



# The System of Operative Determination of the Level of Tension in Society Based on Data from Social Networks

**Maksym Shchoholiev** (✉), **Violeta Tretynyk**

*National Technical University of Ukraine "Igor Sikorsky Kyiv Polytechnic Institute," Kyiv, Ukraine, <https://kpi.ua/en>*

## ABSTRACT:

The system presented in this article was developed for state structures with responsibilities for ensuring security and stability in the society. The main task of the system is to determine the impact of information stored on various Internet resources on society. It serves to assess the level of tension on the basis of comments of social networks' users on certain news and posts. The system takes into account the number of news' reposts and the average emotionality of comments on this news. The methods TF-IDF and word2vec are used for keyword determination from the text and their transformation into numerical values. The level of emotionality is defined by an artificial neural network. To get the resulting estimate, a data aggregation block is developed and used to reduce various independent quantitative and qualitative characteristics into one value.

## ARTICLE INFO:

RECEIVED: 09 MAY 2019

REVISED: 30 AUG 2019

ONLINE: 23 SEP 2019

## KEYWORDS:

social tension, keyword, TF-IDF, word2vec representation, context, data aggregation, social networks, hybrid war



Creative Commons BY-NC-SA 4.0

## Introduction

Today the world is going through a phase of aggravation of information wars. Due to the high level of accessibility of information, millions of people, regardless of their moral and socio-political views, are the target of information attacks. The

basis for a successful information attack is in conflicts and divergences among people and authorities or among certain groups of people. The purpose of many attacks is to increase tensions between these groups. In particular, it is advisable to single out information attacks aimed at cultivating a sense of threat among people. These threats include threats to life of an individual or lives of family members, property loss threats, and the so on.

The system is designed to detect real-time information attacks targeting Ukraine, in particular by Russia, and to determine the impact of these attacks on Ukrainian society.

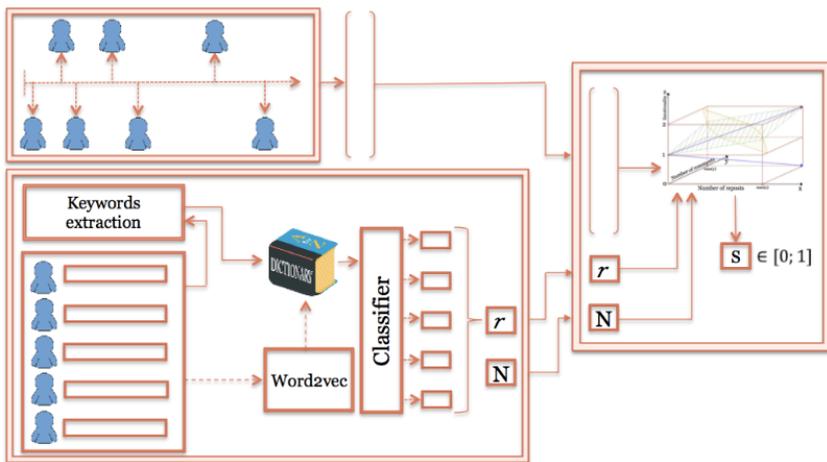
The purpose of work is to develop a system for the dynamic determination of the level of social, informational, political or other tension in the society according to the data from social networks.

The first section of the article presents the architecture of the proposed system. The basic structural blocks, their mathematical interpretation and interconnections are described. The second section describes software implementation, the structure of training and test data, and test results. The final section presents conclusions and plans for future research.

## System Architecture

Schematic representation of the developed system is shown in Figure 1. It follows the main idea of the scheme proposed by Ruchansky, Seo and Liu.<sup>1</sup>

The system consists of three blocks. The first and the second blocks work autonomously. Data in each of them is collected and counted parallelly and independently. The third block receives data from the first two blocks and reduces it to a certain value.



**Figure 1:** The scheme of the system of operative determination of the level of tension in society.

### Subsystem for Tracking the Dynamics of Sharing Particular News

In most social networks the number of news' reposts is counted. The moment of publication of the article fixes on the timeline. The system then automatically calculates the increase of its shares per unit time  $\Delta t$ , for example, per 5 minutes. In this way it is possible to explore the dynamics of readers' interest in the title and content of the post, as well as to determine the exact moment when the frequency of shares is maximal. At the output of this block a numeric vector is obtained:

$$\Delta \bar{p} = (\Delta p_1, \Delta p_2, \dots, \Delta p_n) \in \mathbb{N}^n, \quad (1)$$

$\Delta p_i = p_i - p_{i-1}$ ;  $i = \overline{1, \dots, n}$ ;  $p_i \in \mathbb{N}$  – current number of shares;  $p_0 = 0$ . Each component of this vector shows the increase of number of reposts during a given time. Numbers in the vector are arranged in order of their fixation over time.

### Subsystem for Evaluating the Emotionality of Comments on News

This subsystem is designed to evaluate the content of comments on the news by keywords. In order to use word groups to classify text corpuses, the words are presented as numeric vectors. This representation is created using the word2vec method with the CBOW (continuous bag-of-words) learning algorithm.<sup>2</sup> In this way a dictionary is formed consisting of ordered pairs  $(w_s, l_s) \in (\mathbb{N}^n \times \mathbb{R}^m)$ ;  $s = \overline{1, \dots, q}$ ;  $q$  – current number of words in the dictionary;  $w_s \in \mathbb{N}^n$  – words in Unicode encoding;  $l_s \in \mathbb{R}^m$  – numerical presentation of these words in context, having a fixed length.

For keywords extraction the TF-IDF (term frequency-inverse document frequency) method is used.<sup>3</sup> From each comment on news  $M$  keywords are extracted which theoretically are the most important in the context of this comment. These  $M$  keywords  $\tilde{w}_j \in \tilde{W} \subseteq \mathbb{N}^n$ ;  $j = \overline{1, \dots, M}$  get into the dictionary. There each word is assigned a numeric vector  $\tilde{l}_j \in \mathbb{R}^m$ ;  $j = \overline{1, \dots, M}$ , using Word Embedding technique. Therefore, each comment is represented as a vector  $\tilde{w} = (w_{1k}, w_{2k}, \dots, w_{Mk})$ ;  $k$  – sequence number of the explored comment. At the output of the dictionary a vector  $\tilde{l} = (l_{1k}, l_{2k}, \dots, l_{Mk})$  is obtained.

This vector is submitted to the input of the classifier, which evaluates the degree of emotionality of the comment associated with this vector. As a classifier, the method of artificial neural networks is used. The classifier can be presented as a function  $F: \mathbb{R}^n \rightarrow \mathbb{R}$ :

$$F(\tilde{l}) = r_k, \quad (2)$$

$r_k \in \mathbb{R}$  – valuation of comment  $k$ ;  $r_k \in [0, 2]$ , where the most emotionally positive comments have value 0, the most emotionally negative have value 2.

It is known that comments can have different nesting levels. For example, comments at level 0 are reviews for news or opinions expressed by readers about the

described event. Comments at level 1 are comments on the comments of the previous level, etc. Obviously, the valuations of the comments of higher levels should have a greater impact on the overall value of the level of tension. Thus, it will be considered that with each subsequent level of nesting the valuation of the comment is reduced by half. The average rating of perception of the news is calculated by the formula of the weighted average:

$$\bar{r} = \frac{\sum_k (r_k - 1) \left(\frac{1}{2^n}\right)_k}{\sum_k \left(\frac{1}{2^n}\right)_k} + 1, \quad (3)$$

where  $n$  is the level of nesting. The average estimate characterizes now not the comment itself, but the news, that was commented.

### **Subsystem for Data Aggregation**

It is designed for counting an estimation of influence of the article on people taking into account the data from first two subsystems.

The resulting block of data aggregation gets such data:

- $\bar{r}$  – average value of perception of the news by readers;
- $N$  – number of explored comments;
- $\Delta\bar{p}$  – vector of increase of the number of news reposts.

Data aggregation and derivation of the resulting value  $S \in [0; 1]$  can be represented in a function  $A: \mathbb{R}^3 \rightarrow \mathbb{R}$ :

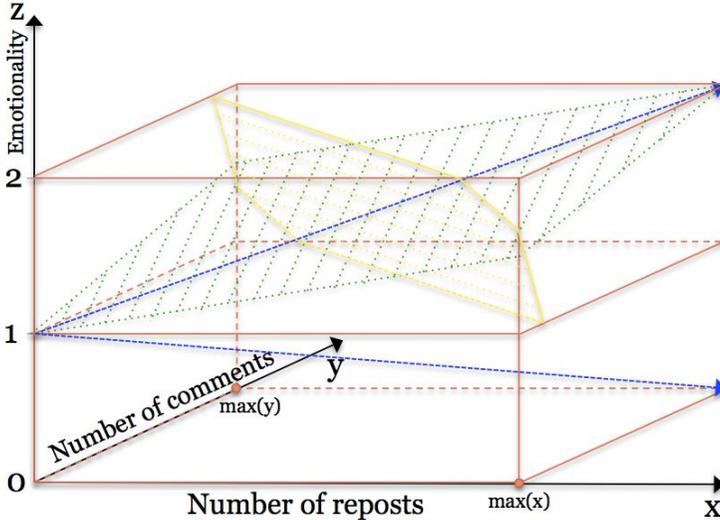
$$A(\bar{r}, N, \Delta p_i) = S. \quad (4)$$

$\bar{r} \in [0; 2] \subseteq \mathbb{R}; N \in \mathbb{N}; \Delta p_i \in \mathbb{N}$ .

It should be noted that regardless of the length of the explored time intervals, the length of the vector of shares, which comes to the data aggregation block, is fixed and predetermined. This vector represents the change in the activity of article shares for a shorter period of time. In this way, the ability to study the change of the popularity of the article in dynamics is achieved.

The resulting three characteristics don't have obvious interdependencies. In order to take into account the influence of each of them equally, the values of these characteristics were laid on the coordinate axes  $Ox$ ,  $Oy$  and  $Oz$ . Drawing the lines parallel to the coordinate axes from the points of the conditional maximum of the characteristics, a right quadrilateral prism is obtained (Figure 2).

The average estimations of emotions are laid on the vertical axis  $Oz$ . It was decided to evaluate the emotionality on the interval  $[0, 2]$ . Value 0 shows maximum positive emotionality of comments, value 2 – maximum negative emotionality. Such marginal estimates are unlikely and will probably indicate unpopularity of the news, the rejection of the seriousness of its content or the author of the post.



**Figure 2: Prism, built on the characteristics, which are aggregated in the resulting block of the system.**

Value 1 is more realistic. It shows the neutrality of the comments or the balance of positive and negative reviews. The maximum number of comments and the maximum number of reposts are determined experimentally.

The prism was divided into two half-prisms by the plane  $z = 1$ , which is parallel to the plane of quantitative characteristics. It's the plane of neutrality. All news which ratings got into the lower half-prism were perceived more positively. The ones which ratings got into the upper half-prism – more negatively.  $(0; 0; 1)$  is the initial point of two vectors indicating the direction of increasing the level of emotionality. The vector in the lower half-prism indicates the direction of increase in positive emotions. The vector in the upper half-prism – in negative emotions.

The vector of the upper half-prism has coordinates  $(\max(x); \max(y); 1)$ ,  $\max(x)$  – maximum number of reposts,  $\max(y)$  – maximum number of comments. Thereby the equation of the set of planes perpendicular to this vector, has the form:

$$\max(x) \cdot x + \max(y) \cdot y + z + d = 0, \tag{5}$$

$d \in [-(\max(x))^2 - (\max(y))^2 - 2; -1]$  – constant term of the equation.

The vector of the lower half-prism has coordinates  $(\max(x); \max(y); -1)$ . The equation of the set of planes perpendicular to this vector has the form:

$$\max(x) \cdot x + \max(y) \cdot y - z + d = 0, \tag{6}$$

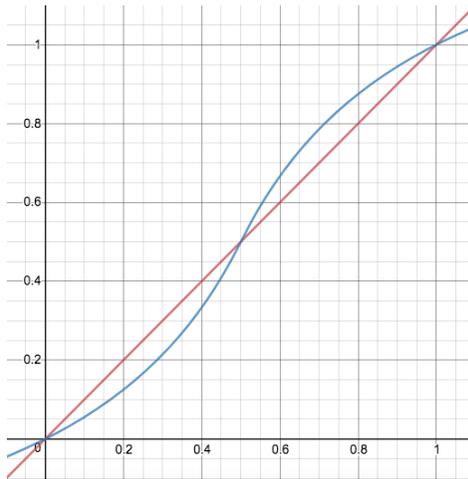
$d \in [-(\max(x))^2 - (\max(y))^2; 1]$ .

To each point inside the prism an estimate on the interval  $[0; 1]$  is given. These points show the level of tension in the society provoked by news, taking into account all quantitative characteristics. Values 0, 0,5 and 1 correspond to the points  $(\max(x); \max(y); 0)$ ,  $(0; 0; 1)$  and  $(\max(x); \max(y); 2)$  respectively. All other possible values on the interval  $[0; 1]$  are obtained on the basis of analysis of the plane (the section of the prism), to which the point belongs. Obviously, the points which get into the lower half-prism have estimates on the interval  $[0; 0,5]$ , and those which get into the upper half-prism have estimates on the interval  $[0,5; 1]$ .

Each plane is determined by its coefficient  $d$ , which is the coefficient of translation. Points which get on one plane have the same resulting estimate. There is a direct relationship between the intervals of the coefficient  $d$  and the intervals of resulting estimates. The estimates will increase from the point  $(0; 0; 1)$  to the point  $(\max(x); \max(y); 2)$  and decrease from the point  $(0; 0; 1)$  to the point  $(\max(x); \max(y); 0)$ .

In order to increase the influence of the estimate of the news, if it tends to negative (greater than 0,5) or to reduce the influence of the estimate, if it tends to positive (less than 0,5), the transformed rational sigmoid function can be applied (Figure 3):

$$f = \frac{(x - 0,5)}{|x - 0,5| + 0,5} + 0,5. \quad (7)$$



**Figure 3: Graph of transformed rational sigmoid function (7). Graph of function  $y = x$  is also shown.**

### Program Realization

The program realization was developed using MATLAB programming language and concerns the subsystem for evaluating emotionality of comments on news. The

neural network was created with the help of the framework Neural Network Toolbox.

For quality check of the program 783 comments were collected from the social network Facebook pages of the 8 most popular news resources in Ukraine. On each of the resources the reaction of people on one news was explored. Comments to each of the news were saved in text files. The structure of all files is the same. Each row keeps data about one comment and consists of 4 elements, separated by a “|” character. The first element is the order number of the comment within one news resource. The second element is the level of comment’s nesting (0 or 1). The third element is the level of emotionality of the comment (on the interval [0, 2]). The fourth element is the text of the comment. The level of emotionality of each comment was determined manually by the experts.

All comments were divided into two groups: 683 comments were prepared for training the classifier, 100 comments – for testing. Accuracy was determined by the mean-square deviations of the results obtained at the output of the neural network from the actual results. For the word2vec method the window with the length 5 was chosen. It’s the longest possible window because there are comments which have only 5 words. The length of the vectors which correspond to each word in the dictionary and the optimal number of neurons on the hidden layer of neural network were chosen experimentally. The feedforward neural network with one hidden layer was used. There were conducted numerous experiments for different combinations of numbers of elements in the word vectors and numbers of neurons on the hidden layer. For each combination 10 experiments were performed and the arithmetic mean of their results was counted. The results of these experiments are presented in the Figure 4.

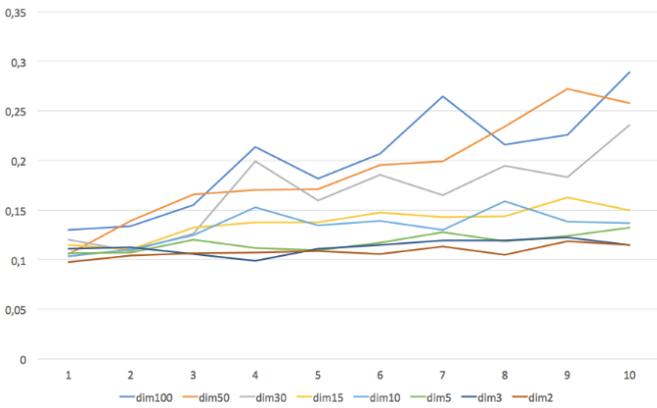


Figure 4: Graphs of the dependencies of mean-squared deviations from the number of neurons on the hidden layer of neural network for the word vectors of different length.

Thus, we have concluded that the increasing number of neurons and the length of word vectors lead to decreasing of accuracy of the results.

## Conclusions and Future Work

The architecture and the mathematical model of the system of operative determination of the level of tension in the society were described in detail. Also, the main schemes, graphs and formulas that allow to understand the principles of the system work were presented. Besides, the brief description of program realization and its results were presented.

Future work will focus primarily on the development of a software implementation of two remaining subsystems and on the integration of these subsystems into one system. Further research and classifier trainings will be conducted with more comments which tension levels will be determined automatically. The work will also be continued to improve the quality of the system by more careful parameter selection.

## Acknowledgement

This study is funded by the NATO SPS Project CyRADARS (Cyber Rapid Analysis for Defense Awareness of Real-time Situation), Project SPS G5286.

## References

- <sup>1</sup> Natali Ruchansky, Sungyong Seo, and Yan Liu “CSI: A Hybrid Deep Model for Fake News Detection,” *Proceedings of the 26th ACM International Conference on Information and Knowledge Management (CIKM)*, Singapore, 2017, available at <https://arxiv.org/abs/1703.06959>.
- <sup>2</sup> Tomas Mikolov, Quoc V. Le, and Ilya Sutskever, “Exploiting Similarities among Languages for Machine Translation,” *arXiv: Computation and Language* (2013), available at <https://arxiv.org/abs/1309.4168>.
- <sup>3</sup> Christopher Manning, Prabhakar Raghavan, and Hinrich Schütze, “An Introduction to Information Retrieval,” *Natural Language Engineering* 16, no. 1 (2010): 100–103.

## About the Authors

Maksym **Shchokoliev** holds B.Sc. (2017) and M.Sc. (2019) in applied mathematics from the National Technical University of Ukraine “Igor Sikorsky Kyiv Polytechnic Institute.” Currently he is data scientist and programmer in the Department of Applied Mathematics of the same university.

Violeta **Tretynyk** holds an M.Sc. degree in theoretical physics from Taras Shevchenko National University of Kyiv and PhD degree in mathematical physics. Currently she is Professor in the Department of Applied Mathematics at the National Technical University of Ukraine “Igor Sikorsky Kyiv Polytechnic Institute.”