

Cluster-Based Ensemble Learning Model for Rapid Detection of Aortic Dissection

Yan Gao

Central South University

Min Wang

Central South University

Guogang Zhang

Central South University

Lingjun Zhou

Central South University

Jingming Luo

Central South University

Lijue Liu (✉ ljliu@csu.edu.cn)

Central South University

Research Article

Keywords: aortic dissection, resampling, imbalanced data, ensemble learning, bagging

Posted Date: July 8th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-653868/v1>

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

[Read Full License](#)

1

2 **Cluster-based ensemble learning model for rapid detection**
3 **of Aortic Dissection**

4

5 Yan Gao¹, Min Wang¹, Guogang Zhang², Lingjun Zhou¹, Jingming Luo², Lijue Liu¹§

6 ¹School of Automation, Central South University, Hunan Province, China.

7 ²Xiangya Hospital, Central South University, Hunan Province, China.

8

9

10 § Corresponding author

11 Lijue Liu: ljliu@csu.edu.cn

12

13

14 Email addresses:

15 Yan Gao: gaoyan@csu.edu.cn

16 Min Wang: wm490602605@csu.edu.cn

17 Lijue Liu: ljliu@csu.edu.cn

18 Guogang Zhang: zhangguogang@csu.edu.cn

19 Lingjun Zhou: csu_0918@163.com

20 Jingming Luo: jmluo0618@139.com

21

Abstract

1

2 **Background**

3 Aortic dissection (AD) is a rare and high-risk cardiovascular disease with dangerous
4 morbidity and high mortality, so it needs rapid diagnosis and timely treatment.
5 However, due to its complex and changeable clinical manifestations and the lack of
6 special symptoms and signs, it is easy to cause missed diagnosis and misdiagnosis.

7 **Methods**

8 The data set used in this paper comes from 53213 patients, which collected from
9 XiangYa Hospital in Hunan Province from 2008 to 2016. The data includes 802
10 patients with aortic dissection and 52411 patients with non-aortic dissection. In order
11 to help clinicians predict AD, we designed an ensemble learning model based on
12 clustering: Cluster Random Under-sampling Smote-Tomek-link Bagging
13 (CRST-Bagging). This model combines the advantages of clustering-based compound
14 resampling (CRST) method and Bagging ensemble classifier. It achieves good results
15 on aortic dissection data sets.

16 **Results**

17 The model validates the effectiveness of the CRST sampling method on the AD data
18 set. We compared the CRST-Bagging model with the classical ensemble models
19 RUSBoost and SMOTE-Bagging on the AD data set. The experimental results show
20 that the CRST-Bagging model has the best performance in the detection of AD.
21 Model's accuracy and recall rate are 83.6% and 80.7% respectively. And the F1 value
22 is 82.1%, which is 4.8% and 1.6% higher than that of RUSBoost and
23 SMOTE-Bagging model.

24 **Conclusions**

25 The model proposed in this paper can be used as an auxiliary diagnostic model of AD
26 to provide reference for clinical medicine. The model can also help doctors to judge
27 whether patients need further imaging examination. Thus help to reduce the rate of
28 clinical misdiagnosis and missed diagnosis of AD.

29

30 **Keywords:** aortic dissection, resampling, imbalanced data, ensemble learning,
31 bagging

32

1 **1. Background and introduction**

2 Aortic dissection (AD) is a medial rupture caused by intramural hemorrhage, which
3 leads to the separation of aortic wall layer, followed by the separation of true and false
4 lumen[1]. AD is a dangerous cardiovascular disease with dangerous morbidity, many
5 complications and high mortality. The clinical manifestations of AD are complex and
6 changeable. They lack of special symptoms and signs. And the location, lesion degree
7 and scale of AD are different. So the clinical manifestations and severity are different.
8 In addition, clinicians tend to observe the common symptoms of AD, such as chest
9 pain and back pain. But for painless patients, atypical symptoms make the diagnosis
10 more difficult. It is easy to cause missed diagnosis and misdiagnosis[2]. According to
11 the clinical statistics, the misdiagnosis rate of AD is more than 1/3 in the actual cases
12 of AD [3][4][5]. Mortality can reach as high as 50% within a week of onset and
13 between 60 and 70% within a month. With the help of scientific methods and effective
14 techniques, the timely diagnosis of AD by clinicians is the most effective means to
15 save patients' lives.

16 In recent years, the application of artificial intelligence in the field of medical and
17 health care has attracted much attention. At present, various auxiliary diagnosis
18 systems in the medical field have emerged one after another. Huang et al. [6] used an
19 enhanced resampling method of electronic medical records to classify and predict
20 Major adverse cardiac events (MACEs) of acute coronary syndrome (ACS). Zhou et
21 al. [7] proposed an interpretable pattern discovery method from the perspective of
22 statistical learning methods to interpret clinical chest data and make classification
23 predictions. Song et al. [8] used deep learning technology to develop a
24 histopathological detection system for gastric cancer detection that can be used for
25 clinical diagnosis. Kamal et al. [9] proposed a novel and interpretable PAVE model
26 for sepsis attack prediction and mortality prediction estimates. Song et al. [10] used
27 Bayesian subgroup analysis to evaluate the conditioned treatment effect of adjuvant
28 treatment for patients with synovial sarcoma and help clinicians choose treatment
29 options. Xia et al. [11] proposed a multi-level hypoglycemia early warning system
30 based on sequential pattern mining based on continuous blood glucose monitoring
31 data. With the accumulation of a large amount of medical data, researchers have used
32 deep learning methods to achieve classification prediction[12][13] in the field of
33 cardiovascular diseases. For example, Guo et al. [14] used the LSTM model to predict
34 cardiovascular health trajectories in time series electronic health records; Cheng et al.
35 [15] used deep learning techniques to classify AD using contrast-enhanced CT images.
36 Among 1,000 CT images from 20 patients, the accuracy rate reached 85.0%. However,
37 in the actual situation, due to the lack of clinician experience or unsupported
38 examination equipment, it's difficult to carry out relevant imaging examination in time
39 resulting in patient AD may not be diagnosed by means of CT image, leading to
40 missed diagnosis and misdiagnosis, which threatens the life safety of patients.
41 Therefore, different from the above methods, we studied the diagnostic and predictive

1 method of AD based on the routine examination data of patients, so as to help doctors
2 judge whether patients need further imaging examination.

3 The rarity of AD leads to the significant imbalance in the data set. If the traditional
4 machine learning technology is applied to the aortic dissection data set, the model
5 tends to fit large samples, showing high accuracy and low recall rate, so the model's
6 generalization ability is low. Therefore, imbalanced learning [16][17] and ensemble
7 learning [18][19][20][21] were combined to predict aortic dissection data.

8 According to the characteristics of AD data, this paper proposes a cluster-based
9 ensemble learning model: Cluster Random Under-sampling Smote-Tomek-link
10 Bagging (Hereinafter referred to as CRST-Bagging) to help clinicians detect AD in
11 clinical practice. This model includes two parts: Cluster-Based resampling (CRST)
12 and Bagging classifier. The resampling method CRST combines the advantages of the
13 over-and-under sampling method. It overcomes the difficulties in the detection of AD
14 caused by imbalanced data. Bagging classifier is used to improve the generalization
15 ability of the learning model. To demonstrate the effectiveness of the CRST-Bagging
16 approach, we compared it with the classic ensemble models RUSBoost and
17 Smote-Bagging on AD datasets. Experimental results show that the proposed model is
18 superior to other models, which prove the effectiveness of the model.

19 The main contributions of this paper can be summarized as follows:

- 20 ● In this paper, missing value processing , feature screening and
21 dimension-reduction visualization were performed in the AD data set. These
22 methods enable us to have a priori knowledge of the distribution of AD data,
23 which can be used in clinical medicine to explore the pathological
24 mechanism of AD.
- 25 ● A new compound resampling method, CRST is proposed. This method
26 combines the advantages of clustering ideology and SMOTE + Tomek-Link
27 sampling methods. This method not only makes the collected samples
28 effectively represent the characteristics of different samples, but also ensures
29 the randomness of sampling, which can effectively reduce the imbalanced
30 ratio of AD data.
- 31 ● CRST-Bagging ensemble learning model is proposed to predict AD disease.
32 Through experimental comparison and analysis, the model shows excellent
33 performance and good generalization ability on AD data sets. Therefore, this
34 model can be used for clinical auxiliary diagnosis.

35 The rest of this article is arranged as follows. In the second section, we introduce the
36 data set used in this paper, our resampling method CRST and the imbalanced
37 algorithm integration model. In the third section we present our experimental results
38 and model performance evaluation. In the fourth section, the experimental results are
39 discussed. Finally, the summarization and the discussion of future work are described
40 in the last section.

2. Data

2.1 Data Overview

The dataset used in this paper comes from the examination indicators of 53213 patients, which collected from XiangYa Hospital in Hunan Province from 2008 to 2016. The data includes 802 patients with AD and 52411 patients without AD. The dataset has 71-dimensional features and 1-dimensional tags (The table of Description of Features is included as additional file 1 of the supplementary material). This dataset has a high imbalanced ratio, with the number of AD samples approximately 67 times that of non-AD samples. In addition, this paper also uses a test set to better verify the classification performance and generalization ability of our model. The test set includes the examination indicators of 235 patients from the same hospital, and the data format is the same as the above data set.

2.2 Data preprocessing

Firstly, some non-numerical indicators are normalized by binary coding, and then standardized. Secondly, we made statistics on the missing rate of samples and features in the original AD dataset as shown in Figure 1 (the abscissa represents the features and the ordinate represents the missing rate). Six features with a deletion rate of more than 50% were found, namely Plasma antithrombin III antigen determination(r_28)(missing rate is 81.5%), Plasma plasminogen antigen determination (r_20)(missing rate is 80.7%), Hypersensitivity thyrotropin (r_64) (missing rate is 75.6%), erythrocyte sedimentation rate (r_59) (missing rate is 63.8%), D-dimer (r_19) (missing rate is 62.6%), free triiodothyronine (FT3) (r_62) (missing rate is 51.9%).

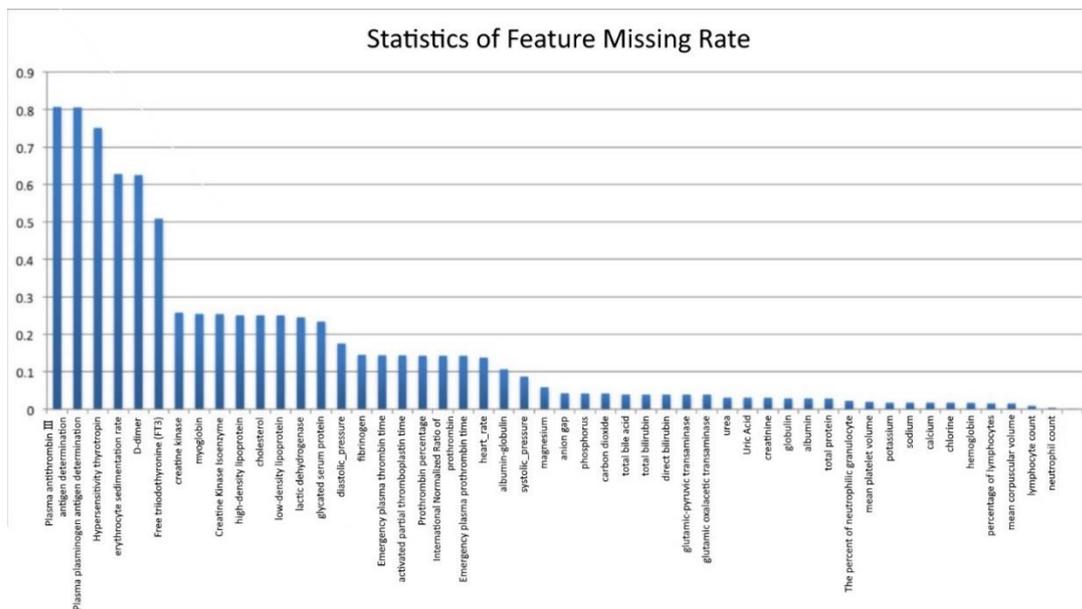


Figure 1. Feature Missing Rate Statistics

24

25

26

27

Due to the high-missing rate of the above-mentioned six-dimensional features, it is difficult to complete them. The general method is to delete them. However, the

1 etiology and related diagnostic indicators of AD are not yet clear, we cannot directly
 2 determine whether the missing feature indicators are key indicators, so they cannot
 3 simply be deleted. Therefore, the XGboost method is used to analyze feature
 4 importance [22][23]. The result is shown in Figure 2, with the horizontal coordinate
 5 as feature numbers and the vertical ordinate as feature importance scores.

6 From figure 2, we find that among the 6-dimensional features with a deletion rate
 7 greater than 50%, the feature importance scores of free triiodothyronine (FT3) and
 8 D-dimer are ranked in the top 10, which indicates that these two features are
 9 important for detecting whether a patient suffers from AD. Therefore, we only remove
 10 the four characteristics of Plasma antithrombin III antigen determination, Plasma
 11 plasminogen antigen determination, Hypersensitivity thyrotropin, erythrocyte
 12 sedimentation rate. Free triiodothyronine (FT3) and D-dimer remain and are
 13 complemented with the remaining features. The adjusted new sample set size is
 14 (53213,67).

15 In this paper, the data set is filled by the method of classified random filling method.
 16 Compared with ordinary random filling, the method of random filling by class is to
 17 fill the positive and negative samples respectively. The missing values of the samples
 18 are randomly filled with the non-null values of the same kind of samples. This filling
 19 method can effectively avoid the intersection of feature values.

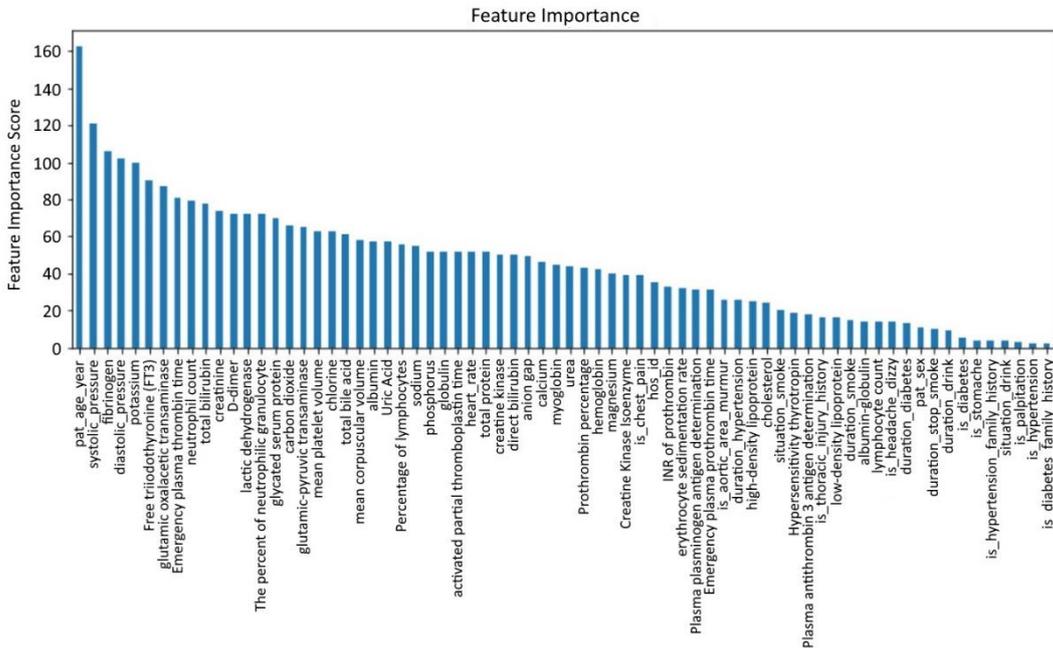


Figure 2. Feature importance analysis

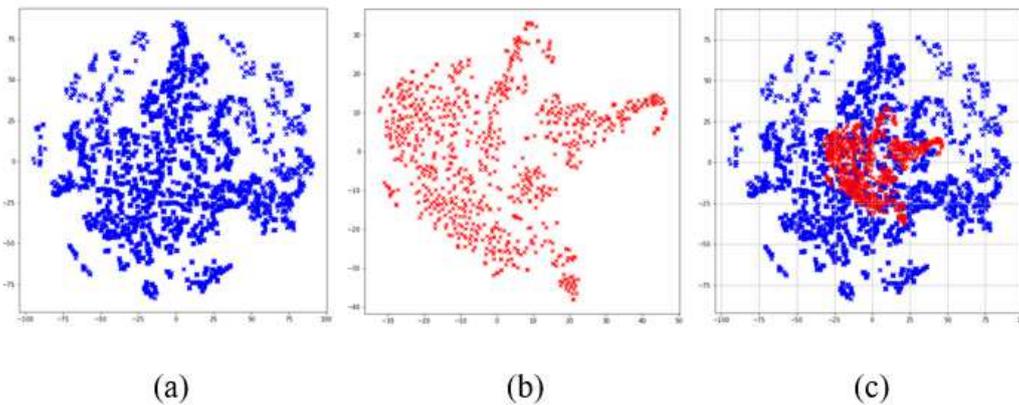
2.3 Dimensionality reduction visualization

In order to have a more systematic and in-depth understanding of the high-dimensional AD data set, this paper uses the method of dimensionality reduction visualization to analyze the data distribution. This allows us to understand the data more intuitively and provide information for the design of AD prediction algorithm.

We analyze the existing methods of dimensionality reduction. T-SNE[24] algorithm

1 can retain both global and local data structures. Therefore, we used the T-SNE
2 algorithm to reduce the dimensionality of the dataset. We analyze the distribution of
3 samples through dimensionality reduction and visualization. By observing the data
4 distribution, it is concluded that the clustering algorithm is feasible to improve the
5 under-sampling method.

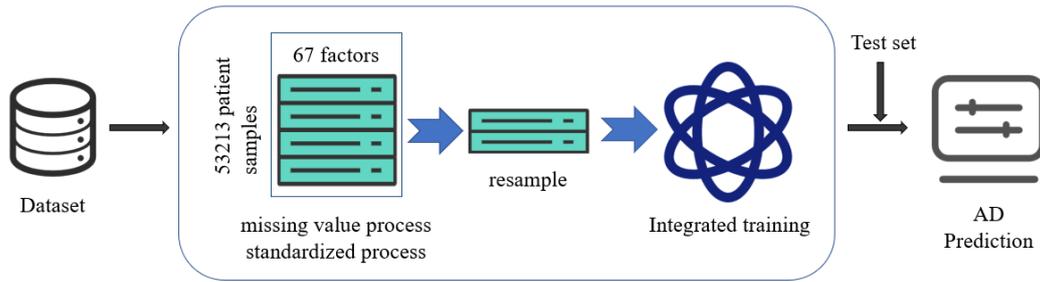
6 The results of the AD dataset using the t-SNE method for dimensionality reduction
7 are shown in Figure 3. The red sample point is the positive sample, and the blue
8 sample point is the negative sample. As can be seen from figure 3, the data
9 distribution of positive samples is agglomerated. This shows that there is a certain
10 similarity between the cases of AD. Therefore, from the analysis of data distribution,
11 the under-sampling method of clustering is effective. At the same time, it can also be
12 seen that there is obvious overlap between the two kinds of samples in space from the
13 visualization of data distribution. Therefore, it is necessary to construct a nonlinear
14 classification model.



15
16 Figure 3. AD dataset t-SNE dimensionality reduction map. (a) The data distribution of negative samples
17 after dimensionality reduction. (b) Data distribution of positive samples after dimensionality
18 reduction. (c) The data distribution of positive and negative samples in the same space after
19 dimensionality reduction.

20 3. Methods

21 According to the data characteristics of AD, this paper proposes an ensemble learning
22 model based on clustering: Cluster Random Under-sampling Smote-Tomek-link
23 Bagging (CRST-Bagging). The model structure is shown in figure 4. Data
24 pre-processing is carried out first (see section 2). Then in the second step, a
25 clustering-based resampling algorithm is proposed to resample the imbalanced data
26 set to reduce the imbalanced ratio of the data. Finally, the Bagging ensemble model is
27 used to construct a powerful nonlinear classifier to predict AD. The methods are
28 described in detail below.



1
2

Figure 4. Model structure diagram

3.1 Cluster Random Under-sampling Smote + Tomek-link Approach (CRST)

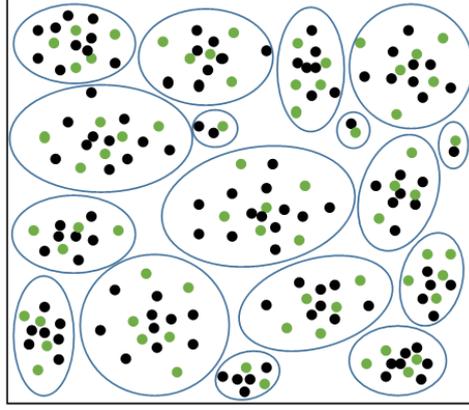
4 Resampling technique is to obtain a balanced data set from the original non-balanced
5 data set by using different sampling methods. From the perspective of data level,
6 resampling methods are mainly divided into three kinds: over-sampling,
7 under-sampling and combined sampling methods. The classic and commonly used
8 methods are: SMOTE and its improved method[25][26][27], Tomek-links method[28],
9 SMOTE + Tomek-Link, SMOTE + ENN[29] combination method, etc.

10 In view of the imbalanced and clustered characteristics of AD data, we proposed
11 Cluster Random Under-Sampling Smote + Tomek-link Approach (hereafter, CRST).
12 This method is an under-sampling method which takes the cluster center as the
13 representative point. It combines the advantages of K-means++ and Smote +
14 Tomek-link sampling method.

15 Firstly, the training samples in majority class were clustered by K-means++ algorithm,
16 in which K is obtained by super-parameter optimization. Then random
17 under-sampling is carried out for each cluster. The degree of sampling $p\%$ can be
18 determined by the actual situation. After the under-sampling is completed, SMOTE +
19 Tomek-Link combined sampling method is used to form a new balanced data set. By
20 iterating the above operations many times, we get several new balanced sub-datasets.

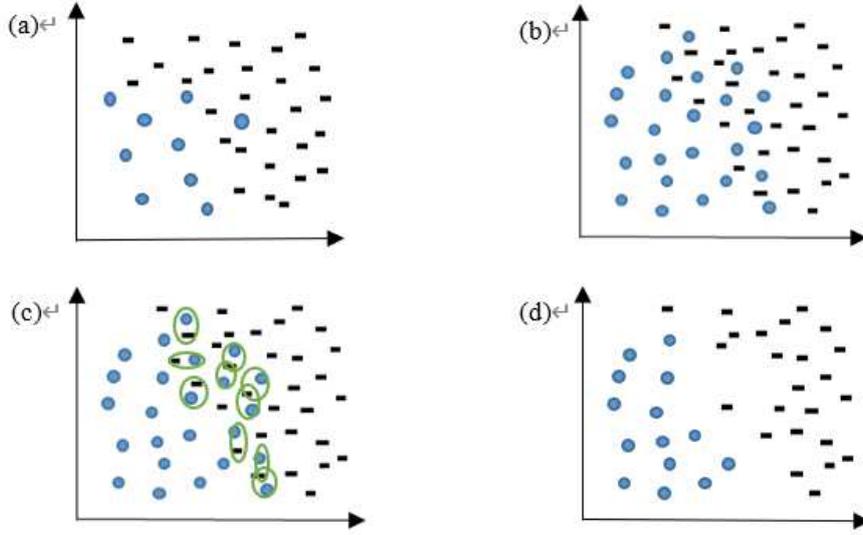
21 The clustering of samples in majority classes is visualized in the figure 5. The green
22 dots are the selected remaining of majority class sample points after $p\%$ random
23 under-sampling for each cluster. It can be seen that in this way the sample points can
24 be uniformly sampled in each cluster by under-sampling, maintaining the original data
25 distribution.

26 Finally, Smote + Tomek-link(S-T) sampling method is applied to generate some
27 minority samples, thus the sample loss caused by the under-sampling method is
28 compensated and the imbalanced ratio is alleviated. As shown in Figure 6, S-T
29 generates minority samples through SMOTE method, while Tomek-Link method is
30 adopted to solve the problem of fuzzy boundary caused by excessive generation of
31 minority samples. This method can reduce the redundancy of samples. The procedure
32 of algorithm is shown in Table 1.



1
2

Figure 5. Schematic Diagram of the Under Sampling Process of majority samples



3
4
5
6

Figure 6. Raw data set (a); Data set after SMOTE over-sampling(b); Tomek-Link recognition process (C); Data set after the boundary and noise samples are removed (D).

Table 1. CRST Sampling Method

Input: All standardized training sample set P ;

Output: The new balanced sub-datasets Z

Method:

1. Divide the input training sample set P into the majority sample set P_{max} and the minority sample set P_{min} according to the labels, and then remove the labels.
2. Use the K-means++ algorithm to cluster the training sample set P_{max} of majority classes to get K clusters. K is obtained by super-parameter optimization. denoted as: $\Omega = \{C_1, C_2, \dots, C_K\}$, wherein, $C_i = \{d_{i1}, d_{i2}, \dots, d_{im}\}$ m represents the sample's number of C_i
3. Take p% samples for each cluster class to obtain a new K cluster class sample set randomly, denoted as $\Omega' = \{C'_1, C'_2, \dots, C'_K\}$, wherein $C'_i = \{d_{i1}, d_{i2}, \dots, d_{in}\}$, $n = m * p\%$
4. Combine the majority class sample set Ω' and the minority class sample set P_{min} to synthesize the sample set Q , S-T method is used for Q to obtain a balanced data set Z .

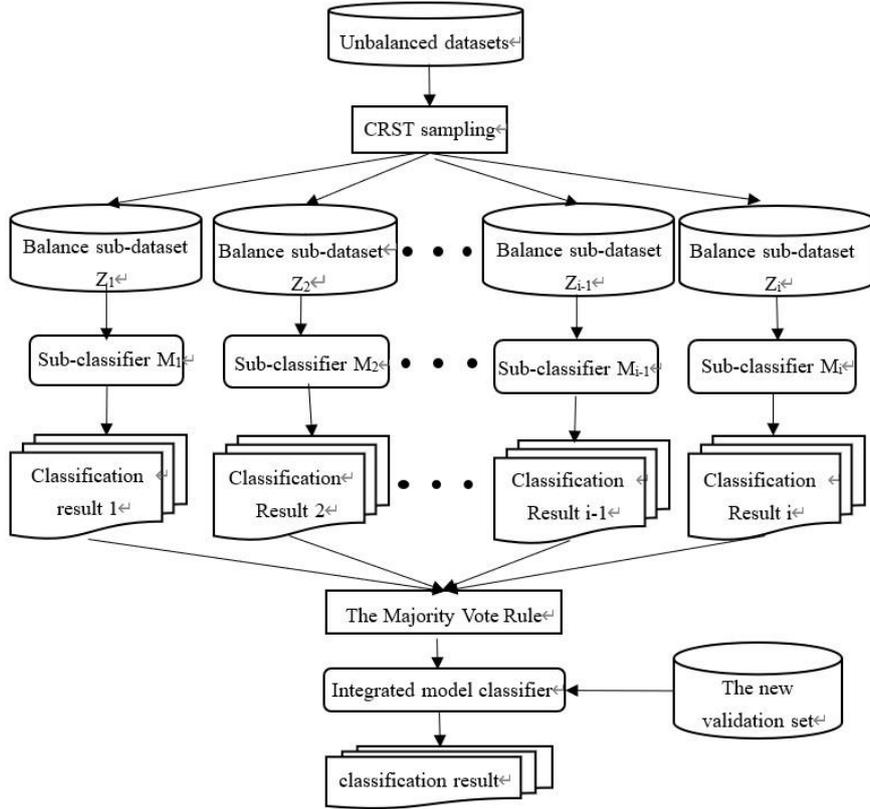
1 **3.2 ensemble model based on CRST**

2 As can be seen from the visualization analysis results of the data in Section 2.1, the
 3 data of aortic dissection has high overlap. So the classification boundary is blurred,
 4 and it is necessary to construct a nonlinear classification model with strong
 5 generalization ability. On the basis of the CRST sampling method proposed in Section
 6 2.2, integrated with the idea of Bagging[30] ensemble learning, CRST-Bagging
 7 ensemble learning algorithm is proposed in this section. It overcomes the limitation of
 8 a single classifier.

9 The CRST-Bagging algorithm is to generate a new sample set $B = \{Z_1, Z_2, \dots, Z_T\}$
 10 by using CRST sampling method iteratively. Then each sub-sample set Z_i is used to
 11 construct a sub-classifier M_i separately. A complete ensemble model classifier can be
 12 obtained by integrating the results of the T sub-classifier. The integration rule used in
 13 the algorithm is Majority Vote rule[31]. For the classifier, if P_{i1} is greater than or
 14 equal to P_{i2} , then R_1 gets one vote; if P_{i1} is less than P_{i2} , then R_2 gets one vote.
 15 R_1 and R_2 represents the sample category. This rule can be expressed by formula
 16 (1)(2). After the construction of the classifier, we send the verification set to the
 17 classifier for verification to evaluate the effect of the model. The model structure is
 18 shown in Figure 7.

19 $R_1 = \sum_{i=1}^T f(P_{i1}, P_{i2}), R_2 = \sum_{i=1}^T f(P_{i2}, P_{i1}), where \quad f(x, y) = \begin{cases} 0, & x < y \\ 1, & x > y \end{cases}$
 20 (1)
 21 C

22 $= \underset{j}{argmax} R_j$ (2)



1
2

Figure 7. CRST-Bagging algorithm structure diagram

3 4. Experiment

4 In this section, the CRST sampling method we proposed is compared with the
5 classical sampling methods Smote and Smote + Tomek-link, to verify its effectiveness
6 on the AD data set. And then, on this basis, we compare the classification
7 performance of CRST-Bagging model, RUSBoost[32] and SMOTEBagging[33] on
8 AD data sets. RUSBoost and SMOTEBagging are ensemble learning algorithm for
9 imbalanced data sets. The effects and advantages of CRST-Bagging method are
10 analyzed.

11 4.1 Validity experiment of the CRST sampling method

12 4.1.1 Experiment settings

13 We apply Smote, Smote + Tomek-link, CCST and CRST sampling methods to the AD
14 data set, and use the seven-fold cross-validation method to measure the sampling
15 performance. In order to ensure the consistency of the experiment, XGboost is used as
16 the classifier for this comparative experiment. The specific practices of the three
17 methods are as follows:

- 18 ● **Smote**, and **Smote + Tomek-link(S-T)**: At first, the samples in majority class
19 are randomly sampled so that the ratio of majority to minority is 2:1. Then
20 Smote and S-T are carried out respectively so that the ratio of positive and
21 negative samples after resampling is 1:1. The number of samples in each
22 category is 1604.

- 1 ● **CCST:** CCST method refers to the clustering of most class training sample
2 sets through K-means++ algorithm. Then, N sample points closest to the
3 center of each cluster were selected. Finally, SMOTE + Tomek-Link method
4 was applied to balance the data. In this experiment, we set K=802, N=2.
5 After under-sampling, the number of majority samples is 1604.
- 6 ● **CRST:** The number of clusters K is obtained by super-parameter selection,
7 which is 30. Randomly select 3.1% of each cluster class for under sampling,
8 and then apply S-T to the samples in minority class. The number of samples
9 in each category is 1624.

10 4.1.2 Experimental results and analysis

11 Table 2. Experimental results for the Smote, S-T, CCST, CRST methods

Method	Precision	Recall	F1
Smote	0.789	0.711	0.748
S-T	0.793	0.723	0.749
CCST	0.778	0.765	0.771
CRST	0.782	0.774	0.778

12 The experimental results of the seven-fold cross-validation on the original dataset
13 (53213,67) are shown in Table 2. Compared with the single over-sampling method
14 Smote and S-T method, CCST and our proposed method CRST has a great
15 improvement in recall rate and F1 value. As CCST selects the samples which is
16 closest to the sample center of each cluster, the sample distance between different
17 clusters after sampling is too far and the sample distribution is uneven, so the effect is
18 inferior to that of CRST. By contrast, the CRST method performed better. It shows
19 that CRST can reduce the occurrence of missed diagnosis in patients with AD to some
20 extent. This is because the CRST is an over-and-under sampling. While clustering and
21 under-sampling the majority samples, the CRST method uses the S-T method to
22 generate the same amount of minority class samples. In this way, we can not only
23 retain the original distribution of majority samples and select representative sample
24 points, but also balance the number of minority samples through over-sampling.

25 4.2 CRST-Bagging Model Effect Comparison Experiment

26 4.2.1 Experimental settings

27 The seven-fold cross-validation method is used in the comparative experiments. The
28 first group of experiments was performed in the original data set (53213, 67). The
29 second group experiment was tested on the test set. The experimental details of the
30 various algorithms of this experiment are as follows:

- 31 ● **RUSBoost:** the base learner type is decision tree C4.5, the number is 100, the
32 depth is 5;
- 33 ● **SMOTEBagging:** Set the number of clusters K = 5, the base learner is
34 decision tree C4.5, the number is 100, the depth is 6;
- 35 ● **CRST-Bagging:** Set the number of clusters K = 50, p%=3.1%.

1 **4.2.2 Experimental results and analysis**

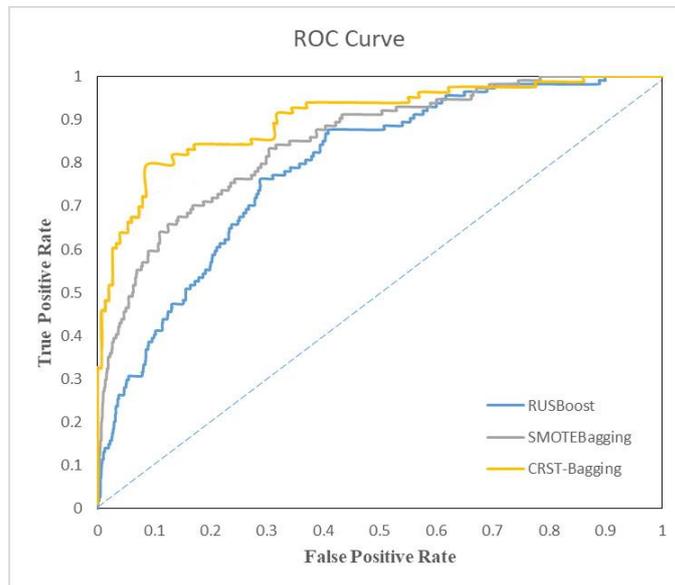
2 The experimental results of the seven-fold cross-validation on the original dataset
 3 (53213,67) are shown in Table 3. Figure 8 shows the ROC curve of the algorithm. It
 4 can be seen that CRST-Bagging performs best on the AD dataset, SMOTEBagging
 5 algorithm is second, and the RUSBoost algorithm performs worst. Compared with
 6 RUSBoost and SMOTEBagging methods, CRST-Bagging improved classification
 7 performance, and significantly improved accuracy and F1 value. It indicates that the
 8 rate of missed diagnosis and misdiagnosis in patients with AD are significantly
 9 reduced by CRST-bagging.

10 Table 3. Average experimental results of RUSBoost, SMOTE Bagging,
 11 CRST-Bagging on the original data set of seven-fold cross-validation

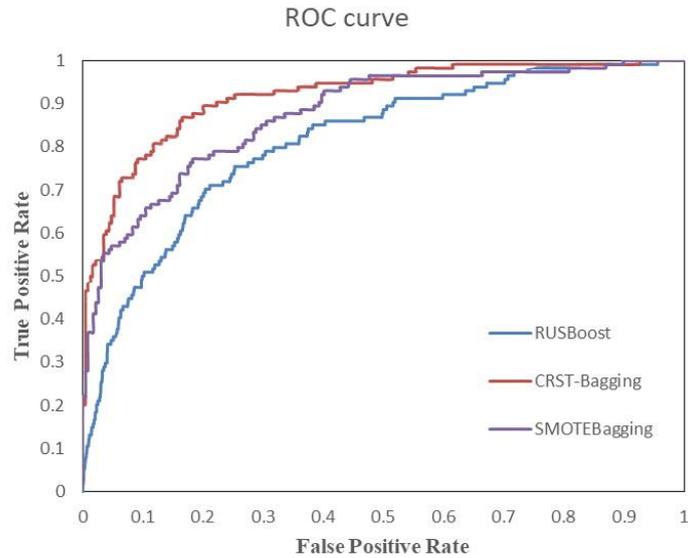
Method	Precision	Recall	F1
RUSBoost	0.774	0.751	0.762
SMOTEBagging	0.791	0.780	0.785
CRST-Bagging	0.841	0.783	0.811

12 Table 4. The average experimental results of the RUSBoost, SMOTEBagging,
 13 CRST-Bagging methods on the test sample set

Method	Precision	Recall	F1
RUSBoost	0.787	0.760	0.773
SMOTEBagging	0.842	0.772	0.805
CRST-Bagging	0.836	0.807	0.821



14
 15 Figure 8. ROC graph of RUSBoost, SMOTE Bagging, CRST-Bagging on the original sample set



1
2 Figure 9. ROC graph of RUSBoost, SMOTEBagging, CRST-Bagging on the test sample set
3 To test the generalization ability of the model, we performed the prediction of AD on
4 the test set (235, 67). The experimental results of the three algorithms are shown in
5 Table 4. Figure 9 shows the ROC curve of the algorithm. As can be found from the
6 table 4 and figure 9, the SMOTEBagging algorithm has the highest accuracy on the
7 test sample set, and the CRST-Bagging algorithm has the highest recall rate and F1
8 value. CRST-Bagging's algorithm performance is significantly improved compared to
9 RUSBoost and SMOTEBagging algorithms on the test sample set. CRST-Bagging
10 algorithm has stronger generalization ability. In other words, CRST-bagging algorithm
11 is more likely to detect potential patients with AD.

12 5. Discussion

13 AD is a rare and high-risk cardiovascular disease. Its complex clinical manifestations
14 and various atypical symptoms lead to serious misdiagnosis and missed diagnosis.
15 The rarity of the disease also leads to a significant imbalance in the data set. This
16 paper studies the misdiagnosis and missed diagnosis of AD and the imbalanced
17 characteristics of data.

18 For the original AD dataset, in the data preprocessing stage, we performed missing
19 value processing, feature screening, data standardization. And at the level of data
20 distribution the data is understood through dimensional reduction visualization. This
21 is different from the general "black box" approach of machine learning algorithms.
22 These methods enable us to have a priori knowledge of the distribution of actual
23 medical data sets. This prior knowledge can inspire clinical medicine to explore the
24 etiology and diagnostic criteria of AD.

25 Aiming at the high imbalance of AD data set, this paper proposes a resampling
26 method CRST based on clustering. This method combines the advantages of
27 traditional sampling methods Smote + Tomek-link and clustering algorithm
28 K-means++. In CRST, a certain percentage of samples are randomly selected from

1 clusters, which not only makes the selected samples effectively represent the
2 characteristics of most kinds of samples, but also ensures the randomness of sampling.
3 Experiments show that CRST scientifically and effectively reduces the imbalanced
4 ratio of rare disease medical data and relieves the obstacles that imbalanced data bring
5 to the construction of classification models.

6 On this basis of CRST, this paper proposes the CRST-Bagging learning model
7 combined with the idea of ensemble learning. After experimental comparative
8 analysis, the CRST-Bagging model presented in this paper shows excellent
9 performance on the AD data set. Not only the accuracy and recall rate of the model on
10 the original AD data set have been improved, but also the generalization ability of the
11 model on the test sample set is also very good. This shows that this model is a good
12 diagnostic model of AD. Clinically, the misdiagnosis rate of AD is close to
13 40%[3][4][5]. The model performance data show that this algorithm can not only
14 reduce the workload of doctors, but also reduce the misdiagnosis rate of AD to save
15 patients' lives effectively.

16 **6. Conclusions**

17 In this paper, a cluster-based ensemble learning model named CRST-Bagging is
18 proposed to assist in diagnosing aortic dissection through the patient's inspection
19 results. Compared with the ordinary classification model, our model pays more
20 attention to the processing of medical data sets with high imbalance ratio. While
21 ensuring high accuracy, CRST effectively improves the recall rate. That is, the missed
22 diagnosis and misdiagnosis rate is reduced. In addition, the algorithm demonstrates a
23 strong ease of use. In many basic hospitals where the equipment is not advanced
24 enough, it is difficult for patients to perform more examination items such as CT,
25 magnetic resonance angiography (MRA), etc. The model proposed in this paper can
26 reduce the burden of doctors and patients to a certain extent and help diagnose the
27 AD.

28 The diagnosis of AD remains one of the most difficult problems in the cardiovascular
29 field. In the future, based on the analysis results of aortic dissection data and the
30 proposed auxiliary diagnostic method for aortic dissection in this paper, we will study
31 the pathological mechanism and key diagnostic indicators of aortic dissection from
32 the perspective of interpretability, and explore whether there is a more definite clinical
33 diagnostic method for aortic dissection.

34 **Abbreviations**

35 AD: aortic dissection

36 LSTM: long short-term memory

37 PAVE: Pattern Attention model with Value Embedding

38 CT: Computed Tomography

39 T-SNE: t-distributed stochastic neighbor embedding

40 CRST-Bagging: Cluster Random Under-sampling Smote-Tomek-link Bagging

- 1 CRST: Cluster Random Under-Sampling Smote + Tomek-link Approach
- 2 S-T: Smote + Tomek-link
- 3 CCST: Cluster-Center Under-Sampling and Smote+Tomeklink Approach

4 **Declarations**

5 **Ethics approval and consent to participate**

6 Ethical approval for this study was obtained from the Ethics Board of Xiangya
7 Hospital, Central South University (201502042).

8 **Consent for publication**

9 Not applicable.

10 **Competing interests**

11 The authors declare that they have no competing interests.

12

13 **Availability of data and materials**

14 Data are provided by Xiangya Hospital and it cannot be shared with other research
15 groups without necessary permission. The data used during the current study is
16 available from the corresponding author on reasonable request. The data description
17 supporting the conclusions of this article is included in the article (and its Appendix).

18

19 **Funding**

20 The study was financially supported by National Natural Science Foundation of China
21 (No. 61502537), and Strategic Emerging Industry Technological Research and Major
22 Technological Achievement Transformation Project, High-tech Development and
23 Industrialization Office (No. 2019GK4013). The funding body had no role in design
24 of the study, collection, analysis, and interpretation of data or in writing the
25 manuscript.

26

27 **Authors' contributions**

28 All of the authors had full access to all of the data in the study and take responsibility
29 for the content of the manuscript. YG designed the model and experiment
30 implementation. MW and LJZ wrote the code. LJL, GGZ and JML contributed to data
31 collection and feature selection. YG, LJL and GGZ perform the results analysis. MW
32 and LJZ drafted the initial manuscript. YG revised the manuscript. All authors read
33 and approved the final draft of the manuscript for publication.

34

35 **Acknowledgements**

36 Not applicable.

37

38 **References**

- 39 [1] Erbel R, Aboyans V, Boileau C, et al. 2014 ESC guidelines on the diagnosis and treatment of
40 aortic diseases: document covering acute and chronic aortic diseases of the thoracic and

- 1 abdominal aorta of the adult. The task force for the diagnosis and treatment of aortic diseases
2 of the European Society of Cardiology (ESC). *Eur Heart J* 2014;35:2873-926.
- 3 [2] Mussa FF, Horton JD, Moridzadeh R, Nicholson J, Trimarchi S, Eagle KA. Acute Aortic
4 Dissection and Intramural Hematoma: A Systematic Review. *JAMA*. 2016;316:754-763.
- 5 [3] Marroush TS, Boshara AR, Parvataneni KC, Takla R, Mesiha NA. Painless Aortic Dissection.
6 *Am J Med Sci*. 2017;354:513-520.
- 7 [4] Yin XB, Wang XK, Xu S, He CY. Type A aortic dissection developed after type B dissection
8 with the presentation of shoulder pain: A case report. *World J Clin Cases* 2021; 9(1): 232-235
9 [PMID: 33511190 DOI: 10.12998/wjcc.v9.i1.232]
- 10 [5] Ahmed, T., et al., Two Intriguing Cases of Stanford Type A Acute Aortic Dissection. *Cureus*,
11 2020. 12(2).
- 12 [6] Zhengxing Huang, Tak-Ming Chan, Wei Dong. MACE prediction of acute coronary
13 syndrome via boosted resampling classification using electronic medical records. 2017,
- 14 [7] Zhou, PY., Wong, A.K.C. Explanation and prediction of clinical data with imbalanced class di
15 stribution based on pattern discovery and disentanglement. *BMC Med Inform Decis Mak* 21,
16 16 (2021).
- 17 [8] Zhigang Song, Shuangmei Zou, Weixun Zhou, et al. Clinically applicable histopathological
18 diagnosis system for gastric cancer detection using deep learning. 2020, 11(1):7-34.
- 19 [9] Kamal, S.A., Yin, C., Qian, B. et al. An interpretable risk prediction model for healthcare
20 with pattern attention. *BMC Med Inform Decis Mak* 20, 307 (2020).
- 21 [10] Seo, S.W., Kim, J., Son, J. et al. Evaluation of conditional treatment effects of adjuvant
22 treatments on patients with synovial sarcoma using Bayesian subgroup analysis. *BMC Med*
23 *Inform Decis Mak* 20, 320 (2020).
- 24 [11] Yu, X., Ma, N., Yang, T. et al. A multi-level hypoglycemia early alarm system based on
25 sequence pattern mining. *BMC Med Inform Decis Mak* 21, 22 (2021).
- 26 [12] Huo D, Kou B, Zhou Z, et al. A machine learning model to classify aortic dissection patients
27 in the early diagnosis phase. *Sci Rep* 2019;9:2701.
- 28 [13] Dwivedi AK. Performance evaluation of different machine learning techniques for prediction
29 of heart disease. *Neural Comput Appl* 2018;29:685-93
- 30 [14] Guo., Beheshti, R., Khan, Y.M. et al. Predicting cardiovascular health trajectories in
31 time-series electronic health records with LSTM models. *BMC Med Inform Decis Mak* 21, 5
32 (2021).
- 33 [15] Junlong Cheng, Shengwei Tian, Long Yu, Xiang Ma, Yan Xing, A deep learning algorithm
34 using contrast-enhanced computed tomography (CT) images for segmentation and rapid
35 automatic detection of aortic dissection, *Biomedical Signal Processing and Control*, Volume
36 62, 2020, 102145, ISSN 1746-8094,
- 37 [16] JAPKOWICZ N. Supervised Versus Unsupervised Binary-Learning by Feedforward Neural
38 Networks[M]. 2001.
- 39 [17] SHENG V S, LING C X B T-N C on A I. Thresholding for Making Classifiers Cost
40 Sensitive[C]//National Conference on Artificial Intelligence. AAAI Press. 2006.
- 41 [18] MIRZA B, KOK S, LIN Z et al. Efficient Representation Learning for High-Dimensional
42 Imbalance Data[M]. 2016.

- 1 [19] XIE Y N, YU L, GUAN G H et al. An Overlapping Cell Image Synthesis Method for
2 Imbalance Data[J]. Analytical cellular pathology (Amsterdam), Hindawi, 2018, 2018(1):
3 1-12.
- 4 [20] HU J, YANG H, LYU M R et al. Online Nonlinear AUC Maximization for Imbalanced Data
5 Sets[J]. IEEE Transactions on Neural Networks and Learning Systems, 2018, 29(4): 882-895.
- 6 [21] SHENG C, HAIBO H, GARCIA E A. RAMOBoost: Ranked Minority Oversampling in
7 Boosting[J]. IEEE Transactions on Neural Networks, 2010, 21(10): 1624-1642.
- 8 [22] CHEN T, GUESTRIN C B T-A S I C on K D & D M. XGBoost: A Scalable Tree Boosting
9 System[C]. Acm Sigkdd International Conference on Knowledge Discovery & Data Mining.
10 2016.
- 11 [23] ZHANG D, QIAN L, MAO B et al. A Data-Driven Design for Fault Detection of Wind
12 Turbines Using Random Forests and XGBoost[J]. IEEE Access, 2018, PP(99): 1.
- 13 [24] VAN DER MAATEN L, HINTON G. Visualizing Data using t-SNE[J]. Journal of Machine
14 Learning Research, 2008, 9: 2579–2605.
- 15 [25] CHAWLA N V, BOWYER K W, HALL L O et al. SMOTE: Synthetic Minority
16 Over-sampling Technique[J], 2002, 16(January): 321–357.
- 17 [26] NGUYEN H M, COOPER E W, KAMEI K. Borderline over-sampling for imbalanced data
18 classification[J]. International Journal of Knowledge Engineering and Soft Data Paradigms,
19 2011, 3(1): 4.
- 20 [27] LUSA L, OTHERS. SMOTE for high-dimensional class-imbalanced data[J]. BMC
21 bioinformatics, 2013, 14(1): 106.
- 22 [28] TOMER I. Two Modifications of CNN[J]. IEEE Transactions on Systems Man and
23 Communications AMC-6, 1976: 769–772.
- 24 [29] BATISTA G E A P A, PRATI R C, MONARD M C. A study of the behavior of several
25 methods for balancing machine learning training data[J]. ACM SIGKDD Explorations
26 Newsletter, 2004, 6(1): 20.
- 27 [30] BREIMAN L. Bagging predictors[M]. 1996.
- 28 [31] KITTLER J, HATEF M, DUIN R P W et al. On Combining Classifiers[J]. IEEE Transactions
29 on Pattern Analysis & Machine Intelligence, 1998, 20(3): 226-239.
- 30 [32] SEIFFERT C, KHOSHGOFTAAR T M, VAN HULSE J et al. RUSBoost: A hybrid approach
31 to alleviating class imbalance[J]. IEEE Transactions on Systems, Man, and Cybernetics Part
32 A:Systems and Humans, 2010, 40(1): 185-197.
- 33 [33] WANG S, YAO X. Diversity analysis on unbalanced data sets by using ensemble models[J].
34 IEEE Symposium on Computational Intelligence&Data Mining, 2009, 1(5):324-331.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile1.docx](#)