



# Speech network analysis and anomaly detection based on FSS model

Xinhui Yan<sup>1</sup>

Received: 9 January 2020 / Accepted: 22 June 2020 / Published online: 30 June 2020  
© Springer Science+Business Media, LLC, part of Springer Nature 2020

## Abstract

Speech network analysis and anomaly detection based on the FSS model is analyzed in this research. Aiming at the problem of parallel detection and processing of massive data in distributed intrusion detection, a data segmentation algorithm based on capability and load is proposed. The algorithm evaluates the capabilities and actual load of the nodes, weighs the actual state of each node and the data distribution relationship in the cluster, and allocates more data to be processed to nodes with strong data processing capabilities and light loads. High-order statistics are used to describe the intrusion characteristics of persistent attacks in the link layer of the speech sensor networks, and vector quantitative decomposition is used to analyze the fusion characteristics of advanced persistent intrusion symbols in mobile terminals. Different terminals are estimated based on the machine learning algorithms, and the FSS model is integrated to achieve the comprehensive analysis of the speech analytic models. The experiment compared with the state-of-the-art methods have proven the efficiency of the framework. The detection accuracy is higher than the latest methodologies.

**Keywords** Speech network · Anomaly detection · FSS model · Speech technology · Data mining

## 1 Introduction

As machine learning algorithms have gradually been widely used in intrusion detection systems, especially neural network-based intrusion detection algorithms, due to their unique advantages, they have been widely used in intrusion detection systems. Network anomalies can significantly damage and reduce the quality of the network services, and they usually appear as anomalies in traffic. Therefore, real-time monitoring and management of network traffic is of great significance for timely detection of the network anomalies and improvement of network reliability and availability. The anomaly detection idea is to first establish a feature model of the normal data, and determine whether an attack occurs based on whether the current data fits the model effectively during the detection. Since anomaly detection does not require an attack to establish a knowledge base, it is only necessary to ensure that most of the training data is normal data, so the actual application has more practical promotion

and engineering application value (Faizin et al. 2018; Sailunaz et al. 2018; Muckenhirn et al. 2018).

Based on the literature review, the existing detection models need to face with the following challenges. (1) Because most of the existing anomaly detection systems use one or more single network feature vectors as the basis for learning and judging, the description of network data flow anomaly is also relatively thin. Secondly, fewer network feature vectors are selected in the cooperative operation of intrusion detection system, which may affect the scalability of the detection system. (2) Because the complex cloud environment is an interconnected network without centralized management of the multiple management domains, but the intrusion detection requires that all detection systems operate cooperatively, so it is very important to provide shared data as the main content of cooperative operation. (3) Since the above intrusion detection method only detects one or several feature vectors in the network data stream, and the selected feature vectors have no specific attack meaning, when the detection system alarms, it only knows that some feature vectors in the network are abnormal, but it cannot judge what kind of attack has occurred. In the Fig. 1, we present the normal system for the references (Savchenko and Savchenko 2019; Ekber et al. 2019; Shadiev et al. 2019).

✉ Xinhui Yan  
christiansen\_orv@yeah.net

<sup>1</sup> Beijing Information Technology College, Beijing 100018, China

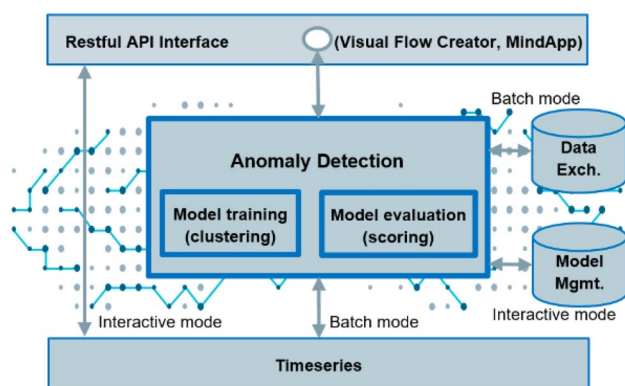


Fig. 1 The normal anomaly detection framework

In recent years, more and more methods have been introduced into the field of network anomaly detection, such as the neural networks, support vector machines, improved Naive Bayes classification, and push-through network anomaly detection methods. Compared with some other methods, the direct push network anomaly detection method has the advantages of high detection rate and low false alarm rate. Early production systems usually used index count matching. The basic idea is to create an index table for all patterns of all the rules. A specific fact in the working memory can be used to find all possible matching pattern instances. For each pattern instance, a count field is given to record in the working memory as the number of facts that can match this pattern; for each rule, also give a count field, record the number of patterns that can match the working memory; the initial value of the count field of each pattern and rule is 0, whenever When the working memory is modified, the index count matching algorithm starts to work. The implementation of the semantic network display algorithm is to better demonstrate the correctness of the knowledge matching algorithm, and it is based on dynamic demonstration. The premise of displaying a semantic network is to decompose a semantic network, first, and then further display.

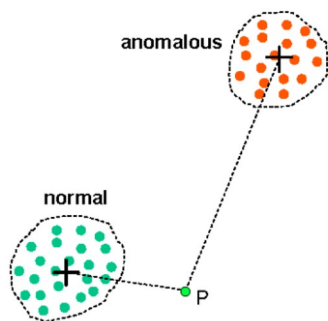
For enhancing the current methodologies, the speech network and the latest anomaly detection framework should be integrated. The research on the sentiment intensity is part of sentiment analysis. Emotion analysis is to divide the user's point of view, whether it is for or against it, and sentiment intensity is an analysis of the strength of criticism. It can reflect subjective information and capture hot events and public opinion. Before the rise of deep learning, the mainstream method for saliency detection was to model low-level features such as the color, texture, and location information. With the introduction of the deep learning, more and more advanced features are introduced into saliency detection, such as shape matching, template matching, deformable templates, active contours, etc. Advanced features can well evaluate target objects in images, but multi-level convolution

and pooling layers will cause the target boundary to become more blurred, the high-level features obtained by the output layer perform poorly, and there is a problem that the target cannot be accurately located. For the comprehensive analysis, the reference (Zheng et al. 2018) used the bidirectional recurrent neural network model for Chinese word segmentation on the premise of sequence labeling. By adding the context information of the word, it can effectively solve the gradient explosion problem, and achieved relatively good word segmentation effect. Proposed in Salamaz et al. (2018) with the method to analyze whether a comma is a clause boundary by extracting features such as syntax, vocabulary, length, etc. However, in this paper's semantic integrity analysis research, if the traditional method is used to determine whether a sentence is semantically complete, on the one hand, the sentence must be syntactic and grammatical analysis, on the other hand, we need to extract appropriate features from the analysis results and analyze the causal relationship between the features and the results. When the problem is more complicated, this method is basically not feasible. Plagens et al. (2018) and Peng et al. (2018) proposed a method based on recurrent nerves. The ancient sentence automatic segmentation method of the network. This method uses the GRU-based two-way recurrent neural network to segment ancient sentences. For visualized understanding, Table 1 gives the core perspectives of the speech networks and analytic framework.

Hence, after considering the speech network analysis, we consider the anomaly detection framework. In order to realize the detection of the advanced persistent intrusion symbols for mobile terminals in wireless sensor networks, the distributed sensor network design of the network attack nodes and also the symbol channel transmission model design of the network intrusion must first be designed, while combined with the intrusion interference node deployment and intrusion symbol amplitude feature extraction method, which performs envelope filtering and feature extraction to improve the ability to detect and identify network intrusion symbols. Our previous inspection methods have exposed a lot of the problems to our modern technology, and he has been unable to meet our needs. For some of the virus and unknown hazards he cannot identify the whole, thus a loophole, affecting our application, at the same time, his bit error rate is abnormal, down to our latest technology to apply it can improve the situation, the genetic algorithm is a good solution to these problems (Wang et al. 2018; Nautsch et al. 2019; Anchalia et al. 2019). We have so many methods in the calculation, the ability to solve is a test of our ability, he does not need very high technology, as long as you have a certain understanding of the problem solved, through a certain way to get the answer. One of the best ways to do this is to calculate the genetic formula, where he has achieved good results and, of course, has some problems. To improve

**Table 1** The core perspectives of the speech networks and analytic framework

Speech networks	Core perspectives characteristics
Composition and cognition of network resources	Because the network entity resources exist in physical space and have typical geo-spatial distribution characteristics, they can also be used as geo-spatial entities. Network virtual resources are composed of a series of digital behavior activities on the physical level of network entity resources, which can be divided into the three categories: application services, virtual entities and primary data resources (Yoshimura et al. 2019; King 2019; Deng et al. 2019)
Symbolic structure and semantic model of network resources	In natural linguistic theory, phonetics, grammar, and semantics are the three elements that make up a linguistic system. Phonetics and semantics are the two parts of the structure of language symbols. Grammar is a language component formed by the combination of both speech and semantics. Therefore, speech and semantics are the basic components of language components
Construction of symbol morphemes	The designed or specified sememes can merge and cluster their semantic meanings to form the semememes database, and the semememes are the product of the combination of the semememes database and the corresponding rule set. Any symbol is composed of a single sign morpheme or a number of sign morphemes. The symbol can also be combined with other symbol morphemes or basic units to form a new symbol



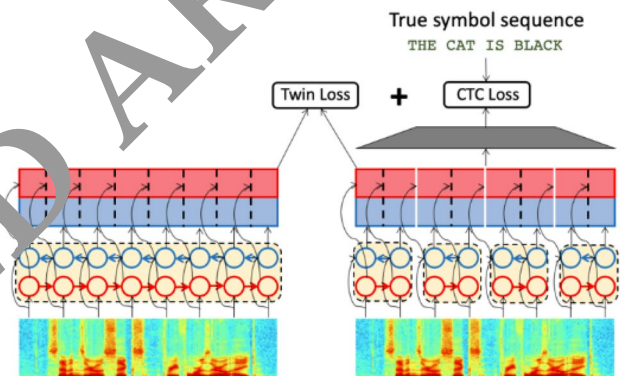
**Fig. 2** The anomaly detection with node information

people’s intrusion technology, we should update the types and types of viruses in a timely manner. We often do not pay attention to some subtle and noticeable viruses in new virus input. Therefore, in the unsupervised learning process, we must pay more attention to these issues. We should pay attention to cooperation in learning. We must learn from the surrounding environment. In the absence of a teacher, we must eliminate interference. Feedback and corrections to our learning and make the right assessment. In the Fig. 2, the detection details are presented. In the following sections, we will introduce the proposed analytic methodology in detail (Fig. 3).

## 2 The proposed model

### 2.1 FSS model framework and details

Computer-aided simulation is an analysis process deduced by a computer when the system structure is sufficiently defined and there are also calculation methods describing



**Fig. 3** The FSS model framework

the expected performance of the system. However, there are many unsatisfactory aspects of the traditional computer simulation technology. First, the mathematical models of complex systems often involve expertise in many fields and are difficult to establish; Second, because the user does not understand the internal mechanism of the system simulation, it is difficult for the simulation results to win the user’s trust. The optimal solution is transferred into Formula 1 (Ren et al. 2019; Ting 2019).

$$(\alpha_1^*, \dots, \alpha_n^*, D^*) = \arg \min_{\alpha_1, \dots, \alpha_n, D} \sum_{i=1}^n J(X_i, \alpha_i, D) \quad (1)$$

In the process of multi-attribute decision-making, in addition to considering the objective factors of decision-making, it is also necessary to consider the subjective factors of the decision-maker, that is, the risk preference of the decision-maker. In order to obtain FSS with the higher quality frequency characteristics, that is, stable center frequency and steep edge characteristics, the design and analysis results of

the dual-screen or multi-screen FSS structure have also been published. However, in practical engineering applications, the process size error exists objectively (Muhammad et al. 2018; Kim and Kweon 2019; Deka et al. 2019).

Fuzzy logic obtains expert experience and knowledge from experts, which uses qualitative knowledge and experience with unclear expression boundaries to make fuzzy reasoning with the help of membership function, thereby solving the problem of regular fuzzy information that is difficult to deal with by conventional methods. This paper introduces the risk preference factor in the fuzzy environment of interval two types to explore the basic influence of different risk preference attitudes on decision-making. Because the risk attitudes of different decision makers are different, some people may prefer the stimulus of big gains and losses, and some people may prefer to "seek stability", which can be divided into the risk aversion, relative risk aversion, risk neutrality, relative risk preference, and risk preference according to the risk attitude of decision makers. Hence, Formula 2 defines the risk function (Vilas et al. 2019). The function is achieved through the curve optimization.

$$\begin{aligned} \min_{\alpha} \quad & \lambda_C \left\| \mathbf{C}\mathbf{P} - \sum_{i=1}^L \alpha_i \mathbf{C}_i \right\|_F^2 + \\ & \lambda_A \left\| \mathbf{P}^{-1} \mathbf{A}\mathbf{P} - \sum_{i=1}^L \alpha_i \mathbf{A}_i \right\|_F^2 + \\ & \lambda_B \left\| \mathbf{P}^{-1} \mathbf{B} - \sum_{i=1}^L \alpha_i \mathbf{B}_i \right\|_F^2 + \lambda_{\alpha} \|\alpha\|_1 \end{aligned}$$

This function reflects the risk preference score of investors in the case of the risk preference attitude towards risk. In the next sections, we will consider using this core optimal function as the kernel parameter.

## 2.2 The speech signal feature extraction

An important step in speaker recognition is to sample feature parameters from the speech signal that collectively then reflect the personality of the speaker. These speech features have these characteristics: the ability to distinguish speakers is very strong which can clearly show the differences between different speakers, maintain relative stability when the speaker's own voice changes due to various factors; when recognition system when affected by changes in the environment, speech features have a stronger ability to adapt. In pattern recognition, how to effectively select the feature parameters is an important issue. One of the currently used solutions is to transform the original feature parameters into a vector space with a lower dimension through a data-driven linear feature transformation method, and further reduce the

dimensions while retaining important distinguishing components. Before the core analysis, we firstly define the parameters as the follows (Wu et al. 2019; Mukherjee et al. 2020).

$$Y = [y(1), \dots, y(N)]^T \quad (3)$$

$$X = [x(1), \dots, x(N)]^T \quad (4)$$

The discovery of new emotional words is not limited to the words found around emotional dictionaries as many new words are produced with emotional tags. For example, "girlfriend" is a noun, but it can express one's intimate emotions towards another. In addition, the emotion of some new emotional word is not the single classification, but there are multiple emotions in it, and the proportion is different. Therefore, the tendency of new emotional words should be then multi-classified. The speech signal is a non-stationary time-varying signal. When analyzing speech signals and extracting feature parameters, they must be first pre-processed. Preprocessing mainly includes general endpoint detection, pre-emphasis, and windowing and frame processing. The traditional method of new word mining is to segment the text first, assuming that the remaining segments of the text that failed to match are new words, and then extract these segments. However, the accuracy of the segmentation results depends on the integrity of the segmentation lexicon. If there are no new words in the word segmentation thesaurus, the result of the word segmentation may make it difficult to form the new words. Hence, following aspects should be considered (Li et al. 2019).

- (1) The primary criterion for judging whether a word can be formed is the degree of cohesion of the word. A notable problem is that in the process of the calculation there is no a priori knowledge. In other words, the word "designer" may be the "design" and "feeling", also may be "set" and "plan", so at the time of the calculation need to enumerate a variety of combination method of a field, and then take the combination with the highest probability.
- (2) In addition to the degree of cohesion within a word, another criterion is the degree of freedom outside the word. The adjacent entropy statistic uses information entropy to measure the uncertainty of the left and right neighbor characters of the candidate new word  $t$ . The higher the uncertainty, the more confused and unstable the string before and after the candidate new word  $t$ , so its higher the probability of word formation.
- (3) Community is also a special sub structure, that is, the internal connection is more closely related to the external structure. According to this definition of density, the meaning of community means that there are many interrelated groups. Of course, there are different defi-

nitions of community, but it is considered that community is a closely connected group based on some similar goals or commonalities (Madhavaraj and Ramakrishnan 2019).

Therefore, in our research, the model is designed as Fig. 4.

Linear prediction analysis is to gradually approach a speech sampled signal by using several previously sampled speech signals or core linear combination between them. According to a certain criterion, the error value between the actual speech sampling signal and the linear prediction sampling reaches a minimum value to then obtain its unique set of the prediction coefficients and use such a set of prediction coefficients in speech recognition and speech synthesis. To this regard, the question can be transferred into the following optimization issue.

$$\max \sum_{i=1}^n \left( \frac{\sum_{j=1}^m w_j \sum_{q=1}^{q_0} \left( H_{ij}^{\sigma(q)} - \left( H_j^{\sigma(q)} \right)^{-} \right)}{\sum_{j=1}^m w_j \sum_{q=1}^{q_0} \left( \left( H_j^{\sigma(q)} \right)^{+} - \left( H_j^{\sigma(q)} \right)^{-} \right)} \right) \tag{5}$$

s.t.  $w = (w_1, w_2, \dots, w_m)^T \in \Delta$

$$w_j \geq 0, j = 1, 2, \dots, m, \sum_{j=1}^m w_j = 1.$$

The topology semantics of graphs include not only the types of sub-graphs or the local structure semantics of sub-graphs, but also the relationships among sub-graphs. The most important relationship between the sub-graphs is the overlap relationship. When two sub-graphs share a point, it indicates that the common point may have two different identities or roles, or it may be the intermediary connecting the two groups. For obtaining the optimal solution, the following steps should be retained. (1) The text content in the web page is extracted from the text library to be retrieved by analyzing the xml language or html language. (2) The text is divided into several fields with a length less than 6. The degree of agglutination and information entropy of each field is calculated. The new vocabulary is obtained by

filtering based on the optimal threshold obtained from the experiment. (3) Word segmentation is carried out on the text. Generally, word segmentation should be then removed after word segmentation, but word vector learning needs to be based on the context, and word segmentation will also have an impact on word. Therefore, word segmentation is not removed at this step. (4) The segmented file is trained using the model, and the parameters are continuously adjusted to obtain satisfactory results. The word vector of all words is obtained after the training.

### 2.3 The speech and semantic networks

In Chinese text, the smallest unit of text is a word, but it is usually not practical to study a single Chinese character. Therefore, most studies on the similarity of the Chinese text are based on the words, sentences, or paragraphs. By calculating the number of words, sentences, or paragraphs in the text, the similarity between the different texts, and the corresponding similarity between different texts can be then obtained after corresponding processing. From the calculation steps, text similarity is mainly divided into three processes: text representation, feature extraction and the similarity calculation. Text after pre-processing is designed to get the corresponding formalized representation, items tend to get a lot of characteristics, and also some characteristics of a text message not much use for not only still can make text mining work multifarious, so some scholars from the improved feature extraction method of basic irrelevant information was lost in the guarantee only the realization of text dimension reduction at the same time. In co-occurrence network, the relationship between node features is determined by the distance between words in the same sentence. However, linguistics shows that there is a syntactic structure between components in a sentence, so it is more practical to determine the related word pairs by dependency syntactic analysis. We define the model as Formula 6.

$$Q(h, h^{(j)}) = E_{X|Y,h} [\log P(Y, X|h)] = \sum_X P(X|Y, h^{(j)}) \log(Y, X|h) \tag{6}$$

Supervised learning method regards relation extraction as a basic classification problem, and extracts features from training data to build a classification model. Although the supervised learning relation extraction method has high accuracy and recall, it depends on the relationship type system and labeled sequence set to be quite high. It requires a lot of training data to get a good classification model. The selection of the initial subset and the interference in the iteration process often affect the experimental results. The unsupervised learning method is essentially a clustering process that uses the context information of each entity pair to express the semantic

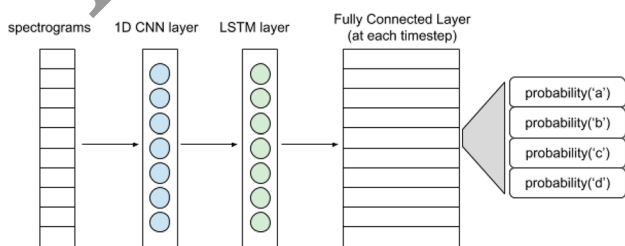


Fig. 4 The designed framework for the analysis

relationship of the core entity pair. Because the unsupervised learning method is difficult to determine the clustering threshold, the basic accuracy of the extracted results is low. Therefore, we define the operation as the follows.

$$Y(k) = X(k)H(k) + W(k) \quad (7)$$

In this model, the following procedures are integrated.

- (1) A semantic density clustering algorithm is introduced to globally cluster predicates in the corpus using the concept of the density, replacing sparse predicates with common predicates in the clustering set to which they belong. Multi-channel semantic synthesis of the CNN, first of all, emotional dictionaries are used to add emotional tendency weights to key words in sentences according to the attention mechanism during data processing. Then the multi-channel approach is adopted to input the text in parallel, and the weight matrix is Shared between the same layer, which greatly reduces the training complexity of the neural network
- (2) The fuzzy mechanism is introduced, and the concept of distance is used to reduce the semantics of the original word vector and improve the correlation with the predicate vector. Choosing a kernel convolution requires multiple convolution operations on the same input. Therefore, while extracting the multi-scale spatial information, it also brings a lot of network parameters and calculations. Therefore, high-precision convolutional neural networks often have a large number of network parameters and calculations, which is not conducive to the application of the model in manned systems.
- (3) Dropout regularization was introduced during the training phase to avoid the problem of neural network overfitting. Finally, the CRF is used to globally normalize the label probabilities to complete the optimal sequence labeling. In the training process, the hidden nodes that are "discarded" are random, that is, the networks used are different in each training process. Because each hidden node used for training is random, not every node can appear in every training process at the same time, which can ensure that the updating of weights does not depend on the joint action of hidden nodes with fixed relation, and ensure the validity and randomness of features to a great extent.

As presented in Fig. 5, the topology of the network is presented. It is the characteristic of language to express the infinite and complex real world through limited symbols. Modern linguistics, which is heavily influenced by structuralism, as forms a complete methodological foundation, which can systematically study the discrete characteristics

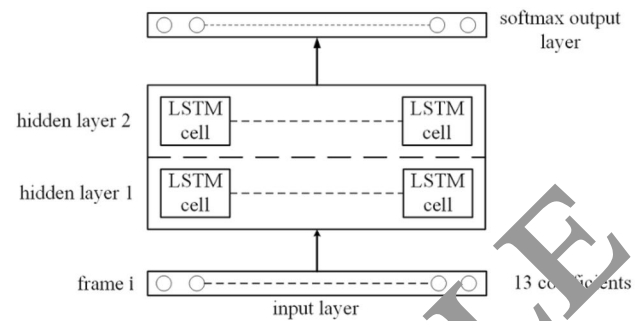


Fig. 5 Speech and semantic network topology and details

of objects from the overall concept. The combination of the geo-spatial information and map semiotics with modern linguistics methodology has formed a spatial information linguistic model and a map linguistic model. The remote supervision assumes that there is some correspondence between the entity pairs contained in the training corpus, and therefore all sentences in the corpus that contain the same entity pair express this relationship. The remote supervision will generate a lot of noise during the feature extraction process. To solve this problem, this paper introduces the attention mechanism to de-noise the impossible relations. In Formula 8, we define the model in detail.

$$\sup_{h,h',s,t,h^k=h^k} \left| U_i(h) - U_i(\tilde{h}) \right| \rightarrow 0 \quad (8)$$

When the model is used for prediction, all hidden nodes will be used, which is equivalent to the effective combination of all training models to get perfect model.

## 2.4 The anomaly detection pipelines

Collect network traffic in the transmission channel of the core wireless sensor network, we construct a transmission channel model for network intrusion, combine abnormal feature extraction methods to implement network intrusion detection, and use high-level statistics to describe persistent attacks in the link layer of the wireless sensor network. Invasive characteristic information. On the basis of constructing the distribution model of wireless sensor network's symbol transmission channel and feature extraction, the network intrusion detection algorithm is then improved and designed. Figure 6 presents the framework.

In the distributed intrusion detection, in order to improve system performance, network data are collected in parallel through multiple distributed sensor nodes, and then the collected massive data are segmented and distributed to multiple nodes for parallel detection and processing, so as to

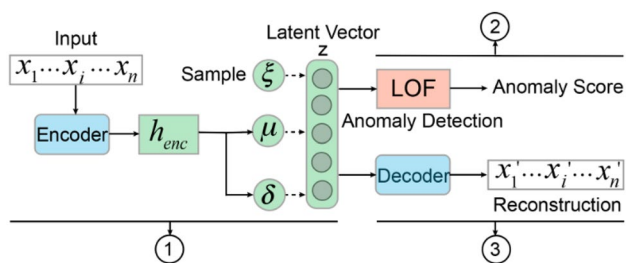


Fig. 6 The designed framework for the detection system

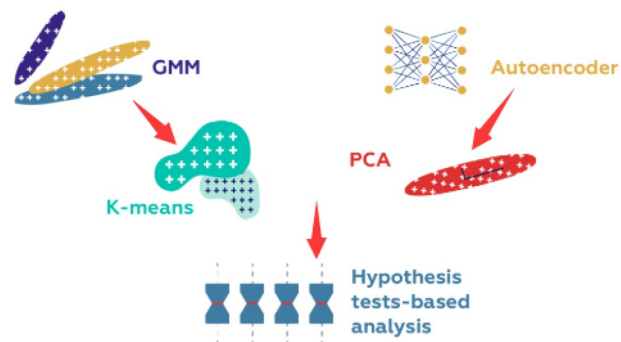


Fig. 7 The comparison of the different models

then improve the system’s data processing capacity. As an essential parameter, the distance is designed as follows:

$$d_{(i,j)} = \sqrt{\left(P_{(1,i)} - P_{(1,j)}\right)^2 + \dots + \left(P_{(n,i)} - P_{(n,j)}\right)^2} \tag{9}$$

The capacity of a node is an important basis for determining the granularity of the data allocated to the node. It is mainly determined by hardware configuration of the node, and relevant performance indicators can be collected through system monitoring. The load of a node is another important factor affecting the granularity of the data allocated to the node. The load evaluation of the data analysis node is completed according to the utilization of the node CPU, memory, and the network bandwidth collected by the system monitoring. Figure 7 provides the comparison.

There are many encoding and decoding methods of VAE. According to the data type of intrusion detection, a relatively simple MLP encoding method is selected. Multivariate Gaussian with diagonal covariance structure is used as the encoder and decoder of general MLP. Quantization noise is always in the digital system, which is determined by the digital system itself. Most of the signals in nature are analog signals, but the storage method of the computer determines that it can only process digital information, so in order to simulate, store, and process signals in nature as much as possible, the analog signals collected by the system are sent to a quantizer and

encoded into digital Signal, each number represents the instantaneous value of the signal obtained by one sampling. During quantization, the quantization level of the signal also brings quantization noise while completing the analog-to-digital conversion. Although the analog signal in nature cannot be fully obtained, the digitized signal is conducive to the transmission and processing of the signal, and it inevitably brings quantization noise. Formula 10 defines the denoise operation:

$$R_{GLS}^{-2} = \frac{n[\hat{\sigma}_\delta^2(0) - \hat{\sigma}_\delta^2(k)]}{n\hat{\sigma}_\delta^2(0)} = 1 - \frac{\hat{\sigma}_\delta^2(k)}{\hat{\sigma}_\delta^2(0)} \tag{10}$$

This paper focuses on improving the performance of the system through parallel detection and processing of massive data. Therefore, in the test, the rationality of the data segmentation algorithm for the massive data segmentation is mainly tested. After the data segmentation, the system load is distributed to the data analysis node for parallel detection and processing balance and improvement of system performance. In the next section, the experiment will be done for the general analysis.

### 3 Experiment

In this section, the simulation will be finalized to obtain the testing on the proposed methodology. On the data analysis node, in order to then ensure the consistency and validity of the detection results, the Snon system is installed and configured. After receiving the divided and distributed data to be detected, the Snort system performs data detection; the system monitoring uses the open source tool Sigar to collect all data analysis node running status; alarm response node performs alarm processing on detected intrusion behavior. Overall testing on the detection accuracy under various scenarios is presented in Fig. 8, the definitions of the axis can be referred to (Deka et al. 2019; Li et al. 2019; Zheng et al. 2018). When the model depth is 4, the detection accuracy is the highest, and when the model depth is 5, the detection accuracy is greatly reduced. This is because when the number of hidden layers is set to 300, the model’s feature learning ability has reached a strong level. Continue Increasing the depth of the model not only greatly increases the training time, but also leads to overfitting. Therefore, for a test set with more "unknown" or even unknown attribute values, the accuracy of model detection will decrease. At the same time, the comparison simulation is validated in Fig. 9.

### 4 Conclusion

Speech network analysis and anomaly detection based on the FSS model is well analyzed in this research. Network intrusion detection method can be divided into general misuse detection and anomaly detection. Most belong to misuse

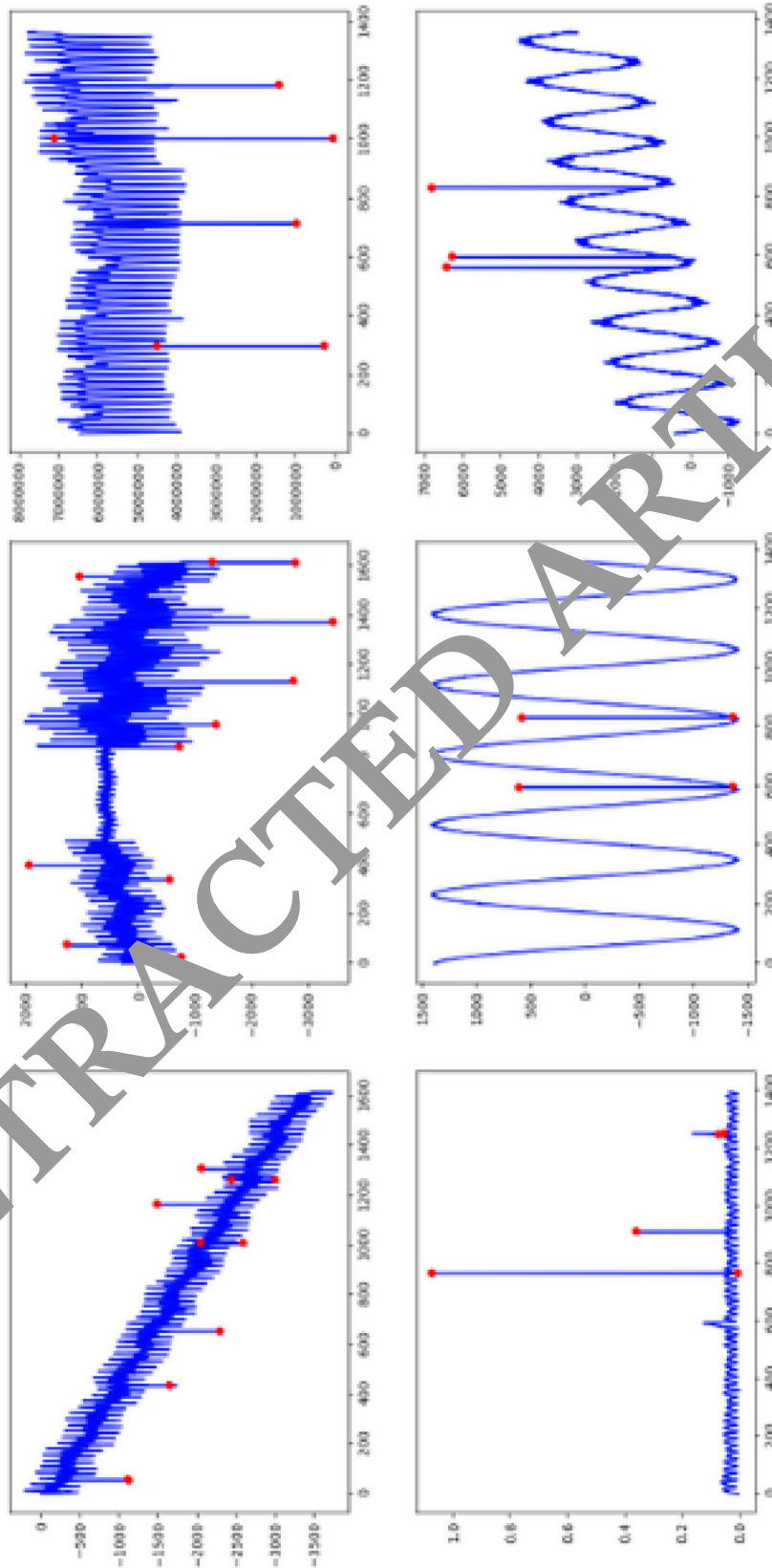
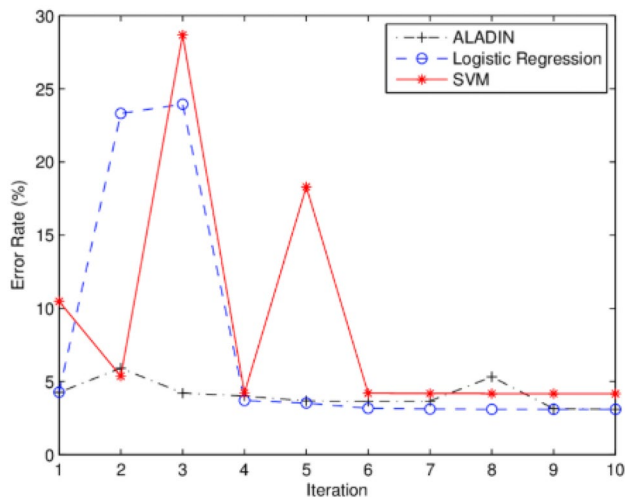


Fig. 8 The overall testing on the detection accuracy under various scenarios





**Fig. 9** The comparison experiment results

detection signature-based attacks, after being identified and analyzed, must be defined by security experts in the attack, after defining signature, all traffic coming into the network to compare the signature to determine whether the service trying to attack. Accordingly, signature-based error detection has better performance for a known type of attack detected. Anomaly detection is suspicious and also normal business operations to cope with a way to compare the attacks, the goal is to stop them before their success in unknown attacks. This paper focuses on the integration of the detection and speech network analysis. The saliency detection and the feature style analysis are integrated to achieve the balance. The experiment has proven the overall performance. In the next stage of the research, we will consider integrating the proposed algorithm into the voice assisted applications.

## References

- Anchalia, U., Reddy, K.P., Modi, A., Neelam, K., Prasad, D., & Nath, V. (2019). Study and design of biometric security systems: fingerprint and speech technology. In *Proceedings of the third international conference on microelectronics, computing and communication systems* (pp. 577–584). Springer, Singapore.
- Deka, A., Garmah, P., Samudravijaya, K., & Prasanna, S.R.M. (2019). Development of assamese text-to-speech system using deep neural network. In *2019 National Conference on Communications (NCC)* (pp. 1–5). IEEE.
- Deng, X., Li, Y., Weng, J., & Zhang, J. (2019). Feature selection for text classification: A review. *Multimedia Tools and Applications*, 78(3), 3797–3816.
- Ekberg, S., Danby, S., Theobald, M., Fisher, B., & Wyeth, P. (2019). Using physical objects with young children in ‘face-to-face’ and telehealth speech and language therapy. *Disability and Rehabilitation*, 41(14), 1664–1675.
- Faizin, B., Ramdhani, M.A., Gunawan, W., & Gojali, D. (2018). Speech acts analysis in Whatsapp status updates. In *International Conference on Media and Communication Studies (ICOMACS 2018)*. Atlantis Press.
- Kim, J. B. & Kweon, H. J. (2019). The analysis on commercial and open source software speech recognition technology. In *International Conference on Computational Science/Intelligence & Applied Informatics* (pp. 1–15). Springer, Cham.
- King, S. O. (2019). How electrical engineering and computer engineering departments are preparing undergraduate students for the new big data, machine learning, and AI paradigm: A three-model overview. In *2019 IEEE Global Engineering Education Conference (EDUCON)* (pp. 352–356). IEEE.
- Li, D., Cai, Z., Deng, L., Yao, X., & Wang, J. H. (2019). Information security model of block chain based on intrusion sensing in the IoT environment. *Cluster Computing*, 22(4), 451–468.
- Madhavaraj, A. & Ramakrishnan, A. G. (2019). Scattering transform inspired filterbank learning from raw speech for better acoustic modeling. In *TENCON 2019–2019 IEEE Region 10 Conference (TENCON)* (pp. 1154–1158). IEEE.
- Muckenhirn, H., Doss, M. M., & Marcell, S. (2018). Towards directly modeling raw speech signal for speaker verification using CNNs. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 4884–4888). IEEE.
- Muhammad, K., Hamza, P., Ahmad, J., Lloret, J., Wang, H., & Baik, S. W. (2019). Secure surveillance framework for IoT systems using probabilistic image encryption. *IEEE Transactions on Industrial Informatics*, 14(8), 3679–3689.
- Mukherjee, H., Obaidullah, S. M., Santosh, K. C., Phadikar, S., & Roy, S. (2020). A lazy learning-based language identification from speech using MFCC-2 features. *International Journal of Machine Learning and Cybernetics*, 11(1), 1–14.
- Nantsch, A., Jiménez, A., Treiber, A., Kolberg, J., Jasserand, C., Kindt, E., et al. (2019). Preserving privacy in speaker and speech characterisation. *Computer Speech & Language*, 58, 441–480.
- Peng, H., Bao, M., Li, J., Bhuiyan, M. Z. A., Liu, Y., He, Y., et al. (2018). Incremental term representation learning for social network analysis. *Future Generation Computer Systems*, 86, 1503–1512.
- Plageras, A. P., Psannis, K. E., Stergiou, C., Wang, H., & Gupta, B. B. (2018). Efficient IoT-based sensor BIG Data collection-processing and analysis in smart buildings. *Future Generation Computer Systems*, 82, 349–357.
- Ren, J., Chen, G., Li, X., & Mao, K. (2019). Striped-texture image segmentation with application to multimedia security. *Multimedia Tools and Applications*, 78(19), 26965–26978.
- Sailunaz, K., Dhaliwal, M., Rokne, J., & Alhadjj, R. (2018). Emotion detection from text and speech: A survey. *Social Network Analysis and Mining*, 8(1), 28.
- Savchenko, A. V., & Savchenko, V. V. (2019). A method for measuring the pitch frequency of speech signals for the systems of acoustic speech analysis. *Measurement Techniques*, 62(3), 282–288.
- Shadiev, R., Sun, A., & Huang, Y. M. (2019). A study of the facilitation of cross-cultural understanding and intercultural sensitivity using speech-enabled language translation technology. *British Journal of Educational Technology*, 50(3), 1415–1433.
- Ting, W. (2019). An acoustic recognition model for english speech based on improved HMM algorithm. In *2019 11th international conference on measuring technology and mechatronics automation (ICMTMA)* (pp. 729–732). IEEE.
- Vilas, A. F., Redondo, R. P. D., Crockett, K., Owda, M., & Evans, L. (2019). Twitter permeability to financial events: An experiment towards a model for sensing irregularities. *Multimedia Tools and Applications*, 78(7), 9217–9245.
- Wang, X., Lorenzo-Trueba, J., Takaki, S., Juvela, L., & Yamagishi, J. (2018). A comparison of recent waveform generation and acoustic modeling methods for neural-network-based speech synthesis. In

- 2018 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 4804–4808). IEEE.
- Wu, J., Hua, Y., Yang, S., Qin, H., & Qin, H. (2019). Speech enhancement using generative adversarial network by distilling knowledge from statistical method. *Applied Sciences*, 9(16), 3396.
- Yoshimura, T., Hashimoto, K., Oura, K., Nankaku, Y., & Tokuda, K. (2019). Speaker-dependent Wavenet-based Delay-free Adpcm Speech Coding. In *ICASSP 2019–2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 7145–7149). IEEE.
- Zheng, Y., Li, Y., Wen, Z., Liu, B., & Tao, J. (2018). Investigating deep neural network adaptation for generating exclamatory and interrogative speech in mandarin. *Journal of Signal Processing Systems*, 90(7), 1039–1052.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

RETRACTED ARTICLE