

Publication Date: 30 June 2024

Archs Sci. (2024) Volume 74, Issue 3 Pages 102-107, Paper ID 2024317.  
<https://doi.org/10.62227/as/74317>

# New Exploration of Legal Research Based on Big Data Technique

Song Zhang<sup>1,\*</sup>

<sup>1</sup>Jilin Police College Law Department, Changchun 130117, Jilin, China.

Corresponding authors: Song Zhang (e-mail: zhangs9797@163.com).

**Abstract** The transformation of the "computing paradigm" of law is conducive to removing the fragmentation problem of "law + information technology". The development of computational law requires the construction of an interdisciplinary academic community to enhance the level of legal scientific research and modern legal capabilities in the era of ubiquitous computing in many countries. Based on the knowledge graph in the judicial field, this paper designs an automatic question answering system. The judicial domain corpus database constructed in this paper can solve the shortage of public data sets in the current judicial domain to a certain extent. Under the big data environment, we build a knowledge graph and automatic question answering to provide reference for the research of Chinese knowledge graph construction technology and automatic question answering technology in other fields. Experiments show that the proposed network can ensure the efficient extraction of entities and relationships in legal text information.

**Index Terms** legal studies, law and technology, converged governance

## I. Introduction

In the Western tradition, law and computing have always been interdependent, and the legal culture is also known as the computing culture. Although "computation" is a cognitive tool discovered and often used by human civilization very early, our scientific understanding of computing has been in the process of continuous deepening [1], [2]. Early "calculations" were mainly mathematical operations such as addition, subtraction, multiplication, and division in the pure mathematical sense that are most common in daily life. This kind of calculation uses rope knots, abacus, calculus or simple calculation experience of human beings to obtain pure mathematics through four algorithms. Mathematical conclusion [3]. With the development of computational science, "computation" began to be applied to the fields of humanities and social sciences during the Renaissance. [4], [5] There are many controversial issues that lack accurate answers in the humanities and social sciences represented by philosophy. Bacon and Descartes reflected on this in their natural philosophy and proposed a scientific method to understand social truth through deductive calculation. This inclusive "computation" mainly referred to computational logic in the legal field at the time. On this basis, Hobbes put forward the classic thesis that "reasoning is calculation" and pursued the precise and scientific rules of social dispute resolution, which is consistent with the final development direction of law in many aspects.

With the continuous advancement of computing tools, people's application and imagination space for "computing" is also expanding. Early manual computing tools, mechanical

computing tools, and electromechanical computing tools are gradually withdrawing from the stage of history [6]. Electronic computers, parallel and distributed computing tools Computing, high-performance cluster computing, and cloud computing have become increasingly popular, and new types of computing, such as quantum computing, social computing, biological computing, and ocean computing, will continue to mature [7]. The production and life of society is evolving from simple computing to complex computing, limited computing to ubiquitous computing. With the expansion of the computable range and the introduction of the theory of ubiquitous computing, famous computer experts in the United States once predicted that we would enter an era of "ubiquitous computing" in which computing is ubiquitous [8], [9]. At this time, a kind of "computationalism" has emerged that points out: the universe is a huge automaton containing computational logic, and the human brain is a kind of Ultra-complex neural network systems, from the universe to the human brain, can be understood and analyzed through cognitive computing [10]. This idea reflects the status of "computing" in the development of modern science. With the development of computing science and technology, the computing society itself has also become a special research object. We need to study the relational structure and behavioral norms of computing space from the perspective of ontology and epistemology, to build an orderly computing society.

At present, the fifth-generation mobile communication technology is in a stage of rapid development, and more and more people choose to use this technology, which has brought

an explosive increase in the amount of information search. Due to the huge amount of information on the Internet, it is not easy for traditional search engines to obtain valuable knowledge quickly and accurately, and the efficiency of users obtaining effective information through traditional search engines is very low. In view of the above problems, the research on automatic question answering technology based on domain knowledge graph becomes more and more important. Currently, most question answering system data comes from open domain knowledge bases. It is not easy to acquire expertise in some specific fields [11]. Facing the need for fast and accurate query, the traditional database storage and retrieval methods in the judicial field have many shortcomings, it is difficult to effectively integrate and correlate massive legal text data, and it is difficult to fully mine the value of legal text data. The emergence of knowledge graph technology provides the possibility for the effective integration and association of massive legal text data. With the development of judicial informatization, more and more legal information is stored in the form of text data. Traditional database storage methods cannot effectively integrate legal text data information. Knowledge graph is convenient for machines to understand data semantic information. Intelligent analysis and mining of data requires machines to understand data semantic information [12]. Therefore, constructing a dataset that meets the requirements and is correct is an urgent problem to be solved. In addition, due to the diversification of legal text features and the low degree of structuring, the entity relationship extraction technology in legal text needs to be paid more attention. It is also one of the research purposes of this paper to store knowledge after entity relation extraction and put knowledge graph into practical application.

Therefore, for the above problems, this paper will study the key technologies in the process of building knowledge graphs in the judicial field, and at the same time study the automatic question answering system based on knowledge graphs to reflect the application value of the research results. In this process, it mainly focuses on legal text named entities Identify the model to optimize, study the rule-based relationship extraction method, construct the domain knowledge graph, construct, and design the automatic question answering system. The construction of judicial data sets can provide data support for the research on knowledge graph related technologies in this field. Exploring the key technologies of knowledge graph construction in the judicial field has certain scientific reference value in the field of knowledge graph research. From the perspective of application value, automatic knowledge graph-based the question answering system can provide a channel for people who need to provide judicial assistance to quickly query and obtain legal information through natural language.

The contributions of this paper are mainly in three aspects:

- 1) we explore the new progress of law from the perspective of computational law;
- 2) we combine big data and law to construct a knowledge graph;

- 3) we establish the correlation coefficient among each law based on the knowledge graph.

## II. Related Work

### A. From Forensic Metrology and Legal Informatics to Computational Law

The research group first adopted legal informatics in 1970 to summarize the new discipline formed by the application of information technology in the field of law [13]. Legal informatics is an independent discipline in Germany, and its content includes not only the application of information technology in the legal field, but also the study of computer-related legal issues. Legal informatics has since made great strides in countries and regions such as Northern Europe and the United States. However, some scholars in the United States searched the relevant research from 1997 to 2005 through the paper database and found that the relevant literature that can be obtained using legal informatics as the key word is very scarce [14]. But this is not to say that there is no relevant research, but that the relevant research does not use the concept of legal informatics. This result is not unrelated to the fact that there are no professional legal informatics associations, no professional legal informatics journals, and only a few law schools in the United States that offer courses in the name of legal informatics [15].

### B. Establishment and dissemination of the concept of computational law

Computational law is not a domestic concept, but a new concept developed because of legal informatics. As early as 1977, papers in the field of legal informatics in Sweden were entitled "computing law", which believed that computational law would become a new discipline, and its main content was legal education, legal information retrieval, legal database, legal information security and related personal rights protection supported by computing technology [16]. Since then, this concept has been cited in some articles on legal research paradigms and information law, but it has not received enough attention for a long time. Stanford University has played an active role in the modern development of the concept of computational law. The International Forum on future law held by Stanford University law school since 2013 has become an exchange center for the theory and practice of computational law, which has widely promoted the dissemination of computational law throughout the world [17], [18]. Since then, international conferences and courses with the theme of computational law have appeared all over the world, and more and more scholars have written monographs on computational law and identified their professional field as "computational law".

## III. Methods

### A. CRF model

Conditional random fields are usually used to solve sequence labeling problems in machine learning. A typical sequence

labeling problem is to label sequences with parts of speech. CRF obtains the global optimal tagging sequence by considering the relationship between adjacent part-of-speech tags. The working principle of CRF is shown in Figure 1.

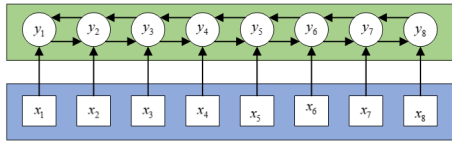


Figure 1: Crf working principal diagram

The mathematical description of CRF is, let  $X, Y$  be random variables,  $Y$  is the corresponding output given  $X$ ,  $P(X, Y)$  is the conditional distribution of  $Y$  given  $X$ , if the random variable  $Y$  constitutes a Markov If it is a random field, then  $P(X, Y)$  is called a conditional random field. CRF is a special case of Markov random field, which is represented by  $G = (V, E)$ , as shown in Eq. (1):

$$P(Y_v | x, Y_w, w \neq v) = P(Y_v | X, Y_w, w \sim v), \quad (1)$$

where  $w - v$  represents all nodes  $w$  that have edge connections with  $v$ ,  $w \neq v$  represents other nodes connected to node  $v$ ,  $Y_w$  represents a random variable with  $v$ , and  $Y_w$  represents a random variable of  $w$ . If  $X = (X_1, X_2, \dots, X_n), Y = (Y_1, Y_2, \dots, Y_n)$ ,  $X$  is the given word and  $Y$  is the output part of speech, as in Eq. (2) shown.

$$\begin{aligned} P(Y_i | X, Y_1, \dots, Y_{i-1}, Y_{i+1}, \dots, Y_n) \\ = P(Y_i | X, Y_{i-1}, Y_{i+1}), \end{aligned} \quad (2)$$

where  $i = 1, 2, \dots, n$ . In the case of random variable  $X$  taking value  $x$ , the conditional probability form of random variable  $Y$  taking value  $y$  is shown in Eq. (3).

$$\begin{aligned} P(y | x) \\ = \frac{1}{Z(x)} \exp \left( \sum_{i,k} \lambda_k t_k(y_{i-1}, y_i, x, i) + \sum_{i,l} \mu_l s_l(y_i, x, i) \right). \end{aligned} \quad (3)$$

Among them,  $z(x)$  represents as shown in Eq. (4):

$$\begin{aligned} z(x) \\ = \sum_y \exp \left( \sum_{i,k} \lambda_k t_k(y_{i-1}, y_i, x, i) + \sum_{i,l} \mu_l s_l(y_i, x, i) \right). \end{aligned} \quad (4)$$

In Eq. (2) and Eq. (3),  $t_k$  represents the transition feature, that is, the feature function of the edge in the undirected graph, which depends on the current position and the previous position,  $s_l$  represents the state feature, that is, the corresponding feature function of the node in the undirected graph, there is a

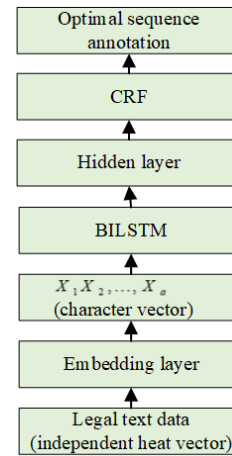


Figure 2: BILSTM-CRF-LER overall model framework

strong dependency on the current position.  $\lambda_k$  and  $\mu_l$  represent the corresponding weights of the transition feature and the state feature, respectively, and  $Zx$  is the function used for specification. It can be seen from Eq. (4) that the function performs a summation operation on all possible output sequences.

### B. BILSTM-CRF-LER model

In this study, a BILSTM-CRF-LER model for the judicial field is established. The model is mainly composed of two sub-modules, namely the BILSTM module and the CRF module. The legal text data is preprocessed and then passed through the embedding layer as the input of the deep learning model. The overall model framework is shown in Figure 2.

The general process of processing legal text data by BILSTM-CRF-LER model is divided into the following three steps:

- 1) Convert one-hot vector to character vector through character embedding matrix.
- 2) The character vector is used as the input of the BILSTM module.
- 3) The various label scores output by the BILSTM module are processed and input into the CRF module.

Since the BILSTM module cannot guarantee that the output tag sequence must be the optimal solution, the CRF module will obtain the optimal tag sequence according to the relationship between adjacent part-of-speech tags. Therefore, this study proposes to add a CRF model on top of the BILSTM model to improve the entity recognition performance of the model.

### C. BILSTM-CRF-LER model design

The BILSTM-CRF-LER model proposed in this study is divided into three layers: embedding layer, BILSTM layer and CRF layer. The character vector processed by the embedding layer is used as the input of BILSTM, and the predicted label corresponding to each character output by the CRF is the output result. The schematic diagram of the model is shown in Figure 3.

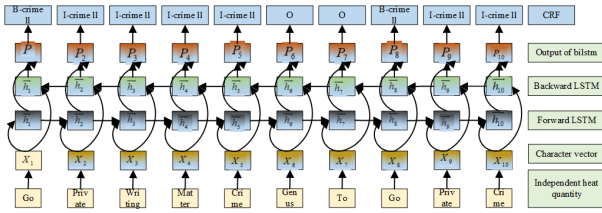


Figure 3: Schematic diagram of the BILSTM-CRF-LER model

The first layer of the BILSTM-CRF-LER model is the embedding layer. a sequence of  $n$  words that represents  $X = (X_1, X_2, \dots, X_n)$ , where  $X_i$  is the number of the  $i$  two word of the sequence in the dictionary map. Initialize the character embedding matrix, and use the pre-trained embedding matrix to convert the one-hot vector into a lower-dimensional vector, Denoted by  $W_i \in R_d$ , the dimension of the character embedding matrix is  $d$ , and the sequence of character vectors  $(X_1, X_2, \dots, X_n)$  is obtained.

Finally, dropout is used to suppress the overfitting problem. The BILSTM layer adds a last linear layer, The purpose is to map the output of the hidden layer from the initial dimension  $m$  to dimension  $k$ , so as to express the meaning of label features. In this study, it represents the total number of labeled labels, so as to obtain the automatically extracted sequence feature matrix.  $P = (P_1, P_2, \dots, P_n) \in p_{k \times n}$ , where  $p_{ji}$  can be considered as the score of the sequence feature  $p$  from the  $i$  to label to the  $j$  to label. The obtained results are directly processed by the logistic regression function, which is equivalent to classifying each position separately, so that the labeled context information cannot be used. Therefore, the sequence feature matrix  $p$  needs to be further processed by the CRF layer. The third layer of the BILSTM-CRF-LER model is the CRF layer, The purpose of setting the CRF layer in this paper is to finally get the optimal sequence annotation according to the contextual semantic information through the sequence annotation obtained in the above steps. In this layer, the  $(k+2) \times (k+2)$  labels are used to transfer the matrix  $A$ , where  $A_{ji}$  represents the transfer score from the  $i$  to label to the  $j$  to label. Since the CRF layer needs to add markers at the beginning and end of the sentence, the  $k$  value is added by 2. If a marker sequence  $y = (y_1, y_2, \dots, y_n)$  The length is consistent with the length of the recorded sequence  $x$ , then the label of the label sequence  $y$  in the model is equal to the label of the sequence  $x$ , as shown in Eq. (5).

$$\text{score}(x, y) = \sum_{i=1}^n p_i, y_i + \sum_{i=1}^{n+1} A_{y_{i-1}, y_i}. \quad (5)$$

It can be seen from Eq. (5) that the result of the label sequence is related to both the second layer model and the third layer, and then the logistic regression function is called to normalize the probability, as shown in Eq. (6).

$$P(y | x) = \frac{\exp(\text{score}(x, y))}{\sum y \exp(\text{score}(x, y'))}. \quad (6)$$

The model uses the maximum log-likelihood function during training to achieve the log-likelihood  $(x, y^x)$  of the training samples, as shown in Eq. (7).

$$\log P(y^x | x) = \text{score}(x, y^x) - \log \left( \sum_{y'} \exp(\text{score}(x, y')) \right). \quad (7)$$

The optimal solution of the model is calculated by the Viterbi algorithm, as shown in Eq. (8).

$$y^* = \text{argmaxScore}(x, y'). \quad (8)$$

## IV. Experiments

### A. Dataset

The experimental dataset uses 20,349 criminal law-related descriptive paragraphs on Hualu.com as the judicial text corpus. Since most of the relevant case data on this website are provided by tourists or volunteers, it is not possible to obtain valid and processable case texts. Therefore, preliminary processing of the obtained data is required. After preliminary data screening cleaning, and fusion, 8105 valid descriptive paragraph data were obtained, of which 80% of the judicial text corpus was used as the training set, and 20% of the judicial text corpus was used as the test set.

### B. Experimental environment and parameter settings

In the experiment, the TensorFlow framework was selected to build the BILSTM-CRF-LER model. The specific experimental environment and configuration are shown in Table 1 below. The parameter settings of the model in this study are shown in Table 2 below.

Category	To configure
Frame	Tensorflow1.13
Operating system	Windows 10
Processor	AMD Ryzen 5 3600 6-core 3.60 GHZ
Memory	16GB
Programing language	Python3.8

Table 1: Configuration of experimental environment

Parameter	Meaning	Parameter setting value
LSTM-dim	Number of LSTM neurons in hidden layer	299
Batch-size	Data size of each batch	17
Epoch	Number of iterations	100
learning rate	Learning rate	0.2
dropout	Copout-1	0.6

Table 2: Model parameter setting

### C. Evaluation indicators

To verify the performance of the model proposed in this study, the general evaluation indicators precision rate, recall rate and F-value is used to evaluate the model performance in this experiment.

**D. Experimental results and comparative analysis**

To verify the validity of the BILSTM-CRF-LER model proposed in this study and the necessity of each module in the model, this paper selects several models used in Chinese named entity recognition related research and the BILSTM-CRF-LER proposed in this paper. The LER models are used for comparative experiments, which are the LSTM model, the BILSTM model and the LSTM-CRF model. Named entities are identified through the LSTM model. This deep learning model can better capture long-distance information and avoid the problem of gradient explosion to a certain extent, but it cannot encode information from the back to the front, nor can it be used for context. association. Aiming at the defects of the LSTM model, the BILSTM model is used to identify entities, which better captures the bidirectional semantics. The LSTM model combined with the CRF algorithm is used to identify military entities, which effectively solves the defect that the LSTM model cannot correlate the context and achieves a high accuracy rate. This paper uses the above three models and the BILSTM-CRF-LER model proposed in this paper to conduct comparative experiments on legal texts. The experimental results of named entity recognition are shown in Table 3. The results of the comparison experiments are shown in Table 4.

Entity	Accuracy	Recall	F value
First class crime	69.59%	75.35%	72.36%
Secondary charge	76.16%	78.22%	77.18%
Concept	73.34%	71.90%	72.61%
Features	80.76%	82.17%	71.46%
Legal provisions	74.37%	76.59%	75.46%
Sentencing standard	72.14%	73.60%	72.86%
All entities	79.74%	81.59%	80.66%

Table 3: Experimental results of BILSTM-CRF-LER model

Entity	Accuracy	Recall	F value
LSTM	58.56%	68.47%	63.12%
BILSTM	77.17%	73.24%	76.70%
LSTM-CRF	64.71%	67.32%	65.98%
BILSTM-CRF-LER	79.74%	81.59%	80.65%

Table 4: Comparison of different models

As can be seen from Table 4, the precision rate of the BILSTM-CRF-LER model proposed in this study is 79.73% and the recall rate is 81.58%. The value of the model is about 14% higher than that of the LSTM-CRF model, and the precision rate value is improved about 15%. It shows that the BILSTM network can greatly improve the accuracy of entity recognition by automatically extracting sequence features and using the context information of words.

The operating system version we use to build the judicial domain knowledge graph is macOS 10.15.7. Due to the performance limitations of the neo4j server and the machine used in the experiment, only a maximum of 1500 nodes can be displayed. The constructed judicial domain knowledge graph is constructed as shown in the Figure 4 shown. The entity part in Table 4 is: green represents the first-level crime node, brown represents the second-level crime node, blue represents the

concept node, pink represents the feature node, brown represents the legal clause node, and red represents the sentencing standard node.

The relationship part: belong to is the affiliation of the first level and second-level crimes, chariot is the relationship between features and second-level crimes, concept of is the relationship between crimes and concepts, law of is the relationship between legal provisions and crimes, and puny of is the corresponding relationship between crimes and sentencing standards.

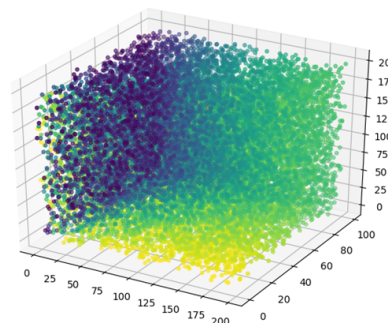


Figure 4: Partial knowledge map in the judicial field

Taking the crime of endangering public security as an example, the knowledge graph of the affiliation between the first and second-level crimes is displayed, as shown in Figure 5.

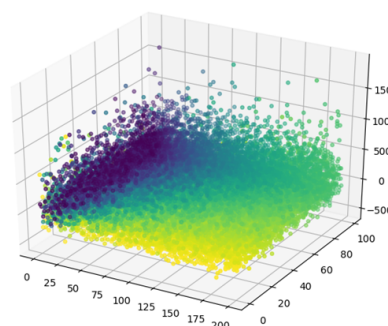


Figure 5: Partial knowledge map of subordination between level I and level II crimes

Taking the crime of sheltering and condoning underworld organizations as an example, the affiliation between this second-level crime and the first-level crime, the corresponding cases of this second-level crime, the concept of this second-level crime, the relevant legal provisions of this second-level crime, and the description of this second-level crime. The knowledge graph constructed by the sentencing standards is shown in Figure 6. The solid part in Figure 6 is: orange represents the second-level crime node, pink represents the feature node corresponding to the second-level crime, green represents the second-level crime concept node, blue represents the first-level crime node to which the second-level crime belongs, and brown represents the second-level crime. The sentencing

standard node of the first-level crime, the red represents the relevant legal provisions of the second-level crime; the relationship part: belong to is the affiliation of the first-level and second-level crimes, chariot is the relationship between the feature and the second-level crime, concept of is the crime-concept relationship, and law of is the legal provisions The relationship with the crime, puny\_ of is the corresponding relationship between the crime and the sentencing standard.

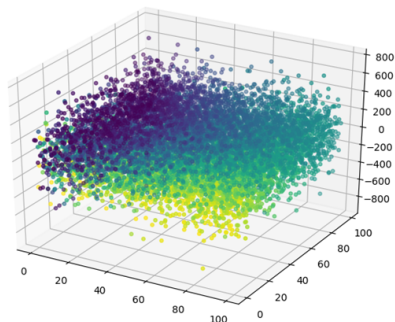


Figure 6: knowledge map display of a secondary charge and its corresponding relationship

## V. Conclusion

In this work, we first introduce the overall design process of building a knowledge graph in the judicial field. Then, we use a rule-based method to extract the relationship between entities. In this way, the knowledge graph in the judicial field can be stored and displayed. Next, we propose a BILSTMCRF-LER model to build the accurate relationship of each sub-unit. Compared with the LSTM model, BILSTM model, LSTM-CRF model, BILSTMCRF-LER model has a good recognition accuracy.

## References

- [1] Xu, Y., Tao, Y., Zhang, C., Xie, M., Li, W., & Tai, J. (2022). Review of digital economy research in China: a framework analysis based on bibliometrics. *Computational Intelligence and Neuroscience*, 2022(1), 2427034.
- [2] Klingemann, J., Witaj, P., Lasalvia, A., & Priebe, S. (2022). Behind the screen of voluntary psychiatric hospital admissions: a qualitative exploration of treatment pressures and informal coercion in experiences of patients in Italy, Poland and the United Kingdom. *International Journal of Social Psychiatry*, 68(2), 457-464.
- [3] Zhang, J., Yang, F., Xu, T., & Zheng, C. (2020). Exploration and practice of training excellent engineers in safety engineering major of Longdong University. *World Scientific Research Journal*, 6(4), 346-352.
- [4] Anderson, C., Pooley, J. A., Mills, B., Anderson, E., & Smith, E. C. (2020). Do paramedics have a professional obligation to work during a pandemic? a qualitative exploration of community member expectations. *Disaster Medicine and Public Health Preparedness*, 14(3), 1-24.
- [5] Milanov, A., & Penchev, G. (2020). The need for establishing a new United Nations body to protect Earth from back contamination and outer space from forward contamination. *International Journal of Criminology and Sociology*, 9(2020), 925-930.
- [6] Cotula, L. (2020). Investment contracts and international law: charting a research agenda. *European Journal of International Law*, 31(1), 353-368.
- [7] Khan, A., Chen, L. R., & Hung, C. Y. (2021). The role of corporate social responsibility in supporting second-order social capital and sustainable innovation ambidexterity. *Sustainability*, 13(13), 6994.
- [8] Ateme, M. E. (2021). Developing marine and coastal resources in Nigeria: prospects and challenges. *Maritime Technology and Research*, 3(4), 335-347.
- [9] Ruhaeni, N. (2020). Direct international responsibility of non-governmental entities in the utilization of outer space. *Padjadjaran Jurnal Ilmu Hukum (Journal of Law)*, 07(1), 102-120.
- [10] Salmivalli, C., Laninga-Wijnen, L., Malamut, S., & Garandeau, C. (2021). Bullying prevention in adolescence: solutions and new challenges from the past decade. *Journal of Research on Adolescence: The Official Journal of the Society for Research on Adolescence*, 31(4), 1023-1046.
- [11] Yan, F., Yang, C., & Hua, L. I. (2022). Methodological research in the field of civil aviation safety. *Asian Agricultural Research*, 14(5), 8.
- [12] Hu, Y., & Sun, P. (2021). Practical exploration of new urbanization in China—a case study of Shichuan River (urban section) comprehensive improvement project. *IOP Conference Series: Earth and Environmental Science*, 692(4), 042054 (7pp).
- [13] Ma, M. (2021). Research on new methods and technologies of engineering and environmental geophysical exploration based on internet of things. *Journal of Physics: Conference Series*, 1744(2), 022058 (5pp).
- [14] Peprah, D., Bangura, J., Vandi, M., Thomas, H., Dea, M., Schneider, A., & Chittenden, K. (2021). Social and Political Dimensions of Disseminating Research Findings on Emerging Zoonotic Viruses: Our Experience in Sierra Leone. *Global Health: Science and Practice*, 9(3), 459-466.
- [15] Ren, X., Ahmed, I., & Liu, R. (2023). Study of Topological Behavior of Some Computer Related Graphs. *Journal of Combinatorial Mathematics and Combinatorial Computing*, 117, 3-14.
- [16] Nazeer, S., Sultana, N., & Bonyah, E. (2023). Cycles and Paths Related Vertex-Equitable Graphs. *Journal of Combinatorial Mathematics and Combinatorial Computing*, 117, 15-24.
- [17] Zeng, Y., & Chu, B. (2023). The Appropriate Scale of Competition Between Online Taxis and Taxis Based on the Lotka-Volterra Evolutionary Model. *Journal of Combinatorial Mathematics and Combinatorial Computing*, 117, 25-36.
- [18] Guo, Q. (2023). Minimizing Emotional Labor through Artificial Intelligence for Effective Labor Management of English Teachers. *Journal of Combinatorial Mathematics and Combinatorial Computing*, 117, 37-46.

...