

Received December 9, 2018, accepted January 11, 2019, date of publication January 21, 2019, date of current version March 1, 2019.

Digital Object Identifier 10.1109/ACCESS.2019.2893877

# Detecting Diseases by Human-Physiological-Parameter-Based Deep Learning

**YULIANG LIU<sup>1,2</sup>, QUAN ZHANG<sup>1,2</sup>, GENG ZHAO<sup>3</sup>, ZHIGANG QU<sup>1,2</sup>,  
GUOHUA LIU<sup>4,5</sup>, ZHIANG LIU<sup>6</sup>, AND YANG AN<sup>1,2</sup>**

<sup>1</sup>College of Electronic Information and Automation, Tianjin University of Science and Technology, Tianjin 300222, China

<sup>2</sup>Advanced Structural Integrity International Joint Research Centre, Tianjin University of Science and Technology, Tianjin 300222, China

<sup>3</sup>Tianjin Medical University Hospital for Metabolic Disease, Tianjin 300070, China

<sup>4</sup>College of Electronic Information and Optical Engineering, Nankai University, Tianjin 300350, China

<sup>5</sup>Electronic Information Experimental Teaching Center, Nankai University, Tianjin 300350, China

<sup>6</sup>School of Electrical Engineering and Automation, Harbin Institute of Technology, Harbin 150001, China

Corresponding authors: Zhigang Qu (zhigangqu@tust.edu.cn) and Guohua Liu (liugh@nankai.edu.cn)

This work was supported by the National Natural Science Foundation of China under Grant 51674176 and Grant 61873187.

**ABSTRACT** The application of artificial intelligence in auxiliary diagnosis diseases has become a current research hotspot. The traditional method of diagnosing diabetes circulatory complication, diabetic peripheral neuropathy hyperlipidemia, diabetes mellitus peripheral angiopathy, and the comprehensive diseases is to distinguish an inspection report by a professional doctor. Its implementation of the clinical decision support algorithm for medical text data faces a challenge with the confidence level and accuracy. We proposed an expanding learning system to detect diseases above in our medical text data, which cover many kinds of physiological parameters of human, such as hematologic parameters, urine parameters, and biochemical detection. First, the raw data were expanded and corrected. Second, the processed data were fed into a 1D-convolution neural network with dropout and pooling. Our algorithm achieves 80.43%, 80.85%, 91.49%, 82.61%, and 95.60% with testing datasets (46 subjects). The effect of data quantification on model performance also had been researched, and the different data quantification methods would affect model performance on distinguishing different diseases. The proposed auxiliary diagnostic systems that have a highly accurate and robust performance can be used for preliminary diagnosis and referral; therefore, it is not only saving many human resources but also resulting in improved clinical diagnostic efficiency.

**INDEX TERMS** Deep learning, automatic diagnosis, physiological parameters of human.

## I. INTRODUCTION

Artificial intelligence (AI) has the potential to promote the development of disease detection by performing classification difficult for human experts and by rapidly processing huge amounts of medical images [1]–[4]. Despite its potential, processing medical text diagnosis of AI remains challenging.

The traditional automatic diagnosis algorithmic based on image recognition technology relied on (1) large number of samples labeling and classification by human experts in the early period, (2) extracting features, (3) using features to classify the samples by traditional neural network or traditional machine learning algorithm. The sample labeling and the data classification required lots of human experts and much time [5]. In addition, medical text is difficult to be used. Besides, traditional neural networks and machine learning algorithms are difficult to achieve a good performance in the task of medical text classification.

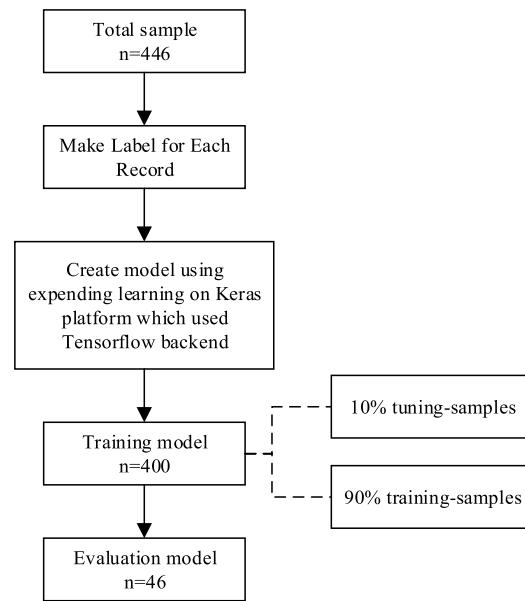
Numerous studies of using medical data to detect disease are reported in recent years. In the international scope, various countries have gradually combined artificial intelligence technology with medical technology and hope that the artificial intelligence technology can make the medical technology further developing [6]–[10]. For example, Liang *et al.* [11] proposed a deep learning model for automated electronic medical records diagnosis of traditional Chinese medicine. Kermany *et al.* [12] proposed a transfer learning deep learning model that uses optical coherence tomography (OCT) images and X-ray images to diagnose diseases and determine treatment effects. Coudray *et al.* [13] classified and predicted mutation from non-small cell lung cancer histopathology images by deep learning. The development of convolution neural network has allowed for significant gains in the accuracy to distinguish pictures and detect objects in image datasets. But it was rarely used to analyze text data.

Motived by the above research problem, 1D-convolution neural network and Long Short Term Memory Network (LSTM) were proposed [14], [15]. Fixed size filter was used to learning data's features in 1D-convolution neural network. As the 2D-convolution neural network, the feature was distinguished by filters. Memory neurons were used in LSTM to learn joint features of data far apart. Extracting the raw data yields the implicit knowledge which underpins the robust decision making process. Therefore, this neural network which replaces the clustering algorithm based on spatial distance can achieve the higher precision [16].

The other challenging of applying deep learning models is the lack of relevant data. When we have too little relevant data, the model cannot learn features enough and classify samples precisely. The other problem is that the model learning features excessively leads poor generalization ability on validation dataset. One method of addressing the lack of data in a given domain is using extended data to supplement raw dataset known as expanding learning. Expanding learning algorithm has been proved to be an effectively method [17], particularly when it faces with domains with limited data [18]. Rather than training a completely blank traditional neural network by complex coding to find relationship between input and output, the model can recognize the implicit knowledge in the raw data. For example, diabetic people have higher blood glucose level and glycosylated hemoglobin value. This model uses significantly fewer training examples and less labor force.

In this study, we sought to propose an effective expanding learning algorithm to process medical text data to provide a timely and effective diagnosis for each sample. We use human hematolgy data and urine data to illustrate the algorithm, and this algorithm's performance also was compared with other models.

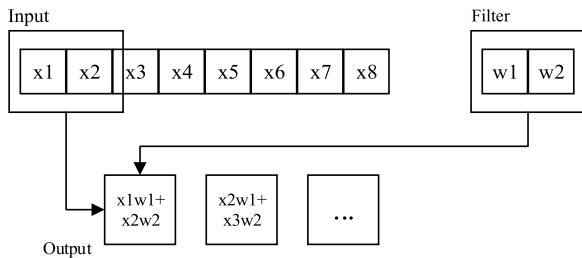
The expanding learning algorithm has been initially applied to process text data of human hematolgy and urinology. The parameters of blood glucose, glycosylated hemoglobin, blood routine, urine routine and biochemical test were used in this study. Different physiological parameters are often obtained by different methods, it can reflect physical condition of human [19], [20]. Hematology parameters are now one of the criteria for the diagnosis and treatment of diabetes and its complications in the worldwide [21], [22]. And urine parameters report metabolic status of human [23]–[25]. Therefore, hematological parameters and urinary parameters are used to determine healthy condition of human. It was estimated that in 2017 there are 451 million (age 18–99 years) people with diabetes worldwide. These figures were expected to increase to 693 million by 2045. In 2017, approximately 5 million deaths worldwide were attributable to diabetes. Moreover, diabetic complications decreased the quality of life of patients [26]. Diabetes has become an important factor threatening human life and health in the worldwide. Therefore, in the worldwide, researchers have done a lot of research on the treatment and detection of diabetes [27]–[34]. Fortunately the advent and wide-spread utilization of insulin



**FIGURE 1.** Workflow diagram (Workflow diagram showing overall experimental design describing the flow of medical text data through the labeling and grading process followed by creation of the expanding learning model, which then the model underwent testing. And the training model step include optimize the network weights and optimize hyperparameters).

medications and food therapy has revolutionized the treatment of diabetes and its complications. It can prolong life and improve quality of life of patients. Hematology parameters and urine parameters are critical for guiding the administration of insulin therapy by reflecting a clear physiological parameter of the patients. Different results representative different healthy condition, which made it possible for automatic clinical examination by these parameters. This method was shown in Fig. 1. The expanding learning algorithm can transfer patients to different departments with less medical resources and less labor cost.

We aim to reduce workload of the clinicians and to enable improving hospital process by providing a robust diagnosis support system for above diseases in a new dataset. In order to improve the performance of the model, two kinds of physiological parameters were used to diagnosis diseases, rather than single blood parameters. Therefore, compared with convention methods, it is comprehensive. This auxiliary diagnostic system gets all the information within the selected data. There is no information reduction through feature extraction. We also proposed a data extension method. When faced with the problem of lacking data, this method also can achieve a better performance. Moreover, we also found that different data quantification method would affect model performance. So we proposed a quantification method of negative test results to improve the accuracy of auxiliary diagnostic system. The auxiliary diagnostic system scored 80.43% 80.85% 91.49% 82.61% 95.60% with testing dataset. This performance figures make us confident that the system will perform well for unknown data in a



**FIGURE 2.** Schematic of 1D-Convolutional neural network.

practical environment. Our work opens up a viable path to auto-diagnostic. Furthermore, the results justify our method to use expanding learning algorithm for diseases detection.

## II. METHODS

The deep learning model in the experiment used 1D-convolution neural network which combined with dropout [35] and pooling technology. The datasets were first processed by convolution and pooling layers in order to extract potential characteristics in the data, and the feature was processed by classifier to classification samples. The gradient descent algorithm was used to update network parameters.

### A. 1D-CONVOLUTION NEURAL NETWORK

Convolution neural network is the core algorithm of learning model, which can effectively process sequence data of certain regularity. The principle of 1D-convolution algorithm is shown in Fig. 2.

The convolution layer learned features by weights sharing method, which maps relationship between input and hidden layers. Each filter learns one feature, and the principle is shown in Equation (1).

$$O = \sigma(b + \sum_{m=1}^N w_m a_{k+m}) \quad (1)$$

where  $O$  is the output of the  $n$ -th neuron,  $\sigma$  is activation function,  $b$  is shared bias,  $w_m$  is the  $m$ -th weight in the weight matrix,  $a_k$  is the input of the  $k$ -th neuron,  $N$  is the length of weight matrix.

### B. STOCHASTIC GRADIENT DESCENT

The stochastic gradient descent method was used to update parameters in the neural network. The principle is shown in Equation (2) and Equation (3).

$$C = \frac{1}{n} \nabla_{\theta} \sum_{k=1}^n L(I^{(k)}, O^{(k)}, \theta) \quad (2)$$

$$\theta' = \theta - \eta C \quad (3)$$

where  $C$  is the gradient of loss function,  $n$  is the number of mini-batch data,  $L$  is the loss function of each sample,  $I^{(k)}$  is the input of  $k$ -th sample,  $O^{(k)}$  is the output of  $k$ -th sample,  $\theta$  is the parameters of network,  $\eta$  is the learning rate,  $\theta'$  is updated parameters.

### C. LSTM

LSTM has good performance when it was used to process text data that has the dependency relationship, the neurons are connected to each other. The principle of updating parameters in the LSTM is shown in Equation (4).

$$s_i^{(t)} = f_i^{(t)} s_i^{(t-1)} + g_i^{(t)} \sigma(b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)}) \quad (4)$$

where  $s$  is state,  $b$  is bias,  $U$  is input weight,  $W$  is forgetting gate weight,  $f$  is control function of forgetting,  $g$  is input gate from external.

### D. SUPPORT VECTOR MACHINE (SVM)

The main idea of the SVM algorithm is to finish linear clustering using kernel functions [36]. It can classify data according to positive or negative of the decision function. The principle is shown in Equation (5).

$$f(x) = b + \sum_j \alpha_j k(x, x^j) \quad (5)$$

In Equation (5),  $f(x)$  is decision function,  $\alpha$  is coefficient vector,  $x^j$  is training data,  $k$  is kernel function. SVM only outputs categories instead of traditional probabilities.

## III. EXPERIMENT

### A. DATA DESCRIPTION

The experiments were conducted based on a new dataset from Tianjin Medical University Hospital for Metabolic Disease, which contained 446 physiological records from different patients. Each record contains 49 parameters which covered blood glucose, glycosylated hemoglobin, blood routine, urine routine and biochemical test. The specific data categories included in each class are shown in Table 1. These data come from medical test results in hospital.

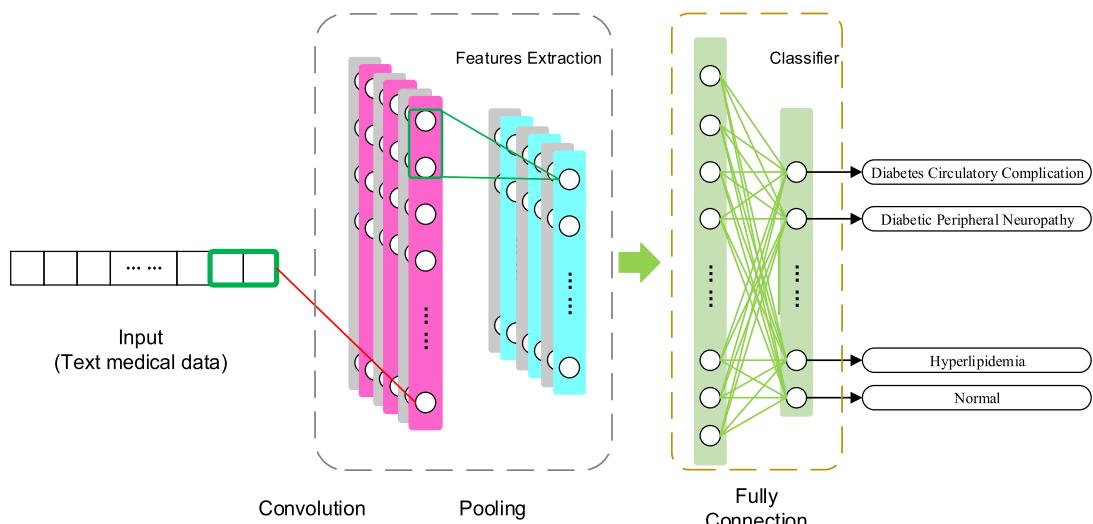
We labeled each data health condition by clinical diagnosis made by professional doctors. Data from 400 patients has been used for training the model and tuning of the model hyper-parameters, with the data from the remaining 46 patients completely withheld for evaluating model. This step ensures that the proposed method generalizes not only to unseen data, but also to unseen patients as well.

All of the subjects in this experiment were recruited from the patients who were going to accept health examinations. Before the experiment, we obtained the permission of the Tianjin Ethics Committee and obtained the consent of the patients and the paper works of informed consent was obtained of each patient. Patient records/information had been anonymized and de-identified before analysis.

There is a phenomenon in original data that a person has multiple diseases at the same time. Table 2 lists the decision of diagnosis and the corresponding number of each category in the raw data.

**TABLE 1.** Specific data categories for each item.

Item Name	Specific Data Categories
Blood Routine	Leukocyte, Neutrophils Percentage, Lymphocyte Percentage, Monocyte Percentage, Eosinophil Percentage, Basophilic Percentage, Erythrocyte, Hemoglobin, Hematocrit, Mean Corpuscular Volume, Average Hemoglobin Amount, Mean Cell Hemoglobin Concentration, Platelet, Thrombocytocrit, Mean Platelet Volume, Platelet Distribution Width, Large Platelet Ratio
Urine Routine	Color, Specific Gravity, PH, Leukocyte, Nitrite, Urine Protein, Glucose, Ketone Body, Urobilinogen, Bilirubin, Urine Erythrocyte
Biochemical Test	Total Protein, Albumin, Globulin Ratio, Total Bilirubin, Indirect Bilirubin, Direct Bilirubin, Alkaline Phosphatase, $\gamma$ -Glutamyltranspeptidase, Alanine Transaminase, Aspartate Aminotransferase, Urea, Creatinine, Uric Acid, Triglyceride, Total Cholesterol, High Density Lipoprotein, Low Density Lipoprotein, Very Low Density Lipoprotein
Blood Glucose	Fasting Venous Blood Glucose
Glycosylated Hemoglobin	Glycosylated Hemoglobin

**FIGURE 3.** The schematic of an expanding learning algorithm.

### B. DATA EXPANDING

Expanding learning algorithm was used in this experiment to solve the problems of the lack of data. We expanded training matrix by adding random disturbance matrix into the training matrix. The principle is shown in Equation (6).

$$Train = A_{m \times n} + D_{m \times n} \quad (6)$$

In Equation (6),  $Train$  is training matrix,  $A_{m \times n}$  is raw training matrix,  $D_{m \times n}$  is random disturbance matrix. The element value of the random perturbation matrix should be much smaller than the original training matrix. And the size of the random perturbation matrix should be same with the original training matrix.

### C. DESIGN OF THE DEEP LEARNING MODEL

Pooling and one-hot approaches were used in this experiment. They can improve the generalization ability and the

**TABLE 2.** List of diagnostic decisions and frequency.

Diagnostic	Frequency
Diabetes Circulatory Complication	62
Diabetic Peripheral Neuropathy	68
Hyperlipidemia	86
Diabetes Mellitus Peripheral Angiopathy	85

recognition accuracy of the model. The schematic of an expanding learning algorithm is shown in Fig. 3.

### D. TRAINING AND EVALUATION MODEL

The raw dataset was divided into 2 parts: (1) training: 400 samples of the data were used to training model

**TABLE 3.** Decision accuracy with different models trained by expanding dataset.

Diagnostic Decision	Learning Model			
	M1	M2	M3	M4
Diabetes Circulatory Complication	<b>80.43</b>	74.47	76.09	73.91
Diabetic Peripheral Neuropathy	<b>80.85</b>	76.60	76.09	73.91
Hyperlipidemia	<b>91.49</b>	89.36	91.30	89.13
Diabetes Mellitus Peripheral Angiopathy	<b>82.61</b>	63.83	73.91	73.91
Four comprehensive	<b>95.60</b>	94.95	60.87	4.35

**TABLE 4.** Decision accuracy with different models trained by raw dataset.

Diagnostic Decision	Learning Model			
	M1	M2	M3	M4
Diabetes Circulatory Complication	73.91	80.43	76.09	76.09
Diabetic Peripheral Neuropathy	71.74	68.48	76.09	73.91
Hyperlipidemia	84.78	82.61	91.30	89.13
Diabetes Mellitus Peripheral Angiopathy	73.91	78.26	73.91	73.91
Four comprehensive	95.38	94.70	60.87	4.35

and (2) evaluation: 46 samples were used to evaluate the model. The training part also has two parts: (1) training-samples: 90% of the dataset were used to optimize the network weights and (2) tuning-samples: 10% of the dataset were used to optimize hyperparameters (such as early stopping for training). The 10 fold-cross-validation algorithm was used to optimize hyperparameters. In the expanding learning algorithm the training-samples were expanded and the tuning-samples were not been expanded. The expanded training dataset was used to train AI system. And the AI system was evaluated in the accuracy of diagnosing diseases in the testing dataset. The early stop algorithm was used to preventing over fitting phenomenon (patience = 5). Besides, the gradient of cross entropy cost function was used to update global parameters. The principle of cross entropy cost is shown in Equation (7).

$$\text{loss} = -\frac{1}{n} \sum [y \ln a + (1 - y) \ln(1 - a)] \quad (7)$$

In Equation (7),  $n$  is the number of samples,  $a$  is the actual output,  $y$  is the target output.

The condition that one person has multiple diseases at the same time was also considered in this paper. In this condition, the labels were made by permutation and combination of diseases, such as the sample of Diabetes Circulatory Complication labeled 1, the sample of patient at the same time suffering from diseases Diabetes Circulatory Complication and Diabetic Peripheral Neuropathy labeled 2 and so on in order to diagnosis comprehensive illness. When judging multiple diseases at the same time, the task of classification layer is to classify multiple situations. This was why this paper used the softmax function as the classification function. Each neuron in output layer represents a kind of health situation.

The softmax function is shown in Equation (8).

$$S(x)_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} \quad (8)$$

The AI system can identify Diabetes Circulatory Complication, Diabetic Peripheral Neuropathy Hyperlipidemia, and Diabetes Mellitus Peripheral Angiopathy. Under this condition, this AI system can conduct preliminary referrals decisions of patients to save medical resources. We also trained LSTM model, SVM and traditional neural network to compare with the expanding learning algorithm performance.

This paper not only diagnosed single disease, but also comprehensive illness by expanding learning algorithm. The performance for all model are shown in Table 3.

M1 is expanding learning algorithm, M2 is LSTM model, M3 is SVM algorithm, M4 is traditional neural network model. When the raw data were used to train different model, the performance of each model are shown in Table 4.

The effect of data quantification on model performance in experiments also studied. Instead of using zero, the physiological parameter that the test result is negative were quantified with a small value closing to zero, this method is called data correction. The results above have been corrected. The performance of different models which trained with uncorrected data is shown in Table 5.

#### IV. DISCUSSION

In this study, an AI system was proposed to diagnosis and referral four common illness causes of diabetes by hematology parameters and urine parameters. By exploiting an expanding learning algorithm, our model demonstrated competitive performance on medical text data analysis without the

**TABLE 5.** Decision accuracy with different models trained by uncorrected and expanded dataset.

Diagnostic Decision	Learning Model			
	M1	M2	M3	M4
Diabetes Circulatory Complication	78.26	78.26	76.09	73.91
Diabetic Peripheral Neuropathy	76.09	71.74	76.09	73.91
Hyperlipidemia	91.30	91.30	91.30	89.13
Diabetes Mellitus Peripheral Angiopathy	80.43	76.09	76.09	73.91
Four comprehensive	95.11	94.84	60.87	4.35

need for large amount of raw data and a large amount of labor. Compared with other models' performance in processing text data, the results were shown in Table 3 and Table 4. When the same expanding training dataset were used to train system, the performance of M1 and M2 had been improved, thereby illustrating the power of the expanding learning system to make higher accuracy of classification and improving model robustness, even with a limited training dataset. The performance of LSTM model trained with extended data was not as well as the expanding learning algorithm. We speculate that this phenomenon is caused by the LSTM over-learning data relationship when the sample size is small. Therefore, when the model lacks raw data, the expanding learning algorithm also performs well in accuracy of diagnosis and referral, and it has powerful and effective diagnostic ability even when it has less raw data.

As shown in Table 3 and Table 4, expanded training dataset cannot improve recognition accuracy of SVM and traditional neural network. Therefore, the expanded learning algorithm is a better choice when the raw data is missing. Specially, the medical text data is difficult to collect in the large amounts necessary to train an absolute blank convolution neural network.

In experiment, a phenomenon was found that the accuracy of Diabetes Circulatory Complication and Diabetic Peripheral Neuropathy were similar and lower than the accuracy of Hyperlipidemia and four comprehensive diagnoses. The reason was speculated that this phenomenon may be that the features of Diabetes Circulatory Complication and Diabetic Peripheral Neuropathy are similar in hematology and urine, and the correlation between the disease and the data is low. When a single disease was diagnosed in this paper, the generalization ability of the model on the test dataset is not powerful because of the interference of similar features.

When the various illnesses were distinguished, the similar features phenomenon was alleviated, because of the increase in the number of types of diseases, which led the model to perform better in diagnosis various diseases. The model achieved the best performance in identifying hyperlipidemia. The reason was speculated that hyperlipidemia has a high correlation with parameters of hematological and urine. Therefore, it is difficult to accurately judge the other three diseases except hyperlipidemia only by hematological parameters and urine parameters.

The method of quantifying raw data also affects the performance of the model. In the experiment, the original data were quantified in different forms. As shown in Table 3 and Table 5, data correction can improve the accuracy of the model on the detecting diseases. The reason may be the corrected data that makes more inputs becoming validate input, it means that the model can learn the negative parameters better. We strive to achieve higher generalization capabilities by exploring more quantitative methods of medical data in the future researches. The data correction method does not improve performance of SVM algorithm and traditional neural network algorithm, which can illustrate that the expanding learning algorithm has more possibilities to be optimized.

Moreover, the model represents a generalized method that can potentially be applied to a wide range of medical text data to make a clinical diagnostic decision. This point was demonstrated by comparing the performance of expanding learning algorithms with other traditional machine learning algorithms and discussing the possibility of the model which could be further optimized in the future. Medical text data has played such a crucial role in diagnosing diseases, guiding treatment and proving medicine recommendations. Therefore, effective analysis medical text data is of great significance to the development of medical diagnosis. The experimental results show that the AI system has the ability to use less raw data to effectively identify more complex data. By demonstrating efficacy with showing the accuracy of different models and with discussing the possibility of the model which further optimized in the future, this expanding learning framework presents a compelling system for further exploration and analysis in medical text data and more generalized application to an automated community-based AI system for the diagnosis and triage of common human diseases. This model could create a referral systems and diagnosis diseases without a large amount of labor, particularly in lacking data areas. Look beyond the narrow area of detecting disease with Deep Learning applied to medical text data, we found similar works. Liang *et al.* [11] used deep generative learning to detect traditional Chinese medicine by electronic health record (EHR) with an accuracy of 87.26%. Choi *et al.* [37] used recurrent neural network to early detect of heart failure onset with an accuracy of 88.3%. Their methods used EHR to detect different diseases. All these studies have achieved good validation results with a conventional processing structure of

feature extraction followed by machine learning and there is no effective and unified method to evaluate the quality of EHR data, which results in the limitation of the accuracy of diagnosis diseases by data from EHR. The method we proposed does not need feature extraction, all the information available in the training dataset is passed on to the deep learning system. And the human physiological parameters used in this paper do not need artificial description, it has uniform standards. Hence, the model proposed by us can extract feature which enables it to make good decisions even for unknown data.

From what has been discussed above, we proposed a new method to detect diseases automatically and precisely even it has a small amount of data. This method can extract all the available information to create the implicit knowledge which underpins the subsequent decision-making processes. Because this method does not artificially reduce the original data, so the new pathogenic factors can be researched by this system. Besides, two kinds of physiological parameters were used to diagnosis diseases instead of single blood parameters, so it has the potential to diagnose complex diseases through multiple dataset. Data extension method and data correction can improve the performance of the model in diagnosis diseases. As a consequence, the auxiliary diagnostic system delivers accurate and robust results (It achieved 80.43% 80.85% 91.49% 82.61% 95.60% with testing dataset, respectively).

Despite its potential, it also has limitations. A limitation of our study comes from the fact that we have used data only contain hematological parameters and urine parameters. Some diseases cannot be judged only by blood and urine parameters, it also needs other information such as the description of symptoms. Another shortcoming is the model used early stop algorithm rather than plot the validation dataset accuracy curve. Although this algorithm can prevent over fitting, the performance may rises after multiple epoch. In future, we plan to investigate the impact of different methods on over fitting by validation accuracy curve.

## V. CONCLUSION

In this paper, a deep learning model was proposed to diagnosis Diabetic Peripheral Neuropathy, Hyperlipidemia, Diabetes Mellitus, and Peripheral Angiopathy by expanding learning algorithm. This paper studied different models performance on diagnostic different diseases. And this paper also researched the effect of different data quantification methods on models performance.

Having an accurate and robust auxiliary diagnosis system using medical text data is a prerequisite for intelligent medical diagnosis. The goal for such intelligent medical diagnosis is to reduce the workload of clinicians. Furthermore, such computer aided diagnosis reduces the risk of inter- and intra-observer variability. These systems have the potential to benefit patients by accelerating the process of medical diagnosis. For clinical studies, new pathogenic factors may be found

because this method proposed by us does not artificially reduce raw data.

The future work is centered on the auxiliary diagnosis of the human different diseases. In order to identify more types of diseases further work is planned to collect more types of medical text data to expand existing data and research the performance of model on diagnosis different diseases. Apart from hematological parameters and urine parameters, there are other physiological message that can be used to diagnose different diseases, such as heart rate signals et al [38]. We will use other reference method on our data after our dataset is expanded. Because the auxiliary diagnostic system gets all the information within the selected data, there is no information reduction through feature extraction. Therefore, we will study the new factors which can affect disease by this system. Apart from addressing these medical engineering aspects, further work is planned to conduct a full ablation study to investigate the importance of individual architectural and algorithmic components of the proposed model (such as preprocessing methods, preventing over fitting method, computational resource consumption, learning rates and other hyper parameters). In the experiment, we found different data quantification methods will affect the performance of model on diagnostic diseases. Therefore, it is necessary to research the preprocessing methods of raw data. In this paper, early stop algorithm was used to preventing over fitting phenomenon. But it may miss the best model's performance, thereby the different methods of preventing over fitting will be researched by plotting training accuracy and validation accuracy curve in further work.

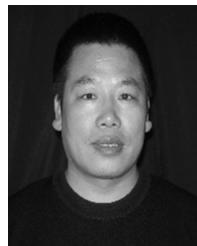
## DATA AVAILABILITY

The data that support the findings of this study are available from the corresponding author upon reasonable request.

## REFERENCES

- [1] H. Choi and K. H. Jin, "Predicting cognitive decline with deep learning of brain metabolism and amyloid imaging," *Behav. Brain Res.*, vol. 344, pp. 103–109, May 2018.
- [2] D. Wong and S. Yip, "Machine learning classifies cancer," *Nature*, vol. 555, no. 7697, pp. 446–447, 2018.
- [3] A. I. Shahin, Y. Guo, K. M. Amin, and A. A. Sharawi, "White blood cells identification system based on convolutional deep neural learning networks," *Comput. Methods Programs Biomed.*, vol. 168, pp. 69–80, Jan. 2019.
- [4] J. G. A. Barbedo, "Factors influencing the use of deep learning for plant disease recognition," *Biosyst. Eng.*, vol. 172, pp. 84–91, Aug. 2018.
- [5] A. D. Hoover, V. Kouznetsova, and M. Goldbaum, "Locating blood vessels in retinal images by piecewise threshold probing of a matched filter response," *IEEE Trans. Med. Imag.*, vol. 19, no. 3, pp. 203–210, Mar. 2000.
- [6] U. R. Acharya, O. Faust, N. A. Kadri, J. S. Suri, and W. Yu, "Automated identification of normal and diabetes heart rate signals using nonlinear measures," *Comput. Biol. Med.*, vol. 43, no. 10, pp. 1523–1529, 2013.
- [7] S. Webb, "Deep learning for biology," *Nature*, vol. 554, no. 7693, pp. 555–557, 2018.
- [8] M. Taddeo and L. Floridi, "How AI can be a force for good," *Science*, vol. 361, no. 6404, pp. 751–752, 2018.
- [9] O. Faust, Y. Hagiwara, T. J. Hong, O. S. Lih, and U. R. Acharya, "Deep learning for healthcare applications based on physiological signals: A review," *Comput. Methods Programs Biomed.*, vol. 161, pp. 1–13, Jul. 2018.

- [10] D. M. Camacho, K. M. Collins, R. K. Powers, J. C. Costello, and J. J. Collins, "Next-generation machine learning for biological networks," *Cell*, vol. 173, no. 7, pp. 1581–1592, 2018.
- [11] Z. Liang, J. Liu, A. Ou, H. Zhang, Z. Li, and J. X. Huang, "Deep generative learning for automated EHR diagnosis of traditional Chinese medicine," *Comput. Methods Programs Biomed.*, to be published.
- [12] D. S. Kermany et al., "Identifying medical diagnoses and treatable diseases by image-based deep learning," *Cell*, vol. 172, no. 5, pp. e9-1122-e9-1131, 2018.
- [13] N. Coudray et al., "Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning," *Nature Med.*, vol. 24, pp. 1559–1567, Sep. 2018.
- [14] F. A. Gers, J. Schmidhuber, and F. Cummins, "Learning to forget: Continual prediction with LSTM," *Neural Comput.*, vol. 12, no. 10, pp. 2451–2471, 2000.
- [15] Y. LeCun, "Generalization and network design strategies," in *Connectionism in Perspective*. 1989, pp. 143–155.
- [16] S. Ahlawat and R. Rishi, "Off-line handwritten numeral recognition using hybrid feature set—A comparative analysis," *Procedia Comput. Sci.*, vol. 122, pp. 1092–1099, 2017.
- [17] W. Yang, M. Zhao, Y. Huang, and Y. Zheng, "Adaptive online learning based robust visual tracking," *IEEE Access*, vol. 6, pp. 14790–14798, 2018.
- [18] P. Y. Simard, D. Steinkraus, and J. C. Platt, "Best practices for convolutional neural networks applied to visual document analysis," in *Proc. Int. Conf. Document Anal. Recognit.*, Aug. 2003, pp. 958–963.
- [19] D. B. Denicola, "Advances in hematology analyzers," *Topics Companion Animal Med.*, vol. 26, no. 2, pp. 52–61, 2011.
- [20] C. Li, Z.-M. Ni, L.-X. Ye, J.-W. Chen, Q. Wang, and Y.-K. Zhou, "Dose-response relationship between blood lead levels and hematological parameters in children from central China," *Environ. Res.*, vol. 164, pp. 501–506, Jul. 2018.
- [21] Y. Kachekouche, M. Dali-Sahi, D. Benmansour, and N. Dennouni-Medjati, "Hematological profile associated with type 2 diabetes mellitus," *Diabetes Metabolic Syndrome, Clin. Res. Rev.*, vol. 12, no. 3, pp. 309–312, 2017.
- [22] E. Letsky, "Haematology of pregnancy," *Medicine*, vol. 32, no. 5, pp. 42–45, 2004.
- [23] H. Chen, C. Tan, Z. Lin, and T. Wu, "The diagnostics of diabetes mellitus based on ensemble modeling and hair/urine element level analysis," *Comput. Biol. Med.*, vol. 50, no. 4, pp. 70–75, 2014.
- [24] L. Duvnjak, M. N. Perković, and K. Blaslov, "Dipeptidyl peptidase-4 activity is associated with urine albumin excretion in type 1 diabetes," *J. Diabetes Complications*, vol. 31, no. 1, pp. 218–222, 2017.
- [25] E. B. Fram, S. Moazami, and J. M. Stern, "The effect of disease severity on 24-hour urine parameters in kidney stone patients with type 2 diabetes," *Urology*, vol. 87, pp. 52–59, Jan. 2016.
- [26] N. H. Cho et al., "IDF diabetes atlas: Global estimates of diabetes prevalence for 2017 and projections for 2045," *Diabetes Res. Clin. Pract.*, vol. 138, pp. 271–281, Apr. 2018.
- [27] U. R. Acharya et al., "An integrated diabetic index using heart rate variability signal features for diagnosis of diabetes," *Comput. Methods Biomed. Eng.*, vol. 16, no. 2, pp. 222–234, 2013.
- [28] M. A. Pfeifer et al., "Quantitative evaluation of cardiac parasympathetic activity in normal and diabetic man," *Diabetes*, vol. 31, pp. 339–345, Apr. 1982.
- [29] A. C. Flynn, H. F. Jelinek, and M. Smith, "Heart rate variability analysis: A useful assessment tool for diabetes associated cardiac dysfunction in rural and remote areas," *Austral. J. Rural Health*, vol. 13, no. 2, pp. 77–82, 2005.
- [30] Z. Trunkvalterova et al., "Reduced short-term complexity of heart rate and blood pressure dynamics in patients with diabetes mellitus type 1: Multiscale entropy analysis," *Physiol. Meas.*, vol. 29, no. 7, pp. 817–828, 2008.
- [31] O. Faust, U. R. Acharya, F. Molinari, S. Chattopadhyay, and T. Tamura, "Linear and non-linear analysis of cardiac health in diabetic subjects," *Biomed. Signal Process. Control*, vol. 7, no. 3, pp. 295–302, 2012.
- [32] R. P. Nolan, S. M. Barry-Bianchi, A. E. Mechetiu, and M. H. Chen, "Sex-based differences in the association between duration of type 2 diabetes and heart rate variability," *Diabetes Vascular Disease Res.*, vol. 6, no. 4, pp. 276–282, 2009.
- [33] R. B. Pachori, P. Avinash, K. Shashank, R. Sharma, and U. R. Acharya, "Application of empirical mode decomposition for analysis of normal and diabetic RR-interval signals," *Expert Syst. Appl.*, vol. 42, no. 9, pp. 4567–4581, 2015.
- [34] E. B. Schroeder, D. Liao, L. E. Chambliss, R. J. Prineas, G. W. Evans, and G. Heiss, "Hypertension, blood pressure, and heart rate variability: The Atherosclerosis Risk in Communities (ARIC) study," *Hypertension*, vol. 42, no. 6, pp. 1106–1111, Dec. 2003.
- [35] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [36] A. Ukil, "Support vector machine," *Comput. Sci.*, vol. 1, no. 4, pp. 1–28, 2002.
- [37] E. Choi, A. Schuetz, W. F. Stewart, and J. Sun, "Using recurrent neural network models for early detection of heart failure onset," *J. Amer. Med. Inform. Assoc.*, vol. 24, no. 2, pp. 361–370, 2016.
- [38] O. Faust, A. Shenfield, M. Kareem, R. S. Tan, H. Fujita, and R. Acharya, "Automated detection of atrial fibrillation using long short-term memory networks with RR interval signals," *Comput. Biol. Med.*, vol. 102, pp. 327–335, Nov. 2018.



**YULIANG LIU** received the B.E. degree from the College of Mechanical Engineering, Hebei University of Technology, Tianjin, China, in 1995, the M.S. degree from the Institute of Modern Measurement and Control Technology, Hebei University of Technology, in 2002, and the Ph.D. degree from the College of Precision Instrument and Opto-Electronics Engineering, Tianjin University, in 2007. He is currently an Associate Professor with the College of Electronic Information and Automation, Tianjin University of Science and Technology, Tianjin. His research interests include biomedical engineering and medical electronic instrument.



**QUAN ZHANG** received the B.E. degree from the College of Electronic Information and Automation, Tianjin University of Science and Technology, Tianjin, China, in 2017, where he is currently pursuing the M.S. degree. His main interests include intelligent clinical diagnosis, deep Learning, and biomedical engineering.



**GENG ZHAO** received the B.E. degree from Tianjin Medical University, Tianjin, China, in 2004. He is currently the Chief Technician with Tianjin Medical University Hospital for Metabolic Disease. His research interest includes clinical test method.



**ZHIGANG QU** received the Ph.D. degree in instrument science and technology from Tianjin University, Tianjin, China, in 2008. He is currently a Professor with the College of Electronic Information and Automation, Tianjin University of Science and Technology, Tianjin. His research interests include non-destructive testing, structural integrity evaluation, and signal processing.



**ZHIANG LIU** is currently pursuing the B.E. degree with School of Electrical Engineering and Automation, Harbin Institute of Technology, Harbin, China. His main interest includes deep learning.



**GUOHUA LIU** received the M.S. and Ph.D. degrees from the Department of Microelectronic Engineering, Nankai University, Tianjin, China, in 1990 and 2004, respectively, where he is currently a Professor and the Dean of the Department of Microelectronics Engineering, School of Electronic Information and Optical Engineering. His research interests include artificial intelligence medical treatment, sensor technology, intelligent systems, and biomedical testing.



**YANG AN** received the Ph.D. degree in instrument science and technology specialty from Tianjin University, China. He is currently with the College of Electronic Engineering and Automation, Tianjin University of Science and Technology, China, as a Lecturer. His research interests include non-destructive testing and signal processing.

• • •