Contents lists available at ScienceDirect

# Neurocomputing

journal homepage: www.elsevier.com/locate/neucom

# A two-stage ensemble method for the detection of class-label noise

## Maryam Sabzevari\*, Gonzalo Martínez-Muñoz, Alberto Suárez

Computer Science Department, Escuela Politécnica Superior, Universidad Autónoma de Madrid, C/ Francisco Tomás y Valiente, 11, Madrid 28049, Spain

#### ARTICLE INFO

Article history: Received 3 August 2017 Revised 3 November 2017 Accepted 5 November 2017 Available online 10 November 2017

Communicated by Prof. Dianhui Wang

Keywords: Noise detection Ensemble learning Subsampling Robust classification Random forest

### 1. Introduction

The presence of noise, which is often unavoidable in real-world settings, is an important nuisance factor that needs to be taken into account in the design of learning algorithms [5,20]. Two types of noise can be found in the data used for automatic induction: class-label and feature noise [20]. Except when the latter is very strong, the presence of erroneous class labels is generally more harmful for learning and generalization [20]. According to its statistical properties label noise falls into three categories [5]. In the Noisy Completely at Random (NCAR) model the probability of a mislabeled instance is uniform in feature and class spaces. In the Noisy at Random (NAR) model, the probability of mislabeled instances is independent of the feature values, but depends on the class. Finally, in the Noisy Not at Random (NNAR) models there are dependencies between the class-label noise, and both feature and class values. In this work, we assume that the noise is completely at random. Nevertheless, the method proposed can be readily adapted to take into account other types of class-label noise.

One way of alleviating the effects of noisy data is to incorporate some regularization mechanism into the design of the learning algorithm [2,6,18]. Another approach is to carry out a preprocessing step in which noisy instances are identified. These instances are then either discarded (filtering), or their class label corrected (cleaning) so that they do not interfere in the learning process [5].

https://doi.org/10.1016/j.neucom.2017.11.012 0925-2312/© 2017 Elsevier B.V. All rights reserved.

### ABSTRACT

The properties of bootstrap ensembles, such as bagging or random forest, are utilized to detect and handle label noise in classification problems. The first observation is that subsampling is a regularization mechanism that can be used to render bootstrap ensembles more robust to this type of noise. Furthermore, appropriate values of the sampling rate can be estimated using out-of-bag data. A second observation is that the ensemble classifiers tend to make more errors in incorrectly labeled instances. Thus, instances for which a sufficiently large fraction of ensemble predictors err are marked as noisy. Suitable values of this threshold, which are problem dependent, are determined by cross-validation using a wrapper method. Instances identified as noisy can then be either filtered (i.e. discarded for training), or cleaned by correcting their class labels. Finally, an ensemble is built afresh on these cleansed training data. Extensive experiments in classification problems from different areas of application show that this procedure is effective to build accurate ensembles, even in the presence of high levels of class-label noise.

© 2017 Elsevier B.V. All rights reserved.

In this work, we combine both strategies using randomized ensembles. Our goal is to take advantage of the robustness of bootstrap ensembles, such as bagging and random forest, to handle noise. Specifically, we will use subsampling as a regularization mechanism, which allows one to build ensembles that are robust to class-label noise [13,14]. Another observation is that the individual classifiers in the ensemble tend to make more prediction errors on incorrectly labeled instances. Therefore, the fraction of incorrect ensemble predictions can be used as an indicator to detect noisy instances. Standard ensemble-based approaches to this problem remove instances for which more than half of the votes are incorrect (majority filtering) or where all votes are incorrect (consensus filtering). However, these choices need not be optimal. In this work, we propose to use a wrapper method to determine a near-optimal value for the percentage of incorrect votes necessary to identify a noisy instance.

The structure of the work is as follows: In Section 2, we present a review of previous work on learning from incorrectly labeled data. Particular emphasis is given to the use of subsampling to improve the robustness to class-label noise of bootstrap ensembles. The conclusions of this analysis are applied in Section 3 to the design of a method for noise detection. In Section 4, we present the results of an empirical evaluation of the proposed procedure and a comparison with other state-of-the-art noise detection methods. Finally, the conclusions of this work are summarized in Section 5.







<sup>\*</sup> Corresponding author. E-mail address: maryam.sabzevari@uam.es (M. Sabzevari).

### 2. Previous work

The problem of induction from noisy data has been extensively addressed in the area of ensemble learning [2,5,6,13,14,18]. Furthermore, as illustrated by recent implementations [8], these methods can be scaled to deal with class-label noise in large problems. Early on in the literature on ensembles, it was noted that the accuracy improvements of an ensemble with respect to a single learner are generally smaller when the training data are contaminated with class-label noise [1]. However, ensembles are generally more robust to this type of noise than single learners. In [3], an ensemble of three classifiers (a univariate decision tree, a k-nearest neighbor and a linear machine) is used to identify noisy instances. The noise detection protocol in this method is as follows: The training set is partitioned into 4-folds. Then, an ensemble is trained using three of these folds. The ensemble is then used to identify noisy instances in the fold not used for training. In this study, two filtering strategies are considered: majority filtering and consensus filtering. In majority filtering, a particular instance of the left-out fold is identified as noisy when the predictions of more than half of the learners are erroneous. The consensus filtering rule is more conservative. An instance is categorized as noisy when all the members of ensemble misclassify it. This process is repeated for each fold to identify noisy instances in all 4 folds. In this study, majority filtering exhibits better overall performance than consensus filtering. In [9] a similar approach is used, albeit with an ensemble of 25 classifiers of different types. The increased diversity of this heterogeneous ensemble reduces the risk of false positives in the noise detection process. In addition, different intermediate strategies between majority and consensus filtering are explored. In the problems investigated, the optimal threshold of erroneous ensemble predictions used to identify an instance as noisy was close to consensus filtering.

In [17], an ensemble of Top-down Induction of Logical Decision Trees (Tilde) is used to filter the training data. In this study either cross-validation or bootstrap sampling are used to generate the training sets on which the base learners are built. Majority and consensus filtering are then applied to identify noisy instances. The best results are obtained with majority filtering. In addition, the use of boosting in noise detection problems is explored. The idea is to filter out instances whose weights are above a threshold after a predefined number of iterations of the boosting algorithm. In the problems investigated, this strategy does not perform well.

In [21], a distributed method for large datasets is presented. In this method, the original training set is partitioned into small subsets. A classifier is induced from each of the these subsets. These classifiers are then used to classify the instances in the complete training set. The local error of an instance is defined as the fraction of classifiers whose training data included that particular instance and misclassified it. The global error for an instance is computed using the classifiers that did not include that instance in their training sets. Majority and consensus filtering are used to identify noisy instances. A necessary condition for an instance to be identified as noisy is that it be misclassified by the base learners whose training sets induced it. The justification is that a classifier has usually higher accuracy on instances that are in its training set. In this work, majority filtering is found to yield better results than consensus filtering.

In [19], a noise detection strategy that combines ideas from bagging and boosting is investigated. In this method, all training instances are initially assigned equal weights. Then, bootstrap samples from the original training data are used to build different classifiers. For each instance a noise count is computed as the number of base learners that have misclassified the instance. The noise count is used in a boosting-like iterative process in which the weight of the instances with higher noise counts is decreased. This process is iterated for a predefined number of rounds. Finally, instances whose noise count values are higher than a threshold are labelled as noisy. Ensemble methods based on ranking have been explored in [15]. In this work, instances are ranked according to the number of base learners that make incorrect predictions. In the medical dataset analyzed, a domain expert singles out the instances with highest ranks for further analysis. Expert knowledge is then used to determine the type of anomaly of the instances with the highest ranks (labeling error, outlier, complex medical case, etc.).

In most of the noise detection methods based on ensembles, majority or consensus filtering are used. Majority filtering was found to be better than consensus filtering in most studies in which homogeneous ensembles were used [3,16,17,21]. By contrast, in small heterogeneous ensembles, agreement rates close to consensus filtering seem to be more effective [9,16]. A qualitative explanation of this observation can be given: In homogeneous ensembles the base learners are of the same type. Diversity is commonly obtained using randomization strategies, such as bootstrap sampling (e.g. in bagging and random forest). Thus, the decision boundaries of the individual classifiers are similar to each other. In fact, the decision boundary is a refinement of the decision border of the individual base learners [12]. Therefore, noisy examples close to this decision boundary can be detected only if majority voting is used. By contrast, in heterogeneous ensembles, the decision borders are more varied because the individual classifiers are of different types. For this reason the dispersion of class labels assignments may simply reflect this variability. In consequence, higher agreement rates should be used to mark an instance as noise. In summary, the optimal agreement rates for noise detection depends on the classification problem and type of ensemble. However, as far as we are aware, the influence of the agreement rate on the effectiveness of the noise detection method has not been systematically evaluated hitherto.

#### 3. A wrapper method for class-label noise detection

In this work, ensembles are used to detect and handle noise in the class labels of the training instances. We focus on randomized ensembles, such as bagging and random forest, in which the individual classifiers are trained on bootstrap samples of the original data. The design of the algorithm leverages two observations: The first one is that reducing the sampling rate in the bootstrapping step can make the ensemble more robust to class-label noise [13,14]. A second observation is that the fraction of erroneous predictions by the ensemble classifiers necessary to detect noisy instances depends on the classification task at hand. Therefore, the value of the threshold  $\theta$  used to mark an instance as noisy should be estimated for each problem during the training phase. The method proposed proceeds in two stages: First, the optimal sampling rate for the bootstrapping step in the construction of the ensemble is determined using out-of-bag data. Then, the optimal threshold for disagreement between the ensemble predictions and the actual class of the instance is determined by means of a wrapper method. Finally, instances labeled as noise are either removed (filtering) or relabeled (cleaning) and an ensemble is built afresh from the cleansed training data.

The pseudocode of the method is detailed in Algorithm 1. In the first stage of the algorithm (steps 1–8 in Algorithm 1), ensembles are built using different values of the sampling rate: 0.1, 0.2, 0.4, 0.6, 0.8, 1.0, and 1.2. From these,  $H_T^*$ , the ensemble that is expected to generalize best is selected and kept for the next phase. In our implementation, out-of-bag instances are used to estimate the generalization error and carry out this selection. In this way, an unbiased selection is achieved, independently of the amount of noise present in the dataset. According to empirical evidence

Algorithm 1: Noise detection using bootstrap ensembles.				
<b>Input:</b> $\mathcal{D}_{train} = \{(\mathbf{x}_n, y_n)\}_{n=1}^{N_{train}}$	% Training set			
T	% Ensemble size			
I Wirapper learner	% Elisellible Size			
cleanse type	% whappen learner method			
Cieunse_type				
<b>Output</b> : $\mathcal{D}_{cleansed}$	% Cleansed set			
1 $min\_error \leftarrow \infty$ % determine	e optimal sampling rate			
2 foreach sampling_rate in [0.1	, 0.2, 0.4, 0.6, 0.8, 1.0, 1.2] <b>do</b>			
3   $H_T \leftarrow bootstrap\_ensemble$	$(\mathcal{D}_{train}, sampling\_rate, T)$			
4 error $\leftarrow$ estimate_error( $H_1$	$\mathcal{D}_{train}$ )			
s if error < min_error then				
$6 \qquad min\_error \leftarrow error$				
$sampling_rate^* \leftarrow sampling_rate$				
8 $H_T^* \leftarrow H_T$				
9 min_error $\leftarrow \infty$ % determin	e optimal disagreement rate			
10 foreach $\theta$ in [0.5, 0.6, 0.7, 0.8	10 foreach $\theta$ in [0.5, 0.6, 0.7, 0.8, 0.9, 1.0] do			
$\mathcal{D}_{cleansed} \leftarrow cleanse\_with\_oob(\mathcal{D}_{train}, cleanse\_type, H_T^*, \theta)$				
12 <i>error</i> $\leftarrow cv\_error(wrapper\_learner, D_{cleansed}, K = 3)$				
<b>if</b> error < min_error <b>then</b>				
$14 \qquad \qquad$				
<b>15</b> $\mathcal{D}_{cleansed} \leftarrow cleanse_with_oob(2)$	$D_{train}, cleanse_type, H_T^*, \theta^*)$			

[13,14], the higher the level of class-label noise, the lower the value of the subsampling rate that is selected. Although the possibility that the presence of noisy instances lead to an incorrect selection of the best ensemble cannot be ruled out, we have not observed this effect in the experiments carried out. This is probably due to the fact that random forests are fairly robust to class-label noise, even without cleansing.

In the second stage (steps 9–14 in Algorithm 1), the predictions of the base classifiers of the selected ensemble,  $H_T^*$ , are used to identify noisy instances. Different values of the threshold  $\theta$  for the disagreement rate are considered; namely, 0.5, 0.6, 0.7, 0.8, 0.9, and 1.0. For a given value of  $\theta$ , instances for which the fraction of incorrect predictions given by the classifiers in the selected ensemble is above this threshold are marked as noise (step 11). Since the data being cleaned and the data used to train  $H_T^*$  are the same, only the predictions of base classifiers in  $H_T^*$  whose training set does not include that particular instance are used (i.e. out-of-bag instances). Cleansing consists in either correcting the label of noisy instances (cleaning), or eliminating them (filtering). Then, we use a wrapper method and compute estimates of the generalization error of a learner built on the cleansed training data. In our implementation, this error is estimated using K-fold cross-validation. The optimal value of the disagreement threshold,  $\theta^*$ , is the one that minimizes this estimate of the generalization error. The details of this selection are as follows: In each of the K iterations of the cross validation procedure, one of the K folds is set apart for validation. Then, a classifier is trained on remaining K - 1 folds cleansed using the corresponding value of  $\theta$ . The classifier is then evaluated on the left-out fold, which is not cleansed. Finally, the cross validation error is calculated by averaging the errors of the validation folds in each of the *K* iterations. The value selected,  $\theta^*$ , is the threshold that minimizes this cross-validation estimate of the generalization error. This optimal threshold value,  $\theta^*$ , and the ensemble  $H_T^*$  are then used to clean or filter the training data. Note that the type of classifier used in the wrapper step is a parameter of the algorithm. In general, we think it is preferable to use the same classifier as the one that is eventually trained with the cleansed data. In the following section we carry out an extensive empirical

ladie I	Та	ble	1	
---------	----	-----	---	--

Characteristics of the classification problems used in the empirical evaluation.

Dataset	Training	Test	Attrib.
Australian	460	230	14
Blood transfusion	499	499	5
Boston	337	169	14
Breast	466	233	9
Chess	2130	1065	37
Crx	460	230	15
Diabetes	512	256	8
German	667	333	20
Heart	178	92	13
Horse-colic	246	122	21
Ionosphere	234	117	34
Liver	230	115	6
Parkinsons	130	65	22
Ringnorm	300	2000	20
Sonar	491	208	60
Spambase	3067	1534	58
Threenorm	300	2000	20
Tic-Tac-Toe	639	319	9
Twonorm	300	2000	20

analysis of the accuracy and resilience to noise of ensembles built in this manner.

#### 4. Empirical evaluation

To assess the effectiveness of the proposed noise detection procedure extensive experiments have been carried out in 19 binary classification problems from the UCI repository [11]. The characteristics of the datasets are summarized in Table 1. The implementation makes use of the R *randomForest* package [10] and the R *adabag* package for AdaBoost [7]. In both packages the default settings are used. Specifically, for random forests, the number of variables considered for splitting at the inner nodes of the random trees is the square root of the total number of attributes in the problem. The minimum size of a terminal node is set to 1. Trees in the forest are grown to their maximum possible size. For AdaBoost, weighted resampling is used. The coefficient that controls the weight update is  $\alpha = 1/2\ln((1 - err)/err)$ .

In all the classification tasks, with the exception of *Ringnorm*, *Threenorm* and *Twonorm*, which are synthetic problems, the labeled instances are randomly assigned to the training and test sets. The sizes of these sets are 2/3 and 1/3 of the original set, respectively. For synthetic problems, 300 examples are used for training and 2000 for testing. In all cases, stratified sampling is used. The results reported are averages over 50 executions. In these executions, different random partitions of the data into training and test sets, or, in synthetic problems, independent realizations of the data are used. The following protocol is followed for each classification task and execution:

- 1. First, noise is injected in the data by randomly switching the class label of a random subset of the training instances. Different noise rates are considered: 0 (no noise is injected), 10, 20, 30 and 40%. This type of class-label noise is known as completely at random noise (NCAR) [5].
- 2. For each noise level, ensembles composed of 501 trees are trained using the following bootstrap sampling ratios: 10, 20, 40, 60, 80, 100(standard) and 120%. Then, the out-of-bag error is computed for each ensemble. The ensemble with the best out-of-bag accuracy  $(H_T^*)$  is kept for the next step. Both random forests and bagging ensembles of unpruned CART trees have been considered. According to the results of our experiments, they are both equally effective to identify noisy instances. Therefore, since random forest are generally more ac-





**Fig. 1.** Comparison of the average ranks of the different types of ensembles for different levels of class-label noise: (a) without injected noise; (b) 10%; (c) 20%; (d) 30%; and (e) 40% noise. Horizontal lines connect methods whose average ranks are not significantly different according to a Wilcoxon signed-ranks test (*p*-value < .05).

curate than bagging ensembles, only results for the former are reported.

- 3. For each instance in the training set, we compute a tally of votes (class label predictions). In this tally, only the predictions of those classifiers in  $H_T^*$  whose training sets do not include that particular instance are used.
- 4. Instances are tentatively marked as noisy if the percentage of incorrect predictions by the individual ensemble classifiers is above a specified threshold  $\theta$ . Noisy instances are either cleaned (i.e. their class labels are corrected by assigning the majority label in the out-of-bag predictions) or filtered (i.e. removed from the training set). The following values of the threshold  $\theta$  are tested: 0.5 (majority filtering), 0.6, 0.7, 0.8, 0.9 and 1.0 (consensus filtering). As described in the previ-

ous section, the optimum value  $\theta^*$  is selected by *K*-fold crossvalidation within the training set by means of a wrapper method. From the results of exploratory experiments with different values of *K*, reliable estimates are obtained with K = 3, which is the value used in the experiments. Random forest composed of 501 trees is used in this wrapper stage.

- 5. The training instances for which the fraction of incorrect outof-bag predictions is above  $\theta^*$  are definitively marked as noise. These noisy instances are then either corrected or removed from the training set.
- 6. Finally, a random forests of 501 random trees, is built on the cleansed training set. The labels *Fl\_rf* for random forest with filtering and *Cl\_rf* for random forest with cleaning will be used throughout the experiments.



# (a) Sampling rate



(b) Threshold for cleansing

Fig. 2. Optimal hyperparameters for random forest with filtering (Fl\_rf).

Table 2           Test errors of the different methods for different noise levels (I).					
	Noise				
Dataset	(in %)	Fl_rf	Cl_rf	Fl_maj_rf	
Australian	0	$133 \pm 19$	$135 \pm 19$	$138 \pm 20$	

Dataset	(in %)	Fl_rf	Cl_rf	Fl_maj_rf	Boosting	RF
Australian	0 10 20 30 40	$\frac{13.3 \pm 1.9}{13.6 \pm 1.9}$ $\frac{14.1 \pm 1.8}{16.3 \pm 2.9}$ $\frac{26.0 \pm 6.3}{2}$	$\begin{array}{c} 13.5\pm1.9\\ 13.7\pm1.9\\ \underline{14.0}\pm\underline{2.0}\\ \underline{16.1}\pm\underline{2.8}\\ \textbf{24.5}\pm\textbf{5.2}^* \end{array}$	$13.8 \pm 2.0 \\ 13.7 \pm 2.0 \\ 13.8 \pm 2.1 \\ 15.8 \pm 2.6 \\ 26.1 \pm 5.1 \\ 15.8 \pm 2.1 \\ 15.1 \pm 5.1 \\ 15.$	$\begin{array}{c} \textbf{13.1} \pm \textbf{1.9} \\ 18.2 \pm 2.2 \\ 23.8 \pm 3.1 \\ 33.6 \pm 3.6 \\ 41.9 \pm 3.2 \end{array}$	$\begin{array}{c} \textbf{13.1} \pm \textbf{1.9}^* \\ \textbf{13.5} \pm \textbf{1.9} \\ 15.3 \pm 2.2 \\ 21.7 \pm 3.0 \\ 34.3 \pm 3.5 \end{array}$
Blood transfusion	0 10 20 40	$\frac{\underline{21.9} \pm \underline{1.6}}{\underline{22.2} \pm \underline{1.5}}$ $\underline{\underline{23.3} \pm \underline{2.3}}{\underline{34.0} \pm \underline{4.6}}$	$\begin{array}{c} 21.8 \pm 1.9 \\ 22.0 \pm 1.7 \\ 23.2 \pm 2.4 \\ 32.6 \pm 4.9^* \end{array}$	$\begin{array}{c} 22.1 \pm 1.8 \\ 22.9 \pm 2.0 \\ 24.9 \pm 2.6 \\ 38.5 \pm 4.5 \end{array}$	$\begin{array}{c} 25.7 \pm 2.3 \\ 28.0 \pm 2.8 \\ 31.2 \pm 3.2 \\ 43.2 \pm 4.3 \end{array}$	$\begin{array}{c} 24.7 \pm 2.1 \\ 26.4 \pm 2.1 \\ 29.8 \pm 2.8 \\ 42.1 \pm 4.5 \end{array}$
Boston	0 10 20 30 40	$13.2 \pm 2.2 \\ 13.7 \pm 2.6 \\ \underline{15.5} \pm \underline{2.8} \\ \underline{19.2} \pm \underline{3.8} \\ \underline{26.8} \pm \underline{6.1} \\ \end{array}$	$13.5 \pm 2.2$ $13.8 \pm 2.6$ $15.5 \pm 2.4$ $18.9 \pm 3.7$ $26.1 \pm 6.5$	$14.7 \pm 2.7 \\ 14.6 \pm 2.2 \\ 15.2 \pm 2.8 \\ 20.2 \pm 4.0 \\ 28.8 \pm 6.5 \\ \end{array}$	$12.4 \pm 2.4 \\ 16.3 \pm 2.8 \\ 23.6 \pm 3.9 \\ 32.5 \pm 4.2 \\ 41.1 \pm 4.1$	$\frac{12.8 \pm 2.1}{13.7 \pm 2.6}$ 17.8 ± 2.7 25.1 ± 4.4 36.2 ± 5.1
Breast	0 10 20 30 40	$\begin{array}{c} \textbf{3.1} \pm \textbf{1.1}^* \\ \underline{3.6} \pm \underline{1.2} \\ \textbf{4.3} \pm \textbf{1.5} \\ \textbf{6.6} \pm \textbf{3.4}^* \\ \underline{15.5} \pm \underline{6.4} \end{array}$	$\begin{array}{c} \textbf{3.1} \pm \textbf{1.0}^{*} \\ \underline{3.6} \pm \underline{1.1} \\ \underline{4.4} \pm \underline{1.8} \\ \textbf{6.6} \pm \textbf{3.2}^{*} \\ \textbf{14.8} \pm \textbf{5.8} \end{array}$	$\frac{3.3 \pm 1.0}{3.5 \pm 1.2}$ 4.3 ± 1.4 $\frac{7.6 \pm 3.5}{18.4 \pm 4.8}$	$\begin{array}{c} 3.4 \pm 1.1 \\ 8.7 \pm 1.8 \\ 14.9 \pm 2.5 \\ 24.1 \pm 3.8 \\ 33.8 \pm 4.2 \end{array}$	$\begin{array}{c} \textbf{3.1} \pm \textbf{1.1}^* \\ \textbf{4.1} \pm \textbf{1.3} \\ \textbf{6.7} \pm \textbf{1.9} \\ \textbf{12.7} \pm \textbf{3.8} \\ \textbf{26.4} \pm \textbf{4.7} \end{array}$
Chess	0 10 20 30 40	$\begin{array}{c} 1.7 \pm 0.4 \\ \textbf{2.2} \pm \textbf{0.4}^* \\ \textbf{3.4} \pm \textbf{0.7}^* \\ \textbf{5.5} \pm \textbf{1.2} \\ \textbf{16.6} \pm \textbf{3.0}^* \end{array}$	$\begin{array}{c} 1.7 \pm 0.4 \\ \underline{2.3} \pm \underline{0.4} \\ \underline{3.6} \pm \underline{0.8} \\ 6.0 \pm 1.1 \\ \underline{17.7} \pm \underline{3.2} \end{array}$	$\begin{array}{c} 2.6 \pm 0.4 \\ 3.0 \pm 0.6 \\ 3.7 \pm 0.8 \\ \underline{5.7 \pm 1.0} \\ \underline{17.7 \pm 2.3} \end{array}$	$\begin{array}{c} \textbf{0.4} \pm \textbf{0.2}^{*} \\ 4.1 \pm 0.8 \\ 5.8 \pm 0.8 \\ 11.0 \pm 1.7 \\ 24.6 \pm 2.2 \end{array}$	$\begin{array}{c} \underline{1.6} \pm \underline{0.4} \\ \underline{2.3} \pm \underline{0.5} \\ 4.1 \pm 0.6 \\ 10.0 \pm 1.1 \\ 25.7 \pm 1.7 \end{array}$
Crx	0 10 20 30 40	$\frac{12.9 \pm 1.8}{13.7 \pm 1.8}$ $\frac{14.1 \pm 2.2}{16.9 \pm 3.3}$ $\frac{24.5 \pm 5.9}{100}$	$13.0 \pm 2.1$ $13.6 \pm 1.8$ $14.2 \pm 2.3$ $16.2 \pm 2.9$ $24.3 \pm 5.9$	$13.2 \pm 2.1$ $13.6 \pm 2.0$ $14.0 \pm 2.2$ $17.0 \pm 3.0$ $25.8 \pm 5.0$	$13.8 \pm 1.8 \\ 18.8 \pm 2.5 \\ 25.3 \pm 3.1 \\ 33.4 \pm 3.8 \\ 41.5 \pm 3.4$	$\begin{array}{c} \textbf{12.6} \pm \textbf{1.9} \\ 14.0 \pm 1.9 \\ 16.2 \pm 2.1 \\ 22.3 \pm 3.3 \\ 33.8 \pm 3.7 \end{array}$
Diabetes	0 10 20 30 40	$\frac{24.1 \pm 2.0}{24.6 \pm 2.3}$ $\frac{25.4 \pm 2.5}{27.3 \pm 2.6}$ $32.5 \pm 4.9$	$24.3 \pm 1.8$ $24.5 \pm 2.3$ $25.4 \pm 2.4$ $26.7 \pm 2.7^*$ $32.5 \pm 3.8$	$24.3 \pm 1.9$ $24.3 \pm 2.2$ $25.3 \pm 2.4$ $\underline{27.3} \pm \underline{2.7}$ $\underline{33.4} \pm \underline{4.7}$	$\begin{array}{c} 27.5 \pm 2.1 \\ 30.1 \pm 2.9 \\ 34.4 \pm 3.0 \\ 39.6 \pm 3.3 \\ 44.4 \pm 4.0 \end{array}$	$\begin{array}{c} \textbf{24.0} \pm \textbf{1.9} \\ \textbf{25.4} \pm \textbf{2.2} \\ \textbf{27.8} \pm \textbf{2.5} \\ \textbf{31.7} \pm \textbf{3.0} \\ \textbf{39.3} \pm \textbf{4.5} \end{array}$

\* *p*-value < 0.05.

7. The generalization error of the resulting ensembles is estimated on the unperturbed test set.

As a benchmark, a second cleansed dataset is obtained using majority filtering following the proposals of previous studies [3,17,21], and, in particular, for homogeneous ensembles [16]. In this benchmark, the data are filtered using K-fold cross-validation. At each iteration, an ensemble is trained using data from K-1folds. Then, the ensemble is used to predict the labels of the instances in the remaining fold. The examples of the holdout fold that are incorrectly classified are removed from the dataset. The process is repeated to filter the other folds. We implemented this method using random forest of size 501 and K = 3. Finally, a random forest composed of 501 random trees is built on the training set cleaned with majority filtering. The label *Fl\_maj\_rf* is used to denote this benchmark. As an additional reference, we also present the results of standard random forest (labelled as RF) and AdaBoost (labeled as Boosting) trained on the original uncleaned training set.

#### 4.1. Predictive accuracy

The results of the experiments carried out to compare the accuracies of the different methods considered are summarized in Tables 2–4. The test errors reported correspond to averages over 50 realizations of the training and test sets. These averages followed by their standard deviations after the  $\pm$  sign. To assess the significance of these observations, the results of an overall comparison of the different methods are summarized in Fig. 1 using the methodology introduced in [4]. In these plots, the average ranks of the different methods are compared. The differences between the average ranks of two methods are statistically significant at a level  $\alpha = 0.05$  if they are above a critical distance (CD). Methods whose average ranks are not significantly different are linked by a horizontal line. Average rank plots are presented for experiments with different levels of class-label noise injected: 0, 10, 20, 30 and 40%.

From the results presented in Fig. 1, Tables 2-4, one concludes that random forests with optimal filtering or cleaning (Fl\_rf and *Cl\_rf*) are among the most accurate ensembles at all noise levels. When no noise is injected, boosting and random forest trained on the original data are slightly better than Fl\_rf and Cl\_rf. However, the differences are not statistically significant. Because of its progressive emphasis on incorrectly classified instances, boosting is not robust to errors in the class labels. This is apparent from the degradation of its performance in problems with higher levels of noise. In fact, boosting has the worst average rank for 10-40% noise levels. In these cases, the differences with most other methods are statistically significant. The differences between filtering and cleaning are not statistically significant. Filtering seems to have a slightly better performance at lower noise levels. When either 30% or 40% of the class labels of the training instances are perturbed, cleaning is slightly better than filtering. This is probably a consequence of the loss of information that results form discarding instances in filtering. Finally, we observe that selecting adequate values of the threshold  $\theta$  for noise filtering or cleaning is superior to using standard majority filtering at all noise levels. These improvements are more pronounced at lower noise levels.

Dataset	Noise (in %)	Fl_rf	Cl_rf	Fl_maj_rf	Boosting	RF
German	0 10 20 30 40	$\frac{24.5 \pm 2.2}{25.0 \pm 1.9}$ $26.5 \pm 2.7$ $28.9 \pm 2.7$ $32.2 \pm 4.2$	$\begin{array}{c} 24.6 \pm 1.9 \\ 25.6 \pm 2.4 \\ \underline{26.7} \pm \underline{2.3} \\ \underline{28.6} \pm \underline{2.8} \\ \textbf{31.7} \pm \textbf{3.9} \end{array}$	$26.5 \pm 2.1 26.7 \pm 2.4 27.1 \pm 2.2 28.3 \pm 2.6 32.8 \pm 4.0$	$\begin{array}{c} 25.1 \pm 2.1 \\ 28.1 \pm 2.2 \\ 32.7 \pm 3.1 \\ 38.3 \pm 3.0 \\ 43.3 \pm 2.9 \end{array}$	$\begin{array}{c} \textbf{24.2} \pm \textbf{2.0} \\ \textbf{24.8} \pm \textbf{1.9} \\ \textbf{27.1} \pm \textbf{2.7} \\ \textbf{31.0} \pm \textbf{2.7} \\ \textbf{37.9} \pm \textbf{3.6} \end{array}$
Ionosphere	0 10 20 30 40	$6.9 \pm 2.0$ 7.7 ± 2.0 $9.4 \pm 3.1$ 14.8 ± 4.5 $32.2 \pm 4.2$	$\begin{array}{c} 6.9\pm 2.0\\ \underline{7.8}\pm \underline{2.2}\\ 9.5\pm 3.7\\ \underline{14.1}\pm \underline{4.0}\\ \textbf{31.7}\pm \textbf{3.9} \end{array}$	$7.3 \pm 1.9 \\ 8.0 \pm 2.4 \\ 9.3 \pm 3.0 \\ 13.8 \pm 4.6 \\ 32.8 \pm 4.0 \\$	$6.4 \pm 1.8 \\ 11.3 \pm 3.3 \\ 19.0 \pm 4.2 \\ 28.0 \pm 5.3 \\ 43.3 \pm 2.9$	$\begin{array}{c} \underline{6.8} \pm \underline{1.9} \\ \underline{7.8} \pm \underline{2.1} \\ 11.3 \pm 3.7 \\ 18.4 \pm 3.8 \\ 37.9 \pm 3.6 \end{array}$
Heart	0 10 20 30 40	$\begin{array}{c} \textbf{17.5} \pm \textbf{3.4} \\ \underline{19.6} \pm \underline{4.3} \\ 21.7 \pm 4.0 \\ \underline{24.2} \pm \underline{5.6} \\ \underline{34.2} \pm \underline{7.6} \end{array}$	$18.1 \pm 3.8 \\ 19.7 \pm 4.3 \\ \textbf{21.4} \pm \textbf{4.5} \\ \underline{24.2} \pm \underline{5.5} \\ \textbf{33.7} \pm \textbf{7.2}$	$\frac{17.9 \pm 3.8}{18.9 \pm 4.8}$ $\frac{21.5 \pm 4.2}{24.0 \pm 5.3}$ $34.5 \pm 7.3$	$\begin{array}{c} 20.9\pm3.4\\ 26.3\pm4.1\\ 32.1\pm5.8\\ 37.1\pm5.3\\ 43.2\pm6.6\end{array}$	$18.1 \pm 3.6 \\ 20.3 \pm 3.8 \\ 23.5 \pm 4.3 \\ 28.7 \pm 5.7 \\ 38.4 \pm 6.2$
Horse-colic	0 10 20 30 40	$\frac{15.6 \pm 2.5}{17.5 \pm 3.0}$ $\frac{20.2 \pm 3.6}{26.2 \pm 4.7}$ $\frac{26.0 \pm 6.3}{36.0 \pm 6.3}$	$\begin{array}{c} \textbf{15.5} \pm \textbf{2.7} \\ \textbf{17.5} \pm \textbf{3.1} \\ 20.5 \pm 4.1 \\ 26.5 \pm 4.9 \\ \textbf{35.2} \pm \textbf{6.5} \end{array}$	$17.6 \pm 3.6$ $\underline{18.2} \pm \underline{3.4}$ $\underline{20.3} \pm \underline{4.0}$ $26.0 \pm 5.5$ $36.0 \pm 6.5$	$\begin{array}{c} 15.7 \pm 2.6 \\ 22.3 \pm 3.4 \\ 28.3 \pm 4.9 \\ 35.4 \pm 4.5 \\ 43.6 \pm 4.7 \end{array}$	$\begin{array}{c} \textbf{15.5} \pm \textbf{2.7} \\ \textbf{17.5} \pm \textbf{2.9} \\ 20.7 \pm 3.7 \\ 29.0 \pm 4.9 \\ 39.7 \pm 4.5 \end{array}$
Liver	0 10 20 30 40	$     \frac{28.4 \pm 4.2}{31.4 \pm 4.3^{*}}     34.5 \pm 4.6     \frac{38.0 \pm 5.9}{43.8 \pm 5.6}   $	$\begin{array}{c} 29.7 \pm 3.8\\ \underline{32.7} \pm \underline{4.9}\\ 35.5 \pm 4.2\\ 38.5 \pm 5.3\\ \textbf{43.1} \pm \textbf{5.4} \end{array}$	$31.2 \pm 4.4$ $33.6 \pm 4.8$ $36.0 \pm 4.7$ $38.8 \pm 5.7$ $44.6 \pm 4.6$	$\begin{array}{c} 29.6 \pm 3.4 \\ 34.4 \pm 3.9 \\ 37.7 \pm 4.6 \\ 40.5 \pm 5.5 \\ 44.9 \pm 4.6 \end{array}$	$27.4 \pm 3.8^{*}$ $31.4 \pm 4.4^{*}$ $34.6 \pm 4.5$ $37.7 \pm 5.8$ $43.5 \pm 5.3$
Parkinsons	0 10 20 30 40	$11.9 \pm 3.5 \\ 13.5 \pm 3.8 \\ 17.8 \pm 4.7 \\ \underline{21.0} \pm 5.2 \\ \underline{31.8} \pm \underline{8.2}$	$11.4 \pm 3.8 \\ 13.5 \pm 3.8 \\ \underline{18.1} \pm \underline{5.2} \\ 21.5 \pm 4.7 \\ \underline{31.8} \pm \underline{8.0}$	$\begin{array}{c} 16.3 \pm 4.0 \\ 16.4 \pm 4.0 \\ \underline{18.1} \pm \underline{4.5} \\ \textbf{20.7} \pm \textbf{5.4} \\ \textbf{29.3} \pm \textbf{7.5}^* \end{array}$	$\begin{array}{c} \textbf{8.1} \pm \textbf{3.5}^{*} \\ \textbf{12.8} \pm \textbf{3.9} \\ 22.0 \pm 5.0 \\ 29.8 \pm 7.1 \\ 39.0 \pm 6.2 \end{array}$	$\frac{11.0 \pm 3.5}{13.0 \pm 3.7}$ $17.8 \pm 4.7$ $23.9 \pm 5.5$ $34.6 \pm 7.1$
Ringnorm	0 10 20 30 40	$\begin{array}{c} 6.1 \pm 1.0 \\ \underline{6.8} \pm \underline{1.3} \\ \underline{8.4} \pm \underline{1.8} \\ \underline{11.8} \pm \underline{2.7} \\ 21.1 \pm 5.3 \end{array}$	$6.2 \pm 1.0$ $6.9 \pm 1.1$ $8.3 \pm 1.6$ $11.5 \pm 3.1$ $18.2 \pm 5.3^*$	$\begin{array}{c} 8.3 \pm 1.4 \\ 8.7 \pm 1.5 \\ 9.9 \pm 2.2 \\ 12.2 \pm 3.2 \\ \underline{20.8} \pm \underline{6.4} \end{array}$	$\begin{array}{c} \textbf{4.4} \pm \textbf{0.4}^* \\ 8.6 \pm 1.1 \\ 15.2 \pm 1.7 \\ 24.4 \pm 2.7 \\ 36.3 \pm 3.1 \end{array}$	$\frac{6.0 \pm 1.0}{6.7 \pm 1.1}$ 8.3 ± 1.7 12.8 ± 2.4 24.7 ± 3.5

 Table 3

 Test errors of the different methods for different noise levels (II).

\* *p*-value < 0.05.

#### 4.2. Optimal values of the hyperparameters

In this section, we analyze the trends in the values selected for the sampling rate and the threshold,  $\theta^*$ . We consider first the values of the sampling rate. In plot (a) of Fig. 2 the average sampling rates obtained by the proposed cleaning procedure for each dataset and noise level are displayed. One can observe that in the original unperturbed problems (white bars in the plot) the optimal values for this hyperparameter are strongly problem dependent. For instance, in Chess, Spambase, and Tic-tac-toe these values are above the standard 100% resampling rate. It is likely that for these problems the variability random forests is too large, and that oversampling is an effective mechanism to reduce it. In the remaining problems subsampling, which increases the variability of the ensemble, seems to be more effective. In some cases, such as Blood Transfusion, Breast, Ionosphere, and Twonorm, the optimal values of the sampling ratios are fairly low (around or below 30%, on average). In general, as more class-label noise is injected in the data, the optimal sampling ratios become smaller. This confirms the observation that subsampling becomes more effective as the amount of class-label noise increases [13,14].

In the method proposed, an instance is marked as noisy when it is incorrectly labeled by a fraction of classifiers that exceeds a specified threshold  $\theta^*$ . The optimal value of this parameter is also strongly problem-dependent. However, from the results presented in plot (b) of Fig. 2, it is apparent that as more noise is injected, the values of  $\theta^*$  become smaller and approach 0.5, which corresponds to majority voting. This is consistent with the observation that majority filtering becomes more effective at higher noise levels.

#### 4.3. Noise detection

The value of the threshold used to mark an instance as noisy is determined on the basis of the accuracy of the wrapper classifier trained on cleansed data. The question remains whether this procedure is actually effective for the detection of noisy instances. To investigate this issue, we have recorded the fraction of instances marked as noise for the different classification tasks and with different levels of injected noise. The results of these experiments are presented in the plots of Fig. 3 for the following noise levels: 0% (first row), 10% (second row) and 30% (third row). Each plot of this figure shows the average percentage of instances that are marked as noise with a white bar. From those instances, the ones corresponding to the artificially injected noise are marked in red. The results for the proposed cleaning procedure using filtering are summarized in the plots in the left column of this figure. For reference, the results of majority filtering benchmark are displayed in the plots in the right column.

An analysis of the results for the unperturbed classification tasks (first row of Fig. 3) reveals that the levels of detected noise vary significantly. In the synthetic problems, which, by construction, are noiseless, the proposed cleaning procedure using filtering discards only a small percentage instances. By contrast, in some problems (such as *Blood transfusion, Diabetes, German* and *Liver*) a significant fraction of instances are identified as noise. Filtering these instances does not appear to be detrimental. As a matter







Fig. 3. Percentage of filtered examples (white bars) and filtered examples that correspond injected noise (red part of the bars) for different noise levels: without injected noise (first row), 10% (second row), and 30% (third row) for the proposed cleansing procedure (left column) and majority filtering (right column). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

	Noise					
Dataset	(in %)	Fl_rf	Cl_rf	Fl_maj_rf	Boosting	RF
Sonar	0 10 20 30 40	$18.8 \pm 3.9 \\20.7 \pm 5.2 \\\underline{24.2} \pm 4.8 \\31.7 \pm 7.1 \\38.6 \pm 6.2$	$\begin{array}{c} 19.4 \pm 4.8 \\ 21.5 \pm 4.9 \\ 24.3 \pm 5.4 \\ 31.9 \pm 7.2 \\ 39.7 \pm 7.5 \end{array}$	$\begin{array}{c} 24.5\pm 3.8\\ 25.2\pm 4.6\\ 26.1\pm 5.2\\ 33.1\pm 7.7\\ 39.6\pm 7.6 \end{array}$	$\begin{array}{c} \textbf{14.4} \pm \textbf{3.8}^{*} \\ \underline{20.1} \pm \underline{4.4} \\ 26.9 \pm 5.5 \\ 35.1 \pm 6.6 \\ 42.6 \pm 7.7 \end{array}$	$\frac{18.5 \pm 3.7}{19.9 \pm 4.7}$ 23.7 ± 5.5 31.2 ± 7.6 $\underline{39.4 \pm 7.4}$
Spambase	0 10 20 30 40	$5.1 \pm 0.6$ $5.9 \pm 0.5^{*}$ $6.7 \pm 0.6$ $7.8 \pm 0.8^{*}$ $11.8 \pm 2.2$	$\begin{array}{l} 5.1 \pm 0.6 \\ \underline{6.1} \pm \underline{0.7} \\ \textbf{6.7} \pm \textbf{0.6} \\ \textbf{7.8} \pm \textbf{0.8}^* \\ \textbf{11.1} \pm \textbf{1.5}^* \end{array}$	$\begin{array}{c} 6.3 \pm 0.5 \\ 6.5 \pm 0.5 \\ \underline{6.9} \pm \underline{0.7} \\ \underline{8.6} \pm \underline{0.9} \\ 14.3 \pm 1.5 \end{array}$	$\begin{array}{c} \textbf{4.9} \pm \textbf{0.4} \\ \textbf{9.4} \pm \textbf{0.8} \\ \textbf{13.7} \pm \textbf{1.1} \\ \textbf{20.3} \pm \textbf{1.4} \\ \textbf{32.0} \pm \textbf{1.6} \end{array}$	$\begin{array}{c} \underline{5.0}\pm \underline{0.5}\\ 6.6\pm 0.6\\ 9.3\pm 0.9\\ 14.6\pm 1.1\\ 25.6\pm 1.3\end{array}$
Threenorm	0 10 20 30 40	$17.0 \pm 1.0$ $18.8 \pm 1.5$ $20.9 \pm 2.2$ $25.6 \pm 2.6$ $33.4 \pm 4.5$	$\begin{array}{c} 17.2 \pm 1.3 \\ 19.1 \pm 1.6 \\ 21.2 \pm 2.6 \\ \underline{25.7} \pm \underline{2.5} \\ \textbf{33.0} \pm \textbf{4.7} \end{array}$	$\begin{array}{c} 19.2 \pm 1.4 \\ 20.5 \pm 1.8 \\ 21.7 \pm 2.5 \\ 26.7 \pm 2.5 \\ 34.3 \pm 5.1 \end{array}$	$\frac{16.8 \pm 0.9}{20.7 \pm 1.4}$ $25.5 \pm 1.6$ $32.3 \pm 2.1$ $40.6 \pm 2.9$	$\begin{array}{c} \textbf{16.5} \pm \textbf{0.9} \\ \textbf{18.4} \pm \textbf{1.2}^* \\ \textbf{20.8} \pm \textbf{1.8} \\ \textbf{26.5} \pm \textbf{2.1} \\ \textbf{35.7} \pm \textbf{3.1} \end{array}$
Tictactoe	0 10 20 30 40	$2.5 \pm 0.9$ <b>5.7 ± 2.2</b> $12.5 \pm 2.9$ $22.1 \pm 2.7$ <b>34.4 ± 3.7</b>	$\begin{array}{c} 2.4 \pm 1.1 \\ 6.2 \pm 2.3 \\ 13.8 \pm 2.8 \\ 23.5 \pm 2.6 \\ \textbf{34.4} \pm \textbf{3.5} \end{array}$	$7.2 \pm 2.5 \\12.0 \pm 2.6 \\17.4 \pm 2.8 \\24.1 \pm 3.3 \\35.1 \pm 3.5$	$\begin{array}{c} \textbf{0.6} \pm \textbf{0.5}^{*} \\ 10.4 \pm 1.7 \\ 21.3 \pm 2.3 \\ 30.5 \pm 2.6 \\ 40.5 \pm 3.7 \end{array}$	$\frac{2.3 \pm 1.0}{5.8 \pm 1.6}$ <b>12.4 ± 2.3</b> <b>21.7 ± 2.2</b> $35.7 \pm 3.2$
Twonorm	0 10 20 30 40	$\begin{array}{c} 3.9 \pm 0.5 \\ \underline{4.6} \pm \underline{0.7} \\ 5.6 \pm 1.2 \\ 8.2 \pm 2.9 \\ 17.5 \pm 6.3 \end{array}$	$\begin{array}{c} 3.9 \pm 0.5 \\ \textbf{4.5} \pm \textbf{0.6} \\ \underline{5.5} \pm \underline{0.9} \\ \overline{7.7} \pm \underline{2.8} \\ \underline{16.7} \pm \underline{6.9} \end{array}$	$\begin{array}{c} 4.3 \pm 0.5 \\ \textbf{4.5} \pm \textbf{0.5} \\ \textbf{5.2} \pm \textbf{0.9} \\ \textbf{7.3} \pm \textbf{2.3} \\ \textbf{15.2} \pm \textbf{5.5} \end{array}$	$\begin{array}{c} \textbf{3.7} \pm \textbf{0.3} \\ \textbf{7.7} \pm \textbf{1.0} \\ \textbf{13.7} \pm \textbf{1.6} \\ \textbf{23.7} \pm \textbf{2.4} \\ \textbf{35.7} \pm \textbf{3.2} \end{array}$	$\begin{array}{c} \underline{3.8} \pm \underline{0.4} \\ \overline{4.9} \pm 0.7 \\ 6.5 \pm 1.0 \\ 10.8 \pm 1.8 \\ 23.2 \pm 3.3 \end{array}$

Test errors of the different methods for different noise levels (III).

\* *p*-value < 0.05.

of fact, in these problems the proposed cleaning procedure yields competitive or better accuracy rates with respect to random forest trained on the uncleaned data (see Tables 2 and 3).

Table 4

For the unperturbed classification tasks, majority filtering is much more aggressive and marks many more instances as noise than the proposed procedure. In problems without noise in the class labels, such as Threenorn, Tic-tac-toe, Ringnorm, Twonorm, around or above 5% of the training instances are discarded (see upper right plot on Fig. 3). As a result, there is a significant decrease of the accuracy for random forest trained on these cleansed data (see Tables 3 and 4). In real-world problems, it is not possible to know the level of intrinsic noise. Nonetheless, majority filtering is likely to discard too many instances as well. Specifically, this method marks more than 20% of the instances as noise in five of the datasets analyzed (Blood transfusion, Diabetes, German, Liver and Sonar). The accuracy of random forest trained on the data cleansed by majority filtering in Blood transfusion and in Diabetes is fairly good. However, in Liver and Sonar it is the least accurate among the methods considered.

The results of experiments on classification tasks perturbed with 10% and 30% class-label noise, are displayed in the second and the third rows of Fig. 3, respectively. The red part of the bars is the percentage of instances whose class-label has been switched in the noise injection process that are identified as noisy. In most cases for the proposed cleansing procedure based on optimal filtering, the height of the red bar is well below the level of noise injected. This means that a significant fraction of perturbed instances are not detected. An extreme example is Sonar. We conjecture that this lack of sensitivity is due to the strong overlap of the distributions for the two classes. For this reason, it is difficult to single out noisy instances located in regions where such overlap is high. Using majority filtering, which is more aggressive, it is possible to detect most of the injected noisy instances. In fact, the relative performance of majority filtering improves when the levels of class-label noise are high. Still, even at high noise levels, the precision of the method is rather low: the percentage of instances that are marked as noise by majority voting is significantly larger than the level of noise injected. By contrast, even though the proposed optimal filtering procedure fails to identify some noisy instances, those that are identified by this strategy are more likely to be noise in actuality. In some datasets (*Boston, Breast, Chess, Parkinsons, Ringnorm, Spambase, Tictactoe* and *Twonorm*) the proposed procedure detects a fairly high percentages of the injected noise without removing a significant number of noiseless instances. For these datasets, random forests trained on data cleaned with the proposed procedure are more accurate than those trained based on data cleaned majority filtering, except in *Breast* and *Twonorm*, where the differences are not statistically significant (see Tables 2–4).

In summary, the proposed cleansing procedure achieves high specificity at the expense of not being able to detect some noisy instances. By contrast, majority filtering detects most noisy instances, but also incorrectly discards a high percentage of valid ones. As shown in Fig. 2 (b),  $\theta^*$ , the optimal threshold for filtering, becomes closer to 0.5 (majority filtering) as the level of class-label noise increases.

#### 5. Conclusions

In this paper, we have proposed a two-stage method for the detection of class-label noise based on the robustness to noise of randomized ensembles that use resampling [13,14]. Near-optimal values of the sampling rate can be determined using out-of-bag data. Typically, the selected sampling ratios become smaller as the level of class label noise increases. Another important observation is that the classifiers in the ensemble tend to make more errors on noisy instances. Therefore, the fraction of incorrect predictions can be used as an indicator for noise detection: if this quantity is above a threshold,  $\theta$ , the instance considered is marked as noisy. Standard values for the threshold are  $\theta = 0.5$  (majority) or  $\theta = 1$  (consensus). In this work, we have shown that the best results are obtained at a value  $\theta^*$  that is intermediate between these extremes. A simple wrapper procedure is proposed to determine  $\theta^*$ . This optimal value depends on the problem under consideration and the

amount of class-label noise. Values of  $\theta^*$  closer to majority filtering are generally obtained for noisy problems. However, majority cleansing tends to discard instances that are correctly labeled. This has been shown to be disadvantageous, specially in problems with low levels of noise. In general, adjusting the threshold used to detect noisy instances allows us to build more accurate ensembles at all levels of class-label noise.

Once the noisy instances have been identified, they can be removed from the training data (filtering) or corrected (cleaning). Filtering is slightly superior at low and medium noise levels. Cleaning tends to be more accurate when the noise levels are high. The reason for this behavior is that filtering discards training instances. This involves some information loss, which could be detrimental if too many instances are discarded.

#### Acknowledgments

The authors acknowledge financial support from the *Spanish Ministry of Economy, Industry and Competitiveness,* projects TIN2013-42351-P, TIN2016-76406-P, and TIN2015-70308-REDT, and of the *Comunidad de Madrid,* project CASI-CAM-CM (S2013/ICE-2845).

#### References

- K.M. Ali, M.J. Pazzani, Error reduction through learning multiple descriptions, Mach. Learn. 24 (3) (1994) 173–202.
- [2] J. Bi, T. Zhang, Support vector classification with input data uncertainty, in: L.K. Saul, Y. Weiss, L. Bottou (Eds.), Advances in Neural Information Processing Systems 17, MIT Press, 2005, pp. 161–168.
- [3] C.E. Brodley, M.A. Friedl, Identifying mislabeled training data, J. Artif. Intell. Res. 11 (1) (1999) 131–167.
- [4] J. Demšar, Statistical comparisons of classifiers over multiple data sets, J. Mach. Learn. Res. 7 (2006) 1–30.
- [5] B. Frénay, M. Verleysen, Classification in the presence of label noise: a survey, IEEE Trans. Neural Netw. Learn. Syst. 25 (5) (2014) 845–869.
- [6] Y. Freund, An adaptive version of the boost by majority algorithm, Mach. Learn. 43 (3) (2001) 293–318.
- [7] Y. Freund, R.E. Schapire, Experiments with a new boosting algorithm, in: Proceedings of the Thirteenth International Conference on Machine Learning (ICML), 1996, pp. 148–156.
- [8] D. García-Gil, J. Luengo, S. García, F. Herrera, Enabling smart data: noise filtering in big data classification, Comput. Res. Repos. (2017). http://arxiv.org/abs/ 1704.01770.
- [9] T.M. Khoshgoftaar, S. Zhong, V. Joshi, Enhancing software quality estimation using ensemble-classifier based noise filtering, Intell. Data Anal. 9 (1) (2005) 3–27.
- [10] A. Liaw, M. Wiener, Classification and regression by randomforest, R News 2 (3) (2002) 18–22. http://CRAN.R-project.org/doc/Rnews/.
- [11] M. Lichman, UCI machine learning repository, 2013. http://archive.ics.uci.edu/ ml.
- [12] G. Martnez-Muoz, A. Surez, Switching class labels to generate classification ensembles, Pattern Recognit. 38 (10) (2005) 1483–1494.
- [13] M. Sabzevari, G. Martínez-Muñoz, A. Suárez, Improving the robustness of bagging with reduced sampling size, in: Proceedings of the Twenty Second European Symposium on Artificial Neural Networks, ESANN, Bruges, Belgium, 2014. April 23–25.
- [14] M. Sabzevari, G. Martnez-Muoz, A. Su rez, Small margin ensembles can be robust to class-label noise, Neurocomputing 160 (Supplement C) (2015) 18–33.

- [15] B. Sluban, D. Gamberger, N. Lavrač, Ensemble-based noise detection: noise ranking and visual performance evaluation, Data Mining Knowl. Discov. 28 (2) (2014) 265–303.
- [16] B. Sluban, N. Lavra, Relating ensemble diversity and performance: a study in class noise detection, Neurocomputing 160 (Supplement C) (2015) 120–131.
- [17] S. Verbaeten, A. Van Assche, Ensemble Methods for Noise Elimination in Classification Problems, Springer, Berlin, Heidelberg, 2003, pp. 317–325.
- [18] D. Wang, M. Li, Robust stochastic configuration networks with kernel density estimation for uncertain data regression, Inf. Sci. 412–413 (Supplement C) (2017) 210–222.
- [19] S. Zhong, W. Tang, T.M. Khoshgoftaar, Boosted noise filters for identifying mislabeled data, Technical report, Department of Computer Science and engineering, Florida Atlantic University, 2005.
- [20] X. Zhu, X. Wu, Class noise vs. attribute noise: a quantitative study, Artif. Intell. Rev. 22 (3) (2004) 177–210.
- [21] X. Zhu, X. Wu, Q. Chen, Eliminating class noise in large datasets, in: Proceedings of the Twentieth International Conference on Machine Learning (ICML), 2003, pp. 920–927.



**Maryam Sabzevari** received her B.Sc (2008) degree in Computer Science and M.Sc (2010) in Artificial Intelligence from Azad University of Mashhad, Iran. Later, she received another M.Sc (2015) in Neuroinformatics from Universidad Autónoma de Madrid (UAM), Spain. Currently, she is conducting her Ph.D. in Computer Science in Universidad Autónoma de Madrid (Spain), where she is also a Teaching Assistant. Her current research interests include machine learning, pattern recognition, neural networks, ensemble learning and learning methods in the presence of noise.



**Gonzalo Martínez-Muñoz** received the university degree in Physics (1995) and Ph.D. degree in Computer Science (2006) from the Universidad Autónoma de Madrid (UAM). From 1996 to 2002 he worked in industry. Until 2008 he was an interim assistant professor in the Computer Science Department of the UAM. During 2008/2009, he worked as a Fulbright postdoc researcher at Oregon State University in the group of Professor Thomas G. Dietterich. He is currently a professor at Computer Science Department at UAM. His research interests include machine learning, computer vision, pattern recognition, neural networks, decision trees, and ensemble learning.



Computer Science Institute (Berkeley, CA) and at MIT (Cambridge, MA). He has worked on relaxation theory in condensed media, stochastic and thermodynamic theories of nonequilibrium systems, latticegas automata, and automatic induction from data. His current research interests include machine learning, computational statistics, quantitative finance, time series analysis and information processing in the presence of noise. He is a member of IEEE.