

DISTRIBUTED ANOMALY DETECTION USING SATELLITE DATA FUSION

KANISHKA BHADURI*, KAMALIKA DAS**, AND PETR VOTAVA***

ABSTRACT. There has been a tremendous increase in the volume of Earth science data over the last decade from modern satellites, in-situ sensors and different climate models. All these datasets need to be co-analyzed to find interesting patterns or to search for extremes or outliers. Information extraction from such rich data sources using advanced data mining methodologies is a challenging task not only due to the massive volume of data, but also because these datasets are physically stored at different geographical locations. Moving these petabytes of data over the network to a single location may waste a lot of bandwidth, and can take days to finish. To solve this problem, in this paper, we present a novel algorithm which can identify outliers in the global data without moving all the data to one location. The algorithm is highly accurate (close to 99%) and requires centralizing less than 5% of the entire dataset. We demonstrate the performance of the algorithm using data obtained from the NASA MODerate-resolution Imaging Spectroradiometer (MODIS) satellite images.

1. INTRODUCTION

The interest in Earth sciences has been growing steadily over the last two decades and together with it there has been a tremendous increase in the volume of Earth science data collected and generated by growing number of satellites, in-situ sensors and increasingly complex ecosystem and climate models. This growth in volume and complexity is going to continue because in order for the scientists to better understand and predict the Earth system processes, they will require far more comprehensive data sets spanning many years and more complex models. With the launch of NASA's Terra and Aqua missions, and the expected launches of number of missions recommended by the Decadal Survey, the need for more efficient and scalable data processing system is crucial. The volume of data itself is often a limiting factor in obtaining the information needed by the scientists and decision makers. This data volume will grow from hundreds of terabytes to tens of petabytes throughout the lifespan of the proposed Decadal Survey missions. More data means more information, only if there are sophisticated means of sifting through the data for extracting the relevant information from this data avalanche.

A very interesting task relevant to the Earth science community is especially interested in identification of anomalies within the ecosystems (*e.g.* wildfires, droughts, floods, insect/pest damage, wind damage, logging), so that experts can then focus their analysis efforts on the identified areas. There are dozens of variables that define the health of the ecosystem and both long-term and short-term changes in these variables can serve as early indicators of natural disasters and shifts in climate and ecosystem health. These changes can have profound socio-economic impacts and it is important to develop capabilities for identification, analysis and response to these changes in a timely manner. In order to fully understand the Earth systems, scientists need to be able to analyze together a number of datasets from satellites, ground sensors and models. Every data component has a different observation or predictive capability and therefore a global analysis on a combination of modalities gives better results than studying a particular feature. For example, observing different but related phenomena, predicting climate impacts at different timesteps, or providing observations of the same

*MCT Inc, MS 269-1, NASA Ames Research Center, Moffett Field, CA-94035. Kanishka.Bhaduri-1@nasa.gov

**SGT Inc, MS 269-3, NASA Ames Research Center, Moffett Field, CA-94035. Kamalika.Das@nasa.gov

***CSU Monterey Bay, NASA Ames Research Center, Moffett Field, CA-94035. Petr.Votava-1@nasa.gov.

phenomena through different means, such as ground sensor or a radar are expected to enable better comprehension and more accurate characterization of changes and disturbances in Earth systems.

The situation is greatly complicated by the fact most of the data representing different modalities are stored at geographically distributed archives, such as NASA’s Distributed Active Archive Centers (DAAC), each containing data specific to only a subset of the scientific community and thus it is almost impossible to perform a globally consistent analysis. Given this scenario, the current approach would be for the scientist to look at only a subset of the dataset available at one site (and thereby compromise on the quality of the results) or to bring all the data together in one place and then perform the analysis. While the second approach works for lower data volumes, it is not feasible to centralize all the data when it grows beyond what can be gathered using current network infrastructure in a timely manner. Another reason why complete centralization is not possible is because the research is done in number of different science teams and organizations in different countries. While there is a trend to consolidate more data at fewer data centers, the capabilities to extract vital information from large distributed datasets will continue to be a key for the Earth science community to be able to gather significant results by analyzing the growing data volumes being accumulated world wide.

In this paper we describe a novel and efficient algorithm for anomaly detection in distributed earth science databases. The contributions of this work based on the state of the art in distributed anomaly detection can be enumerated as:

- To the best of the authors’ knowledge, this is the first algorithm that can scale to terabytes of data when the data is distributed across several sites, with a subset of features only at each site. In the distributed data mining literature this is known as the vertically partitioned scenario.
- For the proposed algorithm, the amount of communication required is less than 1% of that required for centralization, yet is 99% accurate compared to a centralized algorithm in finding the outliers. The accuracy is a function of the data percentage communicated and can be tuned based on the requirement.
- The algorithm is potentially capable of detecting new outliers which are missed by using only a subset of features.

The rest of the paper is organized as follows. In the next section (Section 2) we present the work related to this area of research. We discuss the notations and the one class SVM formulation in Section 3. In Section 4 we present details about the proposed algorithm. We discuss the theoretical analysis of the algorithm in Section 5. Performance of the algorithm on NASA satellite data is presented in Section 6. Finally we conclude the paper in Section 7.

2. RELATED WORK

Outlier or anomaly detection refers to the task of identifying abnormal or inconsistent patterns from a dataset. While outliers may seem as undesirable entities in a dataset, identifying them have many potential applications such as in fraud and intrusion detection, financial market analysis, medical research and safety-critical vehicle health management. Broadly speaking, outliers can be detected using either *supervised* or *semi-supervised* or *unsupervised* techniques [10] [7]. Unsupervised techniques, as the name suggest, do not require labeled instances for detecting outliers. In this category, the most popular ones are the distance-based and density based techniques. The basic idea of these techniques is that outliers are points in low density regions or those which are far from other points. In their seminal work, Knorr *et al.* [12] proposed a distance-based outlier detection technique based on the idea of nearest neighbors. The naive solution has a quadratic time complexity since every data point needs to be compared to every other to find the nearest neighbors. To overcome this, researchers have proposed several techniques such as the work by Angiulli and Pizzuti [1], Ramaswamy *et al.* [14], and Bay and Schwabacher [3]. Density-based outlier detection schemes, on the other hand, flag a point as an outlier if the point is in a low density

region. The density of a point can be evaluated using several techniques such as the ones proposed in [9]. Supervised techniques require labeled instances of both normal and abnormal operation data for first building a model (*e.g.* a classifier) and then testing if an unknown data point is a normal one or an outlier. The model can be a probabilistic such as Bayesian inference [8] or deterministic models such as decision trees, support vector machines and neural networks [11]. This task is similar to building predictive modeling with two or more classes. Semi-supervised techniques only require labeled instances of normal data. Hence they are more widely applicable than the fully supervised ones. These techniques build models of normal data and then flag as outliers all those points which do not fit the model.

There exists a plethora of work for outlier detection from spatio-temporal databases. Barua and Alhajj [2] present a technique for outlier detection from meteorological data using a parallel implementation of the well-known wavelet transformation. The authors show that by implementing the algorithm on modern high performance multi-core processors, they achieve both improved speedup and accuracy. Birant and Kut [5] discuss a way of identifying both spatial and temporal outliers in large databases. They argue that existing methods do not identify both these outliers, and hence they propose a new DBSCAN clustering method to first cluster the dataset based on the density of points and then tags as outliers all points which have low density in its neighborhood. Now depending upon the type of outlier detected, either spatial or temporal neighborhood is considered. Both these methods consider outliers as single points. In practice, there may be a group of points which are outliers *e.g.* a tornado or other natural disaster affecting a large area. Zhao *et al.* [19] present an outlier detection method based on wavelet transformation which can detect region outliers. In their approach, they first transform the image to the wavelet domain and then isolate those coefficients which are greater than a threshold. Inverse wavelet transformation on this thresholded pixels are then candidates for outliers which are further filtered by running an outlier detection method. The work which is closest to this work is the land cover change detection by Boriah *et al.* [6] and Poter *et al.* [13]. In [6], the authors proposed a recursive merging algorithm for change point detection. For more details on the recent work on change detection for land cover data, readers are referred to [6] and the references therein. Several other techniques also exist for building classification and prediction models for mining geospatial data such as [17].

However, none of these techniques can handle the data when it is distributed across a large number of sites. Our proposed algorithm can perform anomaly detection without centralizing all the data to one location.

3. BACKGROUND

In this section we first define the notations and then discuss ν 1-class SVM which forms a building block for our distributed anomaly detection technique.

3.1. Notations. Let P_1, \dots, P_p be a set of nodes where P_0 is designated as the master node and the others are denoted as the computational nodes. Let the dataset at node P_i ($\forall i > 0$) be denoted by $D_i = \begin{bmatrix} \overrightarrow{x_1^{(i)}} & \dots & \overrightarrow{x_m^{(i)}} \end{bmatrix}^T$ of m rows where $\overrightarrow{x_j^{(i)}} \in \mathbb{R}^{n_i}$. Here each row corresponding to an observation and each column corresponds to a feature/attribute/sensor measurement. In the distributed data mining literature, this is referred to as the vertically partitioned data distribution scenario. The global set of features (n) is the vertical concatenation of all the features over all nodes and is defined as $n = [n_1 \ n_2 \ \dots \ n_p]$ (using Matlab notations). Hence, the global data D is the $m \times n$ matrix defined as the union of all data over all nodes *i.e.* $D = [\overrightarrow{x_1} \ \dots \ \overrightarrow{x_m}]^T$ with $\overrightarrow{x_j} \in \mathbb{R}^n$. Note that, here we make the implicit assumption that the ℓ -th row of all the sites corresponds to the ℓ -th observation *i.e.* the observations have been cross-matched. There exist techniques in the literature

to perform this cross-matching using domain knowledge *e.g.* the cross-match tool of Sloan Digital Sky Survey¹.

Let \mathcal{O}_i denote the set of local outliers at node P_i , detected by an outlier detection algorithm running on D_i such that $|\mathcal{O}_i| < |D_i|$. We give a precise definition of outlier and an algorithm to detect these in the next section. The global set of outliers found by a centralized algorithm having access to all the data is denoted analogously by the set \mathcal{O}_c . The set of outliers found by the distributed algorithm is denoted by \mathcal{O}_d .

3.2. One class ν -SVM. Given a training dataset containing two classes of examples, one class SVMs, introduced by Schölkopf [15], is a supervised learning method for drawing a separating hyperplane that separates these two classes. In our discussion, we will refer to positive labeled data points as normal and negative label data points as outliers. Instead of using both types of examples from the training data for constructing the hyperplane, one class SVM uses only instances with positive labels to do the same. It also uses a parameter ν which denotes the maximum allowance of outliers in the training data. During the training phase, the SVM algorithm optimizes the placement of the hyperplane in order to maximize the margin between the hyperplane and the origin, which is the lone representative of the second class with negative label.

In many cases, the decision boundary is non-linear in the input space and the trick is to transform the input data to a higher dimension space; the latter allowing for linear separability. This mapping is often made implicit using a well-known kernel function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ which actually computes the inner product between the input vectors in this (possibly) infinite dimensional space. Throughout this paper, we have used Radial Basis Function (RBF) kernel:

$$(1) \quad k(\vec{x}_i, \vec{x}_j) = \exp\left(\frac{-\|\vec{x}_i - \vec{x}_j\|^2}{2\sigma^2}\right)$$

where $\|\cdot\|$ denotes the Euclidean norm and σ defines the kernel width. σ is often needs to be tuned for a particular dataset.

Schölkopf [15] showed that in the high dimensional feature space it is possible to construct an optimal hyperplane by maximizing the margin between the origin and the hyperplane in the feature space by solving the following optimization problem,

$$(2) \quad \begin{aligned} \text{minimize} \quad & Q = \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) + \rho \left(\nu m - \sum_i \alpha_i \right) \\ \text{subject to} \quad & 0 \leq \alpha_i \leq 1, \quad \nu \in [0, 1] \end{aligned}$$

where α_i and β_i are Lagrangian multipliers, ν is an user specified parameter that defines the upper bound on the fraction of the training error and also the lower bound on the fraction of support vectors, and ρ is the offset of the hyperplane from the origin. The optimal solution returns a set of points SV from the training set known as the *support vectors* for which the $0 < \alpha_i \leq 1$ and also the value of the bias term ρ . Now, for any test point \vec{x}_t , not in the training set, the optimal decision is based on the following inner product computation:

$$(3) \quad f(\vec{x}_t) = \sum_{i \in SV} \alpha_i k(\vec{x}_i, \vec{x}_t) - \rho$$

The point \vec{x}_t is deemed as an outlier if $f(\vec{x}_t) < 0$.

3.3. Overview of algorithm. The distributed outlier detection algorithm that we have developed achieves two tasks. First, it finds a correct set of outliers when compared to a centralized execution. Then it tries to reduce the communication cost of centralization. Both these tasks are achieved by using a prune rule which states that a multi-dimensional point is an outlier, if at least one dimension

¹<http://cas.sdss.org/astrodr6/en/tools/crossid/upload.asp>

is out of bounds. This reduces the communication cost dramatically since, from each site, we only need to test those points which are local outliers. The steps of the algorithm are as follows:

- (1) Run a local anomaly detection algorithm at each of the local sites
- (2) Centralize all the local outliers at the master site
- (3) Collect a sample from all the sites to build a global outlier detector
- (4) Test all the local outliers from all the local sites against the global outlier detector

We elaborate on each of these steps in the next section.

4. ALGORITHM DETAILS

4.1. Pruning rule. As stated earlier, the goal of distributed outlier detection is two-fold: (1) compute the correct set of outliers (with respect to a centralized execution) and (2) minimize the cost to communicate all of the data to a central node. Distributed algorithms often define rules based on the data which minimize communication and also guarantees that the global task is accomplished [4][18][16]. These data dependent rules are such that, if satisfied by all nodes independently, then certain global properties of the dataset hold. As a result, each node can stop communicating messages, as soon as the pruning rule is satisfied for that node.

In this paper we use the following observation to prune the number of messages that need to be sent to the master site for determining the global set of outliers:

Pruning rule: *An observation $\vec{x} \in D$ is a global outlier (with respect to all the features) i.e. $\vec{x} \in \mathcal{O}_d$, if it is an outlier with respect to at least one (or a subset) of the features i.e. $\exists j \in \{1 \dots p\}, x^{(j)} \in \mathcal{O}_j$.*

While this statement may not be true in general, it provides us with a way of pruning the number of observations that needs to be sent to the central site. In our experiments with the NASA Earth Sciences climate data, we have found that this simple pruning strategy can detect more than 99% of the outliers that a centralized execution would find with less than 1% of the communication cost required for centralization. Figure 1 points out the intuition in 2-d. In this figure, the green dots represent the normal points while a single red dot represents the anomalous point. As seen, the red dot is quite far from the green dots based on the actual distance between the red point and the nearest green point. We argue that in order for this to happen, the distance along at least one of the axes will be large. In other words, most of the global outliers will be a local outlier in at least one of the distributed sites. We validate this statement in our experiments using the NASA Earth sciences datasets.

4.2. Detailed description. The overall distributed anomaly detection algorithm consists of two stages. The pseudo code for the first step is shown in Alg. 1. In this step, each node computes the local outliers independently. The input to this local step are the dataset at each node D_i , the size of training set T_s , a seed s of the random number generator, and the parameter ν . The algorithm first sets the seed of the random number generator to s . Then it selects a sample of size T_s from D_i and uses it as the training set T_i . The rest is used for the testing phase H_i . It then builds an SVM model using T_i and ν . Once the model has been built, all points in H_i are tested using M as the set of support vectors. All those elements in H_i whose score is negative is returned as the set of outliers \mathcal{O}_i .

In the second phase (Alg. 2), the local outliers are aggregated to the master site P_0 . A sample of size T_s is drawn from each of the local sites D_i such that the same index (observation) is selected from each node. A global SVM model is then learned on this aggregated sample from each sites. Each element of $\bigcup_{i=1}^p \mathcal{O}_i$ is tested against this global model to assign a score. All those elements in $\bigcup_{i=1}^p \mathcal{O}_i$ whose score is less than 0 is then reported as the true set of outliers \mathcal{O}_d by the distributed algorithm.

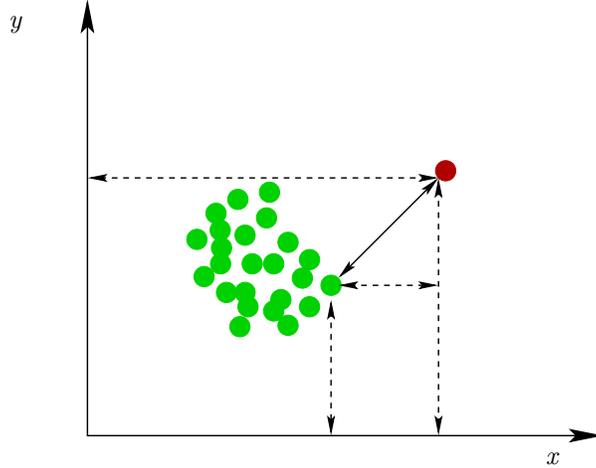


FIGURE 1. This figure shows the basic idea of the pruning rule in 2-d. In this figure, the green dots represent the normal points while a single red dot represents the anomalous point. As seen, the red dot is far away from the green dots. The true distance between the red dot and the closest green dot is show by a bold arrow. The distance along the axes are shown using dotted lines. The observation is that for any true outlier, far away from any of the normal points, the distance along the axes will also be higher. Hence we can only analyze the local outliers from each site.

Algorithm 1: Local outlier detection at each node $P_i, i > 0$

Input: Dataset(D_i), Training sample size(T_s), ν , seed s

Output: Outlier set \mathcal{O}_i

begin

 setseed(s);

$T_i = \text{Sample}(D_i, T_s)$;

 // Training data

$H_i = D_i \setminus T_i$;

 // Test data

$M = \text{SVMTraining}(T_i, \nu)$;

$S = \text{SVMTest}(M, H_i)$;

 // Assign a score to each point in H_i

for $j=1$ to $|H_i|$ **do**

if $S(j) < 0$ **then**

$\mathcal{O}_i(j) \leftarrow [H_i(j) \ S(j)]$;

 Send \mathcal{O}_i to P_0 ;

end

5. ANALYSIS OF ALGORITHM

In this section we provide an analysis of the distributed algorithm.

5.1. Correctness. Correctness of our proposed distributed anomaly detection algorithm is based on the prune rule. A globally correct prune rule guarantees global correctness.

Figure 2 shows a scenario of the algorithm execution in \mathbb{R}^2 . The red and the green dots depict the two classes in the dataset. Hyper-plane A is constructed when both dimension x and y are considered. On the other hand, hyper planes B and C are constructed when only the y and x coordinates are considered separately. Recall that all points that are closer to the origin are denoted

Algorithm 2: Global outlier detection at P_0

Input: $\mathcal{O}_1, \dots, \mathcal{O}_p$, Training sample size(T_s), ν

begin

```
   $T = \text{Sample}(\bigcup_{i=1}^p D_i, T_s);$            // Training data sampled from all the sites
   $H = \bigcup_{i=1}^p \mathcal{O}_i;$                        // Test data
   $M = \text{SVMTraining}(T, \nu);$ 
   $S = \text{SVMTest}(M, H);$                        // Assign a score to each point in  $H$ 
  for  $j=1$  to  $|H|$  do
    if  $S(j) < 0$  then
       $\mathcal{O}_d(j) \leftarrow [H(j) \ S(j)];$ 
  end
```

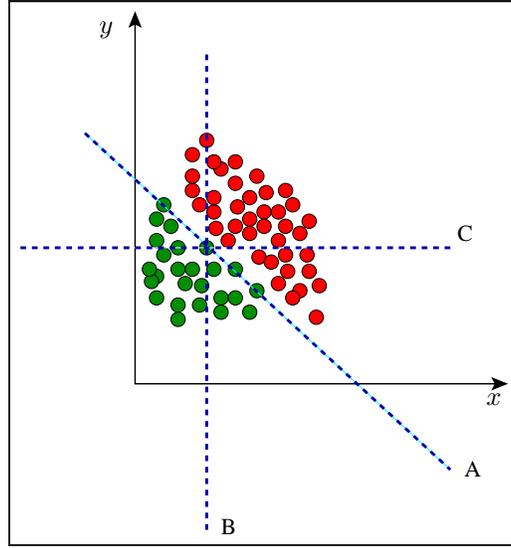


FIGURE 2. This figure shows the different hyper planes drawn by the algorithm when using all the variables (A), only y -dimension values (B) and only x -dimension values (C). Note that different anomalies are found using the different hyper-planes.

as outliers. The sets of outliers that are detected by each of these hyperplanes are not identical. However, the outliers that are closest to the origin (and hence most anomalous points) are detected by all these hyper planes. The missing ones are the boundary outliers, and hence they may offer less value when detected as anomalous points.

5.2. Message complexity. The total number of bytes necessary to centralize all of the data at a single location and run the centralized outlier detection algorithm is:

$$m \times n_1 + m \times n_2 + \dots + m \times n_p = m \times \sum_{i=1}^p n_i$$

For the distributed algorithm, we perform two rounds of communication. First, we centralize the outliers from all the sites and then we gather samples of size T_s from all of them to build a global model and test the outliers found by each of the local sites. The total number of messages is given

by,

$$\underbrace{|\mathcal{O}_1| \times n_1 + |\mathcal{O}_2| \times n_2 + \dots + |\mathcal{O}_p| \times n_p}_{\text{centralizing outliers}} + \underbrace{T_s \times n_1 + \dots + T_s \times n_p}_{\text{centralizing samples}} = \sum_{i=1}^p |\mathcal{O}_i| \times n_i + T_s \sum_{i=1}^p n_i$$

Now since $m \gg \sum_{i=1}^p |\mathcal{O}_i| + T_s$, the distributed algorithm is far more communication efficient than its centralized counterpart. We demonstrate this empirically in Section 6.

5.3. Running time. The running time for the traditional ν -SVM can be written as $O(m^2 \sum_{i=1}^p n_i)$ or $O(m(\sum_{i=1}^p n_i)^2)$, depending on the solution to the primal or the dual problem. In either of these two cases, distributed computing can reduce the running time by splitting n_i across several nodes. Therefore, the load at one node can be reduced from $O(m^2 \sum_{i=1}^d n_i)$ or $O(m(\sum_{i=1}^p n_i)^2)$ to $O(m^2 n_i)$ or $O(mn_i^2)$. When the number of features *i.e.* $\sum_{i=1}^p n_i$ is large, this can provide significant savings in terms of computational complexity at each node. We demonstrate this in the experimental section.

6. EXPERIMENTAL EVALUATION

In this section we describe the experimental evaluation of the algorithm.

6.1. Dataset description. The dataset used in paper is the MODerate-resolution Imaging Spectroradiometer (MODIS) Reflectance product MCD43A4 (version 5) which provides 500-meter reflectance data adjusted using a bidirectional reflectance distribution function (BRDF). The data is stored at intervals of every 8 days as an image file of size 1203×738 where each entry is saved as little-endian 32-bit float value. Each image is saved in 7 separate bands at different wavelengths. Along with the actual reflectance data for each pixel, we also have the latitude and longitude information for them. At the top level, the data is organized by year from 2001 to 2008. Under this top level directory structure are separate files for each band (1 - 7) and each 8-day period of the particular year. Within the period the best observations were selected for each location. Each of the files represent a 2D dataset with the naming conventions as follows:

$$MCD43A4.CA_1KM.005.<YYYYDDD>.<BAND>.flt32$$

where $<YYYYDDD>$ is the beginning year-day of the period and $<BAND>$ represents the observations in particular (spectral) band (band 1 - band 7). The indexing is 0-based, ranging from 0 - 6 (where 0 = band 1, and 6 = band 7). The spectral band frequencies for the MODIS acquisition are as follows (see Table 1):

BAND1	620 - 670 nm
BAND2	841 - 876 nm
BAND3	459 - 479 nm
BAND4	545 - 565 nm
BAND5	1230 - 1250 nm
BAND6	1628 - 1652 nm
BAND7	2105 - 2155 nm

TABLE 1. Spectral band frequencies for MODIS data acquisition.

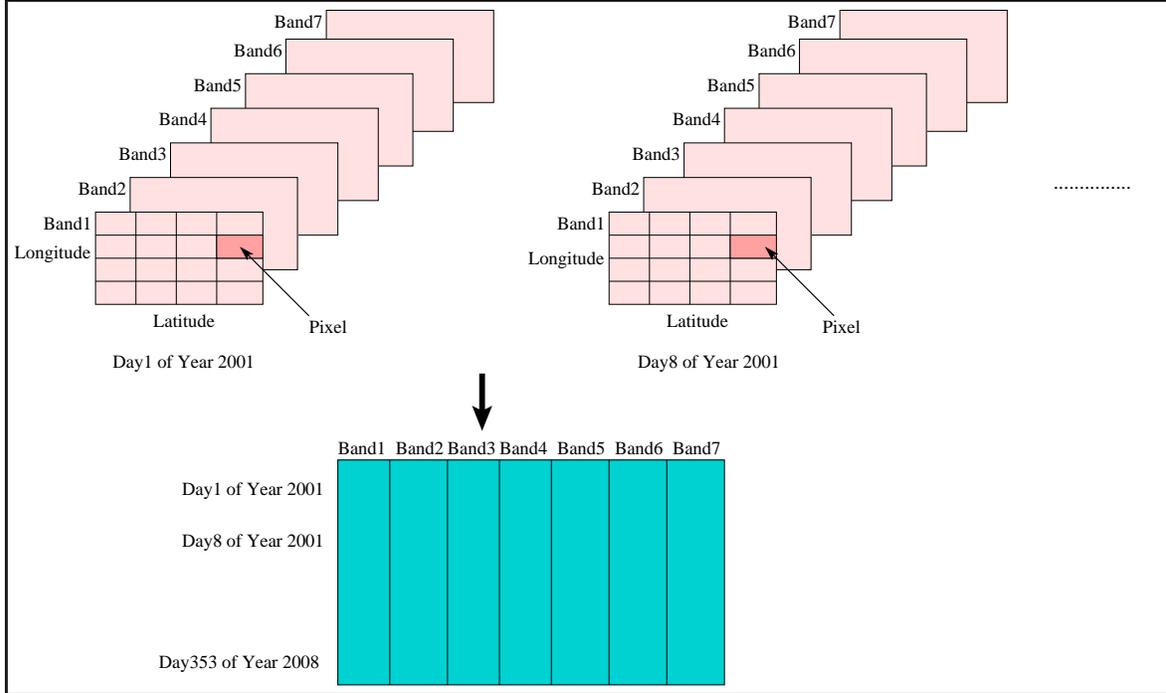


FIGURE 3. This figure shows how the data set is structured. Each file is an image of size 1203×738 . There are seven bands (separate images) for each of the 46 days per year (over 8 years), since data is saved every 8th day. The data contains of both the intensity and the latitude and longitudes of the locations. We arrange the data in a matrix format with each row corresponding to time points and each column corresponds to a band. This technique essentially 'flattens' the 2-D image data into a 1-D representation.

6.2. Dataset preparation. In order to apply our anomaly detection method, we have performed the following preprocessing steps:

- From the 1km dataset, we created a coarser dataset where each pixel covers a $5\text{km} \times 5\text{km}$ region.
- We remove all the pixels which have a fill value of -999.
- For each band, we first aggregate all the data over all the years and days. We then create a single data matrix where each row corresponds to an observation and each column corresponds to a band. The size of this matrix is $12,613,391 \times 7$.
- Along with this, we have also created a latitude and longitude matrix (each of size $12,613,391 \times 7$) for each element in the data matrix.

Figure 3 shows the dataset and the final output of the preprocessing step.

6.3. Measurement metric. In all of our experiments we measured these quantities: (1) the percentage of correct detection, (2) the running time, and (3) the number of outliers detected. By percentage of correct detection we mean the number of common outliers which are found both by our distributed algorithm and a centralized algorithm having access to all of the data but using the same sample size T_s for training as the distributed algorithm. When comparing running time, we plot the running time of our method and the centralized algorithm running on all the features. Note that, for our distributed algorithm since each site can run in parallel, we report the average running time over all the sites. Finally we report the total number of outliers detected by our distributed

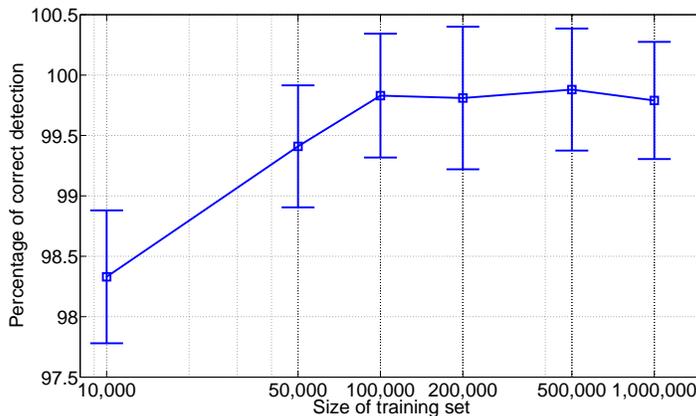


FIGURE 4. Variation of the percentage of correct detection with the size of the training set. The samples are selected at random from the entire dataset. Percentage of correct detection means the number of anomalies detected by the distributed method compared to a centralized SVM algorithm using the entire dataset. As evident, the detection rate increases as the sample size increases.

algorithm, the centralized algorithm, and the unique outliers detected by the distributed algorithm only.

6.4. Performance evaluation. In this section we discuss the performance of the distributed algorithm on the California MODIS dataset. The first figure (Figure 4) shows how the detection rate varies as the size of the training sample (T_s) is varied. By detection rate we mean the number of common outliers found by our proposed distributed algorithm and another centralized algorithm having access to all the data, both using the same training set size T_s drawn uniformly randomly from the entire dataset. We have varied T_s from 10,000 to 1,000,000. Note that the highest value of T_s is only 10% of the entire dataset. For a uniformly selected training set of size 10,000, the percentage of correct detection is 98.33. It remains almost a constant for different sizes of the training set. For approximately 1 million test points, the correct detection rate is close to 99.79%. This shows that our algorithm is extremely accurate and returns the true set of outliers over a different sample sizes.

The next experiment demonstrates the gain of our algorithm with respect to running time. As shown in Figure 5, the running time of our algorithm diverges from the centralized algorithm as T_s is increased. For smaller T_s , the running time is comparable to the centralized algorithm. As T_s increases, our algorithm starts performing better. This is intuitive since with increasing size of training sample, more computation is needed and thus the running time of the centralized algorithm increases sharply. On the other hand, the distributed algorithm exhibits a slower growth in running time since the total processing load is distributed across all the processors. This helps the distributed algorithm to achieve a lower running time.

Message complexity of the algorithm is demonstrated in Figure 6. The x -axis shows the number of samples used for the training and the y -axis refers to the ratio of the bytes transferred by the distributed to the centralized algorithm, expressed in percentage. Note that a value of $y = 100$ means that the distributed algorithm does not provide any communication savings. For all the cases, the percentage message complexity varies between 0.134 and 7.934. This shows that the proposed algorithm is highly communication efficient.

Our final experiment shows the number of outliers detected by our algorithm and the centralized version. Table 2 shows the outliers detected by the various methods. The first column shows

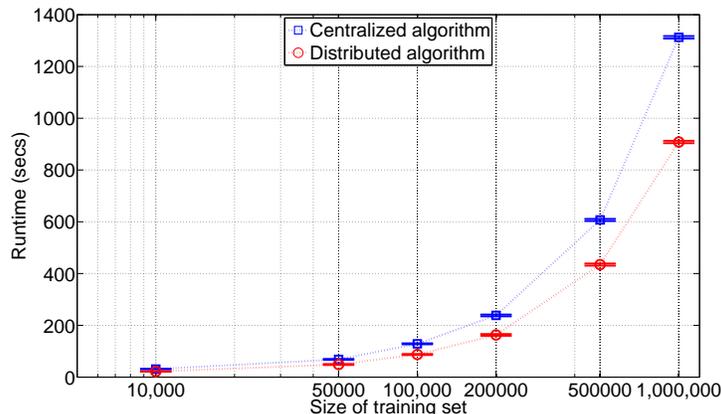


FIGURE 5. Variation of running time with the size of the training set. The samples are selected at random from the entire dataset. Both the running times of our algorithm and the centralized algorithm are shown. Clearly, the distributed algorithm outperforms the centralized one as the sample size increases.

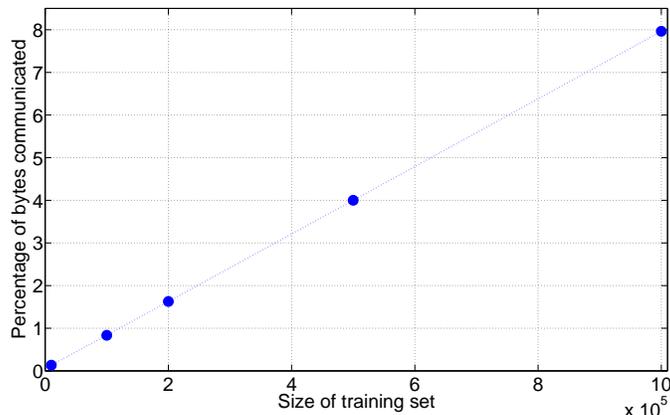


FIGURE 6. Variation of the percentage of bytes communicated with the size of the training set. The samples are selected at random from the entire dataset. The y -axis refers to the ratio of the bytes transferred by the distributed to the centralized algorithm, expressed in percentage. As depicted, the maximum percentage of bytes transferred is close to 8%, proving the excellent scalability of the proposed algorithm.

the number of outliers detected by the distributed algorithm. The second second figure shows the number of outliers detected by the centralized algorithm having access to all the data and using a sample size equal to that of the distributed algorithm. The last column refers to the average number of outliers found by the distributed algorithm and not by any of the sites. For each site, we first compute \mathcal{O}_i . Then we find the distributed outliers \mathcal{O}_d . For each site, then we compute $\{\mathcal{O}_d \setminus \mathcal{O}_i\}$. We take the average over all the sets in order to report the average number of unique outliers.

Figure 7 shows the top 50 outliers for training set size of 100,000. Figure 8 shows the top 50 outliers detected by the distributed algorithm but not detected by the feature at site 1 (*i.e.* Band1) only. Note that the set of outliers in Figure 7 and 8 are different. This is because the top 50 outliers absent in site 1 may be actually ranked lower than the top 50 outliers detected in Figure 7.

Sample size	No of distributed outliers	No of centralized outliers	Average no of unique outliers
10,000	14747	15473	7179
50,000	15382	15473	7284
100,000	13068	13090	6176
200,000	12940	12964	5986
500,000	11033	11046	5090
1,000,000	11221	11197	5233

TABLE 2. Number of outliers detected by the distributed algorithm and the centralized algorithm. The last column shows the unique outliers detected by the distributed algorithm and not detected at any of the local sites (using that feature only).



FIGURE 7. Top 50 outliers detected by the distributed algorithm for $T_s = 100,000$.

7. CONCLUSION

In this paper we have presented a distributed algorithm capable of detecting outliers from distributed data where each site has a subset of the global set of features. To the best of the authors' knowledge, this algorithm is the first which does anomaly detection from vertically partitioned data in a communication efficient manner. Our pruning rule allows us to achieve high accuracy and low communication cost, a must for processing petabytes of data. We have provided a comprehensive theoretical analysis of the algorithm to show its gains. Experimental evaluation is conducted with

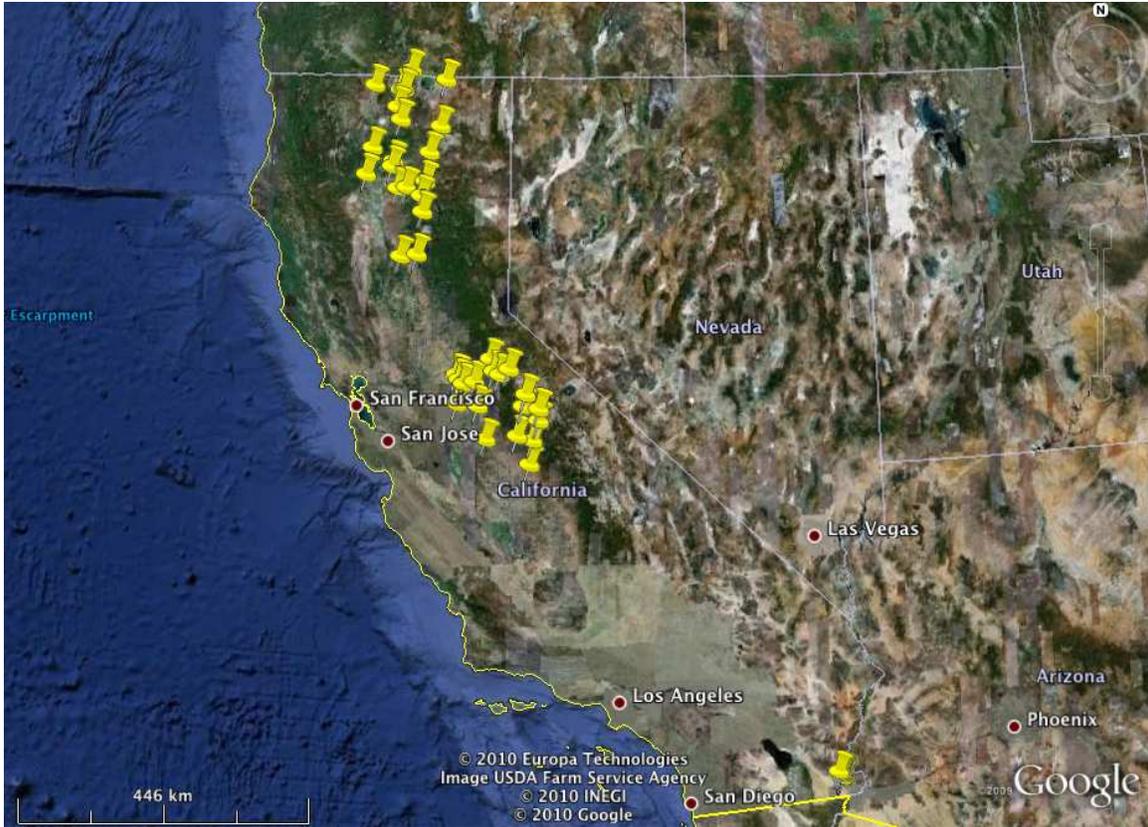


FIGURE 8. Top 50 outliers detected by the distributed algorithm but not detected by the first site using its data only. This is for training set size of 100,000.

the NASA MODIS satellite image dataset. The results show that the algorithm is approximately 99% accurate with only 1% of the communication needed for centralizing all the data.

ACKNOWLEDGEMENT

The data used in this paper are distributed by the Land Processes Distributed Active Archive Center (LP DAAC), located at the U.S. Geological Survey (USGS) Earth Resources Observation and Science (EROS) Center (lpdaac.usgs.gov). This work was supported by the NASA Aviation Safety Program, Intelligent Vehicle Health Management project and the TOPS project. The authors are grateful to the reviewers for their comments and suggestions, which improved this paper.

REFERENCES

- [1] F. Angiulli and C. Pizzuti. Outlier Mining in Large High-Dimensional Data Sets. *TKDE*, 17(2):203–215, 2005.
- [2] S. Barua and R. Alhajj. A Parallel Multi-scale Region Outlier Mining Algorithm for Meteorological Data. In *GIS*, page 51, 2007.
- [3] S. D. Bay and M. Schwabacher. Mining Distance-based Outliers in Near Linear Time with Randomization and a Simple Pruning Rule. In *Proceedings of KDD'03*, pages 29–38, 2003.
- [4] K. Bhaduri and H. Kargupta. A Scalable Local Algorithm for Distributed Multivariate Regression. *Stat. Anal. and Data M. J.*, 1(3):177–194, 2008.
- [5] D. Birant and A. Kut. Spatio-Temporal Outlier Detection in Large Databases. *Journal of Computing and Information Technology*, 14(4):291–297, 2006.

- [6] S. Boriah, V. Kumar, M. Steinbach, C. Potter, and S. Klooster. Land Cover Change Detection: A Case Study. In *Proceeding of KDD'08*, pages 857–865, 2008.
- [7] V. Chandola, A. Banerjee, and V. Kumar. Anomaly Detection: A Survey. *ACM Computing Surveys*, 2008 (to appear).
- [8] K. Das and J. Schneider. Detecting Anomalous Records in Categorical Datasets. In *Proceedings of KDD'07*, pages 220–229, NY, USA, 2007.
- [9] S. Hido, Y. Tsuboi, H. Kashima, M. Sugiyama, and T. Kanamori. Inlier-Based Outlier Detection via Direct Density Ratio Estimation. In *Proceedings of ICDM'08*, pages 223–232, Pisa, Italy, 2008.
- [10] V. Hodge and J. Austin. A Survey of Outlier Detection Methodologies. *Artif. Intell. Rev.*, 22(2):85–126, 2004.
- [11] W. Hu, Y. Liao, and V. R. Vemuri. Robust Anomaly Detection using Support Vector Machines in Computer Security. In *Proceedings of ICML'03*, pages 168–174, 2003.
- [12] E. M. Knorr, R. T. Ng, and V. Tucakov. Distance-based Outliers: Algorithms and Applications. *The VLDB Journal*, 8(3-4):237–253, 2000.
- [13] C. Potter, V. Genovese, P. Gross, S. Boriah, M. Steinbach, and V. Kumar. Revealing Land Cover Change in California With Satellite Data. *Eos Trans. AGU*, 88(26):269–270, 2007.
- [14] S. Ramaswamy, R. Rastogi, and K. Shim. Efficient Algorithms for Mining Outliers from Large Data Sets. *SIGMOD Rec.*, 29(2):427–438, 2000.
- [15] B. Schölkopf, J. C. Platt, J. C. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Estimating the Support of a High-Dimensional Distribution. *Neural Comput.*, 13(7):1443–1471, 2001.
- [16] A. Schuster and R. Wolff. Communication-Efficient Distributed Mining of Association Rules. *Data Mining and Knowledge Discovery*, 8(2):171–196, 2004.
- [17] S. Shekhar, P. R. Schrater, R. R. Vatsavai, W. Wu, and S. Chawla. Spatial Contextual Classification and Prediction Models for Mining Geospatial Data. *IEEE Transactions on Multimedia*, 4(2):174–188, 2002.
- [18] R. Wolff, K. Bhaduri, and H. Kargupta. A Generic Local Algorithm for Mining Data Streams in Large Distributed Systems. *TKDE*, 21(4):465–478, 2009.
- [19] J. Zhao, C. Lu, and Y. Kou. Detecting Region outliers in Meteorological Data. In *Proceedings of GIS'03*, pages 49–55, 2003.