

# Research of insomnia on traditional Chinese medicine diagnosis and treatment based on machine learning

**Yuqi Tang**

Hospital of Chengdu University of Traditional Chinese Medicine <https://orcid.org/0000-0002-4227-1699>

**Zechen Li**

School of Automation,Chongqing University

**Dongdong Yang** (✉ [412409710@qq.com](mailto:412409710@qq.com))

**Yu Fang**

Hospital of Chengdu University of Traditional Chinese Medicine

**Shanshan Gao**

Hospital of Chengdu University of Traditional Chinese Medicine

**Shan Liang**

School of Automation,Chongqing University

**Tao Liu**

Electronic Engineering College,Chengdu University of Information Technology

---

## Research

**Keywords:** TCM, Insomnia, Machine learning, Diagnosis, Association rules, Cluster analysis, Random forest

**Posted Date:** November 17th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-42369/v3>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published on January 6th, 2021. See the published version at <https://doi.org/10.1186/s13020-020-00409-8>.

# **Research of insomnia on traditional Chinese medicine diagnosis and treatment based on machine learning**

Yuqi Tang<sup>1</sup>, Zechen Li<sup>2</sup>, Dongdong Yang<sup>1\*</sup>, Yu Fang<sup>1</sup>, Shanshan Gao<sup>1</sup>

Shan Liang<sup>2</sup>, Tao Liu<sup>3</sup>

\*Corresponding author: Dongdong Yang Email: 412409710@qq.com

1. Department of neurology, Hospital of Chengdu University of Traditional Chinese Medicine, Chengdu, 610072, CN.

2. School of Automation, Chongqing University, Chongqing, 400044, CN.

3. Electronic Engineering College, Chengdu University of Information Technology, Chengdu 610225, CN.

1 **Abstract**

2 **Background:** Insomnia as one of the dominant diseases of traditional Chinese medicine (TCM) has  
3 been extensively studied in recent years. To explore the novel approaches of research on TCM  
4 diagnosis and treatment, this paper presents a strategy for the research of insomnia based on machine  
5 learning.

6 **Methods:** First of all, 654 insomnia cases have been collected from an experienced doctor of TCM  
7 as sample data. Secondly, in the light of the characteristics of TCM diagnosis and treatment, the  
8 contents of research samples have been divided into four parts: the basic information, the four  
9 diagnostic methods, the treatment based on syndrome differentiation and the main prescription. And  
10 then, these four parts have been analyzed by three analysis methods, including frequency analysis,  
11 association rules and hierarchical cluster analysis. Finally, a comprehensive study of the whole four  
12 parts has been conducted by random forest.

13 **Results:** Researches of the above four parts revealed some essential connections. Simultaneously,  
14 based on the algorithm model established by the random forest, the accuracy of predicting the main  
15 prescription by the combinations of the four diagnostic methods and the treatment based on  
16 syndrome differentiation was 0.85. Furthermore, having been extracted features through applying the  
17 random forest, the syndrome differentiation of five zang-organs was proven to be the most  
18 significant parameter of the TCM diagnosis and treatment.

19 **Conclusions:** The results indicate that the machine learning methods are worthy of being adopted to  
20 study the dominant diseases of TCM for exploring the crucial rules of the diagnosis and treatment.

21  
22 **Keywords**

1 TCM; Insomnia; Machine learning; Diagnosis; Association rules; Cluster analysis; Random forest.  
2

### 3 **Background**

4 The application of TCM can be traced back to thousands of years [1]. In spite of the fact that  
5 TCM is still regarded as the complementary and alternative therapy in the field of modern medicine,  
6 it can hardly be ignored that TCM has attracted widespread attention in recent years due to its unique  
7 personalized treatment scheme and the outstanding treatment effect on some dominant diseases [2-3].  
8 Insomnia is one of the dominant diseases of TCM. It has been proven that TCM has been  
9 successfully applied to the treatment of insomnia in the medical field [4-5]. Compared with the  
10 western medicine in the treatment of insomnia, the advantages of TCM treatment are the  
11 personalization of diagnosis and treatment ideas, the non-dependence of treatment drugs and the  
12 diversity of treatment schemes, etc. Unlike the diagnosis and treatment of the western medicine,  
13 which is based on rigorous scientific trials, most of TCM diagnoses are relied on the experience of  
14 doctors to get comprehensive and personalized treatment strategies. Consequently, TCM is  
15 considered as an empirical medicine as well. Nonetheless, it should be noted that a set of core  
16 theories of TCM have been established since the beginning of the TCM development. Subsequently,  
17 the core theories of TCM have been developed into the TCM prescription, acupuncture, meridians  
18 and other theories [6]. Moreover, in the long-term clinical practice, with the constant deepening of  
19 the understanding of the basic theories of TCM, the diagnosis and treatment ideas of TCM have been  
20 promoted tremendously, and the diagnosis and treatment standards have achieved an innovation as  
21 well [7]. Diagnosis and treatment ideas and treatment strategies are the critical points of the clinical  
22 practice. Meanwhile, the medical record data are the embodiment of diagnosis and treatment ideas,  
23 thus worth exploring. The medical record of TCM is composed of four parts, including the basic

1 information, the four diagnoses of TCM, the treatment based on syndrome differentiation and the  
2 main prescription.

3 The concept of wholism and the treatment based on syndrome differentiation are the core  
4 principles for diagnosing and treating disease of TCM. In recent years, many studies have proposed  
5 that the TCM diagnosis and treatment should be integrated and personalized, which was essentially  
6 consistent with the core principles of TCM [8]. Syndrome differentiation and disease treatment are  
7 two inseparable parts in the process of TCM diagnosis and treatment. Syndrome differentiation is the  
8 premise and basis for treatment, and disease treatment is the means and method . The correctness of  
9 syndrome differentiation and treatment can be verified by examining the effect of disease treatment.  
10 The treatment based on syndrome differentiation is the core principle guiding the clinical work of  
11 TCM. In this paper, the four parts of TCM diagnosis and treatment (the basic information, the four  
12 diagnoses of TCM, the treatment based on syndrome differentiation and the main prescription) are  
13 the specific manifestations of TCM diagnosis and treatment process. The whole diagnosis and  
14 treatment process is not only logical, but also indivisible. The diagnosis and treatment of TCM is a  
15 whole from the information collection (including basic information and four diagnoses) to the  
16 treatment based on syndrome differentiation, and then to the establishment of the main prescription.  
17 In the past decades, many efforts have been done to study this process, whereas most researches have  
18 only focused on one part of this process. Zhang et al. applied the data mining technology to explore  
19 the drug rules of pulmonary fibrosis based on TCM medical records [9]. Yu et al. analyzed the dose  
20 data of TCM prescriptions by optimizing the traditional Cheng-Church double clustering algorithm  
21 (CC) [10]. Liu et al. adopted the data mining method to verify the TCM syndrome patterns of PSCI  
22 [11]. These researches have shown some opinions on the diagnosis and treatment process of TCM to

1 some extent. However, their research methods have violated the core principle of integration and  
2 personalization of the TCM diagnosis and treatment, resulting that their conclusions can hardly be  
3 applied in clinical practice. Therefore, for the sake of reliability and comprehensiveness of the  
4 research method adopted in the present paper, the research is carried out logically according to the  
5 sequence of TCM diagnosis and treatment, and the whole will be discussed at last.

6 In recent years, the rapid development of data analysis and artificial intelligence has provided  
7 an innovative research direction for the improvement of the clinical diagnosis and treatment  
8 technology. In the present paper, the medical record data of insomnia are selected as the research  
9 samples. Based on the medical record data, the research method of diagnosis and treatment of  
10 insomnia of TCM is emphatically discussed by applying machine learning methods. Specifically, the  
11 above-mentioned four parts in the process of TCM diagnosis and treatment are analyzed separately  
12 by three analysis methods, including frequency analysis, association rules and hierarchical cluster  
13 analysis. And then, a thorough analysis of the whole four parts is conducted using random forest.  
14 Considering that the data used in each analysis step have unique characteristics, different analysis  
15 schemes are established for different parts of the data.

16

## 17 **1 Data and methods**

### 18 **1.1 Sample data**

19 The sample data are obtained from the Hospital of Chengdu University of Traditional Chinese  
20 Medicine under the confidentiality agreement and the authority approval. According to the  
21 Guidelines for the diagnosis and treatment of insomnia in China(2017) [12] and the International

1 Classification of Sleep Disorders(ICSD-3)(2014) [13], the inclusion criteria are set as follows: the  
2 medical record data should contain one or more symptoms below: ①Sleep latency (SL) is prolonged  
3 and more than 30 minutes; ②Having difficulty in sleep maintenance, mainly manifested by easy and  
4 early to wake up; ③The quality of sleep is decreased, and the patient can hardly get into deep sleep  
5 and have multiple dreams; ④Insufficient sleep duration (less than 6.5 hours); ⑤With daytime  
6 symptoms, including fatigue, emotional problems, memory and attention decline, daytime sleepiness  
7 and work initiative decline, etc. The exclusion criteria are set as: ①The missing of the medical  
8 record data is so severe that it is unable to meet the research requirements; ②The patients have other  
9 serious organic diseases that may cause insomnia.

10 In our preliminary work, 1577 outpatient data (from 2016 to 2020) are collected and screened  
11 from an experienced doctor of TCM according to the above inclusion and exclusion criteria. The  
12 experienced doctor of TCM mentioned here refers to Professor Dongdong Yang. Prof Yang has been  
13 devoted to the clinical diagnosis and treatment of insomnia for decades. In the long-term clinical  
14 practice, a set of unique TCM diagnosis and treatment system has been formed. Finally, only 654  
15 outpatient data are selected as the research samples. Since the selection and analysis of medical  
16 record data of TCM have a high demand for expertise, three professional doctors of TCM (Prof  
17 Yang and the other two professional TCM doctors) are selected to analyze, code and classify the  
18 medical record data information of research samples manually. Meanwhile, the workload is equally  
19 assigned to the three doctors, and the cross-validation is implemented after all work has been  
20 completed, so as to eliminate the impact of subjectivity and artificial errors on the final data. Thus,  
21 there are only a few differences, and mainly in the treatment based on syndrome differentiation part.  
22 For this part, the classifications made by Prof Yang are dominant. Finally, the three doctors would

1 discuss and decide together. And then, the sample database is established. Simultaneously, according  
2 to the TCM diagnosis and treatment ideas, the contents of the sample data are divide into four parts:  
3 the basic information, the four diagnostic treatment, the treatment based on syndrome differentiation  
4 and the main prescription. Each part contains several data, and the specific data processing steps will  
5 be described later. In the light of the characteristics of the data, the machine learning methods,  
6 including frequency analysis, association rules and hierarchical clustering analysis, are adopted to  
7 process and mine the data. Finally, the data of the TCM diagnosis and treatment ideas from the four  
8 diagnoses, the treatment based on syndrome differentiation and the main prescription are integrally  
9 discussed by employing the random forest algorithm. The specific data processing flow designed in  
10 this paper is illustrated in **Figure 1.Flowchart of the data processing designed in this paper.**

11 The code comparative table is compiled by our research team. In the process of coding and  
12 classification, the Guidelines for the diagnosis and treatment of insomnia in China(2017) [12] and  
13 the International Classification of Sleep Disorders(ICSD-3)(2014) [13] are regarded as the basis to  
14 ensure the objectivity and comprehensiveness of the data. In the meantime, based on the personalized  
15 diagnosis and treatment strategy of Prof Yang, a complete code comparative table is shown in **Table**  
16 **1.** Three doctors of TCM are required to complete their work in strict accordance with the code  
17 comparative table.

18

## 19 **1.2 Data processing and machine learning**

### 20 **1.2.1 Data preprocessing**

21 Data preprocessing consists of data alignment, missing value processing and data format

1 conversion, etc. It is worth mentioning that the medical record information is extracted strictly  
2 according to the coding table, and there are a extremely small number of incomplete cases in the  
3 actual medical records. The incomplete items are represented by null values in the process of data set  
4 making. To eliminate the impact of the null value on the research and ensure that the follow-up  
5 research process can be carried out smoothly, the substitute values are selected to fill the null values  
6 of the record data. The substitute values include the course of disease and sleep duration, etc. and  
7 these values are filled with their mean value. The substitute values are specified in **Table 2**.

8 The processed data set are import into Python. The data samples are quantified by programming,  
9 and then analyzed by applying the following machine learning methods.

### 10 **1.2.2 Frequency analysis**

11 Frequency is also known as "time". The total data are divided into groups according to the  
12 preset standards, and then the number of individuals in each group is counted. The relative frequency  
13 is the ratio of the frequency of each group to the total number of data.

### 14 **1.2.3 Association rules**

15 A frequently-used method to study the relationship rules among data is to apply the association  
16 rules of Apriori algorithm [14]. Generally, three indicators, including confidence, support and lift,  
17 can be used to evaluate an association rule. Support is defined as the proportion of the data in the  
18 item set to the data in the data set, thus measuring the frequency of a set appearing in the original  
19 data. For instance, if two sets in the data set are X and Y respectively, then:

$$20 \text{Support}(X \rightarrow Y) = P(X | Y) \quad (1)$$

21 where  $X|Y$  represents the union of X and Y.

1 Confidence is defined for an association rule. The confidence of  $X \rightarrow Y$  can be expressed as  
 2 follows:

$$3 \quad \text{Confidence} = \frac{P\{x|y\}}{P\{X\}} \quad (2)$$

5 Lift can reflect the correlation between X and Y in association rules. As expressed in the  
 6 following function, the lift is defined as the proportion of the probability of the data set containing  
 7 both X and Y to the probability of the data set only containing Y.

$$8 \quad \text{Lift}(X \rightarrow Y) = \frac{P(Y|X)}{P(Y)} \quad (3)$$

10 The higher the lift is ( $\text{lift} > 1$ ), the higher the positive correlation is, and vice versa. The lift  
 11 equal to 1 indicates that there is no correlation.

#### 12 1.2.4 Cluster analysis

13 At present, the cluster analysis is extensively used in the medical field [15]. In general, the  
 14 cluster analysis can be classified into two categories, one is hierarchical clustering algorithm and the  
 15 other is agglomerative clustering algorithm. In the Euclidean space, using hierarchical clustering  
 16 algorithm to analyze small-scale data sets can achieve optimal results. Its basic principle is to  
 17 establish a hierarchical clustering tree by calculating the similarity among different categories of data  
 18 points and adopting the bottom-up aggregation strategy. Each sample set in the data sets is regarded  
 19 as a cluster, and then the clusters with close distance are merged step by step to achieve the expected  
 20 number of clusters.

21 Assuming that there are clusters  $C_i$  and  $C_j$ , the function can be described as follows:

$$22 \quad D_{aug}(C_i, C_j) = \frac{1}{|C_i| |C_j|} \sum_{x \in C_i} \sum_{z \in C_j} dist(x, z) \quad (4)$$

23 where the average distance  $D_{aug}(C_i, C_j)$  is determined by all samples of the two clusters.

1 **1.2.5 Random forest**

2 The random forest algorithm derived from ensemble learning method is composed of multiple  
3 decision trees. The random forest is an extension of the classification tree and the regression tree.  
4 These trees can be used to model the response variables through recursive partition and predict the  
5 final results jointly [16]. The random forest algorithm is commonly employed in data classification  
6 and regression [17]. At present, there are three mainstream decision tree algorithms, including ID3,  
7 C4.5 and CART. In the present paper, the most widely used algorithm, CART, is selected to build  
8 random forest algorithm model. The main function of this algorithm is described below.

9 Suppose that there is a training data set  $D$  with  $k$  classes in total. The Gini index of set  $D$  can be  
10 expressed as follows:

11  
12 
$$Gini(D) = \sum_k \frac{|C_k|}{|D|} (1 - \frac{|C_k|}{|D|}) = 1 - \sum_k (\frac{|C_k|}{|D|})^2 \quad (5)$$

13  
14 where  $C_k$  represents the sample subset of class  $k$ . The  $|C_k|$  and  $|D|$  represent the size of  $C_k$  and  $D$   
15 respectively.

16 In CART algorithm, assuming that feature  $A$  is used to segment the data. If feature  $A$  is a  
17 discrete feature, set  $D$  can be divided into subset  $D_1$  and subset  $D_2$  according to one possible value  $a$   
18 of  $A$ , as shown below.

19  
20 
$$D_1 = \{D | A = a\}; D_2 = \{D | A \neq a\} \quad (6)$$

21  
22 Consequently, the  $Gini(D,A)$  of set  $D$  under the condition of known feature  $A$  can be obtained  
23 by combining the above functions. The Gini index is theoretically similar to entropy, as described  
24 below.

25

$$Gini(D, A) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2) \quad (7)$$

Similar to the principle of entropy, the greater the value of  $Gini(D, A)$  is, the greater the sample uncertainty is. Taking this into account, the value of  $Gini(D, A)$  should be as small as possible when selecting the feature  $A$ .

## 2 Results

### 2.1 Basic information

The basic information mainly consists of the ID, name, clinic time, age and gender of patients. Since the clinic time is not taken as a factor in the screening criteria during the data screening stage, the statistical results may deviate from the actual situation. The ID and name of patients have no impact on the diagnosis and treatment process. As a consequence, the focus of this section is age and gender of patients. Considering that the categories of age and gender data are relatively few, we choose frequency analysis for the data processing. The age distribution of patients shown in **Figure 2. Violin plot of the age distribution of patients.** And **Figure 3. Pie chart of the gender distribution of patients** depicts the gender distribution of patients.

### 2.2 Four diagnostic methods

Four diagnostic methods include inspection (observation), auscultation and olfaction (listening and smelling), interrogation (inquiring or questioning) and palpation (pulse examination). Basically, it is a process of collecting medical history information for doctors of TCM [18]. “Inspection” refers to the observation of patients' external performance, such as tongue picture, expression, reaction and complexion. Moreover, “auscultation and olfaction” is the way that doctors diagnose diseases by

1 hearing and smelling. Additionally, “interrogation” is a sort of diagnostic method for doctors to find  
2 out the occurrence, development, treatment process and past health history of diseases by talking  
3 with patients. Furthermore, “palpation” particularly refers to the method that doctors use index  
4 fingers, middle fingers and ring fingers to touch the special position of radial artery of patients to  
5 check the pathological changes of patients. In the process of collecting medical record information  
6 through the four diagnostic methods, the amount of information obtained by “inspection” and  
7 “auscultation and olfaction” is relatively less than that obtained by “interrogation” and “palpation”.  
8 As a result, in the process of data statistics, the “inspection” and “auscultation and olfaction”  
9 diagnostic data are combined together for further analysis. Meanwhile, the "tongue diagnosis" data,  
10 which is the core of "inspection" diagnosis, are classified and counted separately. In this paper, the  
11 four diagnostic methods are further classified on the basis of the characteristics of the medical record  
12 data of insomnia research samples (shown in **Figure 4. Classification of four diagnostic methods**).

13 Based on the smallest unit of classification, the method of association rules is applied to study  
14 in this section. Considering that the basic information is also a part of TCM interrogation and may  
15 have an effect on the diagnosis and treatment process of the diseases, the basic information is  
16 included in the four diagnostic parts for discussion as well. Taking into account that there are too  
17 many null values in some of the smallest classification units, we attempt to use two methods to  
18 analyze the association rules for the combinations of the smallest units (the combinations items are in  
19 the brackets below), so as to minimize the impact of the null values on the research results. The  
20 results are listed in **Table 3** and **Table 4**.

1 Method 1: four diagnostic methods of TCM: (tongue proper, tongue color and tongue coating),  
2 (sleep duration, sleep status, course of insomnia, concomitant symptoms and emotion), pulse, age  
3 and gender. The results are summarized in **Table 3**.

4 Method 2: (tongue proper, tongue color and tongue coating), pulse, (sleep duration, sleep status  
5 and course of insomnia), emotion, (concomitant symptoms, others), age and gender. The results are  
6 summarized in **Table 4**.

7

### 8 **2.3 Treatment based on syndrome differentiation**

9 Originating from the philosophical culture, the treatment based on syndrome differentiation is  
10 the core of the TCM theories and gradually develops into a complex theoretical framework,  
11 including the yin and yang theory, five elements, eight principles, the Qi and blood theory, the organs  
12 theory and the meridian system [19].

13 The treatment based on syndrome differentiation is a comprehensive analysis by doctors in the  
14 process of diagnosis and treatment of TCM, and its judgment criteria are derived from the objective  
15 medical record information including the four diagnoses. The treatment based on syndrome  
16 differentiation consists of two processes: syndrome differentiation and treatment. It is not only an  
17 essential principle of understanding and treating diseases in TCM, but also a special research and  
18 treatment method of diseases in TCM. In the light of the logic of TCM syndrome differentiation, the  
19 treatment based on syndrome differentiation can be divided into parts, including the eight principal  
20 syndrome differentiation, the organs syndrome differentiation, the meridian syndrome differentiation,  
21 etc. Further, these parts can be separated into several items. The previous studies have either skipped  
22 this process directly or determined the syndrome differentiation category only based on the

1 experience description of doctors, which were too empirical. Based on the characteristics of  
2 insomnia in TCM, this paper focuses on four significant syndrome differentiation points, namely the  
3 syndrome differentiation of asthenia and sthenia, the syndrome differentiation of cold and heat, the  
4 syndrome differentiation of organs and pathogenic factors. The medical record data are extracted by  
5 three professional doctors of TCM, and then classified and coded according to the above four  
6 significant syndrome differentiation points. It is worth mentioning that the organs syndrome  
7 differentiation includes heart, liver, spleen, lung, kidney, gall bladder, stomach, small intestine, large  
8 intestine, bladder and the triple burner; the syndrome differentiation of asthenia and sthenia consists  
9 of asthenia syndrome and sthenia syndrome; the syndrome differentiation of cold and heat is  
10 composed of cold syndrome and heat syndrome; the pathogenic factors include phlegm, fire, blood  
11 stasis and asthenia. The above-mentioned 19 syndrome differentiation factors constitute the section  
12 of treatment based on syndrome differentiation of the insomnia sample data research in this paper. To  
13 ensure the objectivity of each syndrome differentiation factor, the three TCM doctors are supposed to  
14 collect at least two or more kinds of medical record information in the classification and coding stage  
15 of medical record data for determining one syndrome differentiation factor. For instance, the medical  
16 information "wiry pulse" and "irritability" can infer that the syndrome differentiation factor of organs  
17 is liver; the medical information "thin pulse" and "tiredness" can imply that the factor of asthenia and  
18 sthenia syndrome differentiation is asthenia syndrome; the medical information "red tongue"  
19 combined with "tidal fever" and "rapid pulse " indicates that the factor of cold and heat syndrome  
20 differentiation is heat syndrome; the medical information "slippery pulse" combined with "yellow  
21 tongue" and "greasy tongue coating" means that the pathogenic factors is phlegm.

1        Despite that each syndrome differentiation factor in each medical record is relatively  
2 independent, there is a strong correlation among the factors. Therefore, it is reasonable to select  
3 association rules for the analysis. There are two main reasons affecting the confidence. One is that  
4 there are only 654 data selected in this paper, and there are a small number of patients with  
5 incomplete medical records, resulting in sparse data distribution. The other is that a variety of  
6 classification methods are adopted in this paper, leading to complex classification and more  
7 categories of classification combinations. Through a process of trial and error, the confidence is  
8 finally adjusted to 0.7 and the results are summarized in **Table 5**.

9

## 10 **2.4 Main prescription**

11        Basically, the treatment strategy is composed of acupuncture, moxibustion, scraping therapy  
12 and TCM prescription, etc. In the present paper, the TCM prescription is the research focus of this  
13 section. TCM prescription is the embodiment of clinical practice of TCM. Choosing the appropriate  
14 combinations of Chinese medicine under the guidance of the treatment based on syndrome  
15 differentiation not only reflects the typical thoughts of TCM, but also conforms to the treatment  
16 method of drug combinations therapy [20]. Having completed the process from the four diagnoses to  
17 the treatment based on syndrome differentiation, the doctors should determine the main prescription.  
18 And then, on the basis of the main prescription, the doctors should adjust the prescription properly  
19 according to the actual situation of patients. Finally, the treatment prescription can be obtained. Thus,  
20 the determination of the main prescription is particularly significant. The main prescription can not  
21 only prove the personalized treatment advantages of TCM, but also reflect the most core treatment  
22 method in the clinical practice of TCM. The previous studies have achieved some success; however,

1 there are two deficiencies in their researches. First, the previous researches have mainly focused on  
2 the frequency of herb use and interrelation of the herbs. Second, there are few previous researches  
3 concerned about the components of the main prescription [21-22]. Taking the above deficiencies into  
4 account, the less use herbs are removed from the statistics of the herb use frequency, thus reducing  
5 the impact on the research of the main prescription in this paper. **Table 6** shows the herbs in the  
6 prescriptions through the previous data processing. These codes only represent the corresponding  
7 herbs.

8 For the sake of reducing calculation amount and the increasing the code execution efficiency,  
9 all the herbs are replaced with codes, and then the codes are entered into the database. The results of  
10 the analysis of the main prescription using hierarchical clustering analysis are shown in **Figure**  
11 **5.Results of the analysis of the main prescription using hierarchical clustering analysis.** The  
12 frequencies of the main prescription 3 to 13 are less than that of the main prescription 1 and 2, whose  
13 frequencies are 573 and 312 respectively. In order to facilitate comparison and observation, the main  
14 prescription 1 (red line in Figure 5) is determined as Category 1, the main prescription 2 (blue line in  
15 Figure 5) is determined as Category 2, and the sum of frequencies of the main prescription 3 to 13  
16 (green line in Figure 5)(the sum is 577) is determined as Category 3.

17 It is necessary to record the prescription information of the medical record data completely and  
18 accurately, so all the prescriptions that have appeared repeatedly are counted as the main  
19 prescriptions, the main prescriptions of the corresponding serial number are presented in **Table 7**,  
20 and the repeated herb combinations in all main prescriptions are shown in **Table 8**.

21

## 22 **2.5 Diagnosis and treatment idea**

1 In the discussion of the aforementioned four parts, the four parts of TCM diagnosis and  
2 treatment ideas are studied successively, so as to reveal the internal relationship and related research  
3 methods of each part. This section discusses the four parts as a whole. In accordance with the  
4 research process designed in the previous section(in Figure 1), the random forest algorithm is  
5 adopted to establish the model. Simultaneously, the data sets collected from four diagnoses,  
6 treatment based on syndrome differentiation, and the main prescriptions of TCM are put into the  
7 model for cross-validation by k-fold cross-validation method. Consequently, the corresponding  
8 accuracy can be obtained. In the meantime, for the purpose that the internal relationship of TCM  
9 diagnosis and treatment ideas can be explored deeper, this section is divided into two processes for  
10 further discussion. These two processes are illustrated in **Figure 6. Flowchart of the diagnosis and**  
11 **treatment ideas of TCM.**

12 The data set put into the random forest model includes 654 data. According to flowchart Figure  
13 6, two processes are designed and six different random forest models are established to realize the  
14 prediction of TCM diagnosis and treatment ideas. There are 200 CART decision trees for the models  
15 of main prescriptions, and the minimum number of leaf nodes is 5. For the other five models, there  
16 are 200 CART decision trees and the minimum number of leaf nodes is 3. The response variables of  
17 the six models are cold and heat, asthenia and sthenia, five zang-organs combinations, six fu-organs  
18 combinations, pathogenic factors combinations and main prescription combinations.

19 It is worth noting that five zang-organs, six fu-organs and pathogenic factors each contains  
20 several syndrome differentiation factors, which are irregularly combined in the actual medical record  
21 data. In addition, in the actual outpatient service, the prescriptions made by doctors for patients  
22 commonly includes at least one main prescription. Taking the data sample of this paper as an

1 example, 6 independent syndrome differentiation factors labels of the five zang-organs (including a  
2 normal item and 5 independent syndrome differentiation factors) can present 16 different  
3 combinations labels, while 14 independent main prescriptions labels (including an unprescribed item  
4 and 13 independent prescriptions) can present 21 different combinations labels. Therefore, in order  
5 to facilitate data processing, the five zang-organs combinations labels, six fu-organs combinations  
6 labels, pathogenic factors combinations labels and main prescription combinations labels are coded  
7 and loaded into the database. For the sake of presenting the accuracy more intuitively, the method of  
8 confusion matrix is carried out in this paper. The confusion matrix results are shown in **Figure**  
9 **7.Confusion matrix.**

10 As summarized in **Figure 8(A).Accuracy, AUC and Micro-F1 score for each model**, the  
11 accuracy of applying the random forest algorithm model to predict the information of treatment  
12 based on syndrome differentiation through the four diagnostic information is dramatically high.  
13 Simultaneously, the high accuracy is achieved by predicting the main prescription combinations  
14 through the information of the combinations of the four diagnoses and the treatment based on  
15 syndrome differentiation.

16 The ROC curve is introduced to evaluate the discriminating ability the model's discriminating  
17 ability. As depicted in the **Figure 9. ROC curve for each model and Figure 8(B). Accuracy, AUC**  
18 **and Micro-F1 score for each model**, except for the low AUC value of five zang-organs  
19 combinations, the random forest prediction models are effective in verifying the TCM diagnosis and  
20 treatment ideas. The false positive rate is higher than true positive rate of the ROC curves of the  
21 random forest models of five zang-organs combinations, six fu-organs combinations and the main  
22 prescription combinations in some cases. The model of five zang-organs combinations is the most

1 obvious one, leading to the lowest AUC value. There are two reasons for this situation. One is that  
2 the amount of sample data used in this study is small. Second, the data are classified by various  
3 research strategies, and meanwhile, the data distribution is uneven. The distribution of the five  
4 zang-organs data is more uneven than that of other categories data used in other models. The reason  
5 for this is that in the process of data coding, the syndrome differentiation of five zang-organs in most  
6 clinical insomnia medical records is the heart, followed by the liver, while the number of the  
7 syndrome differentiation of spleen, lung and kidney is relatively small. Due to the above reasons, the  
8 output of the models is more inclined to the side with higher accuracy in the process of model  
9 training. In the random division of the data sets, a small number of samples are divided into the test  
10 data sets, thus lacking the training of these samples, resulting in false positive rate [23-24].  
11 Nonetheless, as the amount of data increases, this problem would be alleviated. Similarly, the models  
12 of six fu-organs combinations and the main prescription combinations also show false positive rate in  
13 some cases. Since the data distributions are more even, the false positive cases are less than the true  
14 positive cases. In our future research, we will improve the efficiency of the models by optimizing the  
15 data classification strategies and increasing the amount of data.

16 The Micro-F1 score is selected to evaluate the accuracy of the models established in this paper.  
17 F1 score has been considered as an index used to measure model accuracy in machine learning [25].  
18 Both accuracy and recall of the classification models are taken into account by using F1 score. There  
19 are two evaluation indexes for multivariate classification, namely Micro-F1 score and Macro-F1  
20 score. Micro-F1 score is more suitable for unbalanced data distribution. Due to the various  
21 classification strategies adopted in this paper and the unbalanced data distribution, Micro-F1 score is  
22 selected as the index for evaluating the accuracy of the model. According to the Micro-F1 score

1 (shown in **Figure 8(B). Accuracy, AUC and Micro-F1 score for each model.**), it can be seen that  
2 the accuracy of each model is dramatically high.

3 In process 2 of this section, the random forest algorithm model is applied to extract the  
4 eigenvalues of all data in the data sets. Since the eigenvalues obtained by using the random forest  
5 model are too small to be studied conveniently, the eigenvalues are expanded in the form of  
6 logarithmic transformation to facilitate the observation. The transformed eigenvalues are shown in  
7 **Figure 10. Transformed eigenvalues obtained by using random forest model.**

8

### 9 **3 Discussions**

10 The development of modern medicine mainly has embodied in the continuous improvement of  
11 the basic medical theory and clinical practice. In the meantime, the research on the etiology and  
12 pathogenic mechanism are more inclined to **the** micro research. The machine learning methods have  
13 been widely used in modern medical related fields [26], such as biochemistry, physiology,  
14 microbiology, anatomy, pathology, pharmacology, etc. The concept of wholism and the treatment  
15 based on syndrome differentiation are the core principles for diagnosing and treating disease of TCM.  
16 The more macroscopic characteristics of diagnosis and treatment also lead to the research of TCM  
17 medical record data more complex. Thus, the traditional data analysis methods can hardly be adopted  
18 to comprehensively study the diagnosis and treatment process of TCM [27]. Machine learning, as a  
19 flexible method for processing complex medical data [28], has been employed in the research of  
20 TCM for further progress [29-32]. In the present paper, an innovative research strategy is established  
21 to explore the feasibility of machine learning in the study of TCM diagnosis and treatment data. The  
22 specific research strategy designed in this paper is illustrated in **Figure 11. Flowchart of the**

1 **research strategy designed in this paper.**

2 First of all, the frequency analysis is employed to study the basic information of insomnia so  
3 that the distribution trend of insomnia in gender and age can be analyzed. As illustrated in **Figure 3**,  
4 the proportion of female patients is significantly higher than that of male patients. There is evidence  
5 that more women than men have insomnia, which is related to the complex interaction of biological,  
6 psychological and social factors [33]. The data of insomnia included in this study indicate that  
7 insomnia patients are mostly middle-aged women, and meanwhile, middle-aged women are more  
8 vulnerable to the influence of perimenopausal syndrome [34], resulting in more female insomnia  
9 patients.

10 Subsequently, the association rules are used to study the data obtained from four diagnostic  
11 methods and treatment based on syndrome differentiation. It can be concluded from **Table 3** and  
12 **Table 4** that most of the results are dominantly related to gender and age, while there is no  
13 significant association among the four diagnoses. Based on these results, it can be found that the  
14 information obtained by the four diagnostic methods is complex and relatively independent in the  
15 diagnosis process of TCM. Whether the information without any objective connection can play a  
16 significant role in the treatment based on syndrome differentiation is the answer to be sought in this  
17 paper. At the same time, two innovative research directions can be exploited by concluding the  
18 objective results. On the one hand, this research can be explored deeply through expanding the  
19 sample size and using other methods to find the internal association of the four diagnoses. On the  
20 other hand, there is no obvious external association among the data, but these data have the statistical  
21 significance. These data can be used for the epidemiological study of TCM on condition that the  
22 sample size is large enough. As can be seen from the **Table 5**, besides the associations that can be

1 obtained from the basic theories, such as the associations between fire and sthenia syndrome, fire and  
2 heat syndrome, there are more new-found associations. For example, the complex syndrome of heart,  
3 liver, spleen, asthenia and sthenia → the heat syndrome, the fire stasis syndrome → the heart, liver.  
4 The following conclusions can be drawn by analyzing the treatment based on syndrome  
5 differentiation with association rules. On the one hand, the results can reveal the connections  
6 between complex syndrome differentiation factors and the syndrome differentiation thoughts of  
7 TCM doctors. On the other hand, after applying the above methods to classify the contents of  
8 treatment based on syndrome differentiation, the results can reflect the priority direction of syndrome  
9 differentiation of insomnia to a certain extent, thus having guiding significance for clinical practice.  
10 In the further study, more research methods can be adopted to verify the dominant diseases of TCM  
11 and explore new syndrome differentiation rules.

12 Moreover, due to the complexity of the data classification and the small sample size of this  
13 paper, the Euclidean distance is selected to evaluate the distance [35]. And the hierarchical clustering  
14 algorithm is employed to analyze the small sample data set in the Euclidean distance. According to  
15 the characteristics that the main prescription is composed of a wide variety of herbs, the hierarchical  
16 clustering algorithm is applied to explore the potential classification rules in the data samples of  
17 TCM. There are two purposes of using the hierarchical clustering method to analyze the main  
18 prescriptions in this study. One is to obtain the compatibility of the main prescriptions used by the  
19 attending doctors from a large number of prescription data. The other is to recode the obtained main  
20 prescriptions into the database. The results indicate that the desired purposes can be achieved by  
21 adopting the hierarchical clustering algorithm to analyze the main prescriptions. The rapid  
22 acquisition of the main prescription of TCM is beneficial for the study of the combinations rules of

1 TCM, but also lays a solid foundation for the overall study of the diagnosis and treatment of the  
2 dominant diseases of TCM.

3 Finally, the random forest method is adopted to discuss the whole diagnosis and treatment  
4 process based on the results of the above-mentioned four parts. As illustrated in Figure 10, the most  
5 significant parameter affecting the judgment results is the syndrome differentiation of five  
6 zang-organs combinations, followed by sleep status, pulse conditions, the syndrome differentiation  
7 of asthenia and sthenia and the syndrome differentiation of six fu-organs combinations. Meanwhile,  
8 emotion status, pathogenic factors combinations and tongue picture (including tongue proper, tongue  
9 color and tongue coating) also have a tremendous effect on the judgment results. Nevertheless,  
10 sleeping duration, insomnia course, syndrome differentiation of cold and heat, and other items except  
11 the tongue picture in the inspection and the auscultation and olfaction have less influence on the  
12 selection of the final main prescription. As can be seen from the above results, doctors take the sleep  
13 status, pulse conditions and tongue picture as the most critical indicators when they are obtaining the  
14 four diagnoses information. In the meantime, the emotional status is also taken into account for  
15 understanding the basic situation of the patient's condition. Based on the the syndrome differentiation  
16 of five zang-organs, and combined with the syndrome differentiation of asthenia and sthenia and the  
17 syndrome differentiation of six fu-organs, a comprehensive analysis is conducted to obtain the final  
18 main prescription in the process of syndrome differentiation. Since sleep duration, course of  
19 insomnia and other factors have little impact on the diagnosis and treatment process, they are only  
20 regarded as reference for the diagnosis and treatment. It can be concluded from the above results that  
21 the random forest algorithm model can be applied to quickly and accurately verify the correctness of  
22 TCM diagnosis and treatment ideas.

1 In this paper, in order to explore the feasibility of machine learning processing TCM medical  
2 record data, the comprehensiveness of outpatient medical record data should be taken into account as  
3 the first priority in the data screening phase. On the premise of ensuring the comprehensiveness of  
4 the data, due to the complexity of outpatient medical records, such as the incomplete and  
5 unquantifiable information contained in the data, the therapeutic efficacy can hardly be verified  
6 thoroughly. The preliminary verification of the effectiveness has been carried out by three TCM  
7 doctors including the attending doctor Prof Yang in the screening phase of the medical record data.  
8 Nonetheless, the verification of the effectiveness has mainly relied on the clinical experience of the  
9 three TCM doctors, which may lead to bias on the therapeutic efficacy of single patient with  
10 insomnia. Taking the above reasons into account, the actual therapeutic efficacy of each patient has  
11 not been fully considered in this study. However, the therapeutic effect is one of the significant  
12 evaluation indexes of TCM diagnosis and treatment. In the light of this, it is quite essential to  
13 introduce evaluation methods of therapeutic effectiveness in our future research. In the meantime, it  
14 is also necessary to further standardize the entry methods of outpatient medical records and establish  
15 follow-up records of patients.

16 Meanwhile, the early outpatient data screening work has been carried out by three TCM doctors,  
17 which was extremely time-consuming. In recent studies, various text mining technologies have been  
18 applied to processing medical records [36]. In our further research, we will make an effort to employ  
19 diverse text mining technologies to extract information, so as to tremendously reduce the waste of  
20 resources and improve the efficiency of analysis and processing of TCM medical record data. In  
21 addition, the machine learning methods applied in this paper are limited, especially for the whole  
22 diagnosis and treatment process, only one algorithm model is used, leading to the lack of diversity of

1 methods. In the future study, we will introduce a variety of algorithm models for comparisons  
2 [37-39], and select the optimal model according to the characteristics of different dominant diseases,  
3 so as to further study the feasibility of machine learning methods in TCM diagnosis research.

4 In the future, we will summarize the previous work and establish a simple and easy-to-use TCM  
5 medical record entry and analysis system based on machine learning, which will dramatically  
6 optimize the process of statistics and analysis of TCM diagnosis and treatment data.

7

8

## 9 **4 Conclusions**

10 The results indicate that the machine learning methods can be effectively applied to deeply mine  
11 and analyze the medical record data of the dominant diseases of TCM. The focus of this study is to  
12 analyze the diagnosis and treatment process of the TCM dominant diseases which includes the  
13 acquisition of the patients' condition information through using four diagnostic methods, and the  
14 flexible application of the syndrome differentiation methods to develop the treatment plan and select  
15 the main prescription. And the normalized research strategy established in this paper can efficiently  
16 filter the unessential diagnosis and treatment information, thus helping TCM doctors to quickly and  
17 efficiently obtain valuable information and crucial rules from a substantial number of medical record  
18 data. In the future, it is essential to establish medical record data entry system, introduce more novel  
19 machine learning methods and improve the therapeutic efficacy evaluation of TCM diagnosis and  
20 treatment.

21

## 22 **List of abbreviations**

Traditional Chinese medicine	TCM <sub>1</sub>
Cheng-Church double clustering algorithm	CC
Sleep latency	SL

## **Declarations**

### **Ethics approval and consent to participate**

Informed consent of the study and a statement on ethics approval was waived because of the retrospective nature and the analysis used anonymous clinical data.

### **Consent for publication**

Not applicable

### **Availability of data and materials**

The datasets used and/or analysed during the current study are available from the corresponding author on reasonable request.

### **Competing interests**

The authors declare that they have no competing interests

### **Funding**

Not applicable

### **Authors' contributions**

YT designed the strategy of this research and wrote the manuscript. ZL implemented the strategy with software, and analyzed the data. SG and YF dedicated in experiment results analysis and manuscript revision. DY participated into analysis implementation and organized discussion of the results. SL and TL have polished this paper with professional perspective and put forward some constructive suggestion to the data analysis process. YT and ZL contributed equally to this work. All authors read and approved the final manuscript.

### **Acknowledgements**

Thanks to Mengyu Chen for her professional writing assistance.

## References

- 1、Gao H, Wang Z, Li Y, Qian Z. Overview of the quality standard research of traditional Chinese medicine. *Front Med.* 2011;5(2):195-202.
- 2、Yan D, Liu J, Wang AT, Yang ZR, Yue SJ, Feng XZ. Exploring Research Ideas of Mechanism of Dominant Diseases in Traditional Chinese Medicine Based on Evidence-Based Medicine. *Zhongguo Zhong Yao Za Zhi.* 2018;43(13):2633-2638.
- 3、Chen YB, Tong XF, Ren J, Yu CQ, Cui YL. Current Research Trends in Traditional Chinese Medicine Formula: A Bibliometric Review from 2000 to 2016. *Evid Based Complement Alternat Med.* 2019;2019:3961395.
- 4、Zhang H, Liu P, Wu X, Zhang Y, Cong D. Effectiveness of Chinese herbal medicine for patients with primary insomnia: A PRISMA-compliant meta-analysis. *Medicine (Baltimore).* 2019;98(24):e15967.
- 5、Singh A, Zhao K. Treatment of Insomnia With Traditional Chinese Herbal Medicine. *Int Rev Neurobiol.* 2017;135:97-115.
- 6、Li Z, Xu C. The fundamental theory of traditional Chinese medicine and the consideration in its research strategy. *Front Med.* 2011;5(2):208-211.
- 7、Wang J, Guo Y, Li GL. Current Status of Standardization of Traditional Chinese Medicine in China. *Evid Based Complement Alternat Med.* 2016;2016:9123103.
- 8、Zhou X, Li Y, Peng Y, et al. Clinical phenotype network: the underlying mechanism for personalized diagnosis and treatment of traditional Chinese medicine. *Front Med.* 2014;8(3):337-346.
- 9、Zhang S, Wu H, Liu J, Gu H, Li X, Zhang T. Medication regularity of pulmonary fibrosis treatment by contemporary traditional Chinese medicine experts based on data mining. *J Thorac Dis.* 2018;10(3):1775-1787.
- 10、Yu XW, Gong QY, Hu KF, Mao WJ, Zhang WM. Research on Ratio of Dosage of Drugs in Traditional Chinese Prescriptions by Data Mining. *Stud Health Technol Inform.* 2017;245:653-656.
- 11、Liu Y, Liu D, Zhang Y, et al. Markov Clustering Analysis-Based Validation for Traditional Chinese Medicine Syndrome Patterns of Poststroke Cognitive Impairment. *J Altern Complement Med.* 2019; 25(11):1140-1148.
- 12、Han F, Tang XD, Zhang B. The Guidelines for the diagnosis and treatment of insomnia in China. *Natl Med J China.* 2017;97(24):1844-1856.
- 13、American Academy of Sleep Medicine. International classification of sleep disorders. 3rd ed. Darien, IL: American Academy of Sleep Medicine; 2014.
- 14、Somek M, Hercigonja-Szekeres M. Decision Support Systems in Health Care - Velocity of Apriori Algorithm. *Stud Health Technol Inform.* 2017;244:53-57.
- 15、Xu R, Wunsch DC 2nd. Clustering algorithms in biomedical research: a review. *IEEE Rev Biomed Eng.* 2010;3:120-54.
- 16、Jones FC, Plewes R, Murison L, et al. Random forests as cumulative effects models: A case study of lakes and rivers in Muskoka, Canada. *J Environ Manage.* 2017;201:407-424.
- 17、Sun J, Yu H, Zhong G, Dong J, Zhang S, Yu H. Random Shapley Forests: Cooperative Game-Based Random Forests With Consistency. *IEEE Trans Cybern.* 2020;10.1109/TCYB.2020.2972956.
- 18、Kang H, Zhao Y, Li C, et al. Integrating clinical indexes into four-diagnostic information contributes to the Traditional Chinese Medicine (TCM) syndrome diagnosis of chronic hepatitis B. *Sci Rep.* 2015;5:9395.
- 19、Ma Y, Zhou K, Fan J, Sun S. Traditional Chinese medicine: potential approaches from modern dynamical complexity theories. *Front Med.* 2016;10(1):28-32.
- 20、Jin ZL, Hu JX, Jin HW, Zhang LR, Liu ZM. Analysis of Traditional Chinese Medicine Prescriptions Based on Support Vector Machine and Analytic Hierarchy Process. *Zhongguo Zhong Yao Za Zhi.* 2018;43(13):2817-2823.
- 21、Lin PY, Chu CH, Chang FY, Huang YW, Tsai HJ, Yao TC. Trends and prescription patterns of traditional Chinese medicine use among subjects with allergic diseases: A nationwide population-based study. *World Allergy Organ J.* 2019 Jan

- 26;12(2):100001.
- 22、 Leem J, Jung W, Kim Y, Kim B, Kim K. Exploring the combinations and modular characteristics of herbs for alopecia treatment in traditional Chinese medicine: an association rule mining and network analysis study. *BMC Complement Altern Med.* 2018;18(1):204.
  - 23、 DeVries Z, Hoda M, Rivers CS, et al. Development of an unsupervised machine learning algorithm for the prognostication of walking ability in spinal cord injury patients. *Spine J.* 2020;20(2):213-224.
  - 24、 Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLoS One.* 2015;10(3):e0118432.
  - 25、 Coden A, Savova G, Sominsky I, et al. Automatically extracting cancer disease characteristics from pathology reports into a Disease Knowledge Representation Model. *J Biomed Inform.* 2009;42(5):937-949.
  - 26、 Rajkumar A, Dean J, Kohane I. Machine Learning in Medicine. *N Engl J Med.* 2019;380(14):1347-1358.
  - 27、 Li Z, Xu C. The fundamental theory of traditional Chinese medicine and the consideration in its research strategy. *Front Med.* 2011;5(2):208-211.
  - 28、 Handelman GS, Kok HK, Chandra RV, Razavi AH, Lee MJ, Asadi H. eDoctor: machine learning and the future of medicine. *J Intern Med.* 2018;284(6):603-619.
  - 29、 Zhao C, Li GZ, Wang C, Niu J. Advances in Patient Classification for Traditional Chinese Medicine: A Machine Learning Perspective. *Evid Based Complement Alternat Med.* 2015;2015:376716.
  - 30、 Wang Y, Jafari M, Tang Y, Tang J. Predicting Meridian in Chinese traditional medicine using machine learning approaches. *PLoS Comput Biol.* 2019;15(11):e1007249.
  - 31、 Jafari M, Wang Y, Amirousetfi A, et al. Unsupervised Learning and Multipartite Network Models: A Promising Approach for Understanding Traditional Medicine[J]. *Frontiers in Pharmacology.*2020; 11.
  - 32、 Lin YC, Huang WT, Ou SC, et al. Neural network analysis of Chinese herbal medicine prescriptions for patients with colorectal cancer. *Complement Ther Med.* 2019;42:279-285.
  - 33、 Suh S, Cho N, Zhang J. Sex Differences in Insomnia: from Epidemiology and Etiology to Intervention. *Curr Psychiatry Rep.* 2018;20(9):69.
  - 34、 Hirose A, Terauchi M, Akiyoshi M, Owa Y, Kato K, Kubota T. Subjective insomnia is associated with low sleep efficiency and fatigue in middle-aged women. *Climacteric.* 2016;19(4):369-374.
  - 35、 Yun-Hong S, Zhen-Xiang L I, Lian-Hui S, et al. Contrast between Mahalanobis distance and Euclidean distance in geochemical exploration processing[J]. *Jilin Geology*, 2008.
  - 36、 Sun W, Cai Z, Li Y, Liu F, Fang S, Wang G. Data Processing and Text Mining Technologies on Electronic Medical Records: A Review. *J Healthc Eng.* 2018;2018:4302425.
  - 37、 Heo J, Yoon JG, Park H, Kim YD, Nam HS, Heo JH. Machine Learning-Based Model for Prediction of Outcomes in Acute Stroke. *Stroke.* 2019 May;50(5):1263-1265.
  - 38、 Yang S, Shen Y, Lu W, Yang Y, Wang H, Li L, Wu C, Du G. Evaluation and Identification of the Neuroprotective Compounds of Xiaoxuming Decoction by Machine Learning: A Novel Mode to Explore the Combination Rules in Traditional Chinese Medicine Prescription. *Biomed Res Int.* 2019 Jul 10;2019:6847685.
  - 39、 Yang Y, Ye Z, Su Y, Zhao Q, Li X, Ouyang D. Deep learning for in vitro prediction of pharmaceutical formulations. *Acta Pharm Sin B.* 2019 Jan;9(1):177-185.

1 **The figure legends**

2

3 **Figure 2. Violin Plot of the age distribution of patients. (Page 11)**

4 The average age of patients is 47, the upper quartile of age is 55, the lower quartile of age is 41, and  
5 the mean square deviation of age is 11. Moreover, the maximum age and minimum age are 79 and 14  
6 respectively.

7

8 **Figure 5. Results of the analysis of the main prescription using hierarchical clustering analysis.**  
9 **(Page 16)**

10 The codes in Figure 5 correspond to the codes in Table 6. The line represents the the main  
11 prescription, and the small circle represents the corresponding herb. Frequency of Category 1 is 573,  
12 frequency of Category 2 is 312. Frequency of Category 3 is 577, which is the cumulative frequency  
13 of all green lines.

14

15 **Figure 6. Flowchart of the diagnosis and treatment ideas of TCM. (Page 17)**

16 **Process 1:** the information of the treatment based on syndrome differentiation is deduced from the  
17 data of four diagnoses. The information of the treatment based on syndrome differentiation includes  
18 five parts: cold and heat, asthenia and sthenia, five zang-organs combinations, six fu-organs  
19 combinations and pathogenic factors combinations.

20 **Process 2:** The main prescription combinations is deduced from the four diagnostic information and  
21 the information of the treatment based on syndrome differentiation.

22

1 **Figure 7. Confusion matrix. (Page 18)**

2 **Process 1** is shown in Figure 7(A-E). **Process 2** is presented in Figure 7(F). In the confusion matrix,  
3 the vertical coordinate is the diagnosis made by doctors in the original medical records, and the  
4 horizontal coordinate represents the predicted value made by the random forest. The corresponding  
5 meanings of independent labels are shown in Table 1 and Table 7. Taking the "cold and heat"  
6 confusion matrix (as depicted in Figure 7(A)) in process 1 as an example, the cold and heat  
7 syndrome can be derived from the data of four diagnoses. The total number of medical record  
8 samples is 654, including 73 cases without cold and heat syndrome, 28 cases with cold syndrome,  
9 419 cases with heat syndrome, and 134 cases with cold and heat complex syndrome. As shown in  
10 Figure 7(A), among the predicted values of the random forest model, the numbers of the cases  
11 accurately predicted by the random forest model for the above syndromes are 55, 14, 418 and 94  
12 respectively. In general, a total of 581 cases are accurately predicted, and the prediction accuracy is  
13 0.89. Similarly, the information of asthenia and sthenia (as depicted in Figure 7(B)), five zang-organs  
14 combinations (as depicted in Figure 7(C)), six fu-organs combinations (as depicted in Figure 7(D)),  
15 pathogenic factors combinations (as depicted in Figure 7(E)) and main prescription combinations (as  
16 depicted in Figure 7(F)) can be derived from the the data of four diagnoses, and the numbers of the  
17 cases accurately predicted by the random forest model are 611, 576, 562, 557 and 559 respectively.

1 **Figure 8.Accuracy, AUC and Micro-F1 score for each model. (Page 18/19)**

2 The accuracy of applying the random forest algorithm models to predict the information of treatment  
3 based on syndrome differentiation through the four diagnostic information is shown in Figure 8(A).  
4 AUC and Micro-F1 score for evaluating the effectiveness and accuracy of the random forest  
5 prediction models are shown in Figure 8(B).

6

7 **Figure 9.ROC curve for each model. (Page 18)**

8 The horizontal coordinate represents false positive rate, and the vertical coordinate represents true  
9 positive rate. AUC values of six different models (represented by lines of different colors)  
10 established by using random forest algorithm are shown in the legend.

11

12 **Figure 10. Transformed eigenvalues obtained by using random forest model. (Page 20)**

13 The transformed eigenvalues of each input parameter of the random forest model (refers to the  
14 process 2 of Figure 6).

**Table 1** Code comparative table of various insomnia related symptoms (Page 7)

Item	Content	Code	Item	Content	Code	Item	Content	Code
<b>Tongue color</b>	Normal	0	<b>Insomnia course</b>	≤3 months	1	<b>Sleep duration</b>	Normal	0
	Pale	1		3 months-1 year	2		Pernoctation	1
	Red	2		1-3 years	3		0-1 hours	2
	Dark	3		3-5 years	4		1-2 hours	3
	Others	4		>5 years	5		2-3 hours	4
<b>Cold and heat</b>	Normal	0	<b>Asthenia and sthenia</b>	Normal	0	<b>Sleep duration</b>	3-4 hours	5
	Cold	1		Asthenia	1		4-5 hours	6
	Heat	2		Sthenia	2		5-6 hours	7
	Cold and heat complex	3		Asthenia and sthenia	3		6-7 hours	8
	Others	4		Complex	4		>7 hours	9
<b>Five zang-organs</b>	Normal	0	<b>Pathogenic factors</b>	Normal	0	<b>Emotional status</b>	Taking medicine to sleep	A
	Heart	1		Phlegm	1		Normal	0
	Liver	2		Fire	2		Anxious	1
	Spleen	3		Blood stasis	3		Fear	2
	Lung	4		Asthenia	4		Nervous	3
<b>Tongue proper</b>	Kidney	5	Others	5	<b>Emotional status</b>	Restlessness	4	
	Normal	0	Normal	0		Timidity	5	
	Enlarged tongue	1	Difficult to fall asleep	1		Irritability	6	
	Thin tongue	2	Dysphylaxia	2		Depressed	7	
	Teeth print on tongue	3	Festless sleep	3		Flusteredness	8	
<b>Six fu-organs</b>	Cleft tongue	4	<b>Sleeping status</b>	Hard to fall asleep again after waking up	4	<b>Concomitant symptoms</b>	Others	9
	The vessels of sublingual	5		Dreaminess	5		Normal	0
	The tongue with ecchymosis	6		Others	6		Headache	1
	Others	6		Normal	0		Dizziness	2
	Normal	0		Sallow complexion	1		Lethargic	3
<b>Pulse conditions</b>	Stomach	1	<b>Others</b>	Flushed cheeks	2	Aversion to cold	4	
	Gallbladder	2		Redness of the eyes	3	Aversion to heat	5	
	Large intestine	3		Dark lip	4	Aversion to cold or heat Irregular	6	
	Small intestine	4		Halitosis	5	Tidal fever	7	
	Bladder	5		Eyes bright	6	Night sweat	8	
<b>Pulse conditions</b>	Sanjiao	6	Others	7	Snoring	9		
	Normal	0	Normal	0	Nocturia	A		
	Thin	1	Yellow	1	Fatigue	B		
	Wiry	2	Thin	2	Dry mouth	C		
	Slippery	3	Slimy	3	Bitter taste in the mouth	D		
<b>Pulse conditions</b>	Rapid	4	<b>The tongue coating</b>	White	4	<b>Concomitant symptoms</b>	Abnormal stool and urine	E
	Deep	5		Scanty	5		Others	F
	Floating	6		Thick	6			
	Others	7		Dye	7			
				Others	8			

**Table 2** Comparative table of the null values and the substitute values (Page 8)

Item	Substitute value
Six fu-organs combinations	Normal
Tongue proper	Normal
Pulse conditions	Normal
Emotional status	Normal
Five zang-organs	Heart
Insomnia course	1-3 years
Asthenia and sthenia	Asthenia and sthenia complex
Cold and heat	Cold and heat complex
Tongue color	Others
Pathogenic factors	Others

**Table 3** Summary of the results using Method 1 (Page 12)

Association	Confidence	Support	Lift
39-year old→Female	1.0	0.02	1.37
59-year old→Female	1.0	0.02	1.37
56-year old→Female	1.0	0.02	1.37
29-year old→Female	1.0	0.02	1.37
(Tongue proper: normal/Tongue color: pale/Tongue coating: thin、yellow/Pulse conditions: thin、wiry) →Female	1.0	0.02	1.37
(Tongue proper: normal/Tongue color: pale/Tongue coating: thin、yellow) →Female	1.0	0.04	1.31
(Tongue proper: teeth print on tongue/Tongue color: red、dark/Tongue coating: thin、yellow) →Female	1.0	0.03	1.3
(Tongue proper: normal/Tongue color: pale、dark/Tongue coating: thin、yellow) →Female	1.0	0.03	1.29
(Sleep duration:normal/Sleeping status:Difficult to fall asleep/Insomnia cours >5 years)→Female	0.9	0.03	1.29
(Tongue proper: teeth print on tongue/Tongue color: red/Tongue coating: thin、yellow/Pulse conditions: wiry、rapid) →Female	0.9	0.02	1.28
62-year old→Female	0.9	0.02	1.28
50-year old→Female	0.9	0.05	1.28
30-year old→Female	0.9	0.02	1.25

**Table 4** Summary of the results using Method 2 (Page 12)

Association	Confidence	Support	Lift
(Sleep duration:3-4 hours /Sleeping status:Difficult to fall asleep/Insomnia course:3 months-1 year)→47-year old	1.0	0.01	14.22
(Sleep duration:3-4 hours /Sleeping status:Difficult to fall asleep/Insomnia course:3 months-1 year)→Male	1.0	0.01	3.72
57-year old→Female	1.0	0.01	1.37
38-year old→Female	1.0	0.01	1.37
39-year old→Female	1.0	0.02	1.37
34-year old→Female	1.0	0.01	1.37
59-year old→Female	1.0	0.02	1.37
56-year old→Female	1.0	0.02	1.37
29-year old→Female	1.0	0.02	1.37
Pulse conditions: thin、wiry、rapid、deep→Female	1.0	0.01	1.37
Pulse conditions: wiry、rapid、deep→Female	1.0	0.01	1.37
(Tongue proper: normal/Tongue color: pale/Tongue coating: thin、white) →Female	1.0	0.01	1.37

**Table 5** Summary of the results analyzing the syndrome differentiation factors by adopting association rules (Page 15)

Association	Confidence	Support	Lift
Pathogenic factors: fire→Asthenia and sthenia: sthenia	1	0.08	3.86
Pathogenic factors: fire→Cold and heat: heat	1	0.08	1.56
Pathogenic factors: fire/Cold and heat: heat→Asthenia and sthenia: sthenia	1	0.08	3.86
Asthenia and sthenia: sthenia/Pathogenic factors: fire→Cold and heat: heat	1	0.08	1.56
Pathogenic factors: fire→Asthenia and sthenia: sthenia/Cold and heat: heat	1	0.08	1.45
Pathogenic factors: fire、 blood stasis→Asthenia and sthenia: sthenia	0.96	0.07	3.7
Pathogenic factors: fire、 blood stasis→Cold and heat: heat	0.96	0.07	1.49
Pathogenic factors: fire、 blood stasis、 asthenia→Asthenia and sthenia: asthenia and sthenia complex	0.92	0.08	3.22
Asthenia and sthenia: sthenia→Cold and heat: heat	0.92	0.26	1.43
Pathogenic factors: fire→Five zang-organs: heart、 liver	0.89	0.08	2.8
Pathogenic factors: asthenia→Asthenia and sthenia: asthenia	0.89	0.17	1.97
Pathogenic factors: fire、 blood stasis→Five zang-organs: heart、 liver	0.85	0.07	2.69
Five zang-organs: heart、 liver/Six fu-organs: gallbladder→Cold and heat: heat	0.85	0.13	1.33
Five zang-organs: heart、 liver→Cold and heat: heat	0.85	0.32	1.32
Cold and heat: normal→Asthenia and sthenia: asthenia	0.82	0.11	1.83
Five zang-organs: heart、 liver、 spleen/Asthenia and sthenia: asthenia and sthenia complex→Cold and heat: heat	0.82	0.1	1.28
Five zang-organs: heart、 liver、 spleen/Six fu-organs: normal→Cold and heat: heat	0.8	0.08	1.25
Pathogenic factors: phlegm、 asthenia→Asthenia and sthenia: asthenia	0.79	0.09	1.75

Five zang-organs: heart、 spleen→Asthenia and sthenia: asthenia	0.76	0.21	1.69
Five zang-organs: heart、 spleen/Six fu-organs: normal→Asthenia and sthenia: asthenia	0.75	0.09	1.67
Pathogenic factors: blood stasis、 asthenia→Asthenia and sthenia: asthenia	0.73	0.2	1.63
Five zang-organs: heart、 liver、 spleen→Cold and heat: heat	0.73	0.23	1.14
Asthenia and sthenia: asthenia and sthenia complex→Cold and heat: heat	0.72	0.28	1.12

---

**Table 6** Correspondence between the herbs and the codes (Page 16)

Chinese herbal medicine	Code	Chinese herbal medicine	Code
Spine date seed	1	Pinellia ternate	26
Glycyrrhiza	2	White peony root	27
Anemarrhena	3	Atractylodes Macrocephala	28
Poria cocos	4	Prepared radix rehmanniae	29
Ligusticum wallichii	5	Chinese yam	30
Caulis polygoni multiflori	6	Cornus officinalis	31
Lily	7	Cortex moutan	32
Seed of oriental arborvitae	8	Magnolia officinalis	33
Red peony root	9	Tasteless preserved soybean	34
Gentian	10	Arillus longan	35
Scutellaria	11	Astragalus	36
Gardenia	12	White hyacinth bean	37
Alisma orientale	13	Villous amomum	38
Caulis Aristolochiae Manshuriensis	14	Semen coicis	39
Plantain	15	Os draconis (longgu)	40
Angelica sinensis	16	Oyster	41
Radix rehmanniae	17	Lanceolata	42
Ginseng	18	Radix aucklandiae	43
Polygala	19	Placenta hominis	44
Schisandra chinensis	20	Blighted wheat	45
Coptis chinensis	21	Leonurus japonicus	46
Bamboo shavings	22	Cinnamon	47
Citrus aurantium	23	Eucommia	48
Tangerine peel	24	Jianqu	49
Bupleurum	25		

**Table 7** Correspondence between the main prescriptions and the serial numbers (Page 16)

Serial number	Main prescription	Frequencies
1	Spine date seed, Glycyrrhiza, Anemarrhena, Poria cocos, Ligusticum wallichii, Caulis polygoni multiflori, Lily, Seed of oriental arborvitae, Red peony root	573
2	Bupleurum, Pinellia ternate, Astragalus, Os draconis (longgu), Oyster	312
3	Gentian, Scutellaria, Gardenia, Alisma orientale, Caulis Aristolochiae Manshuriensis, Plantain, Angelica sinensis, Radix rehmanniae	20
4	Angelica sinensis, White peony root, Atractylodes Macrocephala	190
5	Ginseng, Poria cocos, Polygala, Angelica sinensis, Schisandra chinensis, Seed of oriental arborvitae, Radix rehmanniae, Spine date seed	23
6	Coptis chinensis, Bamboo shavings, Citrus aurantium, Tangerine peel	108
7	Bamboo shavings, Citrus aurantium, Tangerine peel	14
8	Bupleurum, Pinellia ternate, Astragalus	108
9	White peony root, Atractylodes Macrocephala, Villous amomum, Ginseng, Chinese yam, Semen coicis	2
10	Atractylodes Macrocephala, Angelica sinensis, Arillus longan, Polygala, Astragalus	30
11	Prepared radix rehmanniae, Chinese yam, Cornus officinalis, Cortex moutan	19
12	Magnolia officinalis	22
13	Tasteless preserved soybean, Gardenia	41

**Table 8** Correspondence between the repeated herb combinations and the serial numbers (**Page 16**)  
combinations

---

Spine date seed, Poria cocos, Seed of oriental arborvitae

Bupleurum, Pinellia ternate, Astragalus

Bamboo shavings, Citrus aurantium, Tangerine peel

White peony root, Atractylodes Macrocephala

Angelica sinensis, Radix rehmanniae

Angelica sinensis, Polygala

---

# Figures

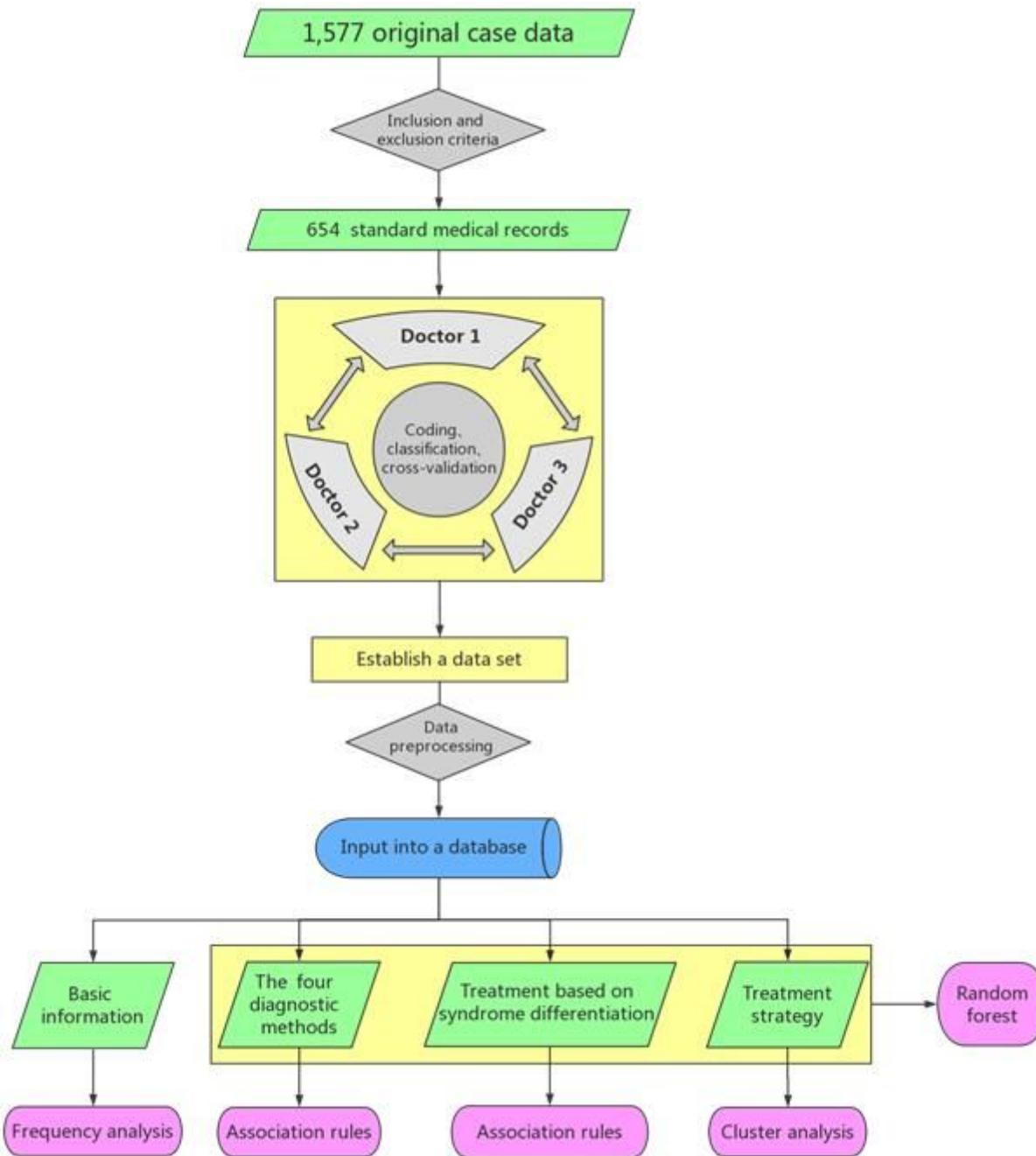
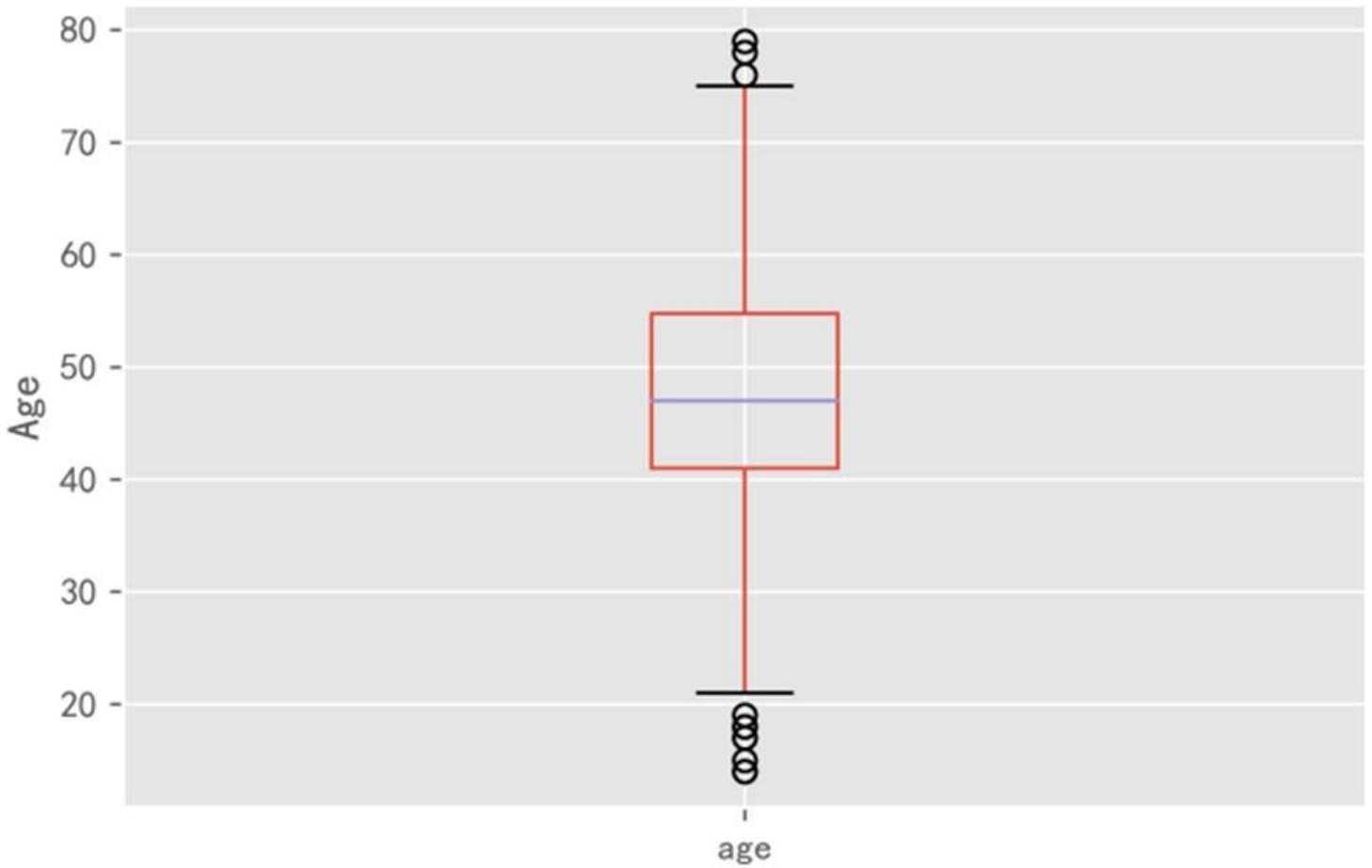


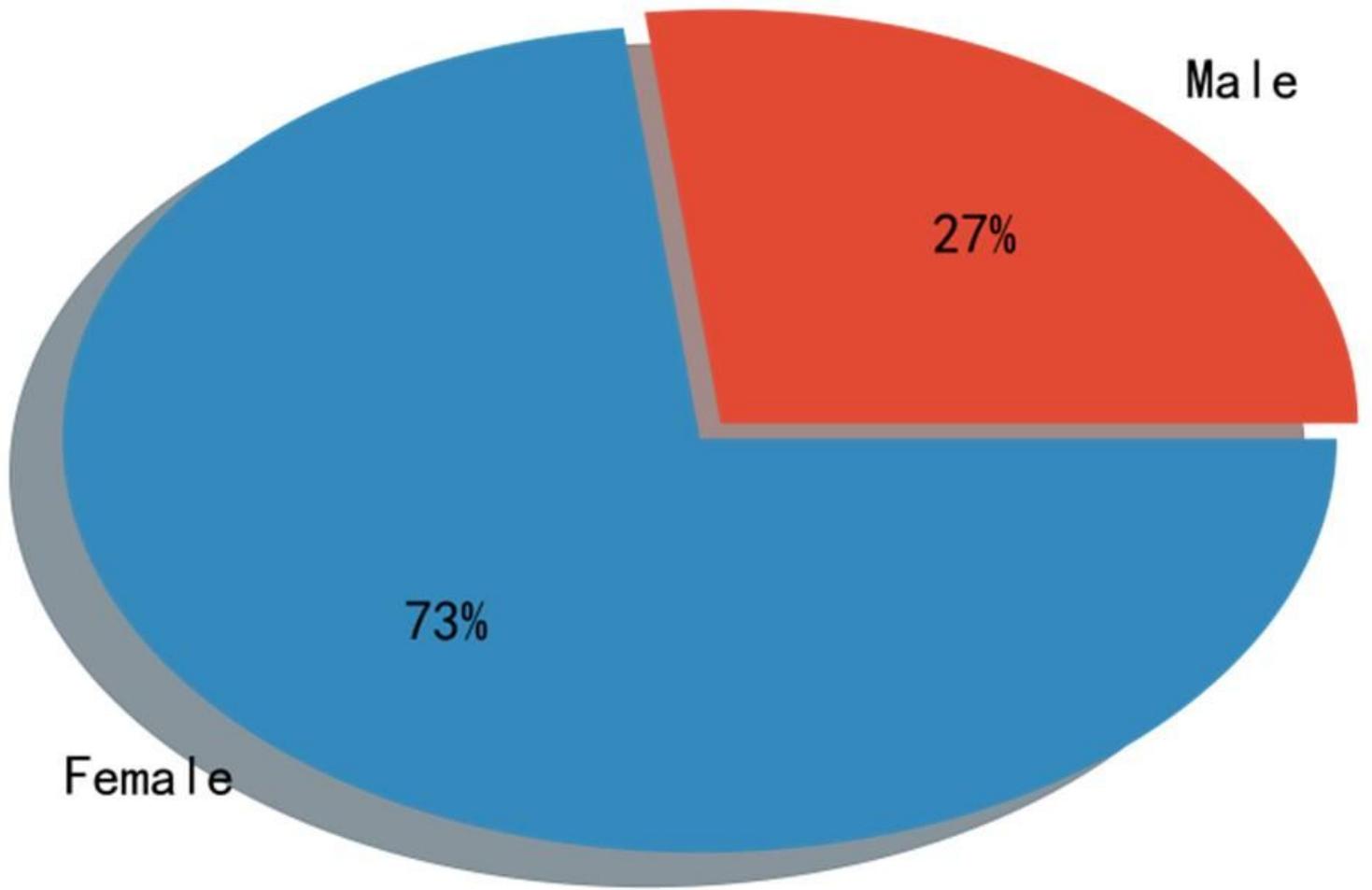
Figure 1

Flowchart of the research strategy designed in this paper.



**Figure 2**

Violin Plot of the age distribution of patients. (Page 11) The average age of patients is 47, the upper quartile of age is 55, the lower quartile of age is 41, and the mean square deviation of age is 11. Moreover, the maximum age and minimum age are 79 and 14 respectively.



**Figure 3**

Pie chart of the gender distribution of patients

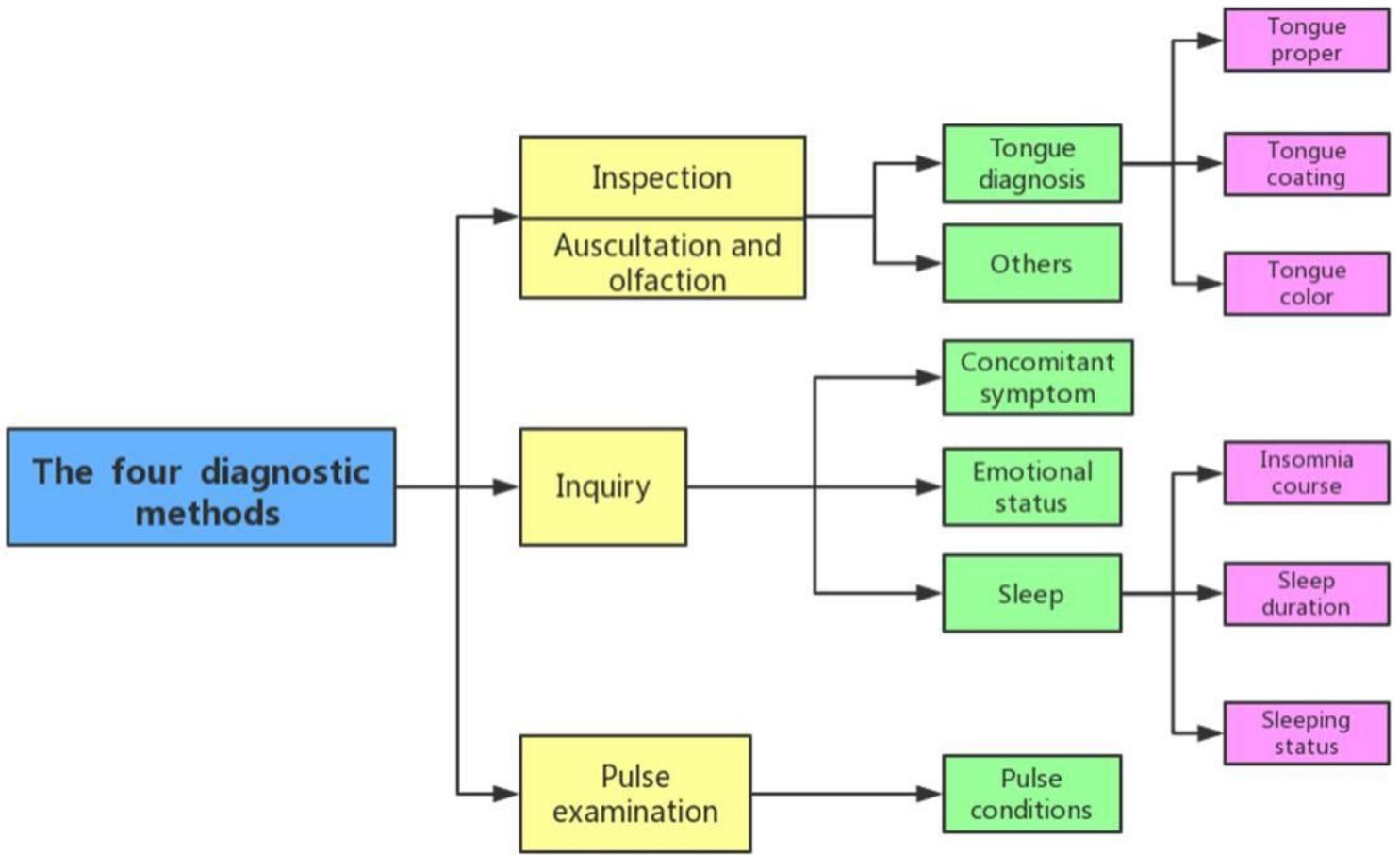
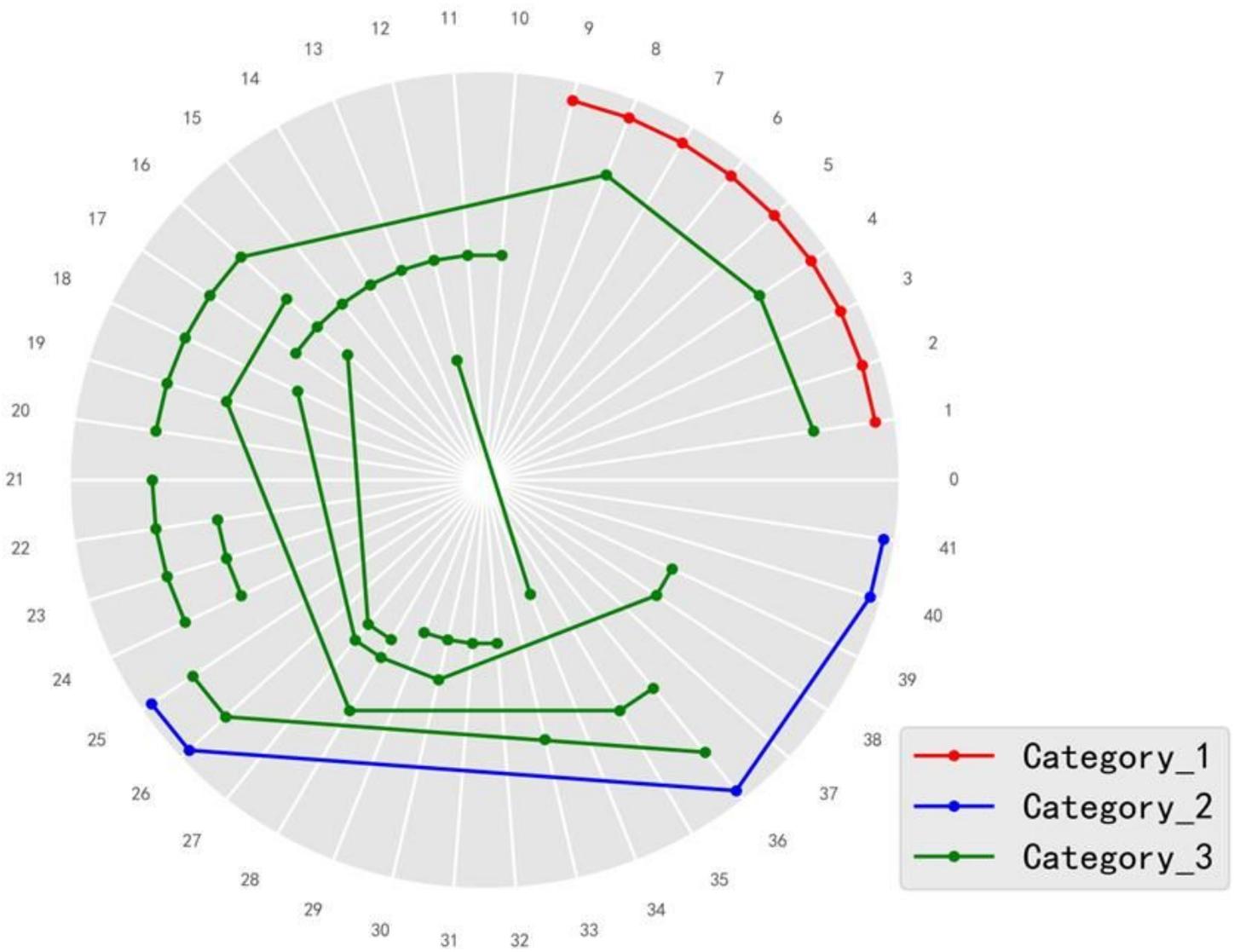


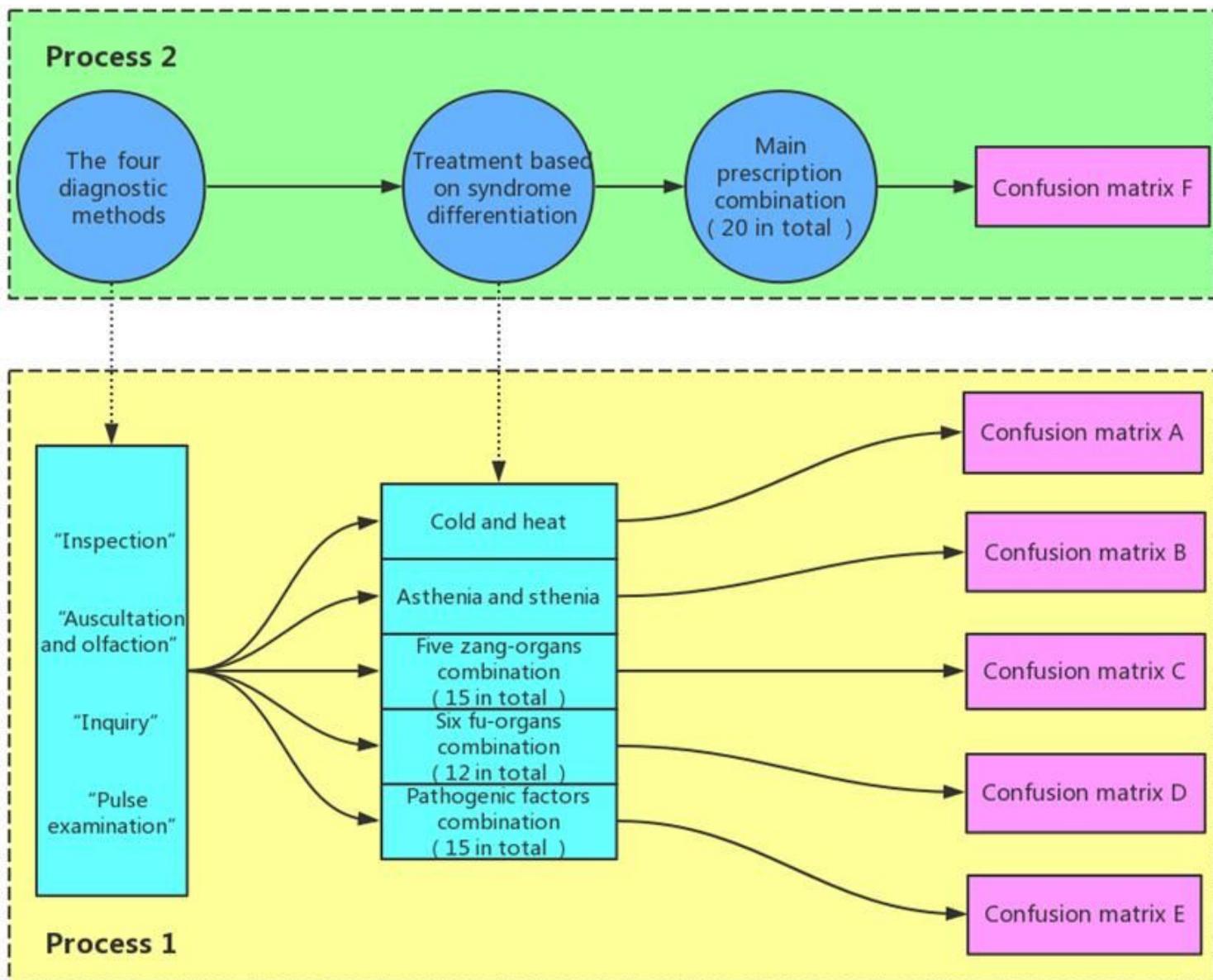
Figure 4

Classification of four diagnostic methods



**Figure 5**

Results of the analysis of the main prescription using hierarchical clustering analysis. (Page 16) The codes in Figure 5 correspond to the codes in Table 6. The line represents the the main prescription, and the small circle represents the corresponding herb. Frequency of Category 1 is 573, frequency of Category 2 is 312. Frequency of Category 3 is 577, which is the cumulative frequency of all green lines.



**Figure 6**

Flowchart of the diagnosis and treatment ideas of TCM. (Page 17) Process 1: the information of the treatment based on syndrome differentiation is deduced from the data of four diagnoses. The information of the treatment based on syndrome differentiation includes five parts: cold and heat, asthenia and sthenia, five zang-organs combinations, six fu-organs combinations and pathogenic factors combinations. Process 2: The main prescription combinations is deduced from the four diagnostic information and the information of the treatment based on syndrome differentiation.

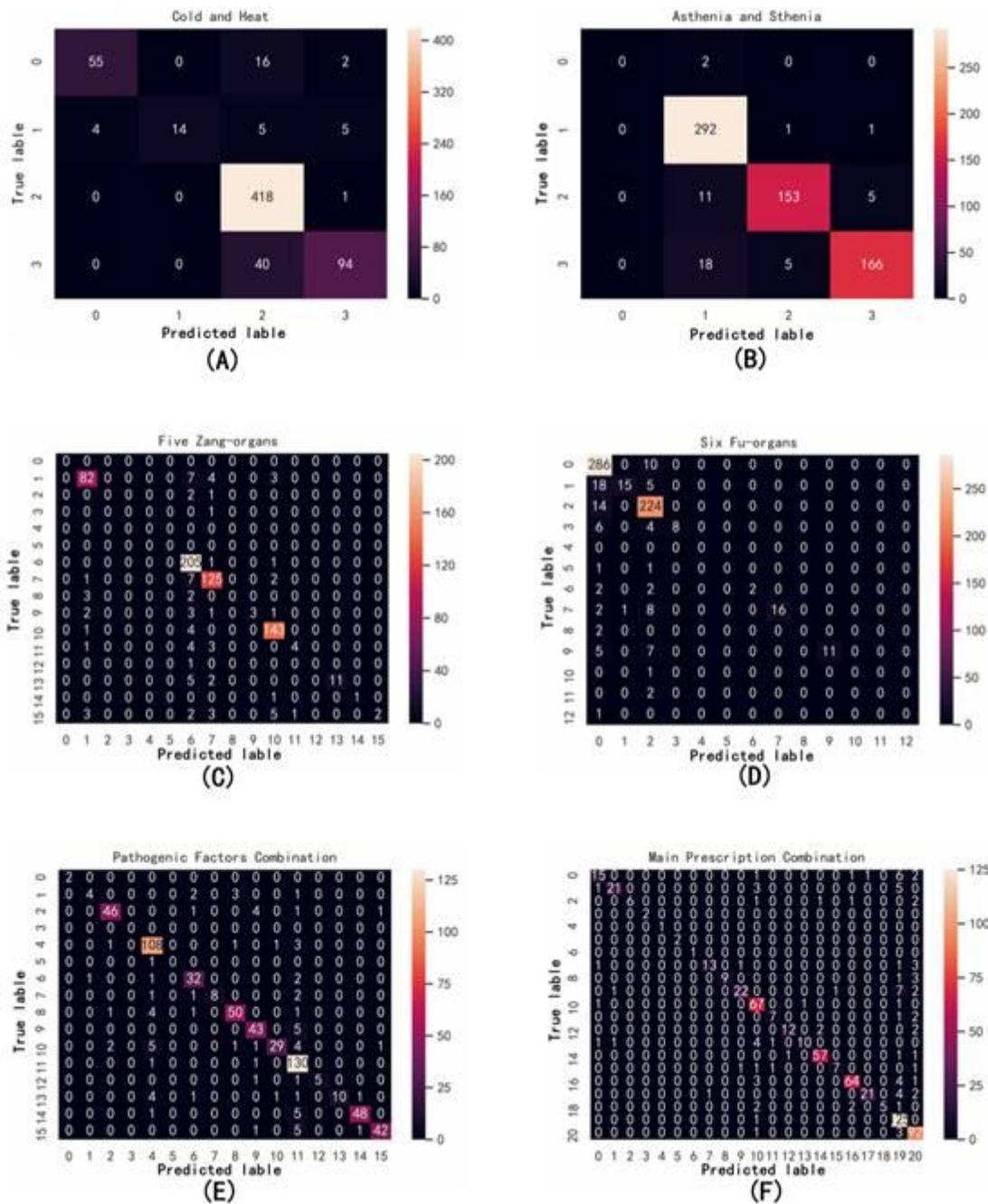
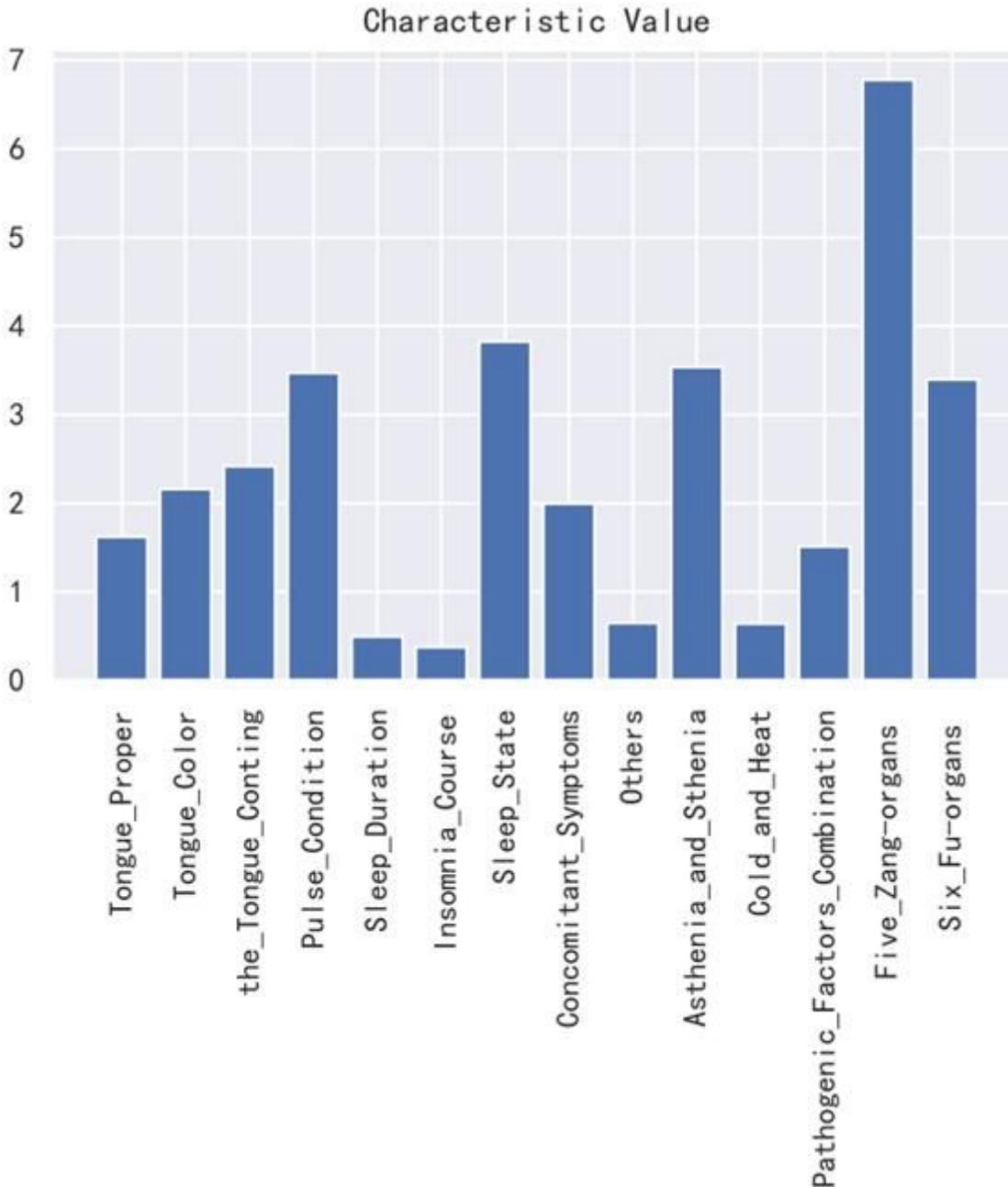


Figure 7

Confusion matrix. (Page 18) Process 1 is shown in Figure 7(A-E). Process 2 is presented in Figure 7(F). In the confusion matrix, the vertical coordinate is the diagnosis made by doctors in the original medical records, and the horizontal coordinate represents the predicted value made by the random forest. The corresponding meanings of independent labels are shown in Table 1 and Table 7. Taking the "cold and heat" confusion matrix (as depicted in Figure 7(A)) in process 1 as an example, the cold and heat syndrome can be derived from the data of four diagnoses. The total number of medical record samples is 654, including 73 cases without cold and heat syndrome, 28 cases with cold syndrome, 419 cases with heat syndrome, and 134 cases with cold and heat complex syndrome. As shown in Figure 7(A), among the predicted values of the random forest model, the numbers of the cases accurately predicted by the

random forest model for the above syndromes are 55, 14, 418 and 94 respectively. In general, a total of 581 cases are accurately predicted, and the prediction accuracy is 0.89. Similarly, the information of asthenia and sthenia (as depicted in Figure 7(B)), five zang-organs combinations (as depicted in Figure 7(C)), six fu-organs combinations (as depicted in Figure 7(D)), pathogenic factors combinations (as depicted in Figure 7(E)) and main prescription combinations (as depicted in Figure 7(F)) can be derived from the the data of four diagnoses, and the numbers of the cases accurately predicted by the random forest model are 611, 576, 562, 557 and 559 respectively.



**Figure 8**

Accuracy, AUC and micro-F1 score for each model. (Page 18/19) The accuracy of applying the random forest algorithm models to predict the information of treatment based on syndrome differentiation

through the four diagnostic information is shown in Figure 8(A). AUC and micro-F1 score for evaluating the effectiveness and accuracy of the random forest prediction models are shown in Figure 8(B).

## Figure 9

Figure 9 not available with this version.

**Figure 10**

Figure 10 not available with this version.

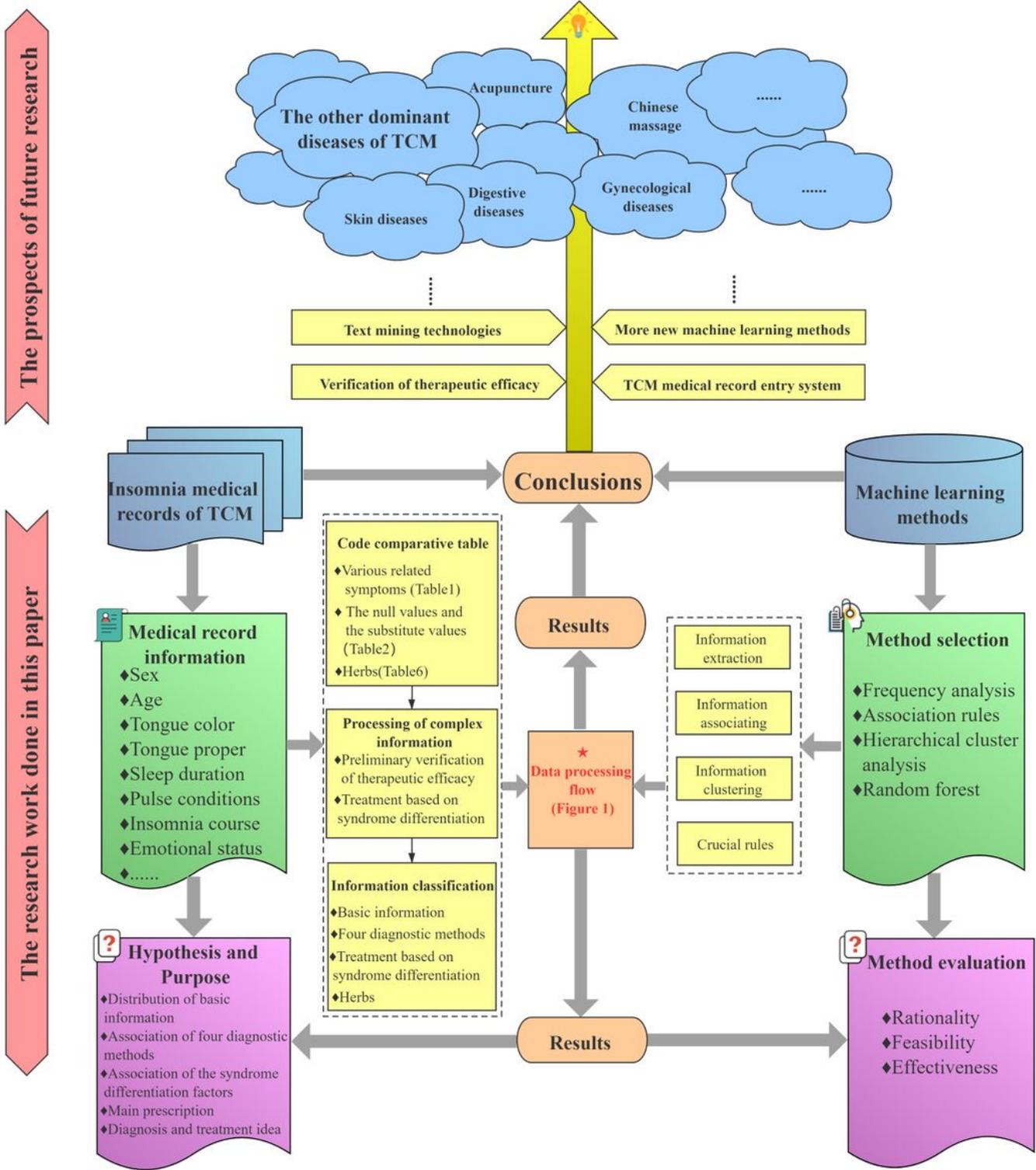


Figure 11

Figure caption not provided with this version.