



# Exploration of classification confidence in ensemble learning



Leijun Li<sup>a</sup>, Qinghua Hu<sup>a,\*</sup>, Xiangqian Wu<sup>a</sup>, Daren Yu<sup>b</sup>

<sup>a</sup> Biometric Computing Research Centre, School of Computer Science and Technology, Harbin Institute of Technology, Harbin 150001, China

<sup>b</sup> School of Energy Science and Engineering, Harbin Institute of Technology, Harbin 150001, China

## ARTICLE INFO

### Article history:

Received 13 February 2013

Received in revised form

20 March 2014

Accepted 23 March 2014

Available online 1 April 2014

### Keywords:

Ensemble learning

Ordered aggregation

Ensemble margin

Classification confidence

## ABSTRACT

Ensemble learning has attracted considerable attention owing to its good generalization performance. The main issues in constructing a powerful ensemble include training a set of diverse and accurate base classifiers, and effectively combining them. Ensemble margin, computed as the difference of the vote numbers received by the correct class and the another class received with the most votes, is widely used to explain the success of ensemble learning. This definition of the ensemble margin does not consider the classification confidence of base classifiers. In this work, we explore the influence of the classification confidence of the base classifiers in ensemble learning and obtain some interesting conclusions. First, we extend the definition of ensemble margin based on the classification confidence of the base classifiers. Then, an optimization objective is designed to compute the weights of the base classifiers by minimizing the margin induced classification loss. Several strategies are tried to utilize the classification confidences and the weights. It is observed that weighted voting based on classification confidence is better than simple voting if all the base classifiers are used. In addition, ensemble pruning can further improve the performance of a weighted voting ensemble. We also compare the proposed fusion technique with some classical algorithms. The experimental results also show the effectiveness of weighted voting with classification confidence.

© 2014 Elsevier Ltd. All rights reserved.

## 1. Introduction

One of the main aims in the machine learning domain has always been to improve the generalization performance. Ensemble learning has gained considerable research attention for more than twenty years [3,8,12,29,33] owing to its good generalization capability. This technique trains a set of base classifiers, instead of a single one, and then combines their outputs with a fusion strategy. Numerous empirical studies and applications show that the combination of multiple classifiers usually improves the generalization performance with respect to its members [1,28,31,47,51].

There are two key issues in constructing an ensemble system: (1) learning a collection of base classifiers and (2) combining them with an effective technique. Various algorithms have been developed for learning base classifiers by perturbing training samples, parameters, or structures of base classifiers [5,6,12,31,50]. For example, Bagging [5] generates different training sets by bootstrap sampling [11], whereas Zhou and Yu proposed a technique of multi-modal perturbation to learn diverse base classifiers [50]. In 2005, a review on

the techniques of learning the diverse members was reported in [6]. The fusion strategy refers to effectively combining the outputs of the base classifiers. Currently available fusion algorithms can be roughly categorized into two schemes: one is to combine all the base classifiers with a certain strategy, such as simple voting [17] and weighted voting [3,13]. However, the investigation in [24] showed that combining part of the base classifiers, instead of all, may lead to better performance. Selective ensembles produced significantly higher accuracies than the original ensembles [41,47,49].

It is well known that a set of diverse and accurate base classifiers is the prerequisite for a successful ensemble. Indeed, effective exploitation of these base classifiers is also an important factor for designing a powerful ensemble. We will focus on the second issue in this work. An ensemble margin is considered an important factor, which has an impact on the performance of an ensemble and is utilized to interpret the success of Boosting [2,34,36,43]. Different boosting algorithms have been developed by constructing distinct loss functions based on the margin [10,12,22,33]. However, the margin defined in [36] just uses the classification decision of the base classifiers and their classification confidences are overlooked. In fact, classification confidence was theoretically proved to be a key factor on the generalization performance [35].

In this work, we want to identify the role of the classification confidence of a base classifier in ensemble learning. We generalize

\* Corresponding author.

E-mail addresses: [lileijun1985@163.com](mailto:lileijun1985@163.com) (L. Li), [hqinghua@hit.edu.cn](mailto:hqinghua@hit.edu.cn) (Q. Hu), [xqw@hit.edu.cn](mailto:xqw@hit.edu.cn) (X. Wu), [yudaren@hit.edu.cn](mailto:yudaren@hit.edu.cn) (D. Yu).

the definition of the margin based on the classification confidence. The weights of the base classifiers are trained by optimizing the margin distribution. This strategy is similar to learning a classification function in a new feature space (meta-learning), just like the stacking technique [39,44]. In stacking, the outputs of the base classifiers are viewed as new features to train a combining function. Here, we show the difference of optimizing different margins from the viewpoint of the stacked generalization, and explain the necessity of incorporating the classification confidence into the margin.

Then, we explore how to utilize the weights and the classification confidences in combining the base classifiers. Four strategies are considered in this work. The weights are used to select the base classifiers. In a selective ensemble, a function should be developed to evaluate the quality of the base classifiers [25]. Similar to feature selection, classifier selection is also a combinational optimization problem. Assume that an ensemble consists of  $L$  base classifiers. Then, there are  $2^L - 1$  nonempty sub-ensembles. Therefore, it is unfeasible to search the optimal solution via exhaustive search. In order to address this problem, several suboptimal ensemble pruning methods were proposed [1,8,16,47,49,51]. In the ordered aggregation technique, the base classifiers are selected based on the order [24–29]. The base classifiers are sorted by a specified rule, and then, they are added into the ensemble sequentially. A fraction of the base classifiers in the ordered ensemble are selected.

How to rank the base classifiers in the aggregation process is the key issue for this technique. In 1997, Reduce-Error pruning and Kappa pruning were proposed [24]. For Reduce-Error pruning, the first classifier is the one with the lowest classification error and the remaining classifiers are sequentially selected to minimize the classification error. Then, in 2004, Reduce-Error pruning without backfitting, Complementarity Measure, and Margin Distance Minimization were proposed to decide the order of the base classifiers [26], respectively. Based on the Complementarity Measure, the classifier incorporated into a sub-ensemble is the one whose performance is most complementary to this sub-ensemble. Recently, ensemble pruning via individual contribution ordering (EPIC) and uncertainty weighted accuracy (UWA) were proposed [23,29]. Moreover, in [25], the performances of some ordered aggregation-pruning algorithms have been extensively analyzed. For the proposed method, the base classifiers are sorted based on their weights in the descending order, which is similar to the method, MAD-Bagging, proposed in [46]. However, MAD-Bagging does not consider the classification confidence. While the major objective of this work is to analyze the influence of the classification confidence in ensemble learning. We try some ordered aggregation techniques to combine the base classifiers. Both the weighted and simple voting strategies are tested after pruning. The objective is to elucidate how to use the weights and the classification confidences of the base classifiers in ensemble optimization.

The main contributions of the work are listed as follows. First, we introduce the classification confidence in defining the ensemble margin and design a margin-induced loss function to compute the weights of the base classifiers. Second, we test several strategies to utilize the weights and the classification confidences in combining the base classifiers. Finally, extensive experiments are conducted to test and compare different techniques, and some guidelines for constructing a powerful ensemble are given.

The rest of the paper is organized as follows. In Section 2, we present some main notations and review the related works. In Section 3, we show how to learn the weights of the base classifiers and reveal the difference of optimizing different margins. In Section 4, we explore how to utilize these weights and the classification confidences to combine the base classifiers and propose a new ordered aggregation ensemble pruning method. Then, we analyze the proposed method in Section 5. Further, we test our algorithm on open classification tasks

and study its mechanism for improving the classification performance in Section 6. Finally, Section 7 presents the conclusions.

## 2. Notations and related works

The main notations used in this paper are summarized as follows:

- $h_j (j = 1, 2, \dots, L)$ : the base classifiers
- $L$ : the total number of the base classifiers
- $X = \{(x_i, y_i), i = 1, 2, \dots, n\}$ : the pruning set
- $y_i$ : the true class label of the sample  $x_i$
- $\hat{y}_{ij}$ : the classification decision of  $x_i$  estimated by the classifier  $h_j$
- $r_{ij}$ : the classification confidence of  $x_i$  estimated by the classifier  $h_j$

In the following, first, we introduce some works related to the classification confidence, margin, and stacked generalization, and then, we present some ordered aggregation pruning methods used in our experiments.

Classification confidence is used in this paper. A classifier  $h_j$  assigns a classification confidence  $r_{ij}$  to its decision  $\hat{y}_{ij}$ . For example, considering a linear real-valued classifier  $h(x) = \psi \cdot x - b$ , the classification decision of the sample  $x$  is 1 if  $h(x) \geq 0$  and  $-1$  otherwise. Then, the value  $|h(x)|$  can be deemed as the classification confidence for its decision. In [35], the bound on the generalization error for this linear classifier was given, and it indicated that the classification confidence was an important factor for generalization. The linear classifier was also generalized to non-linear function and the detailed information can be obtained in [35]. Moreover, the classification confidence has been utilized in certain ensemble learning algorithms [12,30,32,45].

The margin is also considered as an important factor for the generalization performance of ensemble learning [36,43]. In [36], the margin of a sample with respect to an ensemble was introduced. Given a sample  $x_i \in X$ , its margin with respect to  $\{h_1, \dots, h_L\}$  is defined as

$$m(x_i) = \sum_{j=1}^L w_j A_{ij},$$

$$\text{s.t. } w_j \geq 0, \quad \sum_{j=1}^L w_j = 1, \quad (1)$$

where  $w_j$  is the weight of the classifier  $h_j$  and

$$A_{ij} = \begin{cases} 1 & \text{if } y_i = \hat{y}_{ij} \\ -1 & \text{if } y_i \neq \hat{y}_{ij} \end{cases} \quad (2)$$

In this work, a generalized definition of the margin is proposed based on the classification confidence and the weights of the base classifiers are learned through the optimization of the margin distribution. We will discuss the difference of optimizing different margins from the viewpoint of stacked generalization [39,44]. The stacking algorithms learn the weights of the base classifiers by training a function in a new feature space. In [44], the classification decision of the base classifier was used as the input feature. Then, the classification confidence was introduced and the stacking performance was improved [39].

In the classifiers ensemble, we are generally given a set of base classifiers  $\{h_1, \dots, h_L\}$ , which are obtained by certain learning algorithms [5,6,12,31,50]. Then, they are combined with some strategies such as the simple voting or the weighted voting. The simple voting implies that the class that receives the most votes is considered as the final decision. In the weighted voting, the votes are weighted and the final ensemble decision is the class that receives the largest weight coefficients sum of votes.

It was shown that selectively combining some of the base classifiers may lead to a better performance than combining all of

them [24]. Based on this observation, some ensemble pruning techniques were constructed, which select a fraction of the candidate base classifiers from  $\{h_1, \dots, h_L\}$  based on various strategies. As there are  $2^L - 1$  nonempty sub-ensembles, it is impossible to check all of them for obtaining the best solution. Approximate optimal approaches were proposed [40,48].

The focus of this paper is ensemble pruning using the ordered aggregation technique [25]. It selects a sub-ensemble based on modifying the order of the base classifiers, i.e., the random order of  $h_1, \dots, h_L$  should be replaced by a specified sequence  $h_{s_1}, \dots, h_{s_u}$  via a rule. Starting with an empty set  $S$ , the base classifiers are iteratively added into this set based on their order in the specified sequence. Here, the base classifiers are ordered based on the pruning set that can be a training set or a separate set. In [25], it is shown that the performance of using all available data for training as well as pruning is better than that of withholding some data. Thus, the training set is used for pruning in this paper. It can be seen that the key factor for ordered aggregation pruning is sorting the base classifiers and their time complexity is a polynomial in the size of the base classifiers set. In the following, some methods used in our experiments are reviewed briefly.

In [26], the base classifiers are sorted based on margin distance minimization and the corresponding ensemble pruning method is called MDM. In particular, the  $u$ th classifier incorporated into the current sub-ensemble  $S_{u-1}$  is

$$s_u = \arg \min_j d\left(\mathbf{o}, \frac{1}{u} \left(\mathbf{c}_j + \sum_{t=1}^{u-1} \mathbf{c}_{s_t}\right)\right), \quad (3)$$

where  $\mathbf{c}_j$  is the  $n$  dimensional signature vector of  $h_j$  whose  $i$ th component  $(\mathbf{c}_j)_i$  is 1 if the sample  $x_i$  is correctly classified by  $h_j$  and  $-1$  otherwise. The objective point  $\mathbf{o}$  is placed in the first quadrant with equal components  $\mathbf{o}_i = p$  and  $0 < p < 1$ ,  $j$  runs throughout the classifiers outside  $S_{u-1}$ , and  $d(\mathbf{t}, \mathbf{v})$  is the usual Euclidean distance between the vectors  $\mathbf{t}$  and  $\mathbf{v}$ . Then, in [25], the improved version of MDM is proposed, and it uses a moving objective point  $\mathbf{o}$  that allows  $p(u)$  to vary with the size of the sub-ensemble  $u$ . Exploratory experiments show that a value  $p(u) \propto \sqrt{u}$  is appropriate. In this paper, we use the improved version for comparison with our method.

In the Complementarity Measure [26], the  $u$ th classifier incorporated into the current sub-ensemble  $S_{u-1}$  is

$$s_u = \arg \max_j \sum_{i=1}^n I(\hat{y}_{ij} = y_i \wedge H_{S_{u-1}}(x_i) \neq y_i), \quad (4)$$

where  $j$  runs throughout the classifiers outside  $S_{u-1}$ . It can be seen that the selected classifier is one whose performance is complementary to that of sub-ensemble  $S_{u-1}$ .

For EPIC [23], it incorporates the base classifiers based on their contributions to the entire ensemble in the descending order and the contribution of  $h_j$  is defined as

$$IC_j = \sum_{i=1}^n (\alpha_{ij}(2\nu_{\max}^{(i)} - \nu_{\hat{y}_{ij}}^{(i)}) + \beta_{ij}\nu_{\text{sec}}^{(i)} + \theta_{ij}(\nu_{\text{correct}}^{(i)} - \nu_{\hat{y}_{ij}}^{(i)} - \nu_{\max}^{(i)})), \quad (5)$$

where  $\nu_{\max}^{(i)}$  is the number of the majority votes on  $x_i$ ,  $\nu_{\hat{y}_{ij}}^{(i)}$  is the number of predictions  $\hat{y}_{ij}$ , and  $\nu_{\text{sec}}^{(i)}$  is the second largest number of the votes on the labels of  $x_i$ . Further, for a classifier  $h_j$ ,

$$\alpha_{ij} = \begin{cases} 1 & \text{if } \hat{y}_{ij} = y_i \text{ and } \hat{y}_{ij} \text{ is in the minority group;} \\ 0 & \text{otherwise.} \end{cases}$$

$$\beta_{ij} = \begin{cases} 1 & \text{if } \hat{y}_{ij} = y_i \text{ and } \hat{y}_{ij} \text{ is in the majority group;} \\ 0 & \text{otherwise.} \end{cases}$$

$$\theta_{ij} = \begin{cases} 1 & \text{if } \hat{y}_{ij} \neq y_i; \\ 0 & \text{otherwise.} \end{cases}$$

### 3. Learning weights of base classifiers based on margin optimization

It can be seen that the margin defined in the above section does not use the classification confidence. Motivated by the conclusion in [35], we consider the classification confidence and generalize the definition of the margin in ensemble as follows.

**Definition 1.** For  $x_i \in X (i = 1, 2, \dots, n)$ , let  $r_{ij} \in [0, 1]$  be its classification confidence by the classifier  $h_j (j = 1, 2, \dots, L)$ . The margin of  $x_i$  based on the classification confidence is computed as

$$m_{cc}(x_i) = \sum_{j=1}^L w_j A_{ij} r_{ij},$$

$$\text{s.t. } w_j \geq 0, \quad \sum_{j=1}^L w_j = 1. \quad (6)$$

It is a generalization to the margin proposed in [36]. In the following, we show how to learn the weights of the base classifiers,  $w_j (j = 1, 2, \dots, L)$ , through margin distribution optimization. The objective function consists of the classification loss with respect to the margin and a regularization term.

**Definition 2.** For  $x_i \in X$ , its classification loss based on the margin is defined as

$$f(x_i) = (1 - m_{cc}(x_i))^2. \quad (7)$$

Here, the squared loss function is utilized to compute the classification loss. In fact, other loss functions can also be used, such as

$$f_1(x_i) = \log(1 + \exp(-m_{cc}(x_i))). \quad (8)$$

In this work, we will not discuss the other loss functions. The classification loss of  $X$  in terms of the squared loss function is computed as

$$f(X) = \sum_{i=1}^n f(x_i) = \|U - EW\|_2^2, \quad (9)$$

where  $U = [1, \dots, 1]_{n \times 1}^T$ ,  $W = [w_1, \dots, w_L]_{L \times 1}^T$ , and  $E = \{A_{ij} r_{ij}\}_{n \times L}$ .

Now, the optimization function can be written as

$$F = \|U - EW\|_2^2 + \lambda \|W\|_2,$$

$$\text{s.t. } w_j \geq 0, \quad \sum_{j=1}^L w_j = 1. \quad (10)$$

Here, the loss function is regularized with  $l_2$  of the weights for enlarging the margin of the decision function [14]. It is worth noting that we add a constraint to the weights that  $w_j \geq 0$  to guarantee a convex combination of the base classifiers. By minimizing  $F$ , we get  $w_j (j = 1, 2, \dots, L)$ . Several open software packages can be used to determine its solution [21]. The idea of optimizing an objective function based on the margin to compute the weights of the base classifiers was also applied in LPboosting [15]. Its aim was to maximize the minimum margin of an ensemble via linear programming [9]. Experimental results showed that LPboosting could achieve a larger minimum margin than Adaboost; however, LPboosting did not always yield a better generalization performance than Adaboost. According to [36], it is the margin distribution, rather than the minimum margin, that determines the generalization performance. In this work, we are aimed at optimizing the margin distribution of the ensemble, instead of the minimum margin.

We can see from Eq. (10) that the margin is a key factor in the objective function. Naturally, there is a question what is the difference if  $m_{cc}(x)$  is substituted by  $m(x)$  (the new objective function is denoted as  $\bar{F}$ ). We try to answer this question from the viewpoint of the stacked generalization [44].

Consider a two-class supervised learning task and the weighted voting function  $g(x) = \text{sign}(\sum_{j=1}^L w_j \hat{y}_{xj})$ , where  $\hat{y}_{xj} \in \{-1, 1\}$  is the decision of the sample  $x$  by  $h_j$ . The decisions of a sample  $x_i$  can be represented as an  $L$ -dimensional vector  $(\hat{y}_{i1}, \dots, \hat{y}_{iL})$ . Then,  $g(x)$  can be deemed as the classification function in the  $L$ -dimensional feature space, and  $w_j$  ( $j = 1, 2, \dots, L$ ) are the coefficients of this function. It is known that the coefficients can be trained by optimizing an objective function [14]. We use the squared loss function in this work:

$$\sum_{i=1}^n \left[ 1 - y_i * \sum_{j=1}^L w_j \hat{y}_{ij} \right]^2 + \lambda \|W\|_2. \tag{11}$$

It can be seen that  $y_i * \hat{y}_{ij} = \Lambda_{ij}$  and  $m(x_i) = \sum_{j=1}^L w_j \Lambda_{ij} = y_i * \sum_{j=1}^L w_j \hat{y}_{ij}$  when  $\hat{y}_{ij}, y_i \in \{-1, 1\}$ . Thus,

$$\sum_{i=1}^n \left[ 1 - y_i * \sum_{j=1}^L w_j \hat{y}_{ij} \right]^2 + \lambda \|W\|_2 = \sum_{i=1}^n [1 - m(x_i)]^2 + \lambda \|W\|_2 = \hat{F}. \tag{12}$$

It is easy to derive that minimizing  $\hat{F}$  is similar to training  $l_2$ -SVM [37] in the new feature space except that a constraint  $w_j \geq 0$  is considered in ensemble learning.

However, if the sample  $x_i$  is represented as  $(\hat{y}_{i1}, \dots, \hat{y}_{iL})$ , there exists a drawback that its distinguishing ability between different samples is poor. Given two base classifiers, there are only four different cases  $\{(1,1), (1,-1), (-1,1), (-1,-1)\}$  for representing all samples. We can see that most of them will overlap.

In order to overcome this limitation, the classification confidence is used and each sample  $x_i$  is represented with  $(\hat{y}_{i1} * r_{i1}, \dots, \hat{y}_{iL} * r_{iL})$ . Then, the weighted voting can be written as  $g(x) = \text{sign}(\sum_{j=1}^L w_j \hat{y}_{xj} r_{xj})$ . It is known that although two samples have the same classification decision, their classification confidences are generally different. Therefore, compared to  $(\hat{y}_{i1}, \dots, \hat{y}_{iL})$ , the distinguishing ability between different samples is enhanced when  $x_i$  is represented by  $(\hat{y}_{i1} * r_{i1}, \dots, \hat{y}_{iL} * r_{iL})$ . The objective function can be written as

$$\sum_{i=1}^n \left[ 1 - y_i * \sum_{j=1}^L w_j \hat{y}_{ij} r_{ij} \right]^2 + \lambda \|W\|_2. \tag{13}$$

In this case,  $\Lambda_{ij} r_{ij} = y_i * \hat{y}_{ij} r_{ij}$ , where  $\hat{y}_{ij}, y_i \in \{-1, 1\}$ , and  $m_{cc}(x_i) = \sum_{j=1}^L w_j \Lambda_{ij} r_{ij} = y_i * \sum_{j=1}^L w_j \hat{y}_{ij} r_{ij}$ . Thus,

$$\sum_{i=1}^n \left[ 1 - y_i * \sum_{j=1}^L w_j \hat{y}_{ij} r_{ij} \right]^2 + \lambda \|W\|_2 = \sum_{i=1}^n [1 - m_{cc}(x_i)]^2 + \lambda \|W\|_2 = F. \tag{14}$$

Take the UCI data set ‘‘hepatitis’’ as an example. If there are two base classifiers, their outputs form a 2-D decision space, as shown in Fig. 1. Each sample is represented as a 2-dimensional vector. If the classification confidence is not considered, there are only four cases of the decision outputs:  $\{(1, 1), (-1, 1), (1, -1), (-1, -1)\}$ . However, when the classification confidence is considered, the representation ability is enhanced and the samples can be distinguished. Further, their corresponding decision functions are different. In Fig. 1, the solid red and short-dashed green lines denote the decision functions based on  $F$  and  $\hat{F}$ , respectively.

We can obtain the following information from Fig. 1. (1) Some samples are misclassified by these two decision functions; (2) if we consider the classification confidence, the samples are scattered in the 2-D space; otherwise, they are located at four points; (3) there is significant difference between the two classification functions. This difference comes from the second observation. If a sample is misclassified by both classifiers, however, their classification losses are different with respect to the corresponding classifiers as their

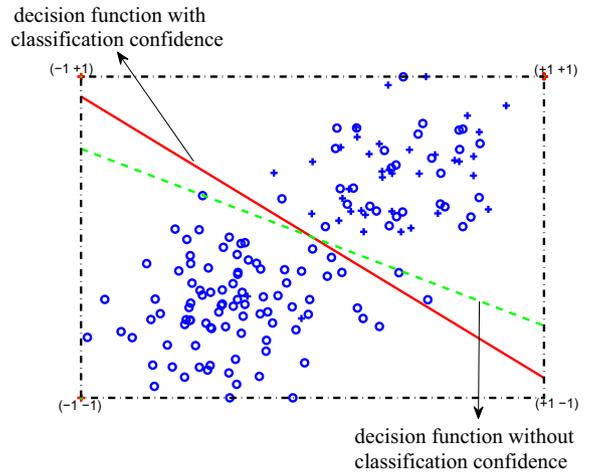


Fig. 1. The difference of decision functions based on different margin definitions. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

classification margins are distinct. When we do not consider the classification confidence, we assign the same classification loss to all the misclassified samples. However, if the classification confidence is taken into account, their losses are different. Thus a finer representation is provided in this case, which leads to the performance improvement in the final ensemble.

#### 4. Ensemble pruning based on classification confidence

Now, we explore how to utilize the classification confidences of the base classifiers and the weights to combine their outputs.

Given  $\{h_1, \dots, h_L\}$  and  $X = \{(x_i, y_i), i = 1, 2, \dots, n\}$ , we minimize  $F$  and obtain  $w_j$  ( $j = 1, 2, \dots, L$ ). The base classifiers are iteratively added into an empty set based on their weights in the descending order, and then, the sub-ensemble with the best performance on the pruning set is selected as the pruned ensemble. The ordered aggregation is a greedy technique and it has been used in some algorithms [7,23,25–29]. In this process, the classification confidence can be utilized in two steps. One is used in learning the weights of the base classifiers via optimizing the margin distribution and the other is used in weighted voting. Therefore, we need to identify whether the classification confidence should be used in both of them. Thus, we try four strategies to use the classification confidences:

1. EP-CC: the classification confidence is used in learning the weights of the base classifiers as well as in weighted voting.
2. EP-CD: the classification decision is used in both learning the weights of the base classifiers and weighted voting.
3. EP-WL-CC: the classification confidence is considered in learning the weights of the base classifiers, but not in weighted voting.
4. EP-WV-CC: the classification confidence is utilized in weighted voting, but not in learning the weights of the base classifiers.

The relationship between these strategies is described in Fig. 2. All the above mentioned strategies can be understood in the framework of the ordered aggregation. Their difference lies in whether classification confidence is utilized in learning the weights of the base classifiers and weighted voting. EP-CC is formulated in Algorithm 1.

#### Algorithm 1. EP-CC.

Input:

- $X = \{x_1, \dots, x_n\}$ : the pruning set;
- $Y = \{y_1, \dots, y_n\}$ : the real labels of the pruning set;

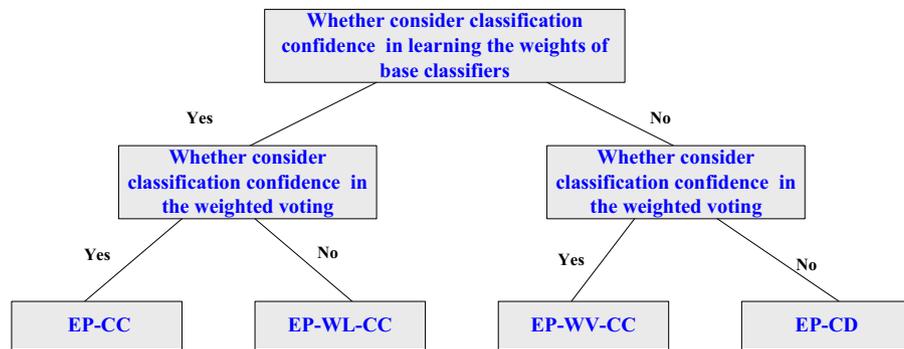


Fig. 2. Relationship between EP-CC, EP-CD, EP-WL-CC and EP-WV-CC.

- $\{h_1, \dots, h_L\}$ : the set of the candidate base classifiers;
- $x$ : a test sample;

Output: the predicted label of  $x$ ;

- 1: Apply  $h_j(j = 1, 2, \dots, L)$  on  $X$  to obtain the classification decision  $H_j = [\hat{y}_{1j}, \dots, \hat{y}_{nj}]$  and the corresponding classification confidence  $R_j = [r_{1j}, \dots, r_{nj}]$ ;
- 2: Minimize  $F$  to obtain the weight coefficients  $w_j(j = 1, 2, \dots, L)$  for the base classifiers;
- 3: Sort the base classifiers as  $\{h_{s_1}, \dots, h_{s_L}\}$  according to their weights in the descending order;
- 4: For  $j = 1, 2, \dots, L$
- 5: Classify  $X$  using the first  $j$  classifiers  $\{h_{s_1}, \dots, h_{s_j}\}$  with the weighted voting based on the classification confidence;
- 6: Compute the classification accuracy  $a_j$ ;
- 7: End for
- 8: Find  $\Gamma$  ( $\Gamma \leq L$ ) with the maximal accuracy  $a_\Gamma$  and select the first  $\Gamma$  classifiers  $\{h_{s_1}, \dots, h_{s_\Gamma}\}$  as the pruned ensemble.
- 9: Use  $\{h_{s_1}, \dots, h_{s_\Gamma}\}$  to classify  $x$  with the weighted voting based on the classification confidence and obtain its prediction.

Table 1

Description of 20 classification tasks.

Data set	Instances	Features	Classes
abalone	4177	8	3
balancescale	625	4	3
crx	690	15	2
derm	366	34	6
ecoli	336	7	8
german	1000	20	2
heart	270	13	2
hepatitis	155	19	2
horse	368	22	2
iono	351	34	2
lung cancer	96	7129	3
movement	360	90	15
mushroom	8124	22	2
pima	768	8	2
satellite	6435	36	7
segmentation	2310	19	7
spam	4601	57	2
wdbc	569	30	2
wdbc	198	33	2
yeast	1484	7	2

Table 2

Classification performances of EP-CC, EP-CD, EP-WL-CC and EP-WV-CC.

Data set	EP-CC	EP-CD	EP-WL-CC	EP-WV-CC
abalone	55.91 ± 0.36	55.54 ± 0.36 •	55.59 ± 0.19 •	55.62 ± 0.31 •
balancescale	89.05 ± 0.65	88.62 ± 0.56 •	88.92 ± 0.53	88.87 ± 0.60
crx	86.61 ± 0.64	85.33 ± 0.35 •	85.71 ± 0.41 •	85.87 ± 0.39 •
derm	98.53 ± 0.55	97.71 ± 0.48 •	97.98 ± 0.46 •	97.92 ± 0.51 •
ecoli	88.01 ± 0.60	87.20 ± 0.64 •	87.36 ± 0.82 •	87.45 ± 0.61 •
german	74.77 ± 0.41	75.58 ± 0.46 ◦	74.69 ± 0.40	74.65 ± 0.50
heart	85.12 ± 0.69	84.87 ± 0.71 •	84.97 ± 0.62	84.95 ± 0.63
hepatitis	90.45 ± 1.49	88.95 ± 1.68 •	89.63 ± 1.54 •	89.54 ± 1.73 •
horse	93.48 ± 0.42	93.05 ± 0.68 •	93.54 ± 0.47	93.23 ± 0.71 •
iono	89.41 ± 0.75	88.73 ± 0.85 •	88.86 ± 0.74 •	88.81 ± 0.71 •
lung cancer	83.14 ± 2.06	81.56 ± 2.15 •	82.05 ± 2.01 •	81.97 ± 2.28 •
movement	79.19 ± 0.95	79.60 ± 1.16 ◦	79.59 ± 0.94 ◦	79.10 ± 1.18
mushroom	98.72 ± 0.07	98.98 ± 0.09 ◦	98.71 ± 0.10	98.75 ± 0.08
pima	78.09 ± 0.49	77.70 ± 0.44 •	77.81 ± 0.41 •	77.87 ± 0.45
satellite	87.08 ± 0.08	87.07 ± 0.12	87.10 ± 0.08	87.05 ± 0.10
segmentation	93.27 ± 0.13	93.24 ± 0.12	93.29 ± 0.11	93.24 ± 0.13
spam	90.85 ± 0.11	90.68 ± 0.13 •	90.71 ± 0.12 •	90.79 ± 0.15
wdbc	98.19 ± 0.26	97.75 ± 0.22 •	97.92 ± 0.24 •	97.83 ± 0.25 •
wdbc	81.98 ± 1.93	80.15 ± 1.97 •	81.05 ± 1.82 •	80.67 ± 2.06 •
yeast	74.89 ± 0.23	74.80 ± 0.14	74.85 ± 0.17	74.87 ± 0.19
Win–Tie–Loss		14–3–3	11–8–1	10–10–0

It should be noted that the weighted voting based on the classification confidence proposed in Algorithm 1 means that  $w_j * r_{xj}$  is computed as the weight of  $\hat{y}_{xj}$ . Here,  $w_j$  is learned as the second step of Algorithm 1,  $\hat{y}_{xj}$  is the classification decision of the sample  $x$  by the base classifier  $h_j$ , and  $r_{xj}$  is its corresponding classification confidence.

In the next section, we will show the classification performances of these four strategies, explore the role of classification confidence, and present the necessity of pruning.

## 5. Algorithm analysis

Table 1 summarizes 20 UCI data sets [4] used in this work. Linear SVM is introduced as the base classification algorithm and the cost of constrain violation is 1. We consider the distance between  $x$  and the hyperplane,  $d_{xj}$ , as the classification confidence. However it takes values in  $[0, +\infty)$ , which is not suitable for Definition 2. Thus we compute the classification confidence of  $x$  with respect to  $h_j$  as  $r_{xj} = 1/(1 + \exp(-d_{xj}))$  [42]. In this case the ensemble margin based on the classification confidence takes value in the interval  $[-1, 1]$ . For a multi-class task, a one-against-one method is utilized to transform the task into multiple two-class tasks. In particular,  $K(K-1)/2$  two-class SVMs are constructed for a  $K$ -class classification task ( $K > 2$ ) and each SVM is trained on only two out of all  $K$  classes. For a given sample  $x$ , every SVM assigns the classification confidence  $r_{xj}$  to its classification decision  $\hat{y}_{xj}(j = 1, \dots, K(K-1)/2)$ . Then, the final classification decision is the class label that receives the most votes. The final

classification confidence is the minimum value in the classification confidences whose classification decision is the final classification decision. In addition, some other learning algorithms can also be used to train the base classifiers as long as an appropriate classification

confidence is given. For example, the classification confidence of the decision tree can be calculated as that in [30].

For the experiments in this paper, 20 runs of 5-fold cross validation are performed. For each trial, the data set is randomly split into

5 subsets: 4 subsets are used for training the base classifiers, optimizing the weights and obtaining the pruned ensemble, 1 for testing. The average accuracy and the standard deviation over 20 runs are calculated to evaluate the performances of the algorithms. In [18], the

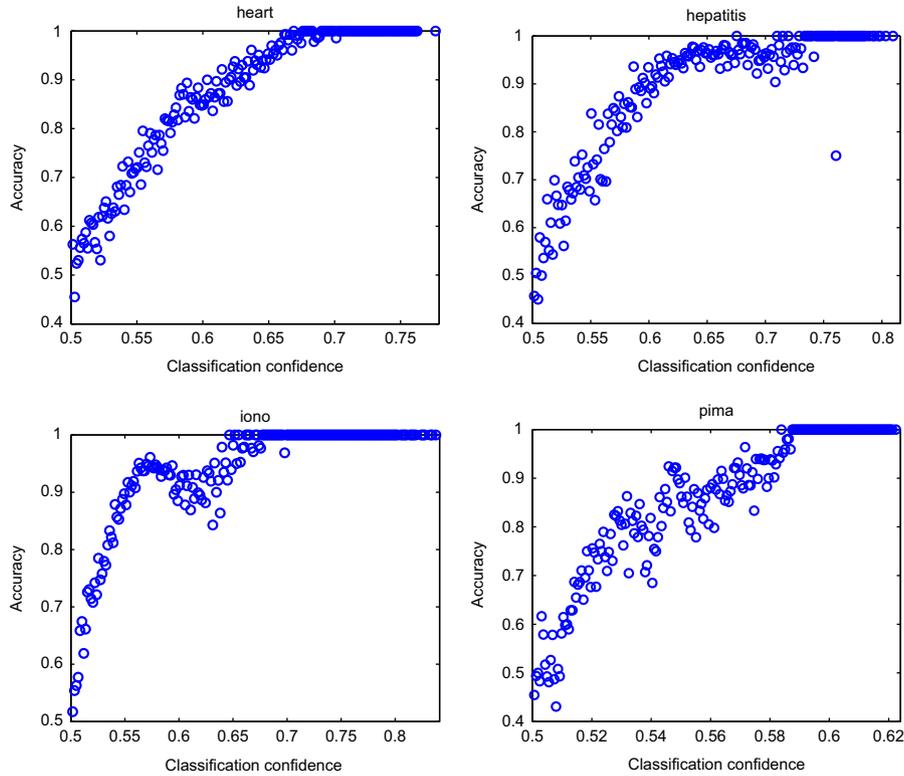


Fig. 3. Variation of classification accuracies with classification confidences.

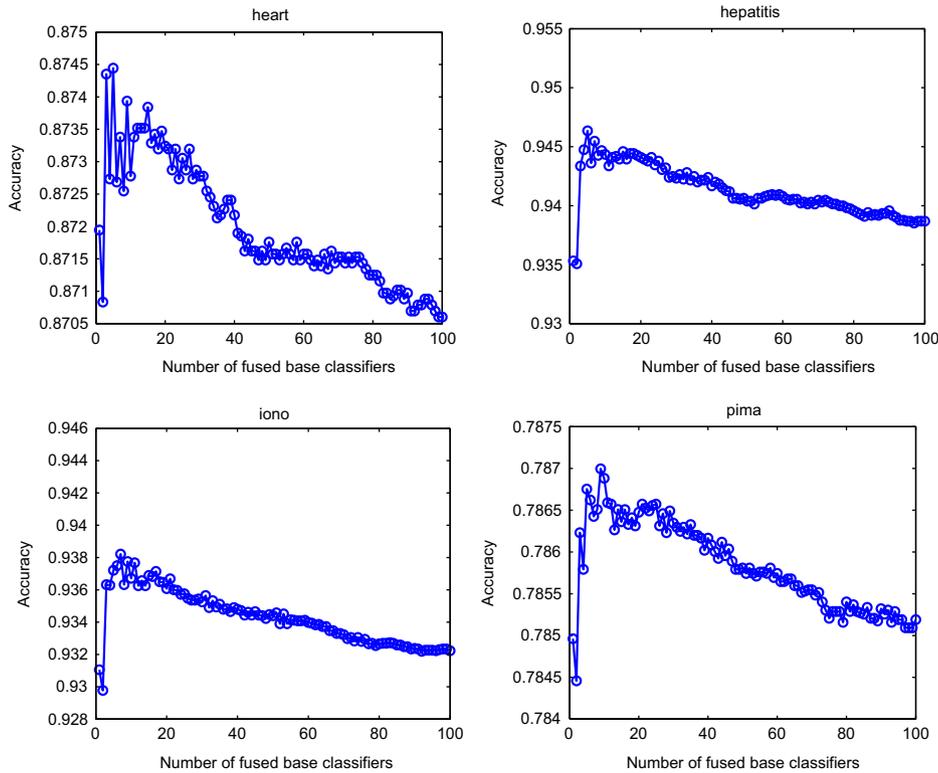


Fig. 4. Variation of classification accuracies with number of base classifiers in pruning set.

bootstrapping technique [11] was utilized to train the individual SVM for ensemble. In this work, we also adopt bootstrapping to train 100 base classifiers and the ratio of bootstrap sampling is set as 0.75.

Further, in order to compare EP-CC with the other methods from the statistical viewpoint, a one-tailed paired *t*-test was performed and the significance level was set as 0.05. A bullet means that EP-CC behaves significantly better than the corresponding method, whereas an open circle indicates that EP-CC is significantly worse than the corresponding method. “Win” means that EP-CC performs significantly better than the corresponding method; “Tie” means that the difference is not significant, and “Loss” indicates that EP-CC behaves significantly worse than the corresponding algorithm.

Table 2 summarizes the average generalization accuracies and the standard deviations for EP-CC, EP-CD, EP-WL-CC, and EP-WV-CC. It can be seen that EP-CC performs significantly better than EP-CD on 14 data sets, worse on 3, and the difference is not significant on the remaining 3 task. Compared to EP-CD, EP-CC utilizes the classification confidence in learning the weights of the base classifiers and the weighted voting. Thus, it validates the necessity of the classification

confidence for improving the generalization performance. Moreover, EP-CC is better than EP-WL-CC and EP-WV-CC. Thus, it can be concluded that the classification confidence should be used in both learning the weights of the base classifiers and the weighted voting. Both of them aid in improving the generalization performance.

We consider the classification confidence in learning the weights of base classifiers due to the assumption that the classification confidence is relevant with the classification performance. Now we use four data sets to show their statistical relation. Fig. 3 shows the variation of classification accuracies with the classification confidences, where the *x*-axis represents the classification confidence and the *y*-axis gives the classification accuracy. In the experiment, we compute the classification confidences of all the test samples in terms of one hundred base classifiers. Then we equally divide the interval between the minimal confidence and the maximal confidence into 200 bins, and compute the classification accuracy of the samples located in each bin with respect to the corresponding base classifiers.

From Fig. 3, we see that statistically the classification accuracy rises if the confidence increases for all the four classification tasks. This results empirically validate the conclusion in [35].

In the pruning process of EP-CC, the base classifiers are sequentially added into a pool according to their weights. Fig. 4 shows the variation of the classification accuracy on the pruning set in this process, where the *x*-axis represents the number of the fused classifiers and the *y*-axis represents the classification accuracies of the ensembles with weighted voting. We see that the classification accuracies rise initially and then drop slowly. It shows that fusion with part of the classifiers can obtain better performance. Finally the subset of base classifiers producing the best performance are selected.

Then, we need to identify whether the obtained ensemble still performs better than combining all the base classifiers with weighted voting (WV) in the test set. We conduct some experiments to answer this question. For comparison, the classification performance of simple voting with all base classifiers (SV) is also given. Table 3 lists the results of EP-CC, SV, and WV. It is easy to see that the proposed weighted voting strategy is better than simple voting when all the base classifiers are combined. Moreover, pruning can boost the generalization performance further.

Finally, the classification performances of the ensembles with fixed ratios of the base classifiers in the test set are listed in Table 4. In this experiment, the base classifiers are ordered according to their weights, and then, the ensembles of the first 5%, 10%, 20%, 40%, 60%, and 100% of the base classifiers are evaluated. The bold accuracy is the

**Table 3**  
Classification performances of SV, WV and EP-CC.

Data set	EP-CC	SV	WV
abalone	55.91 ± 0.36	54.45 ± 0.12 •	55.18 ± 0.37 •
balancescale	89.05 ± 0.65	87.27 ± 0.63 •	87.89 ± 0.82 •
crx	86.61 ± 0.64	85.12 ± 0.16 •	85.59 ± 0.33 •
derm	98.53 ± 0.55	97.26 ± 0.41 •	97.76 ± 0.49 •
ecoli	88.01 ± 0.60	83.58 ± 0.44 •	86.09 ± 0.52 •
german	74.77 ± 0.41	73.65 ± 0.54 •	73.78 ± 0.56 •
heart	85.12 ± 0.69	83.57 ± 0.58 •	83.87 ± 0.61 •
hepatitis	90.45 ± 1.49	88.27 ± 0.90 •	88.43 ± 1.07 •
horse	93.48 ± 0.42	92.01 ± 0.48 •	93.05 ± 0.53 •
iono	89.41 ± 0.75	87.11 ± 0.61 •	88.27 ± 0.69 •
lung cancer	83.14 ± 2.06	79.79 ± 1.57 •	80.29 ± 1.95 •
movement	79.19 ± 0.95	75.82 ± 0.99 •	76.84 ± 1.01 •
mushroom	98.72 ± 0.07	98.40 ± 0.17 •	98.51 ± 0.14 •
pima	78.09 ± 0.49	77.23 ± 0.47 •	77.45 ± 0.43 •
satellite	87.08 ± 0.08	86.70 ± 0.08 •	86.78 ± 0.11 •
segmentation	93.27 ± 0.13	92.68 ± 0.22 •	92.89 ± 0.15 •
spam	90.85 ± 0.11	90.03 ± 0.16 •	90.45 ± 0.12 •
wdbc	98.19 ± 0.26	97.64 ± 0.27 •	97.79 ± 0.22 •
wdbc	81.98 ± 1.93	77.67 ± 0.88 •	80.12 ± 1.82 •
yeast	74.89 ± 0.23	74.21 ± 0.18 •	74.55 ± 0.20 •
Win-Tie-Loss		20-0-0	20-0-0

**Table 4**  
Classification performances of ensembles with fixed ratios of base classifiers.

Data set	<i>r</i> =5%	<i>r</i> =10%	<i>r</i> =20%	<i>r</i> =40%	<i>r</i> =60%	<i>r</i> =100%
abalone	<b>55.59 ± 0.31</b>	55.46 ± 0.36	55.42 ± 0.31	55.29 ± 0.34	55.19 ± 0.39	55.18 ± 0.37
balancescale	<b>88.09 ± 0.64</b>	87.91 ± 0.58	87.87 ± 0.64	87.83 ± 0.70	87.86 ± 0.73	87.89 ± 0.82
crx	85.72 ± 0.52	85.98 ± 0.33	<b>86.03 ± 0.29</b>	85.76 ± 0.34	85.57 ± 0.31	85.59 ± 0.33
derm	97.62 ± 0.43	<b>97.86 ± 0.46</b>	97.65 ± 0.45	97.78 ± 0.47	97.71 ± 0.42	97.76 ± 0.49
ecoli	<b>86.83 ± 0.72</b>	86.37 ± 0.51	86.31 ± 0.49	86.20 ± 0.43	86.16 ± 0.39	86.09 ± 0.52
german	73.68 ± 0.86	73.88 ± 0.47	<b>73.93 ± 0.55</b>	73.87 ± 0.58	73.79 ± 0.62	73.78 ± 0.56
heart	83.81 ± 0.72	<b>84.08 ± 0.59</b>	83.86 ± 0.65	83.99 ± 0.71	83.95 ± 0.73	83.87 ± 0.61
hepatitis	88.36 ± 1.77	<b>88.61 ± 1.57</b>	88.29 ± 1.26	88.52 ± 1.21	88.45 ± 1.16	88.43 ± 1.07
horse	92.81 ± 0.59	93.01 ± 0.44	93.03 ± 0.43	93.01 ± 0.59	92.98 ± 0.56	<b>93.05 ± 0.53</b>
iono	88.21 ± 0.76	<b>88.53 ± 0.73</b>	88.41 ± 0.71	88.31 ± 0.67	88.43 ± 0.64	88.27 ± 0.69
lung cancer	<b>81.43 ± 1.94</b>	81.05 ± 2.04	80.39 ± 2.02	80.45 ± 2.09	80.31 ± 2.12	80.29 ± 1.95
movement	76.56 ± 1.28	77.23 ± 1.42	<b>77.35 ± 1.40</b>	77.00 ± 1.02	77.02 ± 1.08	76.84 ± 1.01
mushroom	98.59 ± 0.12	<b>98.62 ± 0.11</b>	98.56 ± 0.15	98.54 ± 0.10	98.53 ± 0.14	98.51 ± 0.14
pima	77.23 ± 0.45	77.31 ± 0.42	77.39 ± 0.37	<b>77.43 ± 0.41</b>	77.42 ± 0.39	77.45 ± 0.43
satellite	86.85 ± 0.14	86.87 ± 0.09	<b>86.88 ± 0.09</b>	86.84 ± 0.11	86.81 ± 0.11	86.78 ± 0.11
segmentation	93.00 ± 0.15	<b>93.13 ± 0.14</b>	93.07 ± 0.14	93.01 ± 0.16	92.99 ± 0.15	92.89 ± 0.15
spam	90.60 ± 0.10	<b>90.62 ± 0.13</b>	90.53 ± 0.11	90.51 ± 0.15	90.47 ± 0.15	90.45 ± 0.12
wdbc	97.76 ± 0.32	<b>97.85 ± 0.35</b>	97.81 ± 0.33	97.75 ± 0.21	97.79 ± 0.24	97.79 ± 0.22
wdbc	<b>80.25 ± 1.83</b>	80.06 ± 1.79	80.19 ± 1.73	80.02 ± 1.65	80.09 ± 1.76	80.12 ± 1.82
yeast	<b>74.68 ± 0.21</b>	74.67 ± 0.17	74.62 ± 0.17	74.58 ± 0.18	74.57 ± 0.15	74.55 ± 0.20

highest one. From Table 4, we see that the ensembles of the first 5% and 10% achieve the highest accuracy on six and eight data sets, respectively. However, the systems with the first 20%, 40%, and 100% of the candidate base classifiers only achieve the highest accuracy on four, one, and one sets, respectively. These results indicate that most of the candidate base classifiers should be removed.

### 6. Empirical comparison and analysis

In this section, we present the experiments for comparing the performances of EP-CC and some other related algorithms. We compare EP-CC with a single classifier and some classical ordered aggregation pruning algorithms, including EPIC [23], the improved version of MDM [25] and CM [26].

The experimental settings are described in Section 5 and 20 runs of 5-fold cross validation are performed. For EPIC, MDM and CM, we use their default parameters to sort the base classifiers and then the sub-ensemble with the best performance in the pruning set is selected as the final system. Table 5 summarizes the average generalization accuracy of each algorithm. A one-tailed paired *t*-test was performed to compare EP-CC with the other methods. The significance level was set as 0.05. From Table 5, we see that EP-CC performs significantly better than a single SVM classifier on all classification tasks. Compared to EPIC, the statistically significant difference is favorable in 18 tasks over 20, and not significant in 2 sets. Meanwhile, EP-CC outperforms MDM and CM in most of the cases. It validates the effectiveness of the proposed method.

We also summarize the average number of the base classifiers selected by different pruning methods in Table 6. The number in bold is the largest one. We can see that EPIC selects the most base classifiers on 17 data sets. It appears that EPIC tends to select more base classifiers than the other methods.

We want to know why EP-CC can boost the generalization performance compared with the other pruning methods. In what follows, we explore this question.

The fusion strategy used by EP-CC is different from that used by the other pruning methods. EPIC, MDM, and CM use the simple voting to combine the selected base classifiers, whereas EP-CC uses the weighted voting based on the classification confidence. Thus, we wonder whether this fusion strategy is beneficial to the classification performance of EP-CC. In order to elucidate it, we compare the

generalization accuracies of EP-CC and those derived from EP-SV, which uses simple voting to combine the base classifiers selected by EP-CC. The results are given in Table 7. It is easy to see that EP-CC performs significantly better than EP-SV on 16 classification tasks. EP-CC produces the better performances on the other 4 tasks, but the difference is not significant. These results indicate that the weighted voting based on the classification confidence is one factor that aids EP-CC for improving the generalization performance.

However, compared to the other pruning methods, EP-SV generally performs better. Thus, it indicates that other factors should also be considered. It is well known, the accuracies of single base classifiers and their diversity are two important factors for evaluating the ensemble performance [19]. Thus, some experiments were conducted to explore these base classifiers selected by EP-CC and the other pruning methods.

First, we explore the generalization performance of the single base classifier selected by EP-CC. In the EP-CC algorithm, the base classifiers are sorted based on their weights and those with large weights are selected to compose the pruned ensemble. Then,

**Table 6**  
Number of selected base classifiers with different pruning methods.

Data set	EP-CC	EPIC	MDM	CM
abalone	5.18	<b>6.16</b>	3.14	3.35
balancescale	12.75	<b>50.27</b>	26.28	23.03
crx	<b>24.39</b>	4.01	3.22	3.02
derm	3.97	<b>7.07</b>	4.64	5.09
ecoli	3.56	<b>7.95</b>	3.18	4.26
german	10.22	<b>15.67</b>	4.92	5.30
heart	5.01	<b>5.48</b>	2.94	3.88
hepatitis	5.71	<b>7.37</b>	3.27	3.52
horse	7.38	<b>10.54</b>	5.18	4.80
iono	7.81	<b>10.15</b>	4.40	3.57
lung cancer	4.19	3.01	4.25	<b>5.03</b>
movement	12.20	17.75	<b>20.91</b>	16.55
mushroom	8.27	<b>25.12</b>	8.40	10.03
pima	7.52	<b>14.13</b>	3.48	2.97
satellite	10.22	<b>11.92</b>	5.54	6.62
segmentation	7.64	<b>12.06</b>	4.15	4.31
spam	4.63	<b>5.30</b>	4.03	4.25
wdbc	3.89	<b>4.65</b>	3.75	4.02
wdbc	4.93	<b>6.32</b>	4.12	5.03
yeast	4.12	<b>4.34</b>	3.01	4.16

**Table 5**  
Classification performances of SVM and ensemble pruning techniques.

Data set	EP-CC	SVM	EPIC	MDM	CM
abalone	55.91 ± 0.36	54.46 ± 0.15 ●	55.06 ± 0.44 ●	55.12 ± 0.40 ●	55.14 ± 0.39 ●
balancescale	89.05 ± 0.65	87.64 ± 0.28 ●	88.00 ± 0.44 ●	87.75 ± 0.74 ●	87.96 ± 0.63 ●
crx	86.61 ± 0.64	85.07 ± 0.12 ●	85.25 ± 0.36 ●	85.22 ± 0.28 ●	85.19 ± 0.33 ●
derm	98.53 ± 0.55	97.13 ± 0.41 ●	97.34 ± 0.43 ●	97.37 ± 0.59 ●	97.46 ± 0.39 ●
ecoli	88.01 ± 0.60	83.40 ± 0.52 ●	86.00 ± 0.80 ●	86.15 ± 0.55 ●	86.24 ± 0.78 ●
german	74.77 ± 0.41	73.68 ± 0.46 ●	73.96 ± 0.52 ●	73.18 ± 0.55 ●	73.33 ± 0.64 ●
heart	85.12 ± 0.69	83.48 ± 0.56 ●	84.16 ± 1.08 ●	83.96 ± 0.97 ●	84.01 ± 1.08 ●
hepatitis	90.45 ± 1.49	86.63 ± 0.81 ●	87.97 ± 1.96 ●	87.82 ± 2.06 ●	86.54 ± 2.25 ●
horse	93.48 ± 0.42	91.98 ± 0.56 ●	92.20 ± 0.65 ●	92.23 ± 0.73 ●	92.33 ± 0.77 ●
iono	89.41 ± 0.75	86.62 ± 0.67 ●	87.65 ± 1.01 ●	87.29 ± 0.60 ●	87.18 ± 0.90 ●
lung cancer	83.14 ± 2.06	80.10 ± 1.59 ●	81.09 ± 1.72 ●	80.72 ± 1.44 ●	80.47 ± 1.45 ●
movement	79.19 ± 0.95	73.05 ± 0.82 ●	77.46 ± 0.97 ●	76.94 ± 0.89 ●	77.23 ± 1.31 ●
mushroom	98.72 ± 0.07	98.31 ± 0.16 ●	98.76 ± 0.07	98.83 ± 0.09 ○	98.80 ± 0.10 ○
pima	78.09 ± 0.49	77.38 ± 0.50 ●	76.85 ± 0.36 ●	76.96 ± 0.52 ●	77.07 ± 0.53 ●
satellite	87.08 ± 0.08	86.71 ± 0.11 ●	86.83 ± 0.15 ●	86.91 ± 0.12 ●	86.82 ± 0.14 ●
segmentation	93.27 ± 0.13	91.60 ± 0.17 ●	93.06 ± 0.16 ●	92.88 ± 0.19 ●	92.92 ± 0.15 ●
spam	90.85 ± 0.11	90.04 ± 0.11 ●	90.26 ± 0.16 ●	90.18 ± 0.23 ●	90.37 ± 0.17 ●
wdbc	98.19 ± 0.26	97.59 ± 0.25 ●	97.71 ± 0.32 ●	97.53 ± 0.30 ●	97.56 ± 0.34 ●
wdbc	81.98 ± 1.93	76.60 ± 0.97 ●	79.14 ± 1.80 ●	78.84 ± 1.98 ●	78.52 ± 1.42 ●
yeast	74.89 ± 0.23	74.30 ± 0.17 ●	74.71 ± 0.20	74.86 ± 0.15	74.66 ± 0.15 ●
Win-Tie-Loss		20-0-0	18-2-0	18-1-1	19-0-1

we need to identify whether large weight of base classifier means better generalization performance. In other words, we want to know whether the pruned ensembles produce good performance because the selected base classifiers are more accurate than the reduced.

Fig. 5 shows the relationship between the classification accuracies and the weights of the base classifiers in the test set. The *x*-axis represents the ranking of the weights of the base classifiers and the *y*-axis represents its corresponding classification accuracy. On the *x*-axis, “1” means that the weight of the base classifier is the smallest and “100” means that the weight of the corresponding base classifier is the largest. These weights are learned in the second step of

Algorithm 1. Seen from Fig. 5, a large weight does not imply the high generalization accuracy for a base classifier. There is no significant relation between classification accuracies and the weights.

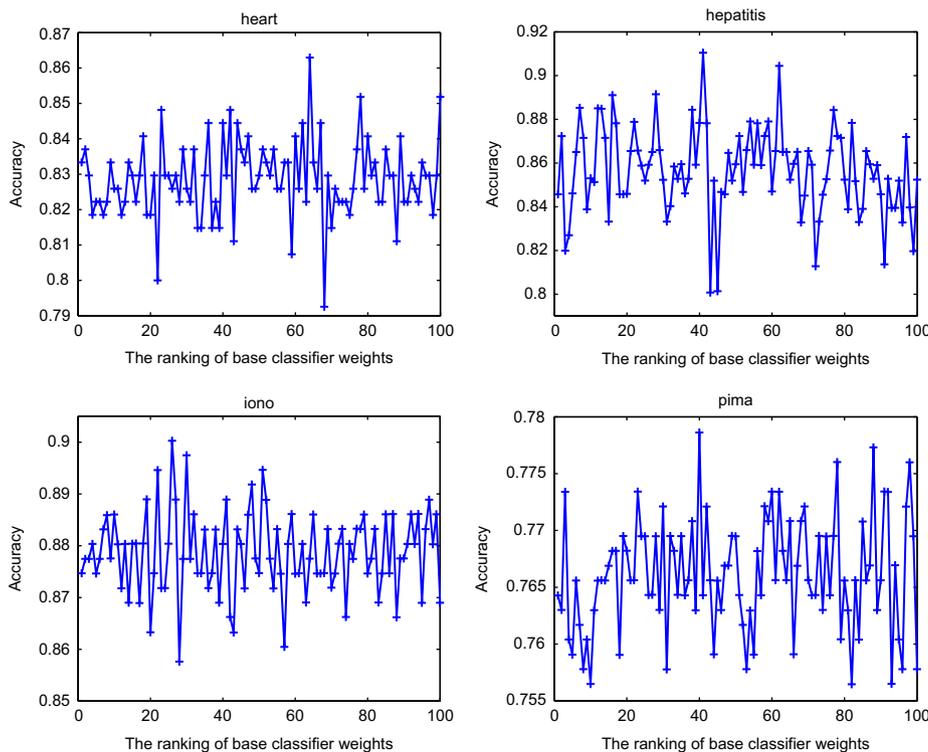
Furthermore, Table 8 summarizes the average accuracies of the base classifiers selected by EP-CC and the other pruning methods in the test set. We can see that the base classifiers selected by EP-CC tend to be more accurate than those selected by the other pruning methods; however, the difference is not significant. From Fig. 5 and Table 8, we know that EP-CC does not always select the classifiers with the high generalization accuracy and the diversity among these base classifiers also affects the performance of the ensembles.

**Table 7**  
Fusion of base classifiers with different strategies.

Data set	EP-CC	EP-SV
abalone	55.91 ± 0.36	55.41 ± 0.28 •
balancescale	89.05 ± 0.65	88.29 ± 0.53 •
crx	86.61 ± 0.64	85.64 ± 0.31 •
derm	98.53 ± 0.55	97.85 ± 0.41 •
ecoli	88.01 ± 0.60	86.93 ± 0.89 •
german	74.77 ± 0.41	74.30 ± 0.35 •
heart	85.12 ± 0.69	84.76 ± 0.74 •
hepatitis	90.45 ± 1.49	88.56 ± 1.65 •
horse	93.48 ± 0.42	93.28 ± 0.69 •
iono	89.41 ± 0.75	88.63 ± 0.81 •
lung cancer	83.14 ± 2.06	81.52 ± 1.79 •
movement	79.19 ± 0.95	78.06 ± 1.05 •
mushroom	98.72 ± 0.07	98.67 ± 0.11
pima	78.09 ± 0.49	77.63 ± 0.51 •
satellite	87.08 ± 0.08	87.02 ± 0.11
segmentation	93.27 ± 0.13	93.23 ± 0.12
spam	90.85 ± 0.11	90.69 ± 0.11 •
wdbc	98.19 ± 0.26	97.94 ± 0.21 •
wdbc	81.98 ± 1.93	80.46 ± 1.67 •
yeast	74.89 ± 0.23	74.78 ± 0.19
Win-Tie-Loss		16-4-0

**Table 8**  
Average accuracies of classifiers selected by different pruning methods.

Data set	EP-CC	EPIC	MDM	CM
abalone	<b>55.33</b>	55.16	55.03	55.06
balancescale	<b>87.83</b>	87.67	87.73	87.64
crx	<b>85.57</b>	85.22	85.23	85.19
derm	<b>97.67</b>	97.09	97.33	97.28
ecoli	85.78	85.64	85.75	<b>85.81</b>
german	<b>73.64</b>	73.17	73.05	73.21
heart	<b>83.96</b>	82.51	83.26	83.44
hepatitis	<b>87.14</b>	85.69	85.67	85.79
horse	<b>92.17</b>	91.42	91.91	91.96
iono	<b>87.79</b>	86.97	87.12	87.18
lung cancer	<b>81.51</b>	81.10	80.68	80.77
movement	<b>72.37</b>	71.19	70.32	71.12
mushroom	98.31	<b>98.34</b>	98.28	98.33
pima	<b>77.05</b>	76.51	76.98	76.89
satellite	<b>86.73</b>	86.65	86.64	86.63
segmentation	92.75	92.55	<b>92.83</b>	92.82
spam	<b>90.46</b>	90.41	90.30	90.26
wdbc	<b>97.63</b>	97.28	97.49	97.42
wdbc	<b>79.51</b>	78.33	78.14	78.01
yeast	<b>74.75</b>	74.56	74.65	74.64



**Fig. 5.** Variation of classification accuracies with the ranking of the base classifiers weights.

We use  $KW$  to measure the diversity among the base classifiers [20].  $KW$  is a symmetrical measurement, computed as

$$KW = \frac{1}{nL^2} \sum_{i=1}^N \phi(x_i)(L - \phi(x_i)), \quad (15)$$

where  $\phi(x_i)$  denotes the number of classifiers which misclassify  $x_i$ .

From Eq. (15), we see that if  $KW$  is large, the diversity among the base classifiers is high. Table 9 presents the  $KW$  values of the base classifiers selected by EP-CC and the other pruning methods in the test set. It can be seen that EP-CC achieves the largest diversity on 14 data sets, whereas EPIC and MDM achieve the highest diversity on 3, respectively.

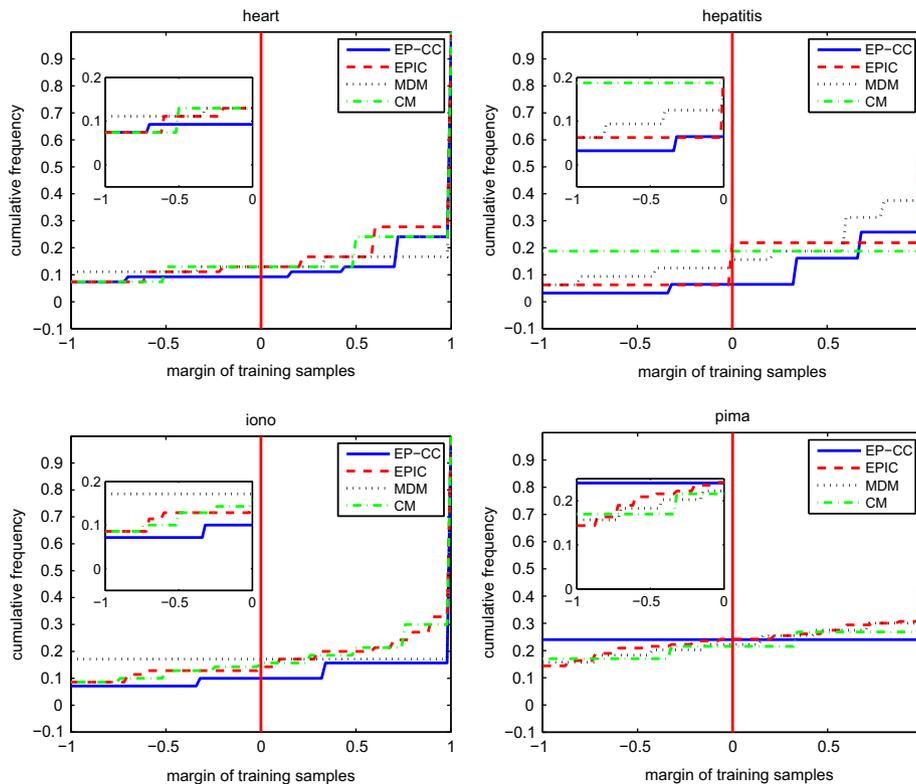
**Table 9**  
KW values computed with base classifiers selected by EP-CC, EPIC, MDM and CM.

Data set	EP-CC	EPIC	MDM	CM
abalone	<b>0.0356</b>	0.0220	0.0159	0.0185
balancescale	0.0142	0.0136	<b>0.0158</b>	0.0148
crx	<b>0.0073</b>	0.0036	0.0010	0.0006
derm	<b>0.0085</b>	0.0068	0.0030	0.0032
ecoli	<b>0.0212</b>	0.0141	0.0066	0.0144
german	0.0377	<b>0.0469</b>	0.0156	0.0113
heart	<b>0.0279</b>	0.0247	0.0051	0.0084
hepatitis	<b>0.0381</b>	0.0366	0.0187	0.0177
horse	<b>0.0211</b>	0.0200	0.0076	0.0082
iono	<b>0.0281</b>	0.0240	0.0084	0.0082
lung cancer	<b>0.0131</b>	0.0075	0.0121	0.0119
movement	0.0890	0.0892	<b>0.0968</b>	0.0910
mushroom	0.0052	0.0058	<b>0.0067</b>	0.0061
pima	0.0219	<b>0.0228</b>	0.0054	0.0062
satellite	<b>0.0126</b>	0.0119	0.0062	0.0089
segmentation	<b>0.0117</b>	0.0104	0.0015	0.0025
spam	<b>0.0055</b>	0.0040	0.0028	0.0050
wdbc	<b>0.0049</b>	0.0034	0.0007	0.0013
wdbc	<b>0.0281</b>	0.0248	0.0173	0.0219
yeast	0.0027	<b>0.0040</b>	0.0007	0.0013

Margin distribution is deemed as an important factor for the generalization performance of ensemble learning. In [36,43], the relationship between the generalization performance and the margin distribution was derived. It indicates that good margin distribution results in a low generalization error. A good margin distribution refers to the fraction of the samples with a small margin is small and most samples have large margins. The detailed information can be obtained in [36,43]. In fact, the diversity and the margin are closely related. In 2006, Tang et al. proved that if the average classification accuracy was set as a constant and the maximum diversity was achievable, maximizing the diversity among the base classifiers was equivalent to maximizing the minimum margin of the ensemble [38]. We now identify if compared to the other pruning methods, EP-CC improves the margin distribution.

Fig. 6 presents the margin distribution of the ensembles generated by EP-CC, EPIC, MDM, and CM in the test set, where the  $x$ -axis represents the value of the margin and the  $y$ -axis represents the fraction of the samples whose margin is not less than the corresponding margin. The small plots inside each graph are used to clearly show the margin distribution in the interval  $[-1, 0]$ . We can see that, compared with the other pruning methods, EP-CC improves the margin distribution, which explains why EP-CC achieves a better classification performance than the other techniques.

In the above experiments, the sub-ensemble with the best performance in the pruning set is selected as the final system. Then, how about their classification performances with the fixed ratios of the base classifiers? Some experiments were conducted to answer this question. In these experiments, the base classifiers are ordered according to the pruning techniques and the ensembles of the first 20% of the original base classifiers are evaluated, respectively. The generalization accuracies and the standard deviations are shown in Table 10. It can be seen that EP-CC performs significantly better than EPIC on 13 data sets, and the difference is not significant on the remained 7 sets. Compared with MDM, the statistically significant difference is favorable on 17 data sets,



**Fig. 6.** Margin cumulative frequency based on EP-CC, EPIC, MDM and CM.

**Table 10**

Classification performances of different pruning methods for ensembles of 20% classifiers.

Data set	EP-CC(20%)	EPIC(20%)	MDM(20%)	CM(20%)
abalone	55.42 ± 0.31	54.99 ± 0.29	54.76 ± 0.23	54.75 ± 0.19
balancescale	87.87 ± 0.64	87.85 ± 0.17	87.82 ± 0.50	87.79 ± 0.37
crx	86.03 ± 0.29	85.49 ± 0.15	85.63 ± 0.14	85.64 ± 0.12
derm	97.65 ± 0.45	97.28 ± 0.32	97.41 ± 0.26	97.39 ± 0.34
ecoli	86.31 ± 0.49	85.59 ± 0.40	85.29 ± 1.00	84.51 ± 0.92
german	73.93 ± 0.55	73.66 ± 0.59	73.61 ± 0.46	73.56 ± 0.43
heart	83.86 ± 0.65	83.70 ± 0.52	83.61 ± 0.60	83.54 ± 0.72
hepatitis	88.29 ± 1.26	87.01 ± 1.34	87.07 ± 1.15	86.63 ± 1.62
horse	93.03 ± 0.43	92.29 ± 0.53	92.31 ± 0.59	92.35 ± 0.65
iono	88.41 ± 0.71	87.62 ± 0.38	87.47 ± 0.54	87.41 ± 0.59
lung cancer	80.39 ± 2.02	80.37 ± 1.92	79.33 ± 1.45	79.59 ± 1.53
movement	77.35 ± 1.40	77.05 ± 1.43	76.91 ± 1.07	77.02 ± 1.21
mushroom	98.56 ± 0.15	98.48 ± 0.16	98.64 ± 0.09	98.35 ± 0.16
pima	77.39 ± 0.37	76.81 ± 0.38	76.93 ± 0.27	77.01 ± 0.31
satellite	86.88 ± 0.09	86.83 ± 0.11	86.76 ± 0.11	86.81 ± 0.12
segmentation	93.07 ± 0.14	93.01 ± 0.21	92.66 ± 0.18	92.69 ± 0.23
spam	90.53 ± 0.11	90.37 ± 0.08	90.15 ± 0.17	90.26 ± 0.12
wdbc	97.81 ± 0.33	97.55 ± 0.22	97.42 ± 0.25	97.45 ± 0.29
wpbc	80.19 ± 1.73	79.09 ± 1.72	78.73 ± 1.26	78.41 ± 1.67
yeast	74.62 ± 0.17	74.41 ± 0.26	74.26 ± 0.21	74.13 ± 0.21
Win-Tie-Loss		13-7-0	17-2-1	15-5-0

unfavorable in 1, and is not significant in 2. Meanwhile, EP-CC also performs better than CM in most of the data sets.

## 7. Conclusions and future work

In this work, we explore the role of the classification confidence in ensemble learning. A generalized definition of the ensemble margin is proposed based on the classification confidence and the weights of the base classifiers are learned through optimizing a margin induced loss function. Then, we try several strategies to utilize the weights and the classification confidences. Some new ensemble pruning and fusion strategies are developed. Extensive experiments are conducted to test the proposed techniques. Some conclusions can be drawn from the study.

(1) The classification confidence should be used in learning the weights of the base classifiers and weighted voting for improving the classification performance.

(2) The proposed weighted voting technique is better than simple voting if all the base classifiers are included in the final fusion.

(3) Pruning via the ordered aggregation can improve the performance of the weighted voting technique further. Moreover, it is better to combine the base classifiers selected by EP-CC with the proposed weighted voting strategy than to combine them with simple voting.

In this work, although the good generalization performance is obtained by considering the classification confidence in ensemble optimization, there are still some questions to be answered. Does there exist relationship between the generalization performance of the voting system and the margin based on the classification confidence? How do we design an appropriate criterion to combine the heterogeneous base classifiers if they are derived with different learning algorithms. We will work on these problems in the future.

## Conflict of interest

The authors declare that there is no conflicts of interest to this work.

## Acknowledgments

This work is supported by the National Program on Key Basic Research Project under Grant 2013CB329304, National Natural Science Foundation of China under Grants 61222210, 60873140, 61170107, 61073125, 61071179, 60963006, and 11078010, National Science Fund for Distinguished Young Scholars under Grant 50925625 and the Program for New Century Excellent Talents in University (No. NCET-12-0399, NCET-08-0155, and NCET-08-0156).

## References

- [1] B. Bakker, T. Heskes, Clustering ensembles of neural network models, *Neural Netw.* 16 (2) (2003) 261–269.
- [2] P.L. Bartlett, For valid generalization, the size of the weights is more important than the size of the network, in: *Advances in Neural Information Processing Systems*, vol. 9, 1997.
- [3] J.A. Benediktsson, J.R. Sveinsson, O.K. Ersoy, P.H. Swain, Parallel consensual neural networks, *IEEE Trans. Neural Netw.* 8 (1) (1997) 54–64.
- [4] C. Blake, E. Keogh, C.J. Merz, UCI Repository of Machine Learning Databases, Depart. Inf. Comput. Sci., University of California, Irvine, CA (Online), Available: (<http://www.ics.uci.edu/mllearn/MLRepository.html>), 1998.
- [5] L. Breiman, Bagging predictors, *Mach. Learn.* 24 (2) (1996) 123–140.
- [6] G. Brown, J. Wyatt, R. Harris, X. Yao, Diversity creation methods: a survey and categorisation, *Inf. Fusion* 6 (2005) 5–20.
- [7] R. Caruana, A. Niculescu-Mizil, G. Crew, A. Ksikes, Ensemble selection from libraries of models, in: *Proceedings of the International Conference on Machine Learning (ICML)*, 2004.
- [8] H. Chen, P. Tino, X. Yao, Predictive ensemble pruning by expectation propagation, *IEEE Trans. Knowl. Data Eng.* 21 (7) (2009) 999–1013.
- [9] V. Chvatál, *Linear Programming*, W. H. Freeman, New York, 1983.
- [10] C. Domingo, O. Watanabe, MadaBoost: a Modification of AdaBoost, in: *Proceedings of Annual Conference on Computational Learning Theory*, 2000, pp. 180–189.
- [11] B. Efron, R. Tibshirani, *An Introduction to the Bootstrap*, Chapman and Hall, New York, 1993.
- [12] Y. Freund, R.E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, *J. Comput. Syst. Sci.* 55 (1) (1997) 119–139.
- [13] G. Fumera, F. Roli, A theoretical and experiment analysis of linear combiners for multiple classifier systems, *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (6) (2005) 942–956.
- [14] T.V. Gestel, J.A.K. Suykens, B. Baesens, et al., Benchmarking least squares support vector machine classifiers, *Mach. Learn.* 54 (1) (2004) 5–32.
- [15] A.J. Grove, D. Schuurmans, Boosting in the limit: maximizing the margin of learned ensembles, in: *Proceedings of the 15th National Conference on Artificial Intelligence*, 1998.
- [16] D. Hernández-Lobato, J.M. Hernández-Lobato, R. Ruiz-Torrubiano, Á. Valle, Pruning adaptive boosting ensembles by means of a genetic algorithm, in: E. Corchado, H. Yin, V.J. Botti, C. Fyfe (Eds.), *Proceedings of 7th International Conference on Intelligent Data Engineering and Automated Learning*, 2006, pp. 322–329.
- [17] J. Kittler, M. Hatef, R.P.W. Duin, J. Matas, On combining classifiers, *IEEE Trans. Pattern Anal. Mach. Intell.* 20 (3) (1998) 226–239.
- [18] H.C. Kim, S.N. Pang, H.M. Je, D. Kim, S.Y. Bang, Constructing support vector machine ensemble, *Pattern Recognit.* 36 (2003) 2757–2767.
- [19] A. Krogh, J. Vedelsby, Neural network ensembles, cross validation, and active learning, in: G. Tesauro, D.S. Touretzky, T.K. Lee (Eds.), *Advances in Neural Information Processing Systems*, vol. 7, MIT Press, Cambridge, MA, 1995, pp. 231–238.
- [20] L. Kuncheva, C. Whitaker, Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy, *Mach. Learn.* 51 (2003) 181–207.
- [21] J. Liu, S.W. Ji, J.P. Ye, SLEP: Sparse Learning with Efficient Projections, Arizona State University, 2010, (<http://www.public.asu.edu/~jye02/Software/SLEP/>).
- [22] H. Lodhi, G.J. Karakoulas, J. Shawe-Taylor, Boosting the margin distribution, in: *Proceedings of the International Conference on Intelligent Data Engineering Automated Learning/Data Mining, Financial Engineering, and Intelligent Agents*, London, UK, 2000, pp. 54–59.
- [23] Z.Y. Lu, X.D. Wu, X.Q. Zhu, J. Bongard, Ensemble pruning via individual contribution ordering, in: *Proceedings of the 16th ACM SIGKDD, KDD*, 2010, pp. 871–880.
- [24] D.D. Margineantu, T.G. Dietterich, Pruning adaptive boosting, in: *Proceedings of the 14th International Conference on Machine Learning*, 1997, pp. 211–218.
- [25] G. Martínez-Muñoz, D. Hernández-Lobato, A. Suárez, An analysis of ensemble pruning techniques based on ordered aggregation, *IEEE Trans. Pattern Anal. Mach. Intell.* 31 (2) (2009) 245–259.
- [26] G. Martínez-Muñoz, A. Suárez, Aggregation ordering in bagging, in: *Proceedings of the International Conference on Artificial Intelligence and Applications*, 2004, pp. 258–263.
- [27] G. Martínez-Muñoz, A. Suárez, Pruning in ordered bagging ensembles, in: *Proceedings of the 23th International Conference on Machine Learning*, 2006, pp. 609–616.

- [28] G. Martínez-Muñoz, A. Suárez, Using boosting to prune bagging ensembles, *Pattern Recognit. Lett.* 28 (1) (2007) 156–165.
- [29] I. Partalas, G. Tsoumakas, I. Vlahavas, An ensemble uncertainty aware measure for directed hill climbing ensemble pruning, *Mach. Learn.* 81 (2010) 257–282.
- [30] J.R. Quinlan, Bagging, boosting, and C4.5, in: *Proceedings of the 13th National Conference on Artificial Intelligence*, 1996, pp. 725–730.
- [31] J.J. Rodríguez, L.I. Kuncheva, Rotation forest: a new classifier ensemble method, *IEEE Trans. Pattern Anal. Mach. Intell.* 28 (10) (2006) 1619–1630.
- [32] R.E. Schapire, Y. Singer, Improved boosting algorithms using confidence-rated predictions, *Mach. Learn.* 37 (1999) 297–336.
- [33] C.H. Shen, H.X. Li, Boosting through optimization of margin distributions, *IEEE Trans. Neural Netw.* 4 (21) (2010) 659–666.
- [34] J. Shawe-Taylor, P.L. Bartlett, R.C. Williamson, M. Anthony, A framework for structural risk minimisation, in: *Proceedings of the 9th Annual Conference on Computational Learning Theory*, 1996, pp. 68–76.
- [35] J. Shawe-Taylor, N. Cristianini, Robust bounds on generalization from the margin distribution, in: *4th European Conference on Computational Learning Theory*, 1999.
- [36] R.E. Schapire, Y. Freund, P. Bartlett, W.S. Lee, Boosting the margin: a new explanation for the effectiveness of voting methods, *Machine Learning: Proceedings of the 14th International Conference*, 1997.
- [37] J. Suykens, J. Vandewalle, Least squares support vector machine classifiers, *Neural Process. Lett.* 9 (1999) 293–300.
- [38] E.K. Tang, P.N. Suganthan, X. Yao, An analysis of diversity measures, *Mach. Learn.* 65 (2006) 247–271.
- [39] K.M. Ting, I.H. Witten, Issues in stacked generalization, *J. Artif. Intell. Res.* 10 (1999) 271–289.
- [40] G. Tsoumakas, I. Partalas, I. Vlahavas, An ensemble pruning primer, *Appl. Superv. Unsuperv. Ensemble Methods* 245 (2009) 1–13.
- [41] G. Tsoumakas, L. Angelis, I. Vlahavas, Selective fusion of heterogeneous classifiers, *Intell. Data Anal.* 9 (2005) 511–525.
- [42] G. Wahba, Support vector machines, reproducing kernel Hilbert spaces and the randomized GACV, in: *Advances in Kernel Methods Support Vector Learning*, MIT Press, MA, 1999, pp. 69–88.
- [43] L.W. Wang, M. Sugiyama, C. Yang, Z.H. Zhou, J.F. Feng, On the margin explanation of boosting algorithms, in: *Proceedings of COLT*, 2008, pp. 479–490.
- [44] D.H. Wolpert, Stacked generalization, *Neural Netw.* 5 (1992) 241–259.
- [45] L. Xu, A. Krzyzak, C.Y. Suen, Methods of combining multiple classifiers and their applications to handwriting recognition, *IEEE Trans. Syst. Man Cybern.* 3 (22) (1992) 418–435.
- [46] Z.X. Xie, Y. Xu, Q.H. Hu, P.F. Zhu, Margin distribution based bagging pruning, *Neurocomputing* 85 (2012) 11–19.
- [47] Y. Zhang, S. Burer, W.N. Street, Ensemble pruning via semi-definite programming, *J. Mach. Learn. Res.* 7 (2006) 1315–1338.
- [48] Z.H. Zhou, *Ensemble Methods: Foundations and Algorithms*, Chapman Hall/CRC, Boca Raton, FL, 2012.
- [49] Z.H. Zhou, J.X. Wu, W. Tang, Ensembling neural networks: many could be better than all, *Artif. Intell.* 137 (1–2) (2002) 239–263.
- [50] Z.H. Zhou, Y. Yu, Ensembling local learners through multimodal perturbation, *IEEE Trans. Syst. Man Cybern., Part B* 35 (2005) 725–735.
- [51] L. Zhang, W.D. Zhou, Sparse ensembles using weighted combination methods based on linear programming, *Pattern Recognit.* 44 (2011) 97–106.

**Leijun Li** got his B.Sc., M.Sc. from Hebei Normal University in 2007 and 2010, respectively. Now he is a Ph.D. candidate with School of Computer Science and Technology, Harbin Institute of Technology. His research interests include ensemble learning, margin theory and rough sets.

**Qinghua Hu** received B.Sc., M.E. and Ph.D. degrees from Harbin Institute of Technology, Harbin, China, in 1999, 2002 and 2008, respectively. He started working with Harbin Institute of Technology from 2006, and was a post doctoral fellow with the Hong Kong Polytechnic University from 2009 to 2011. Now he is a full professor with Tianjin University. His research interests are focused on intelligent modeling, data mining, knowledge discovery for classification and regression. He is a PC co-chair of RSCTC2010 and serves as a referee for a great number of journals and conferences. He has published more than 90 journal and conference papers in the areas of pattern recognition and fault diagnosis.

**Xiangqian Wu** received his B.Sc., M.E. and Ph.D. degrees from Harbin Institute of Technology, Harbin, China, in 1997, 1999 and 2004, respectively. Now he is a full professor with School of Computer Science and Technology, Harbin Institute of Technology. He once visited The Hong Kong Polytechnic University and Michigan State University. His main interests are focused on biometrics, image processing and pattern recognition. He has published more than 50 peer reviewed papers in these domains.

**Daren Yu** received the M.Sc. and D.Sc. degrees from Harbin Institute of Technology, Harbin, China, in 1988 and 1996, respectively. Since 1988, he has been working at the School of Energy Science and Engineering, Harbin Institute of Technology. His main research interests are in modeling, simulation, and control of power systems. He has published more than one hundred conference and journal papers on power control and fault diagnosis.