as a quantitative value within a certain interval. Quantitative factors on which the degree of credibility can be determined may include:

- x_1 — the number of posts, shares, or likes made by users within a short period after the content appears. This helps detect the audience's immediate reaction to certain content. For example, atypical reactions may indicate a high degree of falsehood. If certain news or posts receive many negative reactions or are marked as "fake," this can be a significant indicator;
- x_2 — the number of comments or reactions at specific time intervals, which allows for better understanding of

the speed of information dissemination and the emotional response. For instance, atypical speed of dissemination may indicate a high probability of false content;

- x_3 — the time it takes for information to spread through social networks (for example, how many people interact with the content during the first minutes, hours, or days after publication). This helps to understand the efficiency of the content's spread. Studying how news spreads can point to its unreliability; fake news often spreads rapidly within certain groups or through bots;
- x_4 — the coefficient of viral spread of the content (for example, the average number of shares per user), which helps determine whether the content is prone to rapid dissemination.

These factors describe the profile — that is, the typical or atypical community reaction to published content — and can help determine the degree of credibility or falsehood of that content.

II. PROBLEM FORMULATION FOR IDENTIFYING INTERVAL MODELS OF USER PORTRAITS IN SOCIAL NETWORKS

The influence of uncontrolled other factors on community reaction is proposed to be taken into account in the form of interval estimates of the degree of reliability or unreliability of this content.

The dependence of the indicator $y(\vec{X})$, which characterizes the degree of reliability of certain content on the values of factors $\vec{X} = (x_1, x_2, x_3, x_4)$ will be described in general form by a nonlinear algebraic equation [27]:

$$y(\vec{X}) = f_1(\vec{\beta}, \vec{X}) + f_2(\vec{\beta}, \vec{X}) + \dots + f_m(\vec{\beta}, \vec{X}), \quad (1)$$

where $y(\vec{X})$ — modeled value of content reliability, $\vec{\beta}$ — unknown vector of interval model parameters, $\lambda_{m_s} = \{f_1(\vec{\beta}, \vec{X}), f_2(\vec{\beta}, \vec{X}), \dots, f_m(\vec{\beta}, \vec{X})\}$ — set of basis functions, in general case nonlinear, both in input factors and in model parameters; m_s — number of basis functions of the model, that is, its structural elements.

Data for identifying the mathematical model (1) are obtained in the following form:

$$\vec{X}_i \rightarrow [y_i^-, y_i^+], i = 1, \dots, N, \quad (2)$$

where $[y_i^-, y_i^+]$ — lower and upper bounds of the degree of reliability of the i -th content, $i = 1, \dots, N$, \vec{X}_i — values of factors that describe the portrait, that is, the community's reaction to the published i -th content and can make it possible to determine the degree of reliability or unreliability of this content, N — total number of observations in the experiment.

Based on expressions (1) and (2), we obtain conditions for solving the problem of structural and parametric model identification:

$$y_i^- \leq f_1(\vec{\beta}, \vec{X}_i) + \dots + f_m(\vec{\beta}, \vec{X}_i) \leq y_i^+, i = \overline{1, N}, \quad (3)$$

As we can see, the conditions from which we obtain both the structure and parameters of the model are an interval system of nonlinear algebraic equations. As is known, its solution is a set of mathematical models (corridor). The task of finding a

solution to ISNAR is itself an NP-complete computational problem. To simplify it, we will look for only point estimates of parameters $\vec{\beta}^m$ for one candidate model that satisfies conditions (3).

Under such conditions, structural and parametric identification are based on optimization problems, for solving which methods of multidimensional nonlinear optimization are used [28]. In this case, an optimization problem is solved for model identification (structure and parameters) in the form of such an optimization problem [29]:

$$\delta(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}) \xrightarrow{\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}} \min \quad (4)$$

$$\beta_j^m \in [\beta_j^{low}; \beta_j^{up}], j = 1, \dots, m \quad (5)$$

$$\lambda_{m_s} \in F, \quad (6)$$

$$\alpha_i \in [0, 1], i = 1, \dots, N, \quad (7)$$

where F – set of all possible structural elements of the interval model, λ_{m_s} – number of all possible elements of the s -th structure, α_i – coefficients of linear combination for determining a point within the bounds of the degree of reliability of the i -th content $[y_i^-; y_i^+]$.

For each current structure, as a candidate model, which we form using algorithms of directed enumeration of structural elements, we evaluate model parameters to obtain quantitative estimates of predicted (calculated) values of the degree of reliability of the i -th content in the form of such an expression:

$$\hat{y}_i(\vec{X}_i) = f_1(\vec{\beta}_1, \vec{X}_i) + f_2(\vec{\beta}_2, \vec{X}_i) + \dots + f_m(\vec{\beta}_m, \vec{X}_i), i = 1, \dots, N, \quad (8)$$

and compare with given values using such an objective function:

$$\delta(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}) = \sum_{i=1}^N (\hat{y}_i(\vec{X}_i) - P([y_i^-; y_i^+], \alpha_i))^2, \quad (9)$$

where

$$P([y_i^-; y_i^+], \alpha_i) = \alpha_i \cdot y_i^- + (1 - \alpha_i) \cdot y_i^+, i = 1, \dots, N \quad (10)$$

As an additional stopping criterion for optimization procedures, conditions [8] can be used:

$$\hat{y}_i(\vec{X}) \in [y_i^-; y_i^+], i = 1, \dots, N \quad (11)$$

It should be noted that the obtained optimization problem (4) – (7) can be solved by combining global search methods with gradient methods.

III. HYBRID METHOD FOR IDENTIFYING INTERVAL MODELS OF USER PORTRAITS IN SOCIAL NETWORKS

Considering that the optimization problem (4) – (7) is a nonlinear optimization problem, which also contains two types of constraints:

$$\lambda_m \in F \quad (12)$$

and

$$\beta_j^m \in [\beta_j^{low}; \beta_j^{up}], j = 1, \dots, m \quad (13)$$

$$\alpha_i \in [0, 1], i = 1, \dots, N, \quad (14)$$

It is advisable to transform it to the following form:

$$\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma) \xrightarrow{\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma} \min \quad (15)$$

$$\lambda_{m_s} \in F =$$

$$\{ \varphi_1(\vec{x}), \dots, \varphi_m(\vec{x}), \varphi_{m+1}(\vec{g}), \dots, \varphi_{2m}(\vec{g}) \} \quad (16)$$

As we can see, two unknown coefficients μ, γ are introduced into the objective function by collapsing linear constraints on parameter values $\vec{\beta}^m$ and coefficients $\vec{\alpha}$ using penalty functions in the following form:

$$\varepsilon = \gamma \cdot \sum_{j=1}^N \left(\ln(\hat{\beta}_j^m - \hat{\beta}_j^{low}) + \ln(\hat{\beta}_j^{up} - \hat{\beta}_j^m) \right), \quad (17)$$

$$\Delta = \mu \cdot \sum_{i=1}^N (\ln(\alpha_i) + \ln(1 - \alpha_i)), \quad (18)$$

where γ, μ – given influence coefficients of corresponding penalty functions.

As a result, the objective function in the interval model identification problem takes the following form:

$$\Phi(\lambda_{m_s}, \vec{\beta}, \vec{\alpha}, \mu, \gamma) = \sum_{i=1}^N (\hat{y}_i(\vec{X}_i) - P([y_i^-; y_i^+], \alpha_i))^2 - \gamma \cdot \sum_{j=1}^N (\ln(\hat{\beta}_j^m - \hat{\beta}_j^{low}) + \ln(\hat{\beta}_j^{up} - \hat{\beta}_j^m)) - \mu \cdot \sum_{i=1}^N (\ln(\alpha_i) + \ln(1 - \alpha_i)). \quad (19)$$

Thus, we obtain a barrier objective function, the value of which sharply increases (or decreases) when approaching the boundary values of parameters $\vec{\beta}^m$ or coefficients $\vec{\alpha}$.

Now the interval model identification problem is formulated in the form (15), (16) with objective function (19) in the form of a barrier function. The specified problem is an optimization problem on a discrete set of basis functions of the mathematical model. For its solution, it is necessary to apply a combination of methods of directed enumeration of structural elements with gradient methods.

The work proposes a hybrid method for identifying interval models of user portraits in social networks, which is based on combining a metaheuristic algorithm for model structure synthesis based on the behavioral model of a bee colony and gradient methods for identifying candidate model parameters.

Let us consider the proposed method in detail.

First of all, according to condition (16), the set of basis functions λ_m for a specific candidate model is formed from the general set F , which contains all structural elements from which an interval model can be formed. Increasing the number m of structural elements in the mathematical model can lead to its significant complication. Therefore, in the optimization problem (15), (16), an additional constraint should be added, which concerns the upper bound of the number of structural elements in the model. Let us denote this constraint in the following form:

$$m_s \leq I_{max} \quad (20)$$

where m_s – number of structural elements in the current s -th structure, I_{max} – maximum allowable number of structural elements in candidate model.

Let us introduce constraint (20) in the form of a penalty function into the objective function (19). As a result, we obtain:

$$\begin{aligned} \Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) = & \sum_{i=1}^N (\hat{y}_i(\vec{x}_i) - P([y_i^-; y_i^+], \alpha_i))^2 \\ & - \gamma \cdot \sum_{j=1}^N (\ln(\hat{\beta}_j^m - \hat{\beta}_j^{low}) + \ln(\hat{\beta}_j^{up} - \hat{\beta}_j^m)) \\ & - \mu \cdot \sum_{i=1}^N (\ln(\alpha_i) + \ln(1 - \alpha_i)) - \sigma \cdot \ln(I_{max} - m_s), \end{aligned} \quad (21)$$

where σ – given influence coefficient of the penalty function.

Then, we rewrite the optimization problem (16), (17) in the following form:

$$\begin{aligned} \Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) & \xrightarrow{\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma} \min \quad (22) \\ \lambda_{m_s} & \in F. \quad (23) \end{aligned}$$

Now, we divide the process of solving the optimization problem (22), (23) into two stages:

1. Formation of the current structure based on candidate models (model structure synthesis);
2. Evaluation of its parameters and verification of model adequacy (parametric identification).

At the first stage, we use a metaheuristic algorithm, which is built on behavioral models of a bee colony. The application of this algorithm will make it possible to form and evaluate several different structures of candidate models in parallel. Formally, the bee colony algorithm is based on swarm intelligence, which uses a swarm to search for nectar [30], [31]. Let us consider the main phases of synthesis of candidate model structures, based on analogy with the stages of the behavioral model of a bee colony.

According to the set problem (22), (23), we set initial conditions: $LIMIT$ – number of iterations for exhaustion of the current candidate model structure; S – total number of candidate models within one iteration; I_{max} ; $mcn = 0$ – number of current iteration; MCN total number of iterations; set of structural elements F . We also randomly generate the initial set of structures λ_{m_s} of candidate models, with a total number of S based on the set of all structural elements F .

Worker bees phase. At this phase, we use a number of operators to form a set of candidate model structures [32]. Operator $P(\Lambda_{mcn}, F)$, which transforms each structure λ_{m_s} from the set Λ_{mcn} of interval model in the form (8) to structure λ'_{m_s} , which is similar to structure λ_{m_s} . This operation corresponds to the ABC algorithm in the sense that worker bees explore neighboring nectar sources. As a result, operator $P(\Lambda_{mcn}, F)$ transforms the set of structures Λ_{mcn} to the set of structures Λ'_{mcn} , on the mcn iteration of the structure synthesis algorithm. This transformation occurs by randomly selecting elements from each structure λ_{m_s} and replacing these elements with elements that are randomly selected from the general set F . Using operator $P(\Lambda_{mcn}, F)$ in the current candidate structure, a different number of elements can be replaced, depending on the quality of the specific structure. To determine the number of elements that need to be replaced, it is necessary to evaluate the "quality" of the current candidate model, which

is determined by the value of the objective function $\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma)$. Obviously, the smaller the value of this function, the more accurate the mathematical model and accordingly, the smaller the number of elements that need to be changed in the current structure. It should be noted that the number of elements in the current structure also depends on the total number of elements m_s in this structure. Therefore, the formula for determining the number of elements that need to be replaced in the current structure can be as follows:

$$n_s = \begin{cases} \text{int} \left(\left(1 - \frac{\min\{\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) | s=1..S\}}{\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma)} \right) * m_s \right), \\ \text{if } \Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) \neq \\ \min\{\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) | s=1..S\} \\ \text{and } n_s \neq 0 \\ 1, \text{ if } \Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) = \\ \min\{\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) | s=1..S\} \\ \text{or } n_s = 0 \end{cases} \quad (24)$$

At this same phase, pairwise comparison of generated and current structures is also carried out to select the better one from the pair. Obviously, for this it is necessary to carry out parametric identification for each candidate model using gradient methods. Thus, the operator of pairwise comparison with selection of the better structure has the following form:

$$D(\lambda_{m_s}, \lambda'_{m_s}): \lambda_{m_s}^1 = \begin{cases} \lambda_{m_s}, \text{ if } \Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) \leq \\ \Phi(\lambda'_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) \\ \lambda'_{m_s}, \text{ if } \Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) > \\ \Phi(\lambda'_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) \end{cases} \quad (25)$$

The specified operator (25) performs selection of the best structures $\lambda_{m_s}^1$ from two sets by pairwise comparison of structures from two sets $\lambda_{m_s}, \lambda'_{m_s}$. As a result, we obtain the set of structures of the first iteration $\lambda_{m_s}^1 \in \Lambda_{mcn}^1$.

Scout bees phase. At this phase, a number of other structures are formed around the determined structures from the set Λ_{mcn}^1 . In the context of the behavioral model of a bee colony, scout bees search for new nectar sources around already known ones [33], [34]. This means that for each current structure, it is necessary to generate a certain number R_s of structures. This indicator depends on the quality of the current structure. Therefore, to determine the number of structures that we generate for the current one, we use the following formulas:

$$P_s(\lambda_{m_s}^1) = \frac{1}{\Phi(\lambda_{m_s}^1, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma) \sum_{s=1}^S \frac{1}{\Phi(\lambda_{m_s}^1, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma)}}, \quad (26)$$

$$R_s = \text{ToInt}(P_{s-1}(\lambda_{m_{s-1}}^1), s = 2..S, R_0 = 0) \quad (27)$$

Now, having the known number of structures that need to be formed around the current one, we use operator $P_\delta(\Lambda_{mcn}^1, F)$, which in a similar way to operator $P(\Lambda_{mcn}, F)$ transforms the

set of structures Λ_{mcn}^1 to the set of structures Λ'_s . Again, for each structure $\lambda_{m_s}^1$, R_s structures are formed by randomly replacing the corresponding number n_s of structural elements, randomly selected from set F .

Further, at this phase, using the group selection operator $D_2(\lambda_{m_s}^1, \lambda_{m_s}^1)$ between the current structure and the set of structures formed around it, we choose the best candidate model by selection within the formed group. For this, it is necessary to calculate the value of the objective function (21) for each structure by conducting a parametric identification procedure. As a result, this operator chooses the best structure for those generated in the group. In the selection process, we use formula (25) with the difference that group selection is carried out, not pairwise.

Thus, the specified operator forms structures of the second row of formation Λ_{mcn}^2 on the same mcn iteration.

One of the biggest problems of the above-described algorithm is cycling around the formation of certain structures, that is, at local minima of the optimization problem. To exit local minima, an additional phase of scout bees is provided in ABC.

Scout bees phase. This is a phase of the bee swarm, at which bees randomly choose new nectar sources. In the context of the optimization problem, this means that for some structures it is necessary to form completely new structures randomly [35]. For this, in the computational method for each current candidate structure, a variable $Limit_s$ is introduced. This variable models the exhaustion of the structure and possible complete change of the structure by randomly generating a new set of structural elements. Structure exhaustion occurs when the number of modifications of the current structure exceeds $LIMIT$ without improving its quality. Then we use operator $P_N(I_{max}, F)$, which generates the corresponding new structure.

Thus, the above-described scheme makes it possible to gradually form new structures of candidate interval models and at the same time carry this out in the direction of improving their quality. As already mentioned above, for each fixed current candidate structure, we solve the parametric identification problem for fixed structural elements by solving an optimization problem. Considering the differentiability of the objective function (21), at this stage, we can use gradient methods.

Therefore, at the second stage of interval model identification, the parametric identification algorithm scheme can be as follows:

Step 1. Initialization:

- introduction of experimental data $\vec{X}_i \rightarrow [y_i^-, y_i^+]$, $i = \overline{1, N}$;
- reading the current structure of the model structure λ_{m_s} (a set of structural elements in quantity m_s) from the first stage;
- setting initial values of the components of the parameter vector $\hat{\beta}^m_j \in [\hat{\beta}^{m_{low}}_j; \hat{\beta}^{m_{up}}_j]$, $j = 1, \dots, m$;
- setting initial values for the coefficients-components of vector $\vec{\alpha}$, (usually we will set them as follows: $\alpha_i = 0.5$, $i = \overline{1, N}$);
- introduction of initial values of coefficients μ, γ, σ for penalty functions;
- formation of the objective function $\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma)$;
- introduction of boundary values of stopping criteria;

Step 2. Cycle:

Start of step 2.

While none of the criteria is satisfied, execute:

Step 2.1. Update the barrier function $\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma)$ by formula (21);

Step 2.2. Calculate the gradient of the barrier function (direction of increase in the value of the barrier function):

$$\begin{aligned} \vec{\nabla} \Phi(\vec{\beta}^m, \vec{\alpha}, m_s) = & \vec{\nabla} \left(\sum_{i=1}^N \left(\hat{y}_i(\vec{X}_i) - P([y_i^-, y_i^+], \alpha_i) \right)^2 \right) - \gamma \\ & \cdot \vec{\nabla} \left(\sum_{j=1}^N \left(\ln(\hat{\beta}^m_j - \hat{\beta}^{m_{low}}_j) + \ln(\hat{\beta}^{m_{up}}_j - \hat{\beta}^m_j) \right) \right) - \\ & \mu \cdot \vec{\nabla} \left(\sum_{i=1}^N (\ln(\alpha_i) + \ln(1 - \alpha_i)) \right) - \sigma \\ & \cdot \vec{\nabla} \left(\ln(I_{max} - m_s) \right); \end{aligned}$$

Step 2.3. Determine the descent direction (normalized antigradient vector):

$$-\frac{\vec{\nabla} \Phi(\vec{\beta}^m, \vec{\alpha}, m_s)}{\|\vec{\nabla} \Phi(\vec{\beta}^m, \vec{\alpha}, m_s)\|} = -\vec{\nabla} \tilde{\Phi}(\vec{\beta}^m, \vec{\alpha}, m_s);$$

Step 2.4. Search for the optimal solution at this step with step length s :

$$\begin{aligned} (\vec{\beta}^m_{k+1}, \vec{\alpha}_{k+1}, m_{s_{k+1}}) \\ = (\vec{\beta}^m_k, \vec{\alpha}_k, m_{s_k}) - s \cdot \vec{\nabla} \tilde{\Phi}(\vec{\beta}^m, \vec{\alpha}, m_s); \end{aligned}$$

Step 2.5. Update parameters μ, γ, σ ;

Step 2.6. Check stopping criteria;

End of step 2.

Step 3. Return the parameter vector $\vec{\beta}$.

It should be noted that at step 2.4, we can solve an optimization problem with one parameter s – step length, to ensure greater convergence to the minimum point (such a method is called the steepest descent method), or we can simply reduce the step length value when approaching a local minimum, for example, by halving it at each iteration.

We choose parameters μ, γ, σ in such a way that when approaching the values of variables given in the constraints, the values of the corresponding penalty functions increase significantly and increase the value of this barrier function $\Phi(\lambda_{m_s}, \vec{\beta}^m, \vec{\alpha}, \mu, \gamma, \sigma)$.

We can also note that the parameter of the number of basis functions m_s in the current candidate model can be extracted from the barrier function, since it complicates the parametric identification problem, because it is integer. Instead, use the

procedure of gradual building up of the model structure by directed selection of structural elements and increasing their number. However, such an approach in forming structures can lead to over-complication of the resulting candidate model, since the search for adequate structures will not allow reducing the number of structural elements in the mathematical model.

IV. RESULTS AND DISCUSSIONS

Let us consider the application of the developed method for modeling the portrait of a social network. The specified network belongs to a news network that unites a community oriented towards news content. A feature of this network is a quick reaction to new content. The interaction of community members mainly increases during the first few hours, after which the level of interaction decreases, unless the news causes a long-term resonance. In such a network, anomalies may emerge because of resonant topics, leading to sudden and intense surges in activity. In such cases, the risk of spreading fakes or manipulative content increases due to increased emotionality. As already mentioned, communities oriented towards news content are characterized by rapid dynamics of interaction. Such communities are also characterized by polarized reactions and sensitivity to manipulative content. They have a heterogeneous audience that strives for quick access to relevant information. As a result, a feature of such communities is that their "portrait" changes significantly depending on the nature of the news and the context in which they appear. This feature is used to recognize fake content, which mainly causes unnatural resonance of the user portrait.

For studying the portrait of the specified social network, $N = 20$ cases of posting content of different nature were used, and the community portrait was recorded using such factors: x_1 – the number of posts, shares, or likes made by users within the first 10 minutes after the content appeared; x_2 – dynamics of comments or reactions every 10 minutes (measured by the average number of comment increments); x_3 – the number of unique users for half a day; x_4 – the viral spread coefficient of content. The results of studies of this network are given in Table 1.

Table 1. Results of content studies in the network

N	Number of likes first 10 min	Spread dynamics (average value)	Number of unique users characterizing news spread time (thousands)	Viral spread coefficient	Degree of content reliability
	x_1	x_2	x_3	x_4	$[y_i^-, y_i^+]$
1	52	30,1	2,114	2,5	[0,9; 1]
2	35	24,8	1,878	2,2	[0,95; 1]
3	64	38	2,789	3,4	[0,78; 0,91]
4	92	41	2,830	3,2	[0,54; 0,71]
5	28	20,5	1,500	1,9	[0,96; 1]
6	60	36,7	2,600	3,1	[0,85; 0,93]
7	75	39,9	2,750	3,3	[0,72; 0,84]
8	100	45	3,000	3,8	[0,50; 0,65]
9	110	50,2	3,300	4	[0,35; 0,52]
10	20	15,5	1,300	1,6	[0,97; 1]
11	40	22,3	1,700	2,1	[0,94; 0,99]
12	88	42	2,900	3,5	[0,60; 0,75]

13	99	47,3	3,200	3,9	[0,42; 0,60]
14	120	52,5	3,450	4,2	[0,25; 0,45]
15	18	14,2	1,200	1,4	[0,98; 1]
16	47	27,5	2,000	2,3	[0,88; 0,94]
17	85	43,1	2,950	3,6	[0,66; 0,80]
18	105	49,5	3,350	4,1	[0,38; 0,55]
19	115	53,7	3,600	4,3	[0,21; 0,40]
20	125	55,8	3,800	4,5	[0,12; 0,28]

The mathematical model for evaluating the degree of content reliability is presented in the form of expression (1), where:

$\lambda_{m_s} = \{f_1(\vec{\beta}, \vec{X}), f_2(\vec{\beta}, \vec{X}), \dots, f_m(\vec{\beta}, \vec{X})\}$ – set of basis functions, both in input factors and in model parameters;

$m_s = [3; 5]$ – the number of basis functions of the model, that is, its structural elements is defined in the range;

$\vec{\beta}$ – unknown vector of interval model parameters, we will calculate based on table data using the developed hybrid method for identifying interval models of user portraits in social networks.

Using the data from Table 1, we compose the corresponding interval system (3). Further, based on the generated set of structural elements using the developed hybrid method of structural and parametric identification, we evaluate candidate models, applying the objective function in the form (21).

As a result of implementing the hybrid method, on the 12th iteration of evaluating candidate models, we obtained the following coefficient estimates: $\beta_0 = 0.491103$, $\beta_1 = -0.001786$, $\beta_2 = 0.048772$, and the corresponding mathematical model:

$$y(\vec{X}) = 0,491103 - 0,001786 \cdot x_1 \cdot x_4 + 0,048772 \cdot x_2/x_3 \quad (28)$$

As we can see, the obtained interval model contains three coefficients and accordingly three structural elements. We also see that the mathematical model is nonlinear with respect to the factors on which content reliability depends.

Now we can use the obtained interval model to evaluate the reliability of content posted in the specified network.

Example 1. Let us consider an example of applying the developed mathematical model for predicting content reliability.

News "The European Union has approved a new COVID-19 vaccine adapted to the Omicron variant1" appeared in the network.

The number of likes due to the appearance of the news in the first 10 minutes was $x_1 = 28$; the dynamics (increment) of news spread through the network, recorded over the first six hours every 10 minutes on average was $x_2 = 29$; the number of unique users characterizing the news spread time for half a day was $x_3 = 2,353$ participants; the viral spread coefficient is $x_4 = 2.3$.

Using the obtained mathematical model (28), we get:

$$y(\vec{X}) = 0,491103 - 0,001786 \cdot 28 \cdot 2.3 + 0,048772 \cdot 29/2.353 = 0.98$$

The obtained result means that the news is reliable, which fully corresponds to reality, based on the context of the news.

Example 2. Let us consider the following example of applying the developed mathematical model for predicting content reliability.

News "The European Union mandates all citizens to receive an annual vaccination against COVID-19 with a new "genetic shot."" appeared in the network.

The number of likes due to the appearance of the news in the first 10 minutes was $x_1=138$; the dynamics (increment) of news spread through the network, recorded over the first six hours every 10 minutes on average was $x_2=64$, which means quite fast content spread; the number of unique users for half a day was $x_3=4,372$ participants, which indicates a short spread time; the viral spread coefficient is $x_4=4.2$.

Using the obtained mathematical model (28), we get:

$$y(\vec{X}) = 0,491103 - 0,001786 \cdot 138 \cdot 4.2 \\ + 0,048772 \cdot 64/4.372 = 0.17$$

The obtained result means that the news is unreliable, which fully corresponds to reality, based on the context of the news.

The obtained results indicate that the developed model can be used to evaluate content reliability.

At the same time, it should be noted that the proposed quantitative indicators, as well as the developed mathematical model, are suitable for analyzing and classifying content into two categories – reliable and fake, but to confirm the results of this classification, it is additionally necessary to conduct analysis of content, sources, and context itself. Thus, it is quite difficult to determine that content is fake or 100% reliable only by likes, reposts, or spread speed. Therefore, in the further development of the proposed approach, it is mandatory to take into account additional content verification methods.

V. CONCLUSIONS

Quantitative indicators that reflect the portrait of users in social networks have been analyzed. It has been established that the use of quantitative indicators of community portrait can help identify signs of fake or false content, although this does not guarantee 100% accuracy. Analysis of such data can be useful for recognizing anomalies in audience behavior, which are often characteristic of communities that spread fakes. It has also been shown that the main indicators characterizing audience reaction to certain content are: the number of posts, shares, or likes made by users within a short time after the content appears, which helps to identify instant audience reactions to certain content; the number of comments or reactions at certain time intervals, which allows better understanding of the speed of information spread and emotional response; the time during which information spreads through social networks (for example, how many people interact with content within the first minutes, hours, or days after publication), helps to understand how effective content distribution is; the viral spread coefficient of content, for example, the number of shares from each user, which helps to understand whether content is prone to rapid spread.

For making decisions about the reliability of content posted in social networks, an interval mathematical model is proposed and justified for the first time, which establishes the relationship between the result on which the decision about the reliability or unreliability of content is made and the factors that influence it. The resulting indicator of such a model is the degree of reliability of certain content within the range from 0 to 1. It is proposed and justified to use interval data analysis methods to represent and analyze this indicator based on expert content research. Accordingly, this indicator will not be interpreted as probabilistic, but as a quantitative value on a certain interval.

A hybrid method for identifying interval models of user portraits in social networks is proposed and justified for the first time, which, unlike existing ones, is based on combining a metaheuristic algorithm for model structure synthesis based on the behavioral model of a bee colony and gradient methods for identifying candidate model parameters, which ensured a reduction in the computational complexity of method implementation and the possibility of using standard optimization tools for solving problems of identifying user portrait models in social networks.

Verification of the developed hybrid method for identifying interval models of user portraits in social networks was carried out on the example of modeling user behavior on various types of news in a social network. At the same time, it should be noted that the proposed quantitative indicators are suitable for initial analysis and detection of suspicious content, but to confirm fakeness, it is necessary to conduct additional analysis of content, sources, and context. Therefore, it is quite difficult to determine that content is fake only by likes, reposts, or spread speed. It is possible only to attribute some reliability to the indicated facts. Thus, in the next section, we will consider additional mechanisms that make it possible to improve the reliability of content classification.

References

- [1] O.S. Ulichev, "Research on models of information dissemination and informational influence in social networks," *Systems of Control, Navigation and Communication*, issue 4, pp. 147–151, 2018. [Online]. Available at: http://nbuv.gov.ua/UJRN/suntz_2018_4_31. (in Ukrainian)
- [2] Ye. Ivohin, L. Adzhubey, "On modeling the dynamics of information dissemination based on heterogeneous diffusion hybrid models," *Scientific Bulletin of Uzhhorod University. Series: Mathematics and Computer Science*, pp. 112–118, 2019. [https://doi.org/10.24144/2616-7700.2019.2\(35\).112-118](https://doi.org/10.24144/2616-7700.2019.2(35).112-118). (in Ukrainian)
- [3] O. Ulichev, Ye. Meleshko, D. Sawicki, S. Smailova, "Computer modeling of dissemination of informational influences in social networks with different strategies of information distributors," *Proc. SPIE 11176*, Wilga, Poland, 2019, Article No.: 111761T. <https://doi.org/10.1117/12.2536480>. (in Ukrainian)
- [4] S. Vosoughi, D. Roy, and S. Aral, "The spread of true and false news online," *Science*, vol. 359, no. 6380, pp. 1146–1151, 2018. <https://doi.org/10.1126/science.aap9559>.
- [5] M. Del Vicario, A. Bessi, F. Zollo, F. Petroni, A. Scala, G. Caldarelli, H. E. Stanley, and W. Quattrociocchi, "The spreading of misinformation online," *Proceedings of the National Academy of Sciences*, vol. 113, no. 3, pp. 554–559, 2016. <https://doi.org/10.1073/pnas.1517441113>.
- [6] Y. Wang, F. Ma, Z. Jin, Y. Yuan, G. Xun, L. Jiao, and A. Su, "EANN: Event adversarial neural networks for multi-modal fake news detection," *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 849–857, 2018. <https://doi.org/10.1145/3219819.3219903>.
- [7] J. Ma, W. Gao, P. Mitra, S. Kwon, B. J. Jansen, K. Wong, and M. Cha, "Detecting rumors from microblogs with recurrent neural networks," *Proceedings of the 25th International Joint Conference on Artificial Intelligence*, pp. 3818–3824, 2016.
- [8] K. Shu, D. Mahudeswaran, S. Wang, D. Lee, and H. Liu, "Hierarchical propagation networks for fake news detection: Investigation and exploitation," *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 14, pp. 626–637, 2020. <https://doi.org/10.1609/icwsm.v14i1.7329>.
- [9] J. Zhang, B. Dong, and S. Y. Philip, "FakeDetector: Effective fake news detection with deep diffusive neural network," *2020 IEEE 36th International Conference on Data Engineering*, pp. 1826–1829, 2020. <https://doi.org/10.1109/ICDE48307.2020.00180>.
- [10] M. R. Islam, S. Liu, X. Wang, and G. Xu, "Deep learning for misinformation detection on online social networks: A survey and new perspectives," *Social Network Analysis and Mining*, vol. 10, no. 1, pp. 1–20, 2020. <https://doi.org/10.1007/s13278-020-00696-x>.
- [11] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake news detection on social media: A data mining perspective," *ACM SIGKDD Explorations*

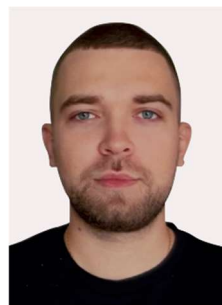
- Newsletter, vol. 19, no. 1, pp. 22–36, 2017. <https://doi.org/10.1145/3137597.3137600>.
- [12] M. Granik and V. Mesyura, "Fake news detection using naive Bayes classifier," 2017 IEEE First Ukraine Conference on Electrical and Computer Engineering, pp. 900–903, 2017. <https://doi.org/10.1109/UKRCON.2017.8100379>.
- [13] S. P. Shary, "Interval methods for data fitting under uncertainty," Journal of Computational and Applied Mathematics, vol. 418, pp. 114–135, 2023. doi: 10.1016/j.cam.2022.114135.
- [14] O. Ulichev, Y. Meleshko, V. Khokh, "The computer simulation method of a social network structure for the research of dissemination processes of informational influences," Scientific and Practical Cyber Security Journal (SPCSJ), 4(3). – Georgia, Tbilisi, 2019, pp. 34–47 (in Ukrainian).
- [15] O.S. Ulichev, Ye.V. Meleshko, "Software modeling of the dissemination of informational and psychological influences in virtual social networks," Collection of Scientific Papers 'Modern Information Systems', Issue 2(2). – Kharkiv: NTU KhPI, 2018, pp. 35–39. <https://doi.org/10.20998/2522-9052.2018.2.06> (in Ukrainian).
- [16] A. Guess, J. Nagler, and J. Tucker, "Less than you think: Prevalence and predictors of fake news dissemination on Facebook," Science Advances, vol. 5, no. 1, pp. 1–8, 2019. <https://doi.org/10.1126/sciadv.aau4586>.
- [17] K. Shu, D. Mahudeswaran, S. Wang, D. Lee, and H. Liu, "FakeNewsNet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media," Big Data, vol. 8, no. 3, pp. 171–188, 2020. <https://doi.org/10.1089/big.2020.0062>.
- [18] M. Zubiaga, A. Aker, K. Bontcheva, M. Liakata, and R. Procter, "Detection and resolution of rumours in social media: A survey," ACM Computing Surveys, vol. 51, no. 2, pp. 1–36, 2018. <https://doi.org/10.1145/3161603>.
- [19] M. J. Metzger, A. J. Flanagan, and R. B. Medders, "Social and heuristic approaches to credibility evaluation online," Journal of Communication, vol. 60, no. 3, pp. 413–439, 2010. <https://doi.org/10.1111/j.1460-2466.2010.01488.x>.
- [20] L. Zhou, D. Zhang, and C. C. Lee, "A survey of opinion mining and sentiment analysis," Mining Text Data, pp. 415–463, 2012. https://doi.org/10.1007/978-1-4614-3223-4_13.
- [21] K. Shu, S. Wang, and H. Liu, "Beyond news contents: The role of social context for fake news detection," Proceedings of the 12th ACM International Conference on Web Search and Data Mining, pp. 312–320, 2019. <https://doi.org/10.1145/3289600.3290994>.
- [22] H. Karimi and J. Tang, "Learning hierarchical discourse-level structure for fake news detection," Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics, pp. 3432–3442, 2019. <https://doi.org/10.18653/v1/N19-1347>.
- [23] O.S. Ulichev, Ye.V. Meleshko, "Modeling of dissemination and neutralization processes of informational influences in a segment of a social network," Scientific Journal 'Information Protection'. – Kyiv: NAU, 2020, pp. 166–176 (in Ukrainian).
- [24] K. Mateńczuk, A. Kozina, A. Markowska, K. Czerniachowska, K. Kaczmarczyk, P. Golec, M. Hernes, K. Lutosławski, A. Koziarkiewicz, M. Pietranik, A. Rot, M. Dyvak, "Financial Time Series Forecasting: Comparison of Traditional and Spiking Neural Networks," Procedia Computer Science, vol. 192, 2021, pp. 5023–5029, <https://doi.org/10.1016/j.procs.2021.09.280>.
- [25] A. K. Jain and B. B. Gupta, "A machine learning based approach for phishing detection using hyperlinks information," Journal of Ambient Intelligence and Humanized Computing, vol. 10, no. 5, pp. 2015–2028, 2019. <https://doi.org/10.1007/s12652-018-0798-z>.
- [26] M. Anderson, P. Wilson, "Interval arithmetic in optimization: Theory and applications," Applied Mathematics and Computation, vol. 456, 2023, pp. 1–18.
- [27] M. Dyvak, I. Spivak, A. Melnyk, V. Manzhula, T. Dyvak, A. Rot, M. Hernes, "Modeling Based on the Analysis of Interval Data of Atmospheric Air Pollution Processes with Nitrogen Dioxide due to the Spread of Vehicle Exhaust Gases," Sustainability, vol. 15, 2023, p. 2163. <https://doi.org/10.3390/su15032163>.
- [28] I. Darmorost, M. Dyvak, N. Porplytsya, T. Shynkaryk, Y. Martsenyuk, V. Brych, "Convergence Estimation of a Structure Identification Method for Discrete Interval Models of Atmospheric Pollution by Nitrogen Dioxide," Proceedings of the 2019 9th International Conference on Advanced Computer Information Technologies (ACIT), Ceske Budejovice, Czech Republic, 2019, pp. 117–120. <https://doi.org/10.1109/ACITT.2019.8779981>.
- [29] M. Dyvak, "Parameters Identification Method of Interval Discrete Dynamic Models of Air Pollution Based on Artificial Bee Colony Algorithm," Proceedings of the 2020 10th International Conference on Advanced Computer Information Technologies (ACIT), Deggendorf, Germany, 2020, pp. 130–135. <https://doi.org/10.1109/ACIT49673.2020.9208972>.
- [30] D. Karaboga, "An idea based on honey bee swarm for numerical optimization," Technical report, Erciyes University, Engineering Faculty, Computer Engineering Department, Erciyes University, 2005, 10 p. [Online]. Available at: https://abc.erciyes.edu.tr/pub/tr06_2005.pdf.
- [31] M. Karaboga, B. Akay, and D. Karaboga, "Artificial bee colony algorithm for optimization problems: A comprehensive review," Applied Soft Computing, vol. 122, 2022, pp. 108–125. doi: 10.1016/j.asoc.2022.108125.
- [32] J. Li, Y. Wang, and H. Chen, "Enhanced artificial bee colony algorithm with adaptive parameter control for global optimization," IEEE Transactions on Cybernetics, vol. 52, no. 8, 2022, pp. 7896–7908. doi: 10.1109/TCYB.2021.3082345.
- [33] R. Sharma, S. Kumar, and P. K. Singh, "Artificial bee colony algorithm for feature selection in machine learning: A systematic review," Expert Systems with Applications, vol. 204, 2022, pp. 117–135. doi: 10.1016/j.eswa.2022.117135.
- [34] L. Zhang, M. Wang, and X. Liu, "Multi-objective artificial bee colony algorithm for interval optimization problems," Information Sciences, vol. 625, 2023, pp. 1–18. <https://doi.org/10.1016/j.ins.2023.01.045>.
- [35] A. Kumar, D. Kumar, "A comprehensive review of artificial bee colony algorithm variants," Swarm and Evolutionary Computation, vol. 44, 2019, pp. 1–15.



M. DYVAK Doctor of Technical Sciences, Vice-Rector for Scientific Research, Professor of the Department of Computer Sciences of the West Ukrainian National University. Scientific interests: mathematical modeling of static, dynamic systems and systems with distributed parameters, structural and parametric identification of systems, calculation methods, parallelization of calculations for interval analysis tasks, interval data analysis, application of interval methods in environmental monitoring and for modeling distributed systems.



TYANDE PAN Master of Software Engineering, PhD student of West Ukrainian National University. Scientific interests: mathematical modeling of a social network based on interval data analysis.



OLEKSANDR KINDZERSKYI Master of Software Engineering, PhD student of West Ukrainian National University, 10 years of experience as a C# developer, lead of software development team. Scientific interests: mathematical modeling of static, dynamic systems, structural and parametric identification of systems, calculation methods, parallelization of calculations for interval analysis tasks.