

## Postprint of Research on Text Sentiment Classification Based on BPSO Random Subspace

**Authors:** Zhang Qingqing, Liu Xilin

**Date:** 2017-11-08T00:00:00+00:00

### Abstract

**Objective:** To address the high dimensionality problem of text feature representation vectors in machine learning-based text sentiment classification research, a selective ensemble algorithm combining BPSO with the random subspace method is proposed. **Method:** Based on the analysis of BPSO and random subspace principles, the model framework and algorithmic flow of BPSO random subspace are presented. After feature representation of Chinese review corpora, experiments are conducted using BPSO random subspace for validation and analysis. **Results:** By varying the subspace rate values in random subspace, the effects of standard random subspace and BPSO random subspace selective ensemble on classification accuracy and system diversity are investigated. The results demonstrate that BPSO random subspace outperforms standard random subspace in both classification accuracy and system diversity. **Limitations:** Validation on English datasets has not yet been conducted. **Conclusion:** Applying BPSO to the random subspace method constitutes a novel selective ensemble model that not only resolves the high dimensionality issue of feature vector space but also improves classification accuracy and generalization capability, providing an effective approach for Chinese text sentiment classification.

### Full Text

### Preamble

**ChinaXiv Cooperative Journal, Issue 5, 2017**

**Text Sentiment Classification Based on BPSO Random Subspace**

Zhang Qingqing<sup>1,2</sup>, Liu Xilin<sup>2</sup>

<sup>1</sup>(School of Management, Xi'an Polytechnic University, Xi'an 710048, China)

<sup>2</sup>(School of Management, Northwestern Polytechnical University, Xi'an 710129, China)

## Abstract

**[Objective]** To address the high dimensionality problem in text feature representation vectors for machine learning-based text sentiment classification research, this paper proposes a selective ensemble algorithm that combines Binary Particle Swarm Optimization (BPSO) with the random subspace method. **[Methods]** Based on an analysis of the principles of BPSO and random subspace, we present the model framework and algorithmic flow of the BPSO random subspace approach. After representing Chinese review corpora as feature vectors, we conduct experimental validation and analysis using the BPSO random subspace method. **[Results]** By varying the subspace rate in the random subspace method, we investigate the impact of standard random subspace versus BPSO random subspace selective ensemble on classification accuracy and system diversity. The results demonstrate that BPSO random subspace outperforms standard random subspace in both classification accuracy and system diversity. **[Limitations]** The approach has not yet been validated on English datasets. **[Conclusion]** Applying BPSO to the random subspace method constitutes a novel selective ensemble model that not only solves the high-dimensional feature vector space problem but also improves classification accuracy and generalization capability, providing an effective method for Chinese text sentiment classification.

**Keywords:** Random Subspace; BPSO; Text Sentiment Classification; Subspace Rate

**Classification Number:** TP391.1

## Introduction

The Internet provides people with abundant information resources, among which subjective texts expressing opinions, views, suggestions, and perspectives—such as technology reviews, product reviews, sports commentary, current affairs commentary, blogs, film and television reviews, news commentary, military reviews, music reviews, and stock reviews—constitute a significant and important portion. These subjective information pieces represent viewpoints, attitudes, opinions, and positions toward specific targets, carrying strong personal emotional coloring. Text sentiment classification is the technology for automatically analyzing, processing, and summarizing such subjective texts, with important application value in e-commerce, e-government, and information forecasting.

Currently, the mainstream approach for text sentiment classification research is machine learning, which primarily focuses on feature representation and the application and improvement of classification models for text sentiment classification tasks. Feature representation aims to identify feature items that can best represent sentence semantics and syntax. Commonly used features for text representation include unigrams, n-grams, part-of-speech (POS), word relationship features, rule-based features, features combined with sentiment lexicons, and

social network features [1-3]. Dependency syntactic relationship features based on dependency grammar have been used in text sentiment classification feature representation and have achieved higher classification accuracy than commonly used features due to their ability to effectively express syntactic structure and modification relationships between words [4].

In terms of classification models, traditional algorithms include Support Vector Machines, Naive Bayes, Maximum Entropy models, k-Nearest Neighbors, and Decision Trees. The superiority of classification algorithms for text sentiment classification tasks remains an open question. Ensemble learning forms a new text classification paradigm based on the principle that the decision results of multiple classifiers are more reliable than those of a single classifier. Using traditional classification algorithms as base classifier training methods, Wang et al. compared three different ensemble learning base classifier generation methods—Boosting, Bagging, and random subspace—training the base classifiers of these three methods with Naive Bayes, Maximum Entropy, Decision Tree, k-Nearest Neighbors, and Support Vector Machine algorithms. Their results showed that the classification accuracy of the integrated system was higher than that of individual classifiers.

## 2.2 Selective Ensemble

Selective ensemble refers to selecting a subset from a batch of trained base classifiers for integration. The concept of selective ensemble was proposed by Zhou et al., who demonstrated that having all classifiers participate does not necessarily guarantee improved generalization capability and provided theoretical analysis showing that selective ensemble outperforms full participation [12]. Selective ensemble is essentially a global optimization process.

Global optimization represents an important research direction for selective ensemble, which can be conveniently transformed into a combinatorial optimization problem. The Particle Swarm Optimization algorithm requires few parameters, executes efficiently, converges quickly, and possesses global search capabilities and advantages in solving combinatorial optimization problems. While particle swarm optimization has been widely applied to other ensemble learning methods such as Bagging and Boosting [13-15], its application to random subspace methods remains limited.

## 2.3 BPSO Algorithm

The BPSO algorithm is a discrete version of the basic Particle Swarm Optimization algorithm based on continuous space, specifically proposed by American social psychologist Kennedy and electrical engineer Eberhart for 0-1 integer programming problems [16].

The basic Particle Swarm Optimization algorithm is described by formula (2), where  $v_{ij}(t+1)$  and  $x_{ij}(t+1)$  represent the velocity and position of the j-th

dimension of the  $i$ -th particle at iteration  $t + 1$ .  $c_1$  and  $c_2$  are acceleration constants, typically valued in  $[0,2]$ , where  $c_1$  regulates the step size for particles moving toward individual optimal positions and  $c_2$  regulates the flight step size toward global optimal positions.  $r_1$  and  $r_2$  are random values in the interval  $[0,1]$ , primarily to increase the randomness of particle flight.  $p_{ij}$  and  $g_{ij}$  represent the individual extremum and global extremum of the  $j$ -th dimension of the particle, respectively. The position information  $x_{ij}(t + 1)$  for the next iteration is obtained by transforming the velocity from the original position.

In BPSO, each position component  $x_{ij}$  takes a value of either 0 or 1, so the velocity component  $v_{ij}$  no longer represents the magnitude of position change but rather reflects the probability of  $x_{ij}$  taking the value 1. When using the velocity update formula, larger values make the particle's position component  $x_{ij}$  more likely to be 1, while smaller values make  $x_{ij}$  tend toward 0. To ensure probability values fall within  $[0,1]$ , BPSO employs a Logistic transformation to process  $v_{ij}$ , as shown in formula (3):

$$S(v_{ij}) = \frac{1}{1 + \exp(-v_{ij})}$$

where  $S(v_{ij})$  represents the probability of position  $x_{ij}$  taking the value 1. The particle changes its position value according to formula (4):

$$x_{ij} = \begin{cases} 1 & \text{if rand() < } S(v_{ij}) \\ 0 & \text{otherwise} \end{cases}$$

where  $\text{rand}()$  is a random number drawn from a uniform distribution in  $[0,1]$ . To avoid  $S(v_{ij})$  approaching 0 or 1, the parameter  $v_{\max}$  is used as the maximum velocity value to limit the range of  $v_{ij}$ .

### 3 BPSO Random Subspace Algorithm

The BPSO random subspace method optimizes the selection of base classifiers using the BPSO algorithm on top of classifiers trained by random subspace. Based on the above analysis of random subspace and BPSO algorithms, the algorithm flow of the BPSO random subspace method is shown in Figure 1 [Figure 1: see original paper].

The key to the BPSO random subspace method lies in the design of particle dimensionality and fitness function in BPSO. In the BPSO selective ensemble learning algorithm, a particle represents a scheme for selecting base classifiers. A particle is a vector whose dimensionality corresponds to the number of base classifiers in a one-to-one relationship, with vector values being either 0 or 1. Assuming there are  $D$  base classifiers, the particle is represented as a  $D$ -dimensional vector. If the value of the  $d$ -th component is 1, it indicates that the  $d$ -th base classifier is selected; conversely, if the value is 0, it indicates that

the  $d$ -th base classifier is not selected. Taking 10 base classifiers as an example, if arranged sequentially, the selection result after optimization might be  $[0, 0, 1, 1, 0, 1, 0, 1, 0, 1]$ , meaning base classifiers numbered 3, 4, 6, 8, and 10 are selected, as illustrated in Figure 2 [Figure 2: see original paper].

In the particle swarm optimization algorithm, the particle's position information indicates whether a base classifier is selected, while the velocity corresponds to the probability of this base classifier being selected. The BPSO algorithm performs global search based on the fitness function. Classification accuracy and system diversity are two metrics for evaluating ensemble learning systems. The most common evaluation function is the current system's prediction error, judging the quality of the ensemble system based on classification accuracy. Another evaluation metric is system diversity, which measures the generalization capability of the ensemble system. This is an indirect method that requires appropriate description to achieve good results. This paper adopts system classification accuracy as the fitness function.

## 4 Experiments

### 4.1 Datasets

The datasets used in this paper come from sentiment analysis corpora provided by Datatang<sup>1</sup>, including hotel review data, book review data, and laptop computer review data, respectively crawled from Ctrip, Dangdang, and JD.com.

Each of the three original datasets contains 4,000 positive texts and 4,000 negative texts, but all exist in paragraph form. This paper studies sentence-level text sentiment orientation and processes the original data accordingly:

1. Sentence segmentation is performed on document data. All documents are segmented using newline characters and Chinese/English question marks (“?” and “?”), periods (“.” and “。”), and semicolons (“;” and “;”) as sentence delimiters. Duplicate removal is then performed on the segmented sentences.
2. Based on the original annotated texts, the newly segmented sentences are screened to delete sentences that do not belong to the original categories.
3. Random sampling is conducted on the three datasets.

This paper studies balanced data for text sentiment classification. For hotel review data, 4,000 sentences are extracted, including 2,000 positive review sentences and 2,000 negative review sentences. For book review data, 2,000 sentences are extracted, including 1,000 positive and 1,000 negative review sentences. For laptop computer review data, 1,000 sentences are extracted, including 500 positive and 500 negative review sentences.

Based on the triple dependency relation feature method proposed in literature [4], Chinese review corpora are transformed into triple dependency relation features. The total number of features obtained from the three datasets is shown in Table 1 .

## 4.2 Evaluation Metrics

The evaluation metrics employed in the experiments are average classification accuracy and system diversity. Average classification accuracy is calculated using formula (5):

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

where TP (True Positive) represents the number of positive sentiment documents correctly classified, TN (True Negative) represents the number of negative sentiment documents correctly classified, FP (False Positive) represents the number of positive sentiment documents incorrectly classified, and FN (False Negative) represents the number of negative sentiment documents incorrectly classified. The sum of TP, TN, FP, and FN equals the total number of documents to be classified. Higher average classification accuracy indicates more accurate judgment of text subjectivity orientation.

Diversity measurement is a unique evaluation criterion for ensemble learning systems. We introduce four commonly used pairwise diversity measures: Q-statistic, correlation coefficient  $\rho$ , disagreement measure *dis*, and double-fault measure DF.

Assuming there are  $L$  base classifiers,  $C_i$  and  $C_j$  ( $i \neq j$ ) are two different classifiers, and  $N$  represents the total number of samples.  $N^{ab}$  represents the number of samples where classifier  $C_i$  classifies correctly and classifier  $C_j$  classifies incorrectly. The specific formulas are shown in Table 2 as formulas (6)-(9).

As shown in Table 2, larger Q-statistic values indicate lower diversity, and the correlation coefficient  $\rho$  follows the same trend. The disagreement measure *dis* focuses on samples where two classifiers produce different classification results—more such samples indicate higher diversity. The double-fault measure DF focuses on samples where both classifiers misclassify; more such samples indicate both lower accuracy and lower diversity.

## 4.3 Experimental Process

To verify the effectiveness of RS\_{BPSO}, the experiments compare and analyze the classification accuracy and system diversity obtained by RS\_{BPSO} and standard RS. The specific experimental process is as follows:

1. Split the review data into training and test sets at a 70:30 ratio.
2. Convert the training and test set texts into structured triple dependency relation feature vectors.
3. Apply bootstrap sampling to the training and test sets for feature subset partitioning.
4. Train base classifiers on the partitioned training sets using Support Vector Machines.

5. Label all base classifiers and use the BPSO algorithm to optimize the selection of base classifiers, determining which base classifiers to retain.
6. Apply the partitioned test sets to the retained base classifiers and fuse the results using majority voting to obtain final classification results.

In the random subspace method, the randomly selected feature subset dimensionality is determined by the subspace rate  $k$ . We investigate the impact of different  $k$  values (0.01, 0.02, 0.03, 0.05) on text sentiment classification accuracy and diversity. The feature subset dimensions obtained for the three datasets under different subspace rates are shown in Table 3 .

As shown in Table 3, as the  $k$  value varies, the feature dimensionality of different datasets is reduced from over ten thousand dimensions to thousands or hundreds. Taking hotel review data as an example, when  $k = 0.01$ , the feature vector space dimensionality for training in base classifiers is 1,409. To ensure all feature items have the possibility of being selected, 50 random samplings are performed on the original feature set, meaning feature items are divided into 50 feature subsets, resulting in 50 base classifiers.

In the particle swarm optimization algorithm, the fitness function uses system classification accuracy. To investigate the impact of particle population size on text classification accuracy and diversity, we compare particle populations of 10, 20, 30, 40, 50, 60, 70, 80, 90, and 100, with iteration count set to 100, learning factor  $c_1 = c_2 = 2$ , and inertia weight using linear iteration  $w_{\min} = 0.4$ ,  $w_{\max} = 0.9$ . Due to the randomness of initial particle swarm values, experimental results are averaged over 20 runs.

## 5 Results and Analysis

### 5.1 Classification Accuracy Results and Analysis

Experiments were conducted following the experimental process for standard RS and RS\_{BPSO} under different subspace rate  $k$  values. The results for the three datasets are shown in Tables 4 -6 . The RS\_{BPSO} classification accuracy column includes the average number of base classifiers after optimization selection in parentheses.

**Table 4** shows the hotel review data classification accuracy results. RS\_{BPSO} achieves higher classification accuracy than the standard random subspace method, with improvements ranging from 3% to 15%, reaching a maximum accuracy of 84.29%. The number of base classifiers selected after BPSO optimization is 17, 14, 13, and 19, respectively, averaging approximately 34 fewer than the original number. Comparing classification accuracy across different  $k$  values, the standard random subspace method shows an increasing trend as  $k$  increases, but after BPSO-based base classifier selection, classification accuracy improves without showing a consistent pattern related to  $k$  values.

**Table 5** presents the book review data classification accuracy results, showing

conclusions consistent with Table 4.  $RS_{\{BPSO\}}$  achieves higher classification accuracy than the standard random subspace method, with improvements ranging from 5% to 14%, reaching a maximum accuracy of 84.34%. The number of base classifiers selected after BPSO optimization is 19, 19, 20, and 21, respectively, averaging approximately 30 fewer than the original number.

**Table 6** shows the laptop computer review data classification accuracy results.  $RS_{\{BPSO\}}$  achieves higher classification accuracy than the standard random subspace method, with improvements ranging from 4% to 7%, reaching a maximum accuracy of 87.62%. The number of base classifiers selected after BPSO optimization is 24, 29, 28, and 22, respectively, averaging approximately 19 fewer than the original number.

From these comparisons across three datasets, we conclude that the BPSO-based random subspace method significantly improves classification accuracy over the standard random subspace method. In the standard random subspace method, classification accuracy increases with larger  $k$  values, but this effect is minimal in  $RS_{\{BPSO\}}$ . The BPSO algorithm substantially reduces the number of base classifiers used, benefiting computational speed and storage efficiency.

## 5.2 System Diversity Results and Analysis

System diversity measures for standard RS and  $RS_{\{BPSO\}}$  are calculated based on the output results of the final selected base classifiers for each test sample, with four different diversity measures computed. Tables 7 -9 present the diversity measurement results for hotel review data, book review data, and laptop computer review data, respectively.

In Table 7, the double-fault measure DF values for  $RS_{\{BPSO\}}$  are higher than those for RS, indicating reduced diversity for  $RS_{\{BPSO\}}$ . However, in the disagreement measure  $dis$ ,  $RS_{\{BPSO\}}$  values are higher than RS when  $k = 0.01, 0.02, 0.03$ . According to the  $dis$  calculation principle, larger  $dis$  values indicate higher system diversity. The Q-statistic and correlation coefficient  $\rho$  show the opposite trend. The data shows that  $RS_{\{BPSO\}}$ 's Q-statistic and correlation coefficient are lower than RS's, and according to their calculation principles, lower values indicate higher system diversity. These results demonstrate that  $RS_{\{BPSO\}}$  achieves higher system diversity.

In Table 8, the double-fault measure DF shows little difference between  $RS_{\{BPSO\}}$  and RS, indicating no significant diversity difference. However, in the disagreement measure  $dis$ ,  $RS_{\{BPSO\}}$  values are 0.02 to 0.04 higher than RS, indicating higher diversity. Similar conclusions are evident in the Q-statistic and correlation coefficient  $\rho$  results, where  $RS_{\{BPSO\}}$  values are lower than RS, clearly showing that  $RS_{\{BPSO\}}$  improves system diversity.

Table 9 yields the same conclusions as Table 8. The double-fault measure DF shows minimal difference between  $RS_{\{BPSO\}}$  and RS. In the disagreement measure  $dis$ ,  $RS_{\{BPSO\}}$  values are higher than RS, indicating higher diver-

sity.  $RS_{\{BPSO\}}$ 's Q-statistic and correlation coefficient  $\rho$  are lower than RS's, demonstrating that  $RS_{\{BPSO\}}$ 's diversity is significantly higher than the standard feature subspace method.

To more intuitively analyze the diversity under different  $k$  values for standard random subspace and discrete binary particle swarm random subspace methods, Figure 3 [Figure 3: see original paper] shows diversity comparison charts for hotel review data, book review data, and laptop computer review data. The figure shows that DF, Q-statistic, and correlation coefficient  $\rho$  exhibit upward trends, while  $dis$  shows a downward trend. All four diversity measures consistently conclude that as  $k$  values increase, diversity decreases in both standard random subspace and  $RS_{\{BPSO\}}$  methods. In pairwise comparisons between RS and  $RS_{\{BPSO\}}$ , the curves for DF, Q-statistic, and correlation coefficient  $\rho$  in  $RS_{\{BPSO\}}$  are above those in RS, while  $dis$  shows the opposite pattern. According to their theoretical foundations, we conclude that  $RS_{\{BPSO\}}$ 's diversity is higher than that of the standard random subspace method. Book review data and laptop computer review data yield the same conclusions as hotel review data.

From this analysis, we conclude that the system diversity of both standard random subspace and  $RS_{\{BPSO\}}$  methods decreases as  $k$  values increase, but  $RS_{\{BPSO\}}$ 's diversity is significantly higher than that of standard random subspace, indicating stronger generalization capability. Comparing these results with classification accuracy, we observe a contradictory trend: in the standard random subspace method, classification accuracy increases with  $k$  values, while diversity decreases. This aligns with Chandra et al.'s proposition that a trade-off exists between classifier accuracy and diversity [17]. Multi-objective optimization ensemble learning, which treats accuracy and diversity as two optimization objectives, has been proposed as a strategy to achieve this balance [17], while other literature has improved existing diversity measures by combining classifier accuracy and diversity into composite diversity functions [18-19].

### 5.3 BPSO Algorithm Optimization Process Analysis

#### (1) BPSO Convergence Analysis

BPSO convergence means that as iteration count increases, the error between algorithm results and true results becomes smaller and approaches a fixed value. Convergence is the opposite of divergence, where the convergence curve cannot stabilize regardless of iteration count.

To better analyze BPSO's optimization performance for base classifier selection in ensemble learning random subspace methods, we analyze BPSO convergence by plotting curves of the best fitness values and average fitness values obtained over 100 iterations. Figures 4 [Figure 4: see original paper]-6 [Figure 6: see original paper] show the variation of fitness values with iteration counts for hotel review data, book review data, and laptop computer review data, respectively.

For hotel review data (Figure 4), the best fitness value for text sentiment classification continuously increases with iteration count, indicating decreasing classification error and gradual convergence toward 0. Although the average fitness value shows local fluctuations, it generally trends toward the maximum classification accuracy, demonstrating excellent convergence of the BPSO algorithm in optimizing base classifier selection for random subspace.

For book review data (Figure 5), the best fitness value continuously increases with iteration count, while the average fitness value, despite local fluctuations, still shows an overall upward trend.

For laptop computer review data (Figure 6), the fitness function shows a continuously rising trend with increasing iteration count, and the average fitness function value demonstrates an overall upward trend.

This analysis demonstrates that BPSO algorithm exhibits excellent convergence when using text sentiment classification accuracy as the fitness function, effectively improving original classification accuracy.

## (2) Impact of Particle Count on Classification Results

All previous results used 50 particles. To explore the impact of different particle counts on text sentiment classification, experiments were conducted with particle counts of 10, 20, 30, 40, 50, 60, 70, 80, 90, and 100. Figure 7 [Figure 7: see original paper] shows the trend of classification accuracy versus particle count for the laptop computer review dataset. Classification accuracy reaches its maximum when particle count is 30. In particle swarm optimization, particle count affects ensemble system performance. Generally, system performance improves with larger particle counts, as more particles can find the global optimal region faster, while fewer particles may require longer search times and risk falling into local optima. However, excessive particle counts increase time costs as each particle must compute fitness values and adjust direction and velocity. In text sentiment classification problems, 30 particles achieve optimal classification accuracy, with slight decreases as particle count increases further. The same conclusions were reached for hotel review data and book review data. Therefore, for the scale of our experimental data, a particle count around 30 is appropriate.

## 6 Conclusion

Text sentiment classification technology has important application value in e-commerce, e-government, and information forecasting. Addressing the diversity and subtlety characteristics of Chinese text sentiment expression, this paper proposes a BPSO-based random subspace ensemble classification method built upon dependency parsing feature representation. The random subspace method reduces the dimensionality of data input in the training model by partitioning features to form training data for individual classifiers. BPSO serves as a selection mechanism that improves both the classification accuracy and generaliza-

tion capability of the ensemble system. Additionally, we comparatively studied the impact of subspace rate on classification accuracy and system diversity for standard random subspace and BPSO random subspace methods, deriving general rules for subspace rate selection in the BPSO random subspace method.

In future research, we will collect more datasets, including English datasets, to further validate our conclusions. Additionally, based on the conclusions obtained in this paper, we will construct more suitable class models for text sentiment classification.

## References

- [1] Agarwal B, Mittal N. Machine Learning Approach for Sentiment Analysis[M]. Springer, International Publishing, 2016: 21-45.
- [2] Vinodhini G, Chandrasekaran R. Sentiment Analysis and Opinion Mining: A Survey[J]. International Journal of Advanced Research in Computer Science and Software Engineering, 2012, 2(6): 282-292.
- [3] Liu B, Zhang L. A Survey of Opinion Mining and Sentiment Analysis[A]// Mining Text Data[M]. Springer US, 2012.
- [4] Zhang Qingqing, Liu Xilin. Sentiment Analysis Based on Dependency Syntactic Relation[J]. Computer Engineering and Applications, 2015, 51(22): 28-32.
- [5] Wang G, Sun J, Ma J, et al. Sentiment Classification: The Contribution of Ensemble Learning[J]. Decision Support Systems, 2014, 57(1): 77-93.
- [6] Wang G, Zhang Z, Sun J, et al. POS-RS: A Random Subspace Method for Sentiment Classification Based on Part-of-Speech Analysis[J]. Information Processing & Management, 2015, 51(4): 458-479.
- [7] Dasarathy B V, Sheela B V. A Composite Classifier System Design: Concepts and Methodology[J]. Proceedings of the IEEE, 1979, 67(5): 708-713.
- [8] Polikar R. Ensemble Based Systems in Decision Making[J]. IEEE Circuits and Systems Magazine, 2006, 6(3): 21-45.
- [9] Dietterich T G. Ensemble Methods in Machine Learning[C]// Proceedings of the 1st International Workshop on Multiple Classifier Systems.2000.
- [10] Ho T K. The Random Subspace Method for Constructing Decision Forests[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 1998, 20(8): 832-844.
- [11] Sun Bo, Wang Jiandong, Chen Haiyan, et al. Diversity Measures in Ensemble Learning[J]. Control and Decision, 2014, 29(3): 385-395.
- [12] Zhou Z H, Wu J X, Jiang Y, et al. Genetic Algorithm Based Selective Neural Network Ensemble[C]// Proceedings of the 17th International Joint Conference on Artificial Intelligence.
- [13] Tama B A, Rhee K H. A Combination of PSO-Based Feature Selection and Tree-Based Classifiers Ensemble for Intrusion Detection Systems [A]// Advances in Computer Science and Ubiquitous Computing[M]. Singapore: Springer, 2015.
- [14] Hedeshi N G, Abadeh M S. Coronary Artery Disease Detection Using

- a Fuzzy-boosting PSO Approach[J]. Computational Intelligence and Neuroscience, 2014, 2014: Article No. 783734. <http://dx.doi.org/10.1155/2014/783734>.
- [15] Tsai C Y, Chen C J. A PSO-AB Classifier for Solving Sequence Classification Problems[J]. Applied Soft Computing, 2015, 27: 11-27.
- [16] Kennedy J, Eberhart R C. A Discrete Binary Version of the Particle Swarm Algorithm[C]//Proceedings of the 1997 Conference on Systems, Man, and Cybernetics. 1997: 4104-4108.
- [17] Chandra A, Chen H, Yao X. Trade-off Between Diversity and Accuracy in Ensemble Generation[A]// Multi-objective Machine Learning[M]. Springer Berlin Heidelberg, 2006.
- [18] Ko A H R, Sabourin R, De Souza Britt Jr A. Combining Diversity and Classification Accuracy for Ensemble Selection in Random Subspaces[C]//Proceedings of the International Joint Conference on Neural Networks.2006.
- [19] Ko A H R, Sabourin R, De Souza Britto Jr A. Compound Diversity Functions for Ensemble Selection[J]. International Journal of Pattern Recognition and Artificial Intelligence, 2009, 23(4): 659-686.

## Author Contributions Statement

Zhang Qingqing: Proposed research ideas, designed research plan, completed experiments.

Liu Xilin: Drafted manuscript and revised final version.

## Conflict of Interest Statement

All authors declare no conflict of interest.

## Supporting Data

Supporting data is self-archived by the authors, E-mail: [suiyue2959@163.com](mailto:suiyue2959@163.com).

- [1] Zhang Qingqing. hotel4000.xls. Hotel review data.
- [2] Zhang Qingqing. book2000.xls. Book review data.
- [3] Zhang Qingqing. notebook1000.xls. Laptop computer review data.
- [4] Zhang Qingqing. useRSTraining.m. Random subspace training program.
- [5] Zhang Qingqing. GetDependencyAndDotlines.java. Triple dependency relation feature extraction program.
- [6] Zhang Qingqing. main.m. BPSO training program.
- [7] Zhang Qingqing. batchRSPSO.m. RS\_{BPSO} training program.

**Received Date:** 2017-03-28

**Revised Date:** 2017-04-24

---

<sup>1</sup><http://www.datatang.com/data/11936>.

*Note: Figure translations are in progress. See original paper for figures.*

*Source: ChinaXiv — Machine translation. Verify with original.*