

Sequential Ensemble Learning for Outlier Detection: A Bias-Variance Perspective

Shebuti Rayana Wen Zhong Leman Akoglu

Department of Computer Science

Stony Brook University, Stony Brook, NY 11794

Email: {srayana,wzhong,leman}@cs.stonybrook.edu

Abstract—Ensemble methods for classification and clustering have been effectively used for decades, while ensemble learning for outlier detection has only been studied recently. In this work, we design a new ensemble approach for outlier detection in multi-dimensional point data, which provides improved accuracy by reducing error through both bias and variance. Although classification and outlier detection appear as different problems, their theoretical underpinnings are quite similar in terms of the bias-variance trade-off [1], where outlier detection is considered as a binary classification task with unobserved labels but a similar bias-variance decomposition of error.

In this paper, we propose a sequential ensemble approach called **CARE** that employs a two-phase aggregation of the intermediate results in each iteration to reach the final outcome. Unlike existing outlier ensembles which solely incorporate a parallel framework by aggregating the outcomes of independent base detectors to reduce variance, our ensemble incorporates both the parallel and sequential building blocks to reduce bias as well as variance by (i) successively eliminating outliers from the original dataset to build a better data model on which outlierness is estimated (sequentially), and (ii) combining the results from individual base detectors and across iterations (parallelly). Through extensive experiments on sixteen real-world datasets mainly from the UCI machine learning repository [2], we show that **CARE** performs significantly better than or at least similar to the individual baselines. We also compare **CARE** with the state-of-the-art outlier ensembles where it also provides significant improvement when it is the winner and remains close otherwise.

1. Introduction

As a significant subject, outlier detection is widely researched in the literature. There exist various approaches for outlier detection such as density based methods [3], [4], [5] and distance based methods [6], [7], which find unusual points by the distance to their k nearest neighbors (kNNs). However, each of these methods can only focus on some specific kinds of outliers based on the datasets collected from different application domains. There exists no known algorithm that could detect all types of outliers that appear in a wide variety of domains. As a result, ensemble learning for outlier detection has become a popular research area more

recently [1], [8], [9], which aims to put together multiple detectors so as to leverage the “strength of the many”.

In contrast to outlier detection ensembles, classification ensembles have been studied for decades. The explosive growth of classification ensemble models provides the new opportunity to design effective methods for other machine learning tasks including outlier mining. One can categorize ensemble methods into two kinds. The first one is the parallel ensemble, where base learners are created independent of each other and their results are combined to get the final outcome; while the second one is the sequential ensemble, where base learners are created over iterations and have dependency among them. Specifically, several outlier ensembles are proposed based on two seminal works of classification ensembles: (i) the parallel ensemble Bagging [10], which creates base components from different subsamples of training datasets parallelly, and (ii) the sequential ensemble AdaBoost [11], which creates base components iteratively. Among those, some try to induce diversity among the base detectors [1], [8], [12], and others selectively combine outcomes from the candidate detectors [9], [13].

Existing outlier ensembles have several limitations, most importantly they avoid discussing the theoretical aspects of outlier detection. Recently, Aggarwal *et al.* [1] argue that although they appear to be very different problems, classification and outlier detection share quite similar theoretical underpinnings in terms of the bias-variance trade-off. Specifically, one can consider the outlier detection problem as a binary classification task where the labels are unobserved, the inliers being the majority class and the outliers the minority class, and the error of a detector can be decomposed into bias and variance terms in a similar way. In existing outlier ensembles, various parallel frameworks combining multiple detector outcomes are designed to reduce variance only, most of which are incapable of overcoming the presence of inaccurate base detectors. On the other hand, it remains challenging to reduce bias in a controlled way for outlier detection or remove inaccurate detectors due to the lack of ground truth to validate the results during the intermediate steps. There exist some successful heuristic approaches to reduce bias. One such commonly used approach is to remove outliers in successive iterations [14] to build more robust outlier models iteratively.

In this paper, we study the feasibility of bias-variance reduction under the unsupervised setting, and propose a

sequential ensemble model called Cumulative Agreement Rates Ensemble (CARE), to reduce both bias and variance for outlier detection. Specifically, each iteration in the sequential ensemble consists of two aggregation phases: (1) in the first phase, we combine the results of feature-bagged base detectors using weighted aggregation, where weights are estimated in an unsupervised way through the Agreement Rates (AR) method by [15], and (2) in the second phase, the result of the current iteration is aggregated with the combined result from the previous iterations cumulatively. These two phase aggregations in each iteration aim to reduce the variance. Furthermore, we use the combined result from the previous iterations to improve the next iteration by removing the top (i.e., most obvious) outliers and perform a variable probability sampling to create the data model to be used for the next iteration. The removal of top outliers in successive iterations aims to reduce the bias.

To the best of our knowledge, this is the first work focusing on reducing both bias and variance for unsupervised outlier detection. In general, this paper offers the following contributions:

- We design a new approach which incorporates weighted aggregation of feature-bagged base detectors, where weights are estimated in an unsupervised fashion (Section 4.3.1 and 4.3.2).
- We devise a sequential ensemble over the weighted combination, which cumulatively aggregates the results from multiple iterations until a stopping condition is met (Section 4.3.3 and 4.3.4).
- We provide a new sampling approach called Filtered Variable Probability Sampling (FVPS) which utilizes the result from the previous iteration to filter the top outliers, and uses variable probability sampling to select points from the original data to create the data model for the next iteration (Section 4.3.3).
- Our sequential ensemble is designed to reduce both bias and variance and improves the overall result. Moreover, we provide experiments with synthetic datasets to support this claim (Section 5).

We evaluate our method on sixteen different real-world datasets majority of which are from the UCI machine learning repository [2]. Our results show that CARE outperforms the baseline (i.e., non-ensemble) detectors in most cases and remains close to the baselines in cases where it falls shorter. We also compare CARE with the existing state-of-the-art outlier ensembles [1], [8], [16]. Similarly, it provides significant improvement when it is the winner, and performs close otherwise. (Section 6)

2. Related Work

2.1. Ensemble Models

Ensemble models for classification have been extensively studied in the literature for several decades. In 1996, Breiman [10] presented a parallel ensemble, today well-known as *bagging*, which consists of multiple predictor

components trained on samples of the original dataset, to infer the label using a plurality vote, enhancing the model through variance reduction. However, bagging omitted the bias term and could only reduce the variance. To make up this deficit, a new sequential ensemble known as *boosting* was devised by Freund *et al.* [11]. Their proposed AdaBoost algorithm assigned larger weights for the misclassified instances to advance the given base classification algorithm and combined weighted sum of multiple weak learners into a boosted classifier to reduce both bias and variance.

After these two important seminal works, a proliferation of ensemble methods followed, aiming to explain and improve over the original methods. A two-stage process called adaptive bagging [17] was proposed to perform bias and variance reduction respectively. It generated an intermediate output in the first stage and the altered first stage output was adopted as new input of the second stage that is bagging. Sun *et al.* [18] analyzed the influence of adding a cost term to AdaBoost and offered the cost-sensitive boosting algorithm for imbalanced data classification.

Our proposed outlier detection approach CARE uses similar insights as in AdaBoost; by sequentially updating the data model on which data points are scored for outlierness as well as by combining multiple detector outcomes parallelly to reduce bias and variance, respectively. Since the above methods can be learned in supervised settings while most anomaly detection tasks provide no labels, our method differs from the existing classifier ensembles in that we focus on reducing both bias and variance in a fully unsupervised setting, which has not been studied before.

2.2. Outlier Ensembles

Outlier ensemble learning, as a rarely explored area, mainly tries to reduce the variance through the combination of different base detectors. A parallel approach called feature bagging, proposed by Lazarevic and Kumar [12], built an ensemble based on randomly selected feature subsets from original features to detect outliers in high-dimensional and noisy datasets. Inspired by random forests [19], Liu *et al.* [16] employed different subsamples of training data to establish an ensemble of trees to isolate outliers on the basis of the path length from the root to the leaves.

In recent years, with more attention focusing on outlier ensembles, several work discussed the theories and emphasized the crucial aspects of ensemble model construction. Aggarwal [20] and Zimek [21] talked about algorithmic patterns, categorization, and the important building blocks of outlier ensembles such as model combination and diversity of base models. Zimek *et al.* [8] analytically and experimentally studied the subsampling technique and improved results through building an ensemble on top of several subsamples without mentioning the subsample selection. Aggarwal and Sathe [1] deduced the bias-variance trade-off theory from classification to outlier detection, clarified some misconceptions about the existing subsampling methods, and proposed more effective subsampling and feature bagging approaches. On base detector combination, Rayana and Akoglu [9] presented unsupervised strategies to select

a subset of trusted detectors while omitting inaccurate ones in an unsupervised way.

Unlike existing outlier ensembles that solely employ a sequential or parallel framework, our proposed method CARE incorporates both of these building blocks to reduce both bias and variance. These two phases respectively involve (i) successively eliminating outliers from the original dataset to build a better data model on which outlierness is estimated (sequentially), and (ii) combining the results from individual base detectors and across iterations (parallelly).

3. Background and Preliminaries

3.1. Outlier Detection Problem

A popular characterization of an outlier is given by Hawkins as “an observation which deviates so much from other observations as to arouse suspicion that it was generated by a different mechanism” [22]. A common approach to outlier detection is to find unusual multi-dimensional points by quantifying a measure of normality relative to their neighboring points. Based on this notion, there exist two major varieties of outlier detectors: (i) distance, and (ii) density based. Specifically, distance based detectors find data points which are far from their nearest neighbors and density based detectors find the points which reside in a lower density region compared to their nearest neighbors. Formally, the problem can be stated as follows:

Given a multi-dimensional data D with n individual points in d dimensions;

Find outliers which are far from the rest of the data (i.e., inliers) or reside in a lower density region.

3.2. Bias-Variance Trade-off for Outlier Detection

The bias-variance trade-off is often explained in the context of supervised learning, e.g., classification, as quantification of bias-variance requires labeled data. Although outlier detection problems are solved using unsupervised approaches due to the lack of ground truth, the bias-variance trade-off can be quantified similarly by treating the dependent variable (actual labels) as unobserved [1].

Unlike classification, most outlier detection algorithms output “outlierness” scores for the data points. We can consider the outlier detection problem as a binary classification task having a majority class (inliers) and a minority class (outliers) by converting the outlierness scores to class labels. The points with scores above a threshold are considered as outliers with label 1 (label 0 for inliers below threshold). After converting the unsupervised outlier detection problem to a classification task with only unobserved labels, we can explain the bias-variance trade-off for outlier detection using ideas from classification. Specifically, the expected error of outlier detection can be split into two main components: reducible error and irreducible error (i.e., error due to noise). The reducible error can be minimized to maximize the accuracy of the detector. Furthermore, the reducible error can be decomposed into (i) error due to squared bias, and (ii) error due to variance. However, there is a trade-off while minimizing both these sources of errors.

Bias of a detector is the amount by which the expected output of the detector differs from the actual label (unobserved) over the training data. While variance of a detector is the expected amount by which the output of a detector over one training set differs from the expected output of the detector over all the training sets. Simply put, the trade-off between bias and variance can be viewed as, (i) a detector with low bias is very flexible in fitting the data well and fits each training set differently with high variance, and (ii) an inflexible detector fits each training set almost similarly yielding high bias but low variance.

3.3. Motivation for Ensembles

Our goal in this work is to improve outlier detection by reducing both bias and variance, and in return decreasing the reducible error. It is evident from the classification ensemble literature that combining results from multiple base algorithms decreases the overall variance of the ensemble [1], [10], [12], which is also true for outlier ensembles. On the other hand, this combination does not provide any evidence for reducing bias, as controlled bias reduction is rather difficult due to lack of ground truth. However, there exist some successful heuristic approaches for reducing bias by removing outliers iteratively to build a more robust successive outlier detector. This iterative approach can be considered as a sequential ensemble. The basic idea here is that the outliers interfere with the creation of the normal data model, and the removal of points with high outlierness scores will be beneficial for the outlier model to produce an output close to the actual (unobserved) labels in expectation.

4. Proposed Approach

4.1. Overview

CARE takes the d -dimensional data, a value for k (nearest neighbor count), and a value for $MAXITER$ as input and outputs an outlierness score list \mathbf{fs} and a rank list \mathbf{r} (ranked based on most to least outlierness) of all the data points. In the experiments we use $k = 5$, which is compatible with the state-of-the-art methods [1], [8]. Moreover, parameter k in subsampling methods is scaled by the inverse of various subsample sizes, as such large k is not required (the smaller the subsample, the larger the relative neighborhood per point becomes for a fixed k). As for $MAXITER$, we set it to 15, a relatively small value. We assume that our approach improves the base detectors over iterations, and the results are stabilized after only a few iterations and the algorithm stops following the stopping criterion.

The main steps of CARE are given in Algorithm 1. Step 3 creates the feature-bagged outlier detectors as base detectors of the ensemble. For the first iteration the sample set S contains the whole data D as shown in step 1. For each base detector, $q \in [d/2, d - 1]$ features are selected randomly to create the corresponding feature bag. We create b ($= 100$) feature-bagged base detectors. Motivated by Platanios *et al.* [15], step 4 calculates the pairwise agreements a_A for all

Algorithm 1: CARE Outlier Detection Ensemble

Input: d -dimensional Data D , NN count $k = 5$,
 $MAXITER = 15$
Output: Score list (\mathbf{fs}) and rank list (\mathbf{r}) of points

- 1: $S = D$ (initially); $E = \emptyset$; $iter = 0$
- 2: **while** $iter \leq MAXITER$ **do**
- 3: Obtain results from (b) feature-bagged base detectors (D, S, k) [Section 4.2]
- 4: Calculate pairwise agreement rates a_A for all base detector pairs in set A
- 5: Estimate detector errors \mathbf{e} ($b \times 1$) based on a_A [Section 4.3.1]
- 6: Compute detector weights using estimated errors [Section 4.3.2]
- 7: Compute pruned weighted outlieriness scores of data points to get combined scores (\mathbf{ws}) [Section 4.3.2]
- 8: $E = E \cup \mathbf{ws}$
- 9: $\mathbf{fs} = average(E)$
- 10: Generate new data sample S from D using FVPS (w/o replacement) on \mathbf{fs} [Section 4.3.3]
- 11: **if** *stopping condition* is TRUE **then**
- 12: **break** [Section 4.3.4]
- 13: **end if**
- 14: $iter = iter + 1$
- 15: **end while**
- 16: $\mathbf{r} = sort(\mathbf{fs})$ (descending order)

possible pairs of base detectors and step 5 estimates the error rates of the individual base detectors in an unsupervised way using a_A . Step 6 calculates weights for the base detectors using their corresponding error rates. Step 7 combines the outlieriness scores from the different base detectors with weighted average combination to get final outlieriness scores \mathbf{ws} . Step 8 stores the outlieriness scores in E at each iteration and step 9 calculates the final outlieriness scores \mathbf{fs} by averaging the results of all previous iterations as well as the current iteration. Based on \mathbf{fs} , step 10 generates the new data sample S (where, $|S| < |D|$) using the FVPS approach w/o replacement (see Section 4.3.3) and step 11 generates the ranked list \mathbf{r} of instances from most to least outlieriness. We repeat steps 3-16 until the stopping condition at step 12 is met or upto the given maximum iteration $MAXITER$.

Unlike existing ensemble techniques, CARE incorporates a two-phase aggregation approach in each iteration; first, it combines the results from the individual base detectors (parallel) and second, it cumulatively aggregates the results from multiple iterations (sequential).

Next we describe the main components of our proposed CARE in detail. In particular, we describe the base detectors in Section 4.2 and consensus approaches in Section 4.3.

4.2. Base Detectors

There exist various approaches for outlier detection based on different aspects of outliers, designed for distinct applications to detect domain-specific outliers. In our work, we are interested in *unsupervised* outlier detection

approaches that assign outlieriness scores to the individual instances in the data, as such, allow ranking of instances based on outlieriness.

4.2.1. kNN based Outlier Detectors. There are a number of well-known unsupervised approaches, e.g., “distance based” and “density based” methods for outlier detection. Distance based methods [6], [23] and their variants try to find the *global* outliers far from the rest of the data based on k nearest neighbor (kNN) distances of the data points. On the other hand, density based methods [3], [4], [5] and their variants try to find the *local* outliers which are located in a lower density region compared to their k nearest neighbors.

In this work, we create two versions of CARE: (1) using the distance based approach AvgKNN (average k nearest neighbor distance of individual data point is used as outlieriness score), and (2) using the most popular density-based approach LOF [3]. We note that CARE is flexible to accommodate any other nearest neighbor based outlier detection algorithm as well.

4.2.2. Feature Bagging. Feature bagging is commonly used in classification ensembles for dimensionality reduction as well as for variance reduction. Like classification ensembles, feature bagging can also be incorporated in outlier ensembles in order to explore multiple subspaces of the data to induce diverse base detectors for high-dimensional outlier detection. As such, in this work we incorporate feature bagging to create multiple base detectors and combine their results to detect outliers with a goal to improve the detection performance by reducing variance. Given a d -dimensional dataset D , for each base detector (either LOF or AvgKNN), we randomly select $q \in [d/2, d - 1]$ features to create b ($= 100$) different feature-bagged base detectors.

4.3. Consensus Approaches

Unlike classification, building an effective ensemble for outlier detection is a challenging task due to the lack of ground truth, which makes it difficult to measure the detector accuracy and combine the results from accurate detectors. Most of the existing approaches either combine outcomes of all the base detectors [12], [24] (hurting the ensemble in presence of poor detectors), or selectively incorporate accurate base detectors in an unsupervised fashion discarding the poor ones [9], [21]. However, the definition of a poor detector varies across different application domains, as some selective approaches are useful for certain applications but not as useful for others. Therefore, in this work we go beyond binary selection and estimate weights for individual base detectors to aggregate their results with a weighted combination. Furthermore, we cumulatively combine the weighted aggregation results for multiple iterations until a stopping condition is satisfied. In the following two sections, we describe the error estimation and weighted aggregation of the base detectors. In Section 4.3.3 the sequential aggregation approach is described and in Section 4.3.4 we introduce a stopping condition for our iterative CARE approach.

4.3.1. Error Estimation. Platanios *et al.* [15] proposed an *unsupervised* approach called Agreement Rates (AR) to estimate errors of multiple classifiers. Motivated by [15], we adapt the unsupervised error estimation of the individual outlier detectors in our work. This estimation is based on the agreement rates for all possible pairs of base detectors in A . Outlier detection can be considered as a binary classification problem with a majority class (inliers = 0) and a minority class (outliers = 1). However, most existing outlier detection algorithms provide outlierness scores for the data points, and not $\{0, 1\}$ labels for them. In order to adapt the AR approach, we calculate the agreement rates for all possible pairs of detectors in A , for which $\{0, 1\}$ labels are needed for the data points. We use Cantelli's inequality [25] to estimate a threshold th_i ($i = 1 \dots b$) with confidence level at 20% to find a cutoff point between inliers (= 0) and outliers (= 1) for each base detector to get a binary list of class labels.

After estimating the class labels, we calculate the agreement rates. As inliers are the majority class, if we take into account all the data points in calculating the agreement rates, it is likely that most values would be large as most detectors often agree on a large number of inliers. Our main goal is to find agreement based on the outliers detected by the base detectors, and ignore a large number of inliers. Therefore, we take the union of all outliers (= 1) across different base detectors to obtain U . Set U contains the important data points (detected as outliers), which we use to calculate the agreement rates for the detector pairs in A .

In the following sections we denote the base detectors as $f_i \in F$ ($i = 1 \dots b, |F| = b$), input data as D , and class labels as Y . The error event E_A of a set of detectors in A is defined as an event when all the detectors make an error:

$$E_A = \bigcap_{i \in A} [f_i(D) \neq Y], \quad (1)$$

where \bigcap denotes set intersection. The error rate of a set of detectors in A is then defined as the probability that all detectors in A make an error together and is denoted as

$$e_A = \mathbb{P}(E_A). \quad (2)$$

The agreement rate of two detectors is the probability that both make an error or neither makes an error. As such, the pairwise agreement rate equation in terms of error rates for the sets in $A : |A| = 2$ can be written as

$$\begin{aligned} a_{\{i,j\}} &= \mathbb{P}(E_{\{i\}} \cap E_{\{j\}}) + \mathbb{P}(\bar{E}_{\{i\}} \cap \bar{E}_{\{j\}}) \\ &= 1 - e_{\{i\}} - e_{\{j\}} + 2e_{\{i,j\}}, \quad \forall \{i,j\} \in A : i \neq j, \end{aligned} \quad (3)$$

where $\bar{\cdot}$ denotes the set complement. On the other hand, the agreement rates for the set of detectors in $A : |A| = 2$ can be directly calculated from the detector output and set U (defined earlier) as follows:

$$a_A = \frac{1}{|U|} \sum_{u=1}^{|U|} \mathbb{I}\{f_i(D_u) = f_j(D_u)\}, \quad \forall \{i,j\} \in A : i \neq j. \quad (4)$$

Provided that one can easily compute the pairwise agreement rates $a_{\{i,j\}}$'s, which can be written in terms of the (unknown) individual and pairwise error rates of the detectors, we can cast the error rate estimation as a constrained optimization problem where the agreement equations in (3) form constraints that must be satisfied as follows:

$$\begin{aligned} \min. \quad & \sum_{\hat{A}:|\hat{A}|\leq 2} e_{\hat{A}}^2 + \epsilon_{\hat{A}} \\ \text{s.t.} \quad & a_A = 1 - e_{\{i\}} - e_{\{j\}} + 2e_{\{i,j\}}, \quad \forall \{i,j\} \in A \quad (5) \\ & 0 \leq e_{\hat{A}} < 0.5 + \epsilon_{\hat{A}}, \\ & 0 \leq \epsilon_{\hat{A}} \end{aligned}$$

where \hat{A} contains individual as well as pairs of detectors (i.e., $\hat{A} = F \cup A$) and $\epsilon_{\hat{A}}$'s denote the slack variables.

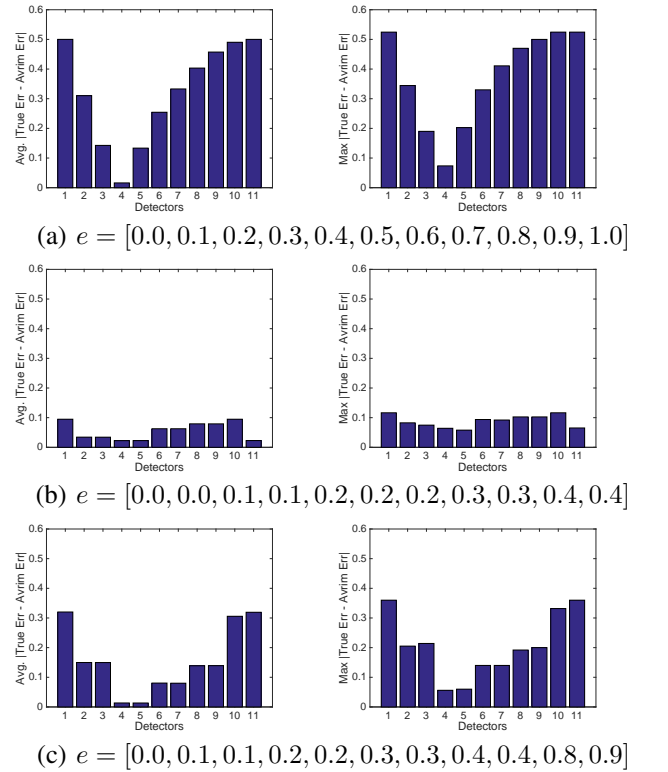


Figure 1. Average (left) and Maximum (right) difference (across datasets) between the true error and the estimated error of different base detectors (e represents true error rates of base detectors). Notice that the differences are high with the presence of many bad detectors, but low otherwise.

In their AR approach, Platanios *et al.* assume that the error rates should be strictly < 0.5 . Different from theirs, we allow the error rates to be above 0.5, for which we introduce a slack variable $\epsilon_{\hat{A}} \geq 0$ in constraints $0 \leq e_{\hat{A}} \leq 0.5 + \epsilon_{\hat{A}}$. In real-world settings, it is possible to have poor base detectors having large errors (i.e., worse than random). Then the question becomes whether the presence of poor detectors (with error ≥ 0.5) hampers the overall estimation of the errors. To answer this question we have designed experiments with synthetic datasets mimicking the real datasets having 1000 data points in total, where 10% of them are outliers and 11 base detectors with different true error rates.

We generate multiple (= 100) snapshots of the synthetic datasets randomly, to analyze results. Figure 1 shows the average and maximum difference between the true error and estimated error of different base detectors. In Figure 1 (a), 6/11 detectors have true error ≥ 0.5 , as a result the average and maximum differences are larger as compared to Figure 1 (b) and (c), where in (b) none and in (c) only 2/11 detectors have errors ≥ 0.5 . We conclude that if all the detectors are good (better than random) or only a few are bad, the optimization in (5) estimates meaningful error rates.

Although the above constrained optimization approach estimates error rates of individual as well as of all possible pairs of base detectors, we only utilize the error rates of the individual detectors to calculate their corresponding weights for aggregation, which we describe next.

4.3.2. Weighted Aggregation. Most commonly used aggregation functions in outlier ensembles are *average* and *maximum*. In most cases, average is preferred over maximum as the latter overestimates the absolute scores. On the other hand, averaging might dilute the final scores with the presence of poor detectors. In CARE, we propose to use *weighted aggregation* to improve the ensemble.

We calculate the weights of the base detectors from their estimated error rates as described in the previous section, such that the weights are positive and inversely proportional to the corresponding errors. Inspired by AdaBoost [11], we employ the error rates of individual detectors to calculate their corresponding weights using the following equation:

$$w_i = \frac{1}{2} \log \left(\frac{2}{e_i} - 1 \right), \quad i = 1 \dots b \quad (6)$$

where $w_i \geq 0$ is the weight of detector i with estimated error $e_i \in [0, 1]$, for $i = 1 \dots b$. Moreover, as we assume that in real-world settings the base detector pool will have poor (i.e., worse than random) detectors, we also incorporate a pruning strategy where we discard the detectors with error $e_i \geq 0.5$. To support the weighted aggregation strategy with pruning, we provide experimental results with synthetic datasets to compare average vs. weighted aggregation as well as pruned vs. un-pruned selection. In Figure 2, we show the distributions of the final ensemble accuracies with different consensus approaches across 1000 samples of a synthetic dataset. Each synthetic data has 1000 points and 4 detectors with true errors $e = [0.2, 0.4, 0.6, 0.8]$ having 5% (left) and 10% (right) of outliers respectively. Here, we calculate the weights and prune the detectors using these estimated errors. We can see from both figures that weighted consensus is better than averaging and pruning is better than un-pruned aggregation as the red curve is more skewed towards the higher accuracy values than the others.

After pruning p detectors with error $e_i \geq 0.5$, we combine the outlieriness scores from the base detectors using weighted aggregation. In order to do weighted aggregation, we need to unify the outlieriness scores, as different base detectors employ different feature sets, hence provide scores with varying range and scale. To standardize, we use Gaussian Scaling [26] to convert the outlieriness scores

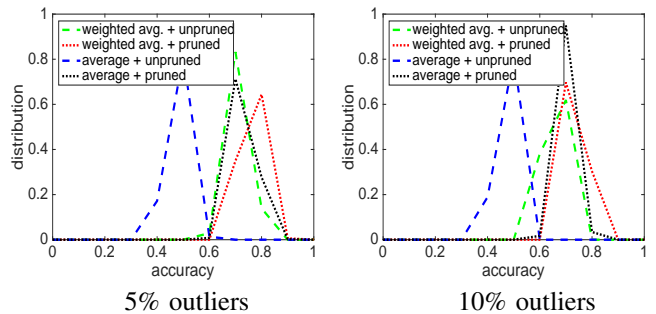


Figure 2. Distribution of accuracies with different consensus approaches. Notice the distribution with pruned weighted aggregation (red curve) is skewed towards higher accuracies.

of AvgKNN or LOF into probability estimates Pr_i ($i = 1 \dots b - p$) $\in [0, 1]$. We calculate the final outlieriness score $ws(x)$ of a data point x using the weighted average of the probability estimates as follows:

$$ws(x) = \frac{\sum_{i=1}^{b-p} w_i \times Pr_i(x)}{\sum_{i=1}^{b-p} w_i} \quad (7)$$

Above, $\sum_{i=1}^{b-p} w_i$ is used to normalize the outlieriness scores. The final scores can be used to sort the data points from most to least outlieriness to produce a ranked list.

Thus far, we described steps 3–7 of Algorithm 1. Next we describe the iterative nature of our sequential ensemble.

4.3.3. Sequential Ensemble. With the weighted aggregation combining multiple feature-bagged base detectors we aim to reduce variance, but our additional goal is to reduce bias. In outlier detection, it is hard to reduce bias in a controlled way, but there exist some successful heuristics to reduce bias. One commonly used approach is to remove outliers in successive iterations [14] in order to build more robust outlier models iteratively. This is a type of sequential ensemble. The basic idea is that the outliers interfere with the creation of a model of normal data, and the removal of points with high outlier scores is beneficial for the model in the following iteration.

As such, we adopt a sequential ensemble approach in CARE where we use the result from the previous iteration to improve the next. In particular, we select a subsample S from the original data D (where $|S| < |D|$) to use it as a *new data model based on which we calculate the outlieriness scores* for all the data points in D . For example, when we need the average k NN distance of a data point $x \in D$, we calculate the distance to its k -nearest neighbors $N_i \in S$. The goal is to construct S that includes as few of the true outliers as possible, such that it serves as a more reliable data model. To do so, we design a sampling approach which we call Filtered Variable Probability Sampling (FVPS). Following are the steps of the FVPS:

- Discard top T outliers detected in previous step from D , where T is the number of outliers selected using Cantelli’s inequality [25] on final outlieriness scores \mathbf{fs} (threshold is selected at 20% confidence level to find the cutoff point between outliers and inliers).

- Select l uniformly at random between $\min\{1 - \frac{T}{n}, \frac{50}{n}\}$ and $\max\{1 - \frac{T}{n}, \frac{1000}{n}\}$, where n is the number of points in the original dataset.
- Build sub-sample S (where $|S| = l \times (n - T)$) by sampling from D' (outliers-discarded) based on the probability of the points being normal (i.e., $(1 - \mathbf{fs})$).

In step 1 of FVPS, we obtain D' by filtering the outliers detected in the previous step to reduce bias and improve the outlier ensemble iteratively. Here, we choose confidence level 20% to get a larger T in order to remove as many outliers as possible. Even though this step might remove some inliers, those should not effect the model as they would have lower probability of being normal points to be removed in the first place. Inspired by Aggarwal and Sathe [1], we adopt variable sampling to select a sample size in step 2. The variable sampling approach has an effect over the parameter choice of the outlier detectors (i.e., k). Varying the subsample size at fixed k effectively varies the percentile value of k in the subsample for different iterations. For some datasets smaller value of k is better, for others larger is better. However, there is no known suitable approach to estimate the correct value of k for a dataset. Therefore, in CARE we select a small value of k (e.g., 5) and employ variable sampling to incorporate the illusion of using different k values in different iterations, which introduces diverse detectors iteratively. After deciding the sample size in step 2, we use probability sampling to create the data model S in step 3. Here, we choose a point from D' to include in S based on its probability of being normal. As a result, we expect to have less interference from the outliers in S as it is mostly built with normal points.

FVPS introduces diverse detectors based on different S in each iteration, hence, we aggregate (e.g. cumulative average) the outlieriness scores \mathbf{ws} over the iterations to compute final scores \mathbf{fs} to further reduce variance and improve the sequential ensemble (step 9 in Alg. 1). Note that \mathbf{fs} is also what FVPS uses for discarding outliers and sampling set S .

4.3.4. Stopping Criterion. We need a proper stopping criterion for the sequential ensemble approach to decide where the iteration should stop and return the final result. As the whole framework is unsupervised, there is no way to use intermediate evaluation to find a stopping point. In CARE, we utilize the pairwise agreement rates a_A between all possible pairs of base detectors to find the stopping point. Experiments reveal a useful strategy: if the distribution of a_A 's is skewed towards higher agreement rates, then the error estimates of the base detectors tend to be more accurate. Intuition is, it is unlikely that most pairs would have high agreement and yet agree on the wrong labels. Therefore, we propose to track the agreement rates and their distribution to decide when to stop iterating CARE.

Specifically, we use the area under the curve (*auc*) of the complementary cumulative distribution function (*ccdf*) of a_A 's as the quantitative measure to decide the stopping point. The *auc* of *ccdf* is large if the distribution of a_A 's is skewed towards higher agreement rates and vice versa. We assume that as CARE sequentially progresses over iterations,

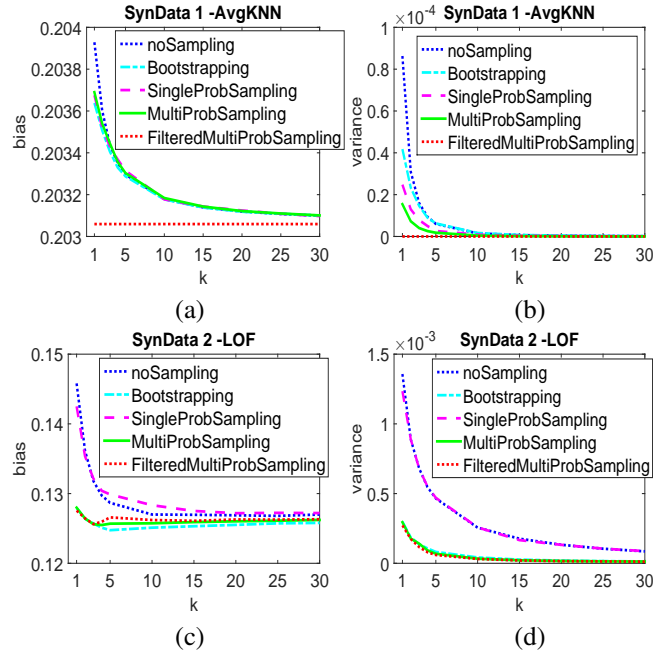


Figure 3. bias (left) and variance (right) vs. k (avg'ed over 10 test datasets) on two synthetic datasets. Notice that our approach (red) w/ probability sampling after top outliers being filtered reduces both bias and variance.

the base detectors improve, and hence the *auc* of *ccdf* for pairwise agreement rates gets larger. However, if at any iteration $t + 1, t \in [0, MAXITER]$, the *auc*($t + 1$) falls below the average by more than the standard deviation of *auc*($0, \dots, t + 1$), the sequential ensemble stops and returns the result at iteration t or otherwise iterates until $MAXITER$ and returns the final result.

5. Reducing Bias and Variance with CARE

According to [1], ensembles with feature-bagged base detectors and with variable sampling tend to reduce variance. In this section, we provide quantitative results through experiments on synthetic datasets to show that filtering top T outliers and probability sampling in our sequential ensemble reduce bias along with variance. To present the bias-variance reduction quantitatively, we design five procedures. For each synthetic dataset, we use a data generation model \mathbb{M} to create R training datasets $D_i, i = 1 \dots R$ of size $m = 210$ (200 inliers and 10 outliers) and 10 test datasets $D_j^{Test}, j = 1 \dots 10$ of size $n = 1000$ by randomly drawing points from \mathbb{M} . Bias and variance of different procedures for different values of k (i.e., # nearest neighbors) for a test data D_j^{Test} are calculated w.r.t. the training data $D'_i, i = 1 \dots R$ sampled from D_i as follows:

$$bias = \sqrt{\frac{\sum_{x=1}^n (f^*(x) - \bar{f}(x))^2}{n}} \quad (8)$$

$$var = \frac{\sum_{x=1}^n \sum_{i=1}^R (f(x, D'_i, k) - \bar{f}(x))^2}{n \times R} \quad (9)$$

Here, $\bar{f}(x) = \frac{\sum_{i=1}^R f(x, D'_i, k)}{R}$, $f^*(x)$ is the actual label

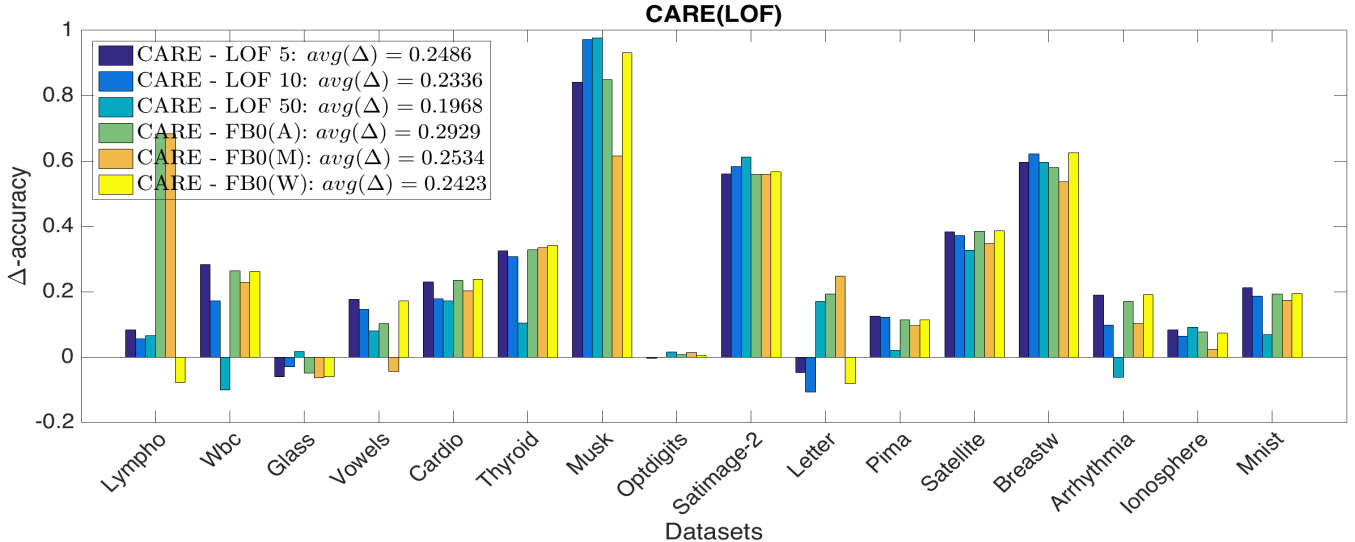


Figure 4. Δ AP (Average Precision) from CARE(LOF) to LOF based baseline approaches on all the datasets. Notice that CARE boosts detection performance significantly for 14/16 datasets over most of the baseline approaches. $avg(\Delta)$ denotes average of Δ AP values across datasets.

of data point $x \in D_j^{Test}$, and $f(x, D'_i, k)$ is the normalized outlierness score of x w.r.t. sampled training set D'_i for k nearest neighbors. For each procedure we design a different approach for sampling D'_i . These five different procedures are briefly described as follows:

(i) *noSampling*: $D'_i = D_i$, (ii) *Bootstrapping*: sampling m times (w/ replacement) from D_i to get D'_i , (iii) *SingleProbSampling*: probability sampling on $f(D_i, D_i, k)$ for a single iteration to get D'_i , (iv) *MultiProbSampling*: probability sampling on $f(D_i, D'_i, k)$ for multiple (i.e. 10) iterations where $D'_i = D_i$ initially, and (v) *FilteredMultiProbSampling*: filtered (top T outliers removed from D_i) probability sampling on $f(D_i, D'_i, k)$ for multiple iterations (i.e. 10) where $D'_i = D_i$ initially (our proposed approach).

In this section, we provide results on only two synthetic datasets (20 dimensional) for brevity, where the inliers are drawn from a mixture of Gaussian distributions and outliers are drawn from (i) power law, and (ii) uniform distribution. Figure 3 shows bias (left) and variance (right) vs. k , where for the top two plots (i.e. (a), (b)) AvgKNN is used to calculate $f(x, D'_i, k)$, and for the bottom two plots (i.e. (c), (d)) LOF is used. We can see from the figure that the curve for MultiProbSampling (green) is below the noSampling (blue) as well as the SingleProbSampling (magenta) curve, showing that probability sampling in multiple iterations helps to reduce both bias and variance. We also see that the FilteredMultiProbSampling (red) reduces bias further thanks to the filtering of top T outliers. Moreover, removing top outliers appears to also reduce variance as the red curve is below all the others in both (b) and (d).

6. Experiments

6.1. Datasets

We evaluate CARE on 16 different real-world outlier detection datasets¹ mostly from the UCI ML repository [2].

1. <http://odds.cs.stonybrook.edu/>

Table 1 provides the summary of the datasets used in this work. The first 9 datasets, Letter dataset, and the following 5 datasets are respectively obtained from [1], [27] and [28].

TABLE 1. REAL-WORLD DATASETS USED FOR EVALUATION, WHERE d IS DATA DIMENSIONALITY, AND % INDICATES THE % OF OUTLIERS.

Dataset	#Pts n	Dim. d	% Outlier Class
Lympho	148	18	classes 1,4 (4.1%)
WBC	278	30	21 sampled malignant class (5.6%)
Glass	214	9	class 6 (4.2%)
Vowels	1456	12	50 sampled class 1 (3.4%), classes 6,7,8, inliers
Cardio	1831	21	176 sampled pathologic (9.6%), normal inliers
Thyroid	3772	6	from [29] (2.5%)
Musk	3062	166	classes 213,211 (3.2%) classes j146,j147,252 inliers
Optdigits	5216	64	150 sampled digit 0 (3%)
Satimage-2	5803	36	71 sampled class 2 (1.2%)
Letter	1600	32	from [27] (6.25%)
Pima	768	8	pos class (35%)
Satellite	6435	36	3 smallest classes (32%)
Breastw	683	9	malignant class (35%)
Arrhythmia	452	274	classes 3,4,5,7,8,9,14,15 (15%)
Ionosphere	351	33	bad class (36%)
Mnist	7603	100	700 sampled digit 6 (9.2%), digit 0 inliers

6.2. Results

6.2.1. CARE vs state-of-the-art baselines. We first compare CARE with simple LOF and AvgKNN based baseline approaches; using $k = \{5, 10, 50\}$, as well as non-sequential feature bagging (FB0) approaches with three types of aggregation; average (A), maximum (M), and weighted (W). Figure 4 shows the Δ Average Precision (AP: area under

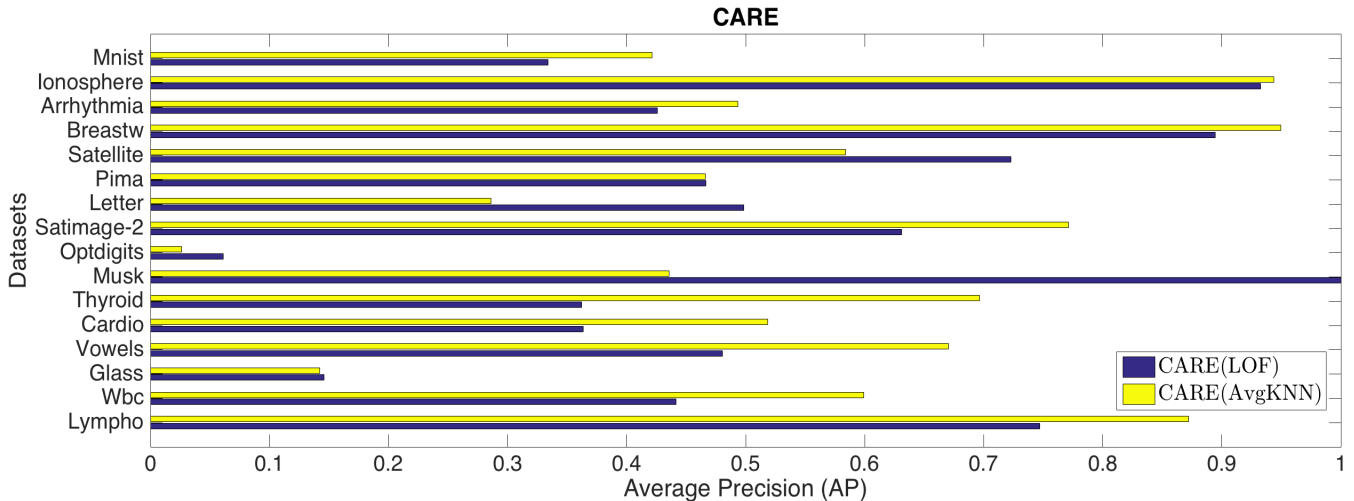


Figure 5. Average Precision (AP) of CARE across datasets for both LOF and AvgKNN based base detectors. CARE(AvgKNN) performs better than CARE(LOF) on 10/16 datasets.

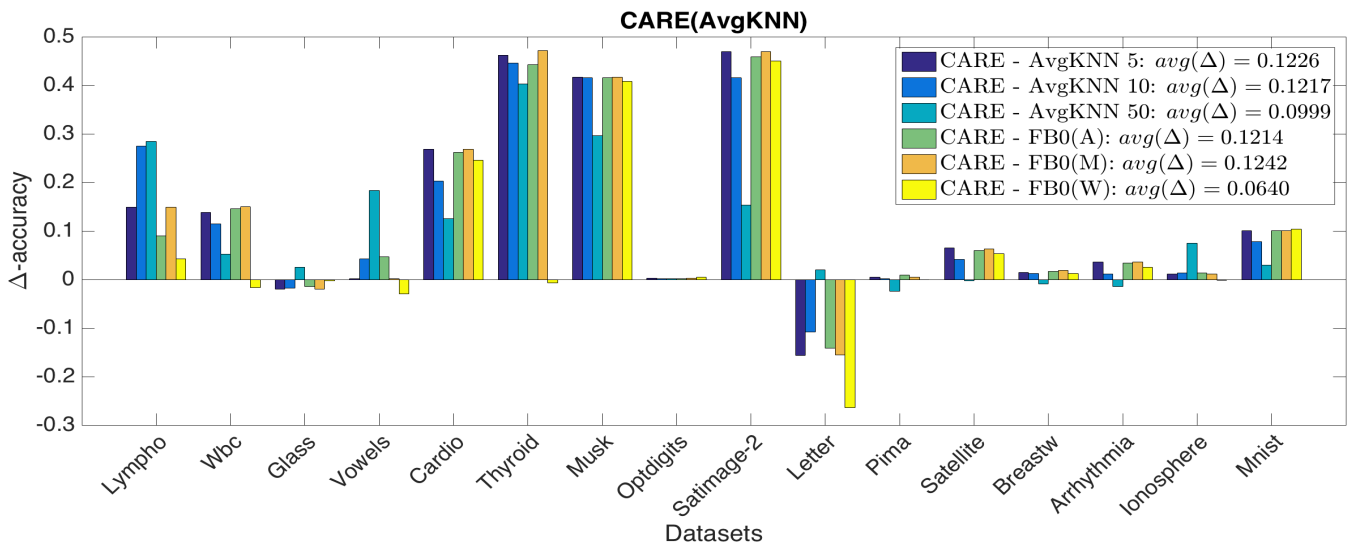


Figure 6. Δ AP values from CARE(AvgKNN) to AvgKNN based baselines. CARE improves over more than half of the baselines on 14/16 datasets.

the precision-recall curve) values from CARE(LOF) to these six state-of-the-art baselines all using the LOF algorithm. That is, the bars depict $AP^{CARE} - AP^{baseline}$. We refer to Figure 5 for the original AP values that CARE(LOF) and CARE(AvgKNN) achieve on the datasets. Results show that CARE outperforms all the base detectors on 9/16 datasets, and more than half of them on 14/16 datasets. Negative Δ values are much smaller as compared to positive ones, which indicates that in cases where CARE is not better than the baselines, it remains close. In the legend of the figure, we provide the overall ΔAP values averaged across all the datasets and positive values indicate that CARE performs better than the individual baselines on average. Similarly, Figure 6 contains the ΔAP values from CARE(AvgKNN) to six baselines, which this time use AvgKNN based subroutines. Again, the average Δ values (in the legend) across different datasets indicate that CARE outperforms the individ-

ual baselines on average. In cases where CARE falls shorter it often remains close to the baselines (note the relatively much smaller negative Δ 's). From these two figures we also conclude that CARE(LOF) provides greater improvement over the baselines compared to CARE(AvgKNN).

6.2.2. CARE vs state-of-the-art ensembles. Next we compare CARE with the existing state-of-the-art outlier ensemble methods, including Aggarwal and Sathé's variable sampling (VS), rotated bagging (RB), and variable rotated bagging (VR) approaches [1], Zimek *et al.*'s subsampling approach [8], as well as the Isolation Forest (iF) ensemble of Liu *et al.* [16]. We employ $b = 100$ base detectors for each of these existing ensemble approaches such that they are comparable with CARE. We present the ΔAP from CARE(LOF) to these six existing state-of-the-art outlier ensembles using the LOF algorithm (except for iF)

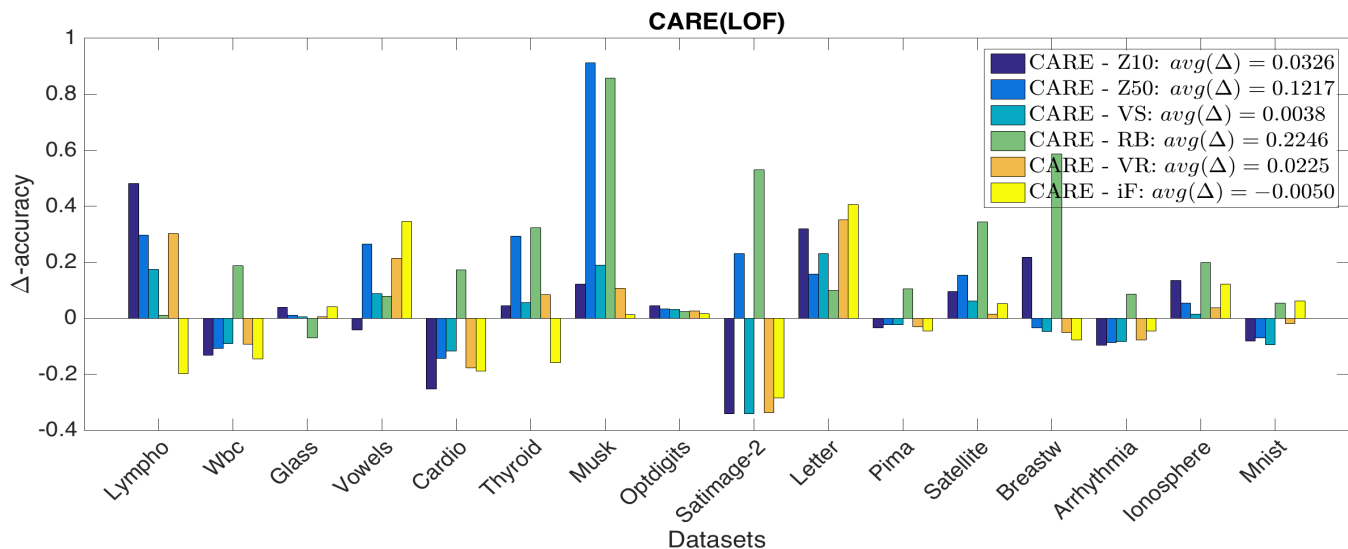


Figure 7. Δ AP from CARE(LOF) to LOF based state-of-the-art ensemble approaches on all the datasets. Notice that CARE outperforms existing ensembles significantly on several datasets and achieves comparable performance otherwise. $avg(\Delta)$'s in the legend denote average of Δ AP values across datasets.

in Figure 7. For Zimek’s subsampling approach we only present the results for sample sizes 10% and 50% (Z10, Z50). In Figure 7, we can see that the performance of CARE(LOF) and VS are close with $avg(\Delta) = 0.0038$ across all the datasets. Notice that CARE mostly improves over Z50 and RB. Although iF is little better than CARE(LOF) with $avg(\Delta) = -0.0050$, for some datasets e.g., Vowels and Letter where iF performs poorly with AP values 0.1341 and 0.0929 respectively, CARE(LOF) provides 2.6 \times improvement with AP value 0.4803 for Vowels, and 4.4 \times improvement with AP value 0.4986 for Letter. Moreover, we note that the magnitude of positive Δ values are larger than the negative ones on average. This indicates that CARE(LOF) provides major improvement in cases when it is the winner and performs similarly to existing ensembles in other cases. Finally, Figure 8 shows the corresponding results for CARE(AvgKNN). Positive average Δ AP values across all datasets show that CARE provides significant improvement when it outperforms an existing ensemble and falls short by a small margin in other cases. iF outperforms CARE significantly on Musk, which has a dense cluster of outliers that avoid detection by nearest neighbor based methods. We also find that none of the existing methods on Optdigits and Glass, where further investigation is needed to understand the type of outliers that they exhibit.

7. Conclusion

In this paper, we proposed CARE, a new sequential ensemble approach for outlier mining with a goal to achieve low detection error through reduced variance and bias. Two main components of CARE are its parallel and sequential building blocks. The former helps reduce variance by a weighted combination of multiple base detectors. Detector weights are derived from their error rates that are estimated through their relation to pairwise agreement rates. On the other hand, the sequential component

is designed to reduce bias. It utilizes results from previous iterations and a new sampling strategy FVPS to weed out top outliers so as to construct a more robust data model based on which outlieriness scores are computed. We evaluate our method on 16 real-world datasets. Extensive experiments validate that CARE provides significant improvement over the baseline methods as well as the state-of-the-art outlier ensembles when it is the winner and performs close enough otherwise. All source codes of our methods are shared openly at <http://shebuti.com/sequential-ensemble-learning-for-outlier-detection/>.

Acknowledgment

This material is based upon work supported by the ARO Young Investigator Program under Contract No. W911NF-14-1-0029, NSF CAREER 1452425, IIS 1408287 and IIP1069147, DARPA Transparent Computing Program under Contract No. FA8650-15-C-7561, a Facebook Faculty Gift, an R&D grant from Northrop Grumman Aerospace Systems, and Stony Brook University Office of Vice President for Research. Any conclusions expressed in this material are of the authors’ and do not necessarily reflect the views, either expressed or implied, of the funding parties.

References

- [1] C. C. Aggarwal and S. Sathe, “Theoretical foundations and algorithms for outlier ensembles.” *ACM SIGKDD Explorations Newsletter*, vol. 17, no. 1, pp. 24–47, 2015.
- [2] M. Lichman, “UCI machine learning repository,” 2013. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [3] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, “Lof: identifying density-based local outliers,” in *ACM sigmod record*, vol. 29, no. 2. ACM, 2000, pp. 93–104.
- [4] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek, “Loop: local outlier probabilities,” in *CIKM*. ACM, 2009, pp. 1649–1652.
- [5] S. Papadimitriou, H. Kitagawa, P. B. Gibbons, and C. Faloutsos, “LocI: Fast outlier detection using the local correlation integral,” in *ICDE*. IEEE, 2003, pp. 315–326.

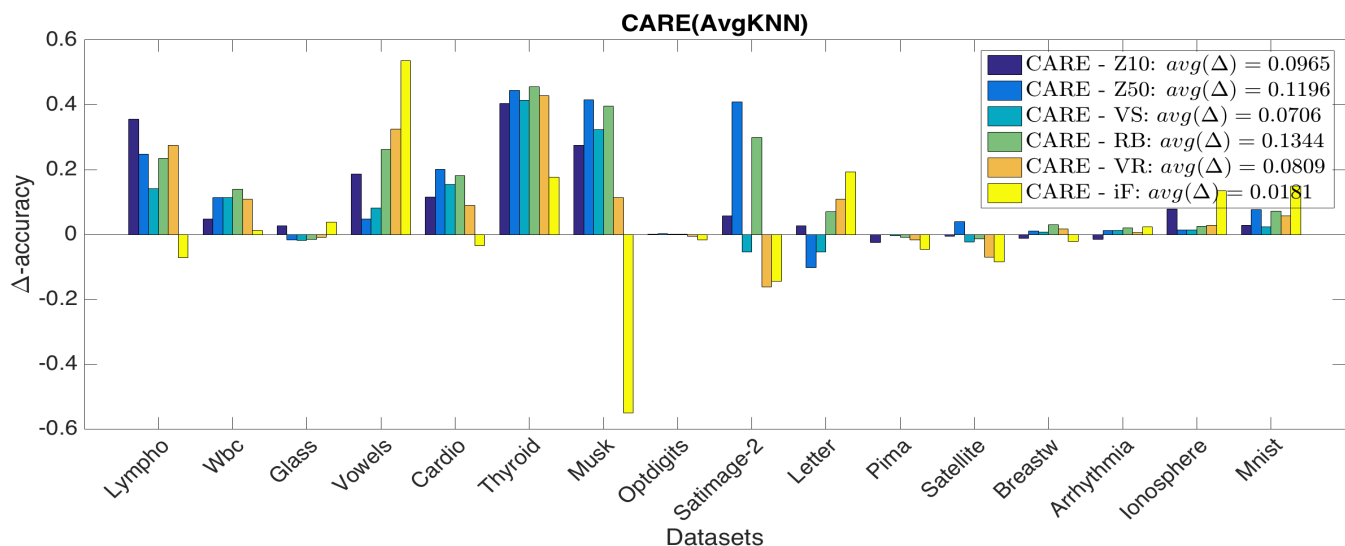


Figure 8. Δ AP values from CARE(AvgKNN) to AvgKNN based state-of-the-art ensemble approaches. Note the generally large positive and otherwise small negative values across datasets.

- [6] K. Zhang, M. Hutter, and H. Jin, "A new local distance-based outlier detection approach for scattered real-world data," in *Advances in Knowledge Discovery and Data Mining*, 2009, pp. 813–822.
- [7] M. E. Otey, A. Ghoting, and S. Parthasarathy, "Fast distributed outlier detection in mixed-attribute data sets," *Data Mining and Knowledge Discovery*, vol. 12, no. 2-3, pp. 203–228, 2006.
- [8] A. Zimek, M. Gaudet, R. J. Campello, and J. Sander, "Subsampling for efficient and effective unsupervised outlier detection ensembles," in *ACM SIGKDD*, 2013, pp. 428–436.
- [9] S. Rayana and L. Akoglu, "Less is more: Building selective anomaly ensembles with application to event detection in temporal graphs." *SDM*, vol. 17, 2015.
- [10] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [11] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *Journal of computer and system sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [12] A. Lazarevic and V. Kumar, "Feature bagging for outlier detection," in *ACM SIGKDD*, 2005, pp. 157–166.
- [13] J. Gao and P.-N. Tan, "Converting output scores from outlier detection algorithms into probability estimates," in *ICDM*, 2006, pp. 212–221.
- [14] C. C. Aggarwal, "Outlier ensembles: position paper." *SIGKDD Explor. Newsl.*, vol. 14, no. 2, pp. 49–58, 2012.
- [15] A. Platanios, A. Blum, and T. M. Mitchell, "Estimating accuracy from unlabeled data," in *In Proceedings of UAI*, 2014.
- [16] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in *ICDM*. IEEE, 2008, pp. 413–422.
- [17] L. Breiman, "Using adaptive bagging to debias regressions," Statistics Dept. UCB, Tech. Rep., 1999.
- [18] Y. Sun, M. S. Kamel, A. K. Wong, and Y. Wang, "Cost-sensitive boosting for classification of imbalanced data," *Pattern Recognition*, vol. 40, no. 12, pp. 3358–3378, 2007.
- [19] L. Breiman, "Random forests," *Machine Learning*, vol. 45, pp. 5–32, 2001.
- [20] C. C. Aggarwal, "Outlier ensembles: position paper," *ACM SIGKDD Explorations Newsletter*, vol. 14, no. 2, pp. 49–58, 2013.
- [21] A. Zimek, R. J. Campello, and J. Sander, "Ensembles for unsupervised outlier detection: Challenges and research questions," *SIGKDD Explor. Newsl.*, vol. 15, no. 1, pp. 11–22, 2013.
- [22] D. M. Hawkins, *Identification of outliers*. Springer, 1980, vol. 11.
- [23] E. M. Knorr and R. T. Ng, "A unified notion of outliers: Properties and computation." in *KDD*, 1997, pp. 219–222.
- [24] S. Rayana and L. Akoglu, "An ensemble approach for event detection in dynamic graphs." in *ACM SIGKDD ODD² Workshop*, 2014.
- [25] G. Grimmett and D. Stirzaker, *Probability and Random Processes.*, 2001.
- [26] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek, "Interpreting and unifying outlier scores." in *SDM*, 2011.
- [27] B. Mícenková, B. McWilliams, and I. Assent, "Learning outlier ensembles: The best of both worlds—supervised and unsupervised," in *ACM SIGKDD ODD² Workshop*, 2014.
- [28] K. Ting, S. Tan, and F. Liu, "Mass: A new ranking measure for anomaly detection," *Gippsland School of Information Technology, Monash University*, 2009.
- [29] F. Keller, E. Müller, and K. Böhm, "Hics: high contrast subspaces for density-based outlier ranking," in *Data Engineering (ICDE), 2012 IEEE 28th International Conference on*. IEEE, 2012, pp. 1037–1048.