



HAL
open science

Generalized isolation forest for anomaly detection

Julien Lesouple, Cédric Baudoin, Marc Spigai, Jean-Yves Tournernet

► **To cite this version:**

Julien Lesouple, Cédric Baudoin, Marc Spigai, Jean-Yves Tournernet. Generalized isolation forest for anomaly detection. *Pattern Recognition Letters*, 2021, 149, pp.109-119. 10.1016/j.patrec.2021.05.022 . hal-03382634

HAL Id: hal-03382634

<https://hal.science/hal-03382634v1>

Submitted on 18 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Open Archive Toulouse Archive Ouverte

OATAO is an open access repository that collects the work of Toulouse researchers and makes it freely available over the web where possible

This is an author's version published in: <http://oatao.univ-toulouse.fr/28086>

Official URL:

<https://doi.org/10.1016/j.patrec.2021.05.022>

To cite this version:

Lesouple, Julien and Baudoin, Cedric and Spigai, Marc and Tourneret, Jean-Yves Generalized isolation forest for anomaly detection. (2021) Pattern Recognition Letters, 149. 109-119. ISSN 0167-8655

Any correspondence concerning this service should be sent to the repository administrator: tech-oatao@listes-diff.inp-toulouse.fr

Generalized isolation forest for anomaly detection[☆]

Julien Lesouple^{a,*}, Cédric Baudoin^b, Marc Spigai^b, Jean-Yves Tournéret^{a,c}

^aTéSA, 7 Boulevard de la Gare, Toulouse 31000, France

^bThales Alenia Space, 26 Avenue Jean-François Champollion, Toulouse 31100 France

^cUniversity of Toulouse/INP-ENSEEIH/IRIT, 2 Rue Charles Camichel, Toulouse, 31071, France

A B S T R A C T

This letter introduces a generalization of Isolation Forest (IF) based on the existing Extended IF (EIF). EIF has shown some interest compared to IF being for instance more robust to some artefacts. However, some information can be lost when computing the EIF trees since the sampled threshold might lead to empty branches. This letter introduces a generalized isolation forest algorithm called Generalized IF (GIF) to overcome these issues. GIF is faster than EIF with a similar performance, as shown in several simulation results associated with reference databases used for anomaly detection.

Keywords:

Anomaly detection

Isolation forest

1. Introduction

Anomaly Detection (AD, Chandola et al. [3]) has gained attention in the past few years, due to the enhancement of modern computers and the increasing interest for machine learning algorithms. AD consists in detecting rare patterns or unobserved samples in data, referred to as anomalies. It is widely used in potentially critical environments, e.g., in credit fraud detection [1], crowd surveillance [9], or in satellite telemetry monitoring [12,15]. AD has received an increasing interest for satellite monitoring in the past few years, with new satellite constellations, resulting in a huge amount of data to be processed at the same time. Time-series resulting from satellite telemetry are of course used for the constellation mission but also for system monitoring and failure prevention.

This letter focuses on unsupervised AD algorithms, which learn the normal behavior of unlabeled data using a so-called training dataset. The performance of the algorithm can then be tested using a labeled dataset called test set. Various AD algorithms have been proposed in the literature including those based on nearest neighbors (Local Outlier Factor, Breunig et al. [2], Local Outlier Probability (LoOP), Kriegel et al. [8] or Neighborhood Construction (NC), Inkaya et al. [7]), support vector machines (Support Vector Data Description, Tax and Duijn [14], One Class Support Vector Machines, Schölkopf et al. [13]), Sparse Coding [4], or Isolation Forest (IF, Liu et al. [10]).

A specific attention is devoted in this letter to IF, which aims at finding anomalies with the idea that in some feature space, anomalies should be “far” from other data. To look for these anomalies, IF generates random isolation trees in order to isolate each data point. The number of branches required to isolate each point is then computed for each tree. The mean of this number of branches defines the expected path length, which is used to isolate a point of interest. The expected path length is generally small for anomalies (contrary to nominal data) since anomalies are far from the majority of nominal data. However, the trees generated by IF are considering a random feature at each node, which can lead to some artefacts in the score map function, as shown in Hariri et al. [6]. In order to improve the isolation of data points, tree branches with random hyperplanes can be considered [6]. Random hyperplanes are not necessarily parallel to one of the components of the feature vector and have been used in the extended IF (EIF) algorithm. Unfortunately, this strategy generates a lot of empty branches, which increases the complexity of the trees belonging to the forest. This letter goes a step further by proposing a new IF construction inspired by the work of [6] leading to the generalized isolation forest (GIF) algorithm. The GIF algorithm generates trees without any empty branch, which significantly improves the execution times when compared to EIF.

This letter is organized as follows: Section 2 recalls the principles of IF and EIF and introduces the proposed GIF algorithm. Section 3 evaluates the performance of GIF using experiments on both synthetic and real benchmark datasets. Conclusion are reported in Section 4.

[☆] Editor: Jose Ruiz-Shulcloper

* Corresponding author.

E-mail address: julien.lesouple@tesa.prd.fr (J. Lesouple).

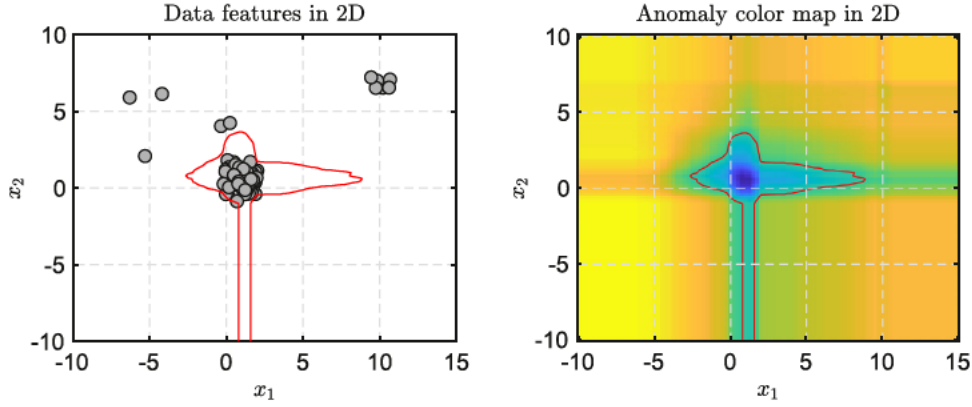


Fig. 1. Illustration of IF problems using artificial 2D data. Training data are depicted in the left figure as well as the curve $s(\mathbf{x}, n) = s_0$ (displayed in red). The right figure shows the heat map of the anomaly score (dark blue corresponds to values next to 0 and light yellow to value close to 1). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

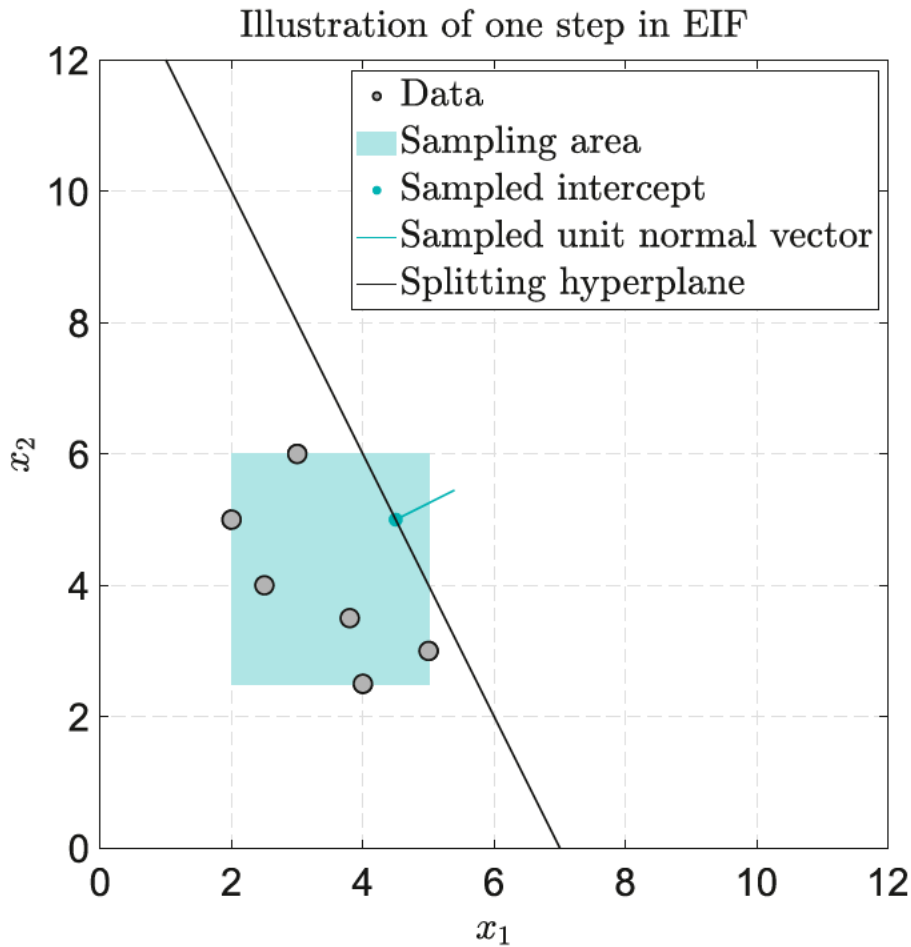


Fig. 2. Illustration of an EIF drawback using artificial 2D data. A splitting hyperplane is created by sampling a random unit vector and a random intercept in the sampling area. As one can see, using this strategy, all the data points are below the hyperplane (for this outcome). Thus the corresponding right branch of the tree will be empty.

2. Isolation forest

2.1. Original formulation

IF generates $t > 0$ random trees to partition the data, and computes for each tree the number of nodes required to isolate each training vector. Anomalies are then detected as the vectors whose average path lengths are the smallest, motivated by the fact that

nominal data are more concentrated than anomalies and thus require more nodes to be isolated.

To create a random isolation tree, assume that we have n training data $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, where $\mathbf{x}_i = [x_{i,1} \ \dots \ x_{i,d}]^T \in \mathbb{R}^d$. We will also use the notation $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ for the matrix gathering all the training data. To create a random node and split the dataset into two subsets, one component of \mathbb{R}^d (denoted as q) is chosen randomly, and a split value p is sampled uniformly in the interval $[\min_{i=1, \dots, n} x_{i,q}; \max_{i=1, \dots, n} x_{i,q}]$. The dataset is then split

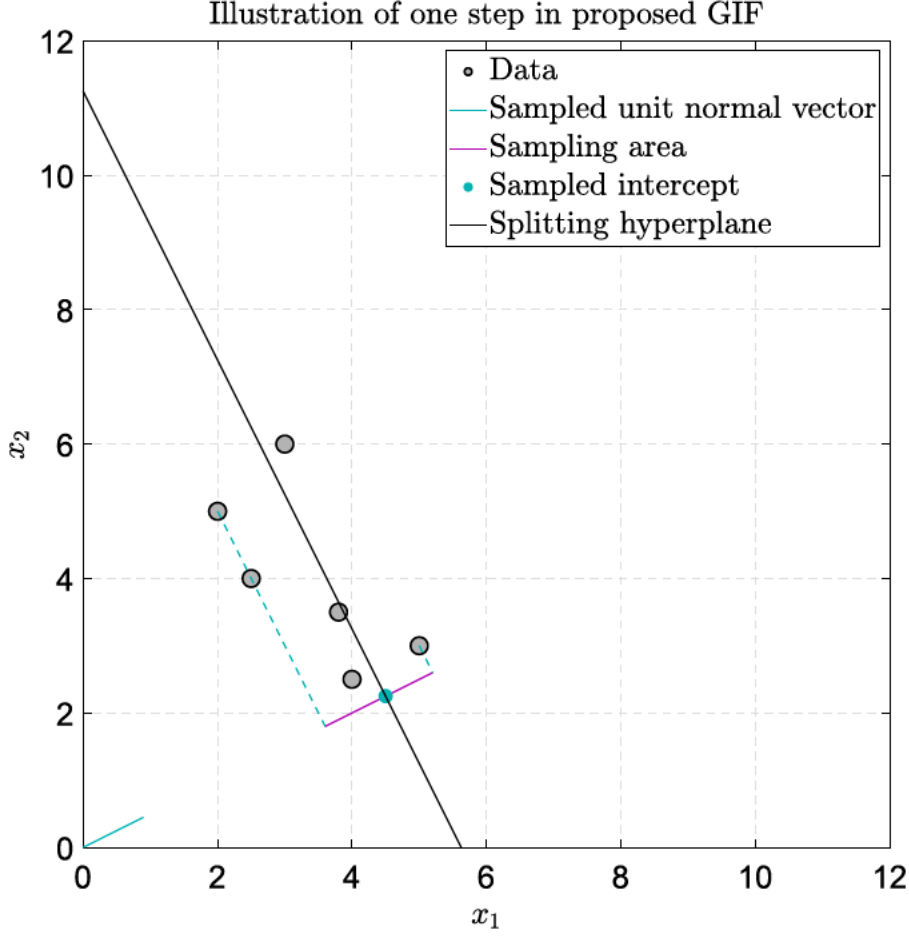


Fig. 3. Illustration of the proposed GIF approach. A splitting hyperplane is created by sampling a random unit vector and a random intercept in the sampling area, which reduces to a line (and not a square). This strategy has the advantage of having data points on each side of the splitting hyperplane.

into two parts: the so-called left branch corresponding to the set $\{\mathbf{x}_i, x_{i,q} \leq p\}$ and the so-called right branch, corresponding to the set $\{\mathbf{x}_i, x_{i,q} > p\}$. The tree is created by applying this procedure iteratively to each branch until a branch contains a unique data point, or until some depth l has been reached. To create an IF, this procedure could be applied several times to the whole learning dataset. However, authors in Liu et al. [10] have shown that for each tree, a sub-sample of the whole dataset of size $\psi > 0$ (chosen to $\psi = 256$ in this letter) can be considered with similar performance and improved computation time.

Once the forest has been created by generating t random isolation trees, the expected path length $h(\mathbf{x})$ to isolate a point \mathbf{x} is computed using the mean of the path lengths required to isolate the point using each generated tree. Finally, an anomaly score is defined as

$$s(\mathbf{x}) = 2^{-\frac{E[h(\mathbf{x})]}{c(\psi)}}, \quad (1)$$

where $c(n)$ is the average value of $h(\mathbf{x})$ for a dataset of size n , which can be computed as

$$c(n) = 2H(n-1) - \frac{2(n-1)}{n}, \quad (2)$$

where $H(n)$ is the n th harmonic number (that can be approximated by $\ln(n) + \gamma$, where $\gamma \approx 0.577$ is the Euler-Mascheroni's constant). Thus, when $E[h(\mathbf{x})] = c(n)$, the anomaly score of \mathbf{x} is $s(\mathbf{x}, n) = 0.5$. When $h(\mathbf{x})$ tends to $+\infty$, i.e., when \mathbf{x} is not an isolated point, the anomaly score tends to 0. Finally, when $h(\mathbf{x})$ is small compared to $c(n)$, i.e., when \mathbf{x} is an isolated point, the corresponding anomaly score tends to 1. Thus we can define an anomaly

Table 1
Proposed values for the various parameters of IF.

Parameters	Meaning	Proposed value
t	Number of trees	100
ψ	Sub-sample size	256
l	Tree maximum depth	$\text{ceil}(\log_2 \psi) = 8$
s_0	Anomaly detection threshold	0.6

threshold $s_0 \in [0, 1]$ such that \mathbf{x} is detected as an anomaly when $s(\mathbf{x}) > s_0$, and as a nominal data when $s(\mathbf{x}) \leq s_0$. Of course, the closer the anomaly score to 1, the more likely \mathbf{x} is an anomaly, and the closer the anomaly score to 0, the more likely \mathbf{x} is a nominal vector. Thus, a trade-off has to be made to determine an appropriate value of s_0 . Authors in Liu et al. [10] have proposed values for the different parameters that are summarized in Table 1.

The resulting IF algorithm is a convenient solution to detect anomalies without assumptions on the data distribution and it is computationally efficient. However, this algorithm suffers from a bias due to the way trees are created. Indeed, by randomly choosing one dimension to split the data, parallel hyperplanes are used (with a normal vector collinear to the selected dimension), and data spread around stripes parallel to the axis and passing through the cluster have a lower anomaly score, as depicted in Fig. 1.

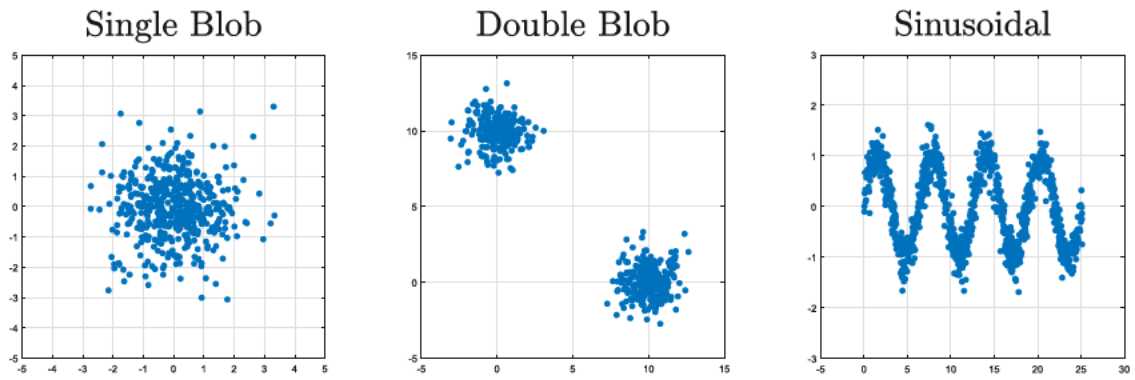


Fig. 4. Synthetic 2D datasets used to visualize the gain of EIF and GIF.

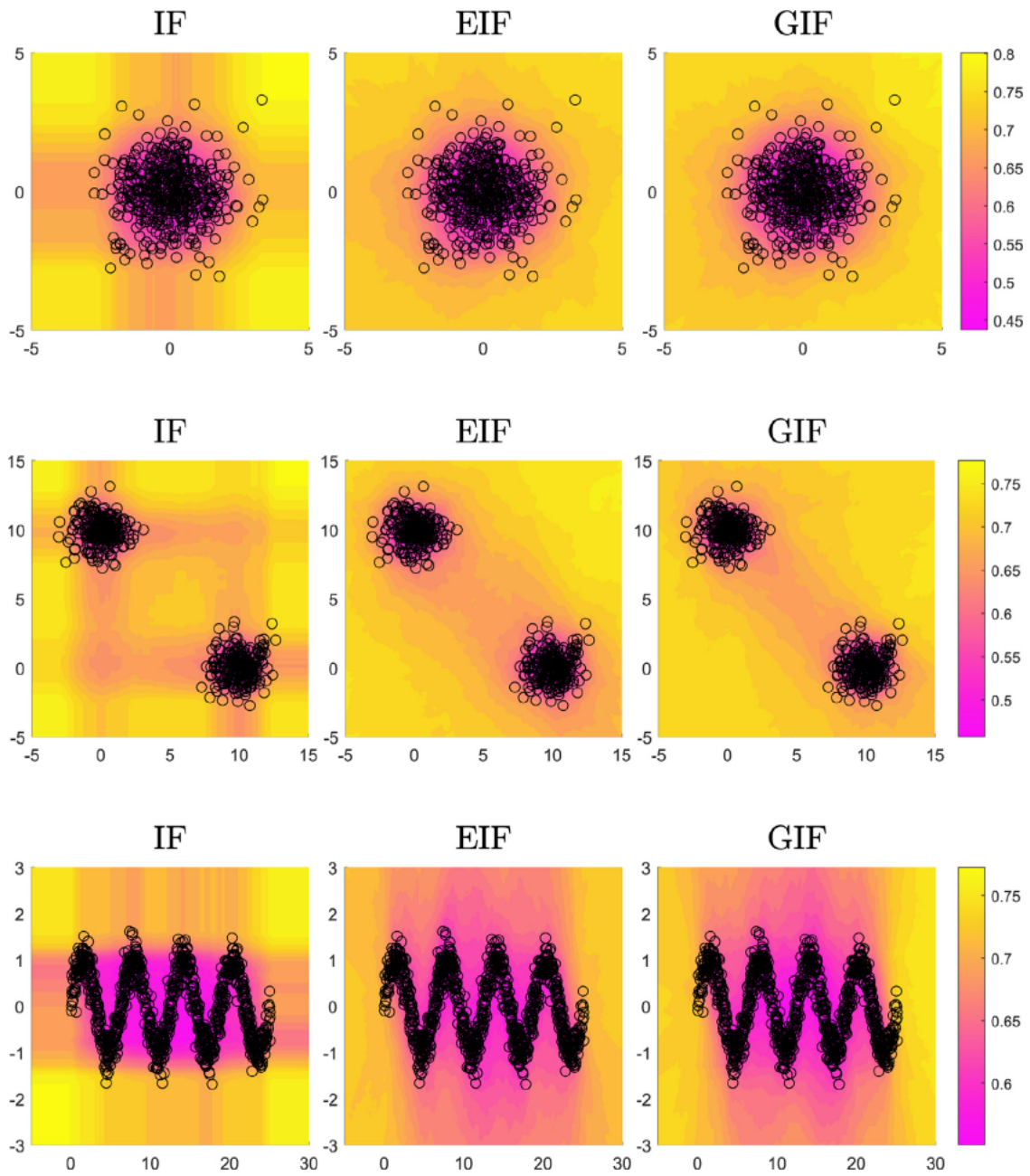


Fig. 5. Heat maps for the three algorithms and the three datasets: IF, EIF, GIF from left to right, and single blob, dual blob and sinusoidal from top to bottom. Pink values correspond to low anomaly scores and yellow to high.

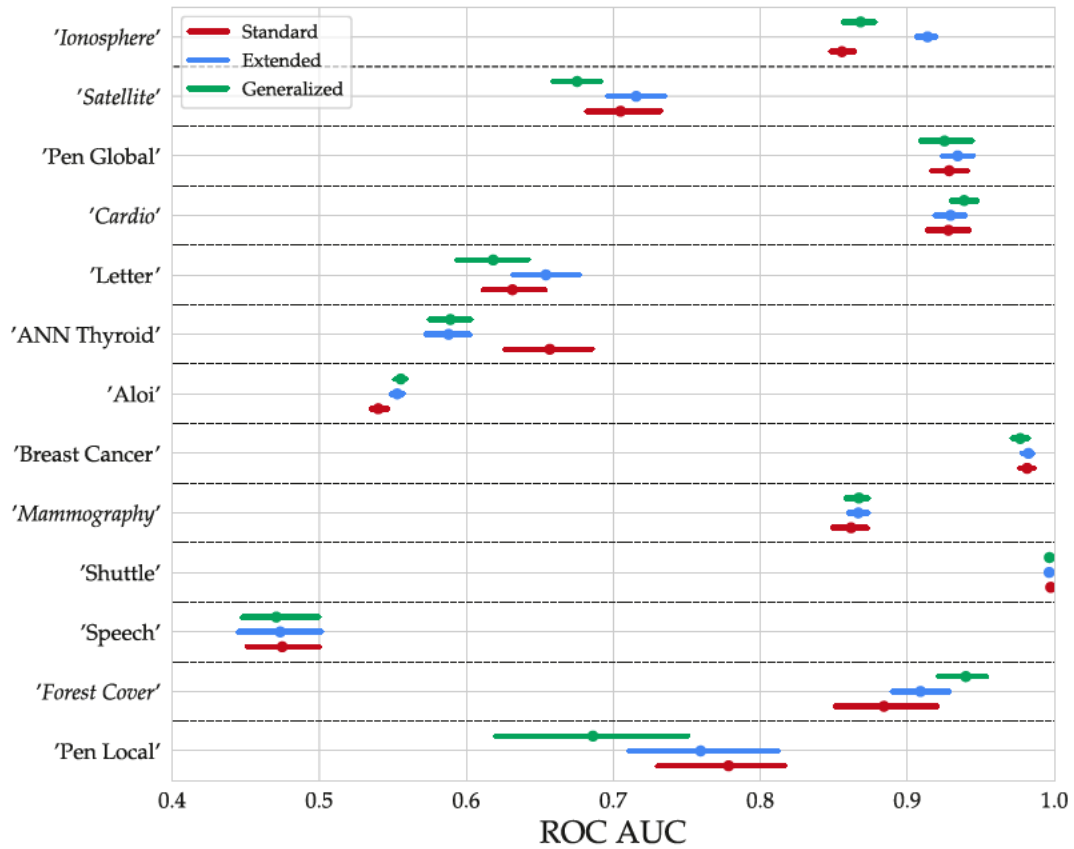


Fig. 6. Comparison of ROC AUC for several datasets (a line represents a 95% confidence interval and a dot the corresponding mean).

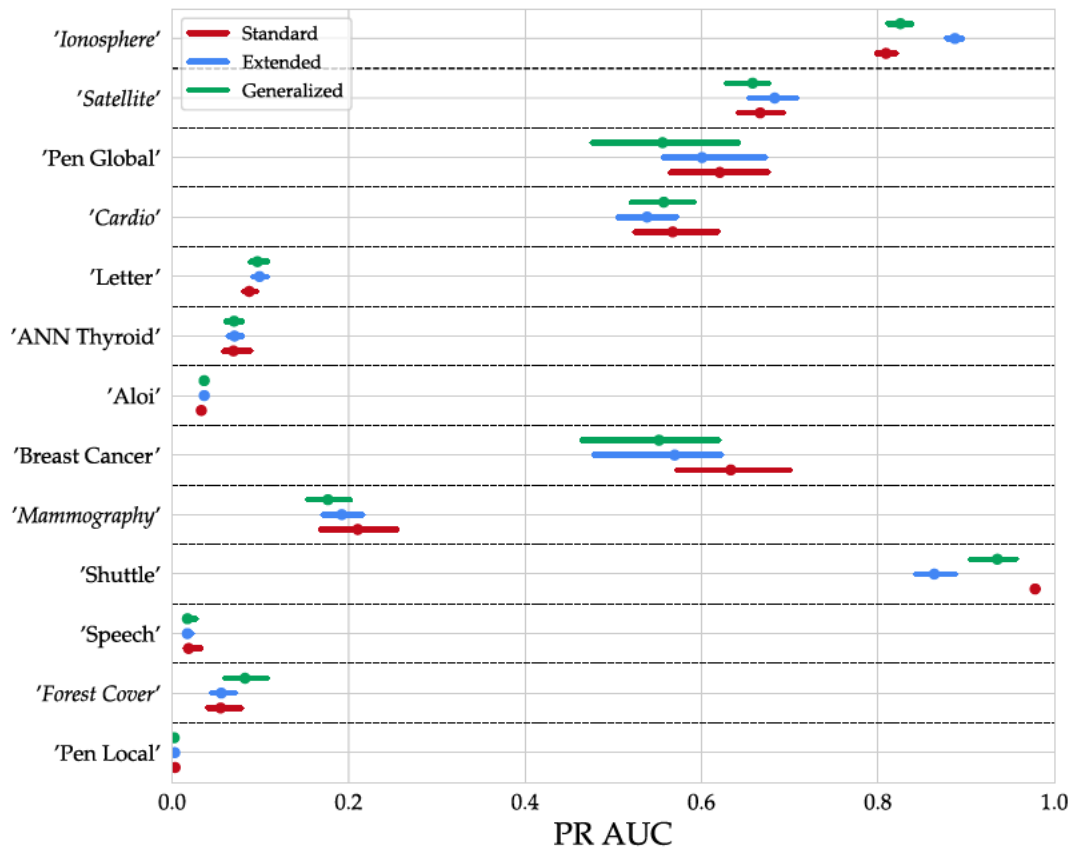


Fig. 7. Comparison of PR AUC for several datasets (a line represents a 95% confidence interval and a dot the corresponding mean).

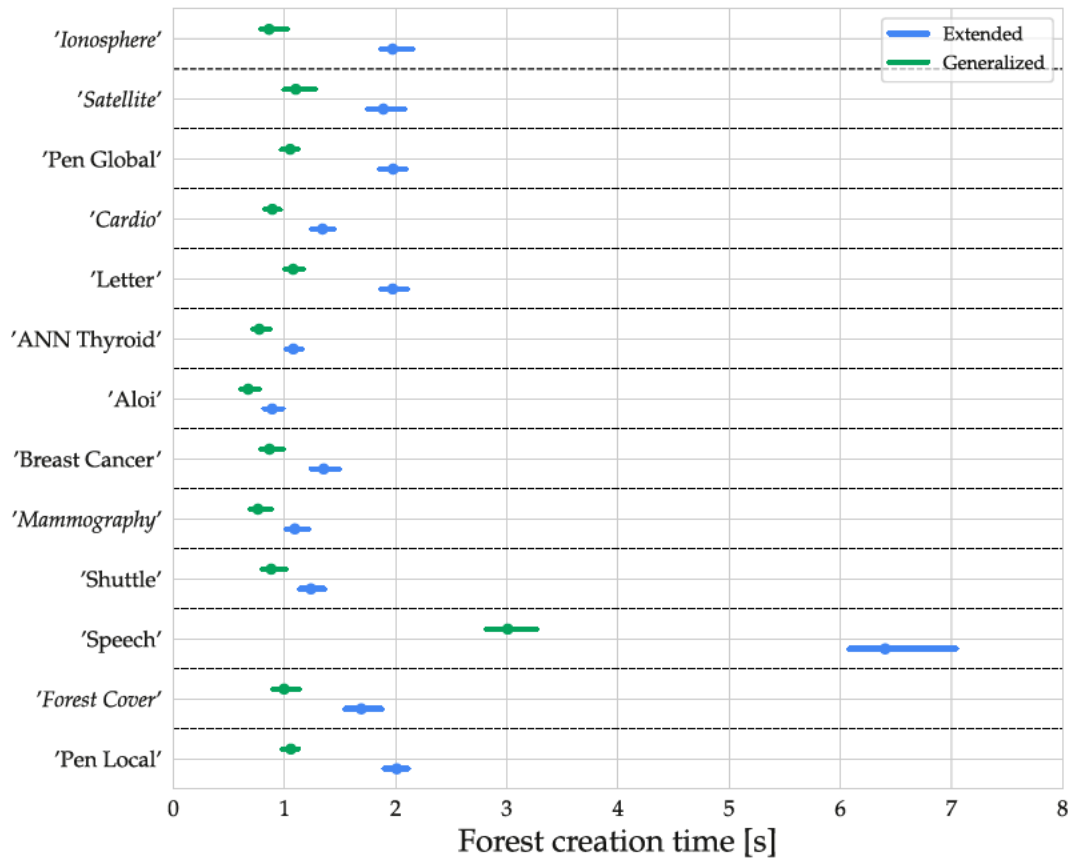


Fig. 8. Comparison of EIF and GIF computation times for several datasets (a line represents a 95% confidence interval and a dot the corresponding mean).

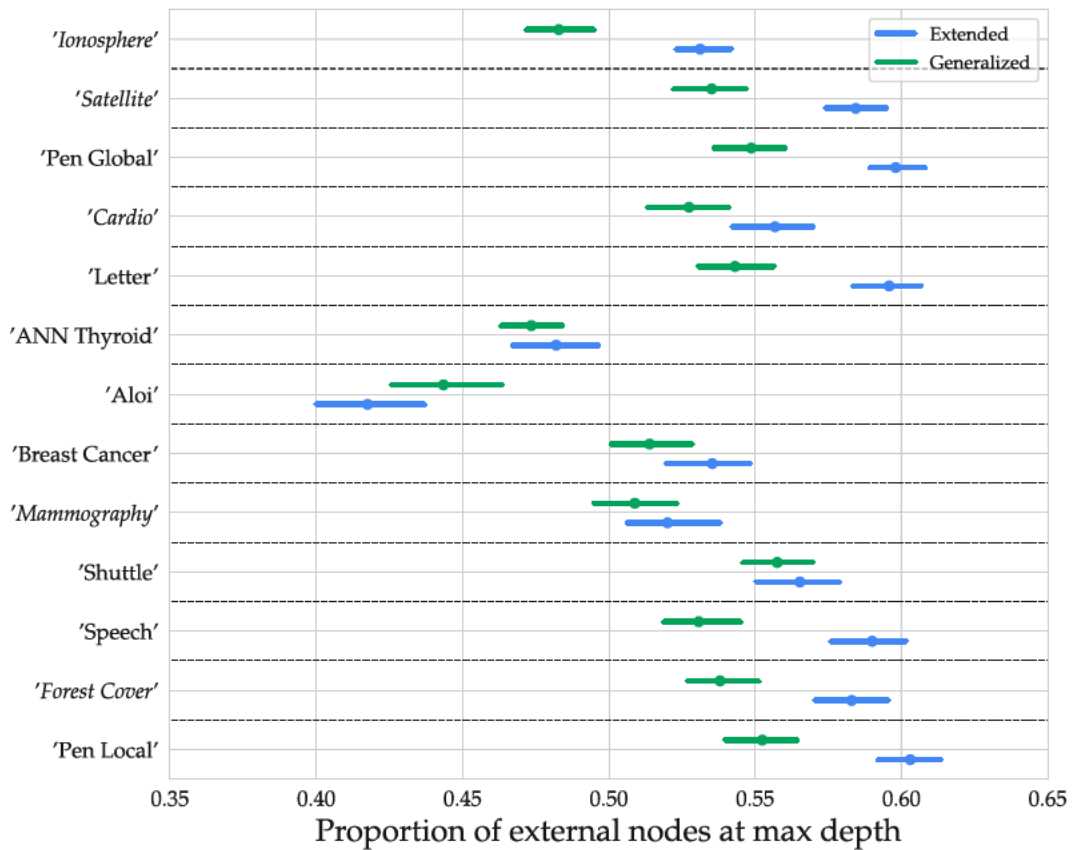


Fig. 9. Comparison of external nodes at maximum depth proportion for several datasets (a line represents a 95% confidence interval and a dot the corresponding mean).

2.2. Extended IF

To avoid artefacts such as those illustrated in Fig. 1, an improved solution was presented in Hariri et al. [6] referred to as EIF. As explained in Hariri et al. [6], the main drawback of IF is due to the way hyperplanes are constructed to split the data. Indeed, since the drawn normal vectors are chosen according to each dimension of \mathbb{R}^d , a discrete set of orthogonal directions is generated, which is at the origin of these vertical lines appearing in the level sets of $s(\mathbf{x})$. To mitigate this problem, a normal vector \mathbf{w} can be sampled for each decision hyperplane randomly chosen in the unit sphere of \mathbb{R}^d [6], i.e., a Gaussian vector is sampled according to $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d) \in \mathbb{R}^d$ and normalized leading to $\mathbf{w} = \mathbf{u}/\|\mathbf{u}\|_2$ [11]. To select the split value, an intercept vector $\mathbf{p} \in \mathbb{R}^d$ is sampled uniformly in the smallest axis-bounding hypercube enclosing all the samples at a branching point (as illustrated in Fig. 2). The two branches of the tree are defined depending on whether $(\mathbf{x} - \mathbf{p})^T \mathbf{w} > 0$ (right branch of the tree) or $(\mathbf{x} - \mathbf{p})^T \mathbf{w} \leq 0$ (left branch of the tree). The hyperplane is thus the one defined by the normal vector \mathbf{w} and containing the intercept point \mathbf{p} . One drawback of this method is that it can lead to empty branches in the tree, which goes against the idea of IF (whose idea is to split the tree until the number of points equals one or until a given maximal depth has been reached in order to efficiently isolate the data). This situation is depicted in Fig. 2 for the previous 2D example.

This letter studies a variation of EIF avoiding empty branches in each tree, referred to as generalized isolation forest (GIF), which is detailed in the next section.

2.3. Generalized isolation forest

In order to avoid empty branches in EIF, we transpose the EIF problem into the original one defining IF. More precisely, we propose to project all the data on the sampled normal unit vector, look for the minimum and maximum values of the projections (identified by the dotted lines in Fig. 3) and sample a split value uniformly between these two values. Note that this sampling ensures that there is at least one data in each branch of a tree: the first branch being defined from the min value and the second branch associated with the max value. This is equivalent to sample an intercept point on the restriction of the line spanned by the normal vector to the segment between the minimum and maximum values of the projected data points as shown in Fig. 3. This strategy ensures that the two branches of a tree are not empty, contrary to EIF. Note that it is equivalent to EIF where the sampling volume has been reduced to the convex hull of the data. Empty branches in EIF are due to intercepts sampled outside the convex hull of the considered samples and inside the axis-bounding hypercube. For EIF, the probability of sampling an intercept leading to an empty branch is therefore the volume between the hypercube and the convex hull, divided by the volume of the hypercube. Conversely, this volume equals 0 for GIF. Note that probability of having an empty branch in EIF increases as the number of dimensions increases, due to the curse of dimensionality, which motivates the need to avoid such situations. Finally, the proposed method can be defined by three algorithms summarized in Algorithms 1–3, inspired by Hariri et al. [6] and Liu et al. [10].

3. Experiments

This section evaluates the performance of the proposed GIF algorithm using synthetic 2D data and some benchmark datasets considered in Hariri et al. [6] and Goldstein [5].

Algorithm 1 Create the forest.

Input: \mathbf{X} - input data, t - number of trees, ψ - subsampling size
Output: *Forest* - a set of *iTrees*

```

1: function iFOREST( $\mathbf{X}, t, \psi$ )
2:   initialize Forest  $\leftarrow$  struct ▷ Empty structure
3:   set  $l = \text{ceil}(\log_2 \psi)$  ▷ Height limit
4:   for all  $i = 1$  to  $t$  do
5:      $\mathbf{r} \leftarrow \text{Sample}(\mathbf{X}, \psi)$  ▷ Subsample of size  $\psi$ 
6:     Forest.Tree( $i$ )  $\leftarrow$  iTREE( $\mathbf{r}, 0, l$ )
7:   end for
8: end function

```

Algorithm 2 Create a tree.

Input: \mathbf{X} - input data, e - current tree height, l - height limit
Output: *iTree* - a tree

```

1: function iTREE( $\mathbf{X}, e, l$ )
2:   initialize Tree  $\leftarrow$  struct ▷ Empty structure
3:   if  $e \geq l$  or  $|\mathbf{X}| \leq 1$  then
4:     Tree.Size  $\leftarrow$   $|\mathbf{X}|$  ▷ Number of remaining data
5:     Tree.Type  $\leftarrow$  'ext' ▷ No nodes after this one
6:   else
7:     draw  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ 
8:      $\mathbf{w} \leftarrow \mathbf{w}/\|\mathbf{w}\|_2$  ▷ Random unit vector of  $\mathbb{R}^d$ 
9:      $p_{\min} \leftarrow \min(\mathbf{X}\mathbf{w})$ 
10:     $p_{\max} \leftarrow \max(\mathbf{X}\mathbf{w})$ 
11:    draw  $p \sim \mathcal{U}([p_{\min}; p_{\max}])$ 
12:     $\mathbf{X}_l \leftarrow \mathbf{X}(\mathbf{X}\mathbf{w} \leq p, :)$ 
13:     $\mathbf{X}_r \leftarrow \mathbf{X}(\mathbf{X}\mathbf{w} > p, :)$ 
14:    Tree.Level  $\leftarrow$   $e$  ▷ Level of the node
15:    Tree.Left  $\leftarrow$  iTREE( $\mathbf{X}_l, e + 1, l$ )
16:    Tree.Right  $\leftarrow$  iTREE( $\mathbf{X}_r, e + 1, l$ )
17:    Tree.Normal  $\leftarrow$   $\mathbf{w}$ 
18:    Tree.Threshold  $\leftarrow$   $p$ 
19:    Tree.Type  $\leftarrow$  'int' ▷ Nodes after this one
20:   end if
21: end function

```

3.1. Synthetic datasets

In order to appreciate the benefits of GIF with respect to EIF and IF, we first consider three datasets of synthetic 2D samples displayed in Fig. 4. For each dataset, IF, EIF and GIF are run on the same data to learn the corresponding isolation forest. After building the isolation forests, a square area containing all the samples is transformed into a 100×100 grid. The anomaly score is computed for each point of this grid in order to build heat maps that are displayed in Fig. 5. The advantages of EIF and GIF with respect to IF, as already highlighted in Hariri et al. [6], are clear: the ‘‘cross’’ on the single blob, the sinusoid, and the ghost blobs for the second example disappear for GIF and EIF. In order to have a quantitative appreciation of the various methods, the next experiments consider several benchmark datasets whose anomalies are detected using the different algorithms.

3.2. Benchmark datasets

This section evaluates the performance of GIF on the datasets investigated in Hariri et al. [6]¹ and Goldstein [5]. Note that the different datasets are described in Table 2 and are ranked in increasing order regarding the anomaly proportion (datasets in italic are those used in Hariri et al. [6]).

¹ The datasets can be downloaded from <http://odds.cs.stonybrook.edu/>

Algorithm 3 Compute isolation score (path length).

Input: x - input vector, $Tree$ - an iTree, e - current path length
1: # e must be initialized to 0 when first called
Output: $Length$ - isolation score
2: **function** PL($x, Tree, e$)
3: **if** $Tree.Type = 'ext'$ **then**
4: **if** $Tree.size > 1$ **then**
5: $Length \leftarrow e + c(Tree.size)$ ▷ see (??)
6: **else**
7: $Length \leftarrow e$
8: **end if**
9: **else**
10: $w \leftarrow Tree.Normal$
11: $p \leftarrow Tree.Threshold$
12: **if** $x^T w \leq p$ **then**
13: $Length \leftarrow PL(x, Tree.Left, e + 1)$
14: **else**
15: $Length \leftarrow PL(x, Tree.Right, e + 1)$
16: **end if**
17: **end if**
18: **end function**

Table 2

Datasets used in the experiments.

Name	Samples n	Features d	Anomalies
Pen Local	6724	16	0.15%
Forest Cover	286,048	10	0.96%
Speech	3686	400	1.65%
Shuttle	46,464	9	1.89%
Mammography	11,183	6	2.32%
Breast Cancer	367	30	2.72%
Aloi	50,000	27	3.02%
ANN Thyroid	6916	21	3.61%
Letter	1600	32	6.25%
Cardio	1831	21	9.60%
Pen Global	809	16	11.12%
Satellite	6435	36	31.64%
Ionosphere	351	33	35.90%

Since IF-based anomaly detectors include some randomness due to the way the trees are built, Monte-Carlo simulations (using 100 iterations) were performed for all the datasets and the three methods (IF, EIF and GIF) to compute the average area under the curve (AUC) for both receiver operational characteristics (ROC) and precision recall (PR) curves, as well as quantiles $\alpha/2$ and $1 - \alpha/2$ where $\alpha = 5\%$ in order to obtain 95% confidence intervals. The results are gathered in Fig. 6 for the ROC and in Fig. 7 for PR curves.

Note that the computations were made using Python, with the IF algorithm from scikit learn², EIF from the author's github³, and our own implementation of GIF. The whole code as long as the datasets are available on the first author's webpage⁴. Note that all datasets have been preprocessed in order to obtain zero mean and unit variance for each feature. This preprocessing is not necessary for IF because splittings are made along a single feature. However, they are useful for EIF and GIF especially in high dimensions since these algorithms are sensitive to scaling.

As one can see, there is not a significant difference between EIF and GIF in terms of ROC AUC, except for the datasets Pen Local, Letter, Satellite and Ionosphere, where EIF seems to give a better result, and datasets Forest Cover and Cardio in favor of GIF. EIF and GIF also provide good results when compared to IF, except for the dataset ANN Thyroid. Regarding PR AUC, EIF and GIF

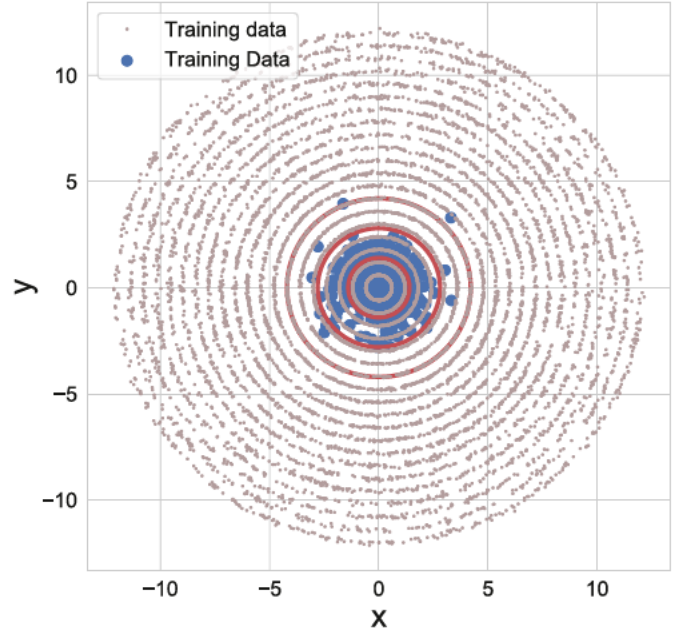


Fig. 10. Dataset for the additional experiments. The red circles represent 1, 2 and 3 data standard deviations. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

seem to have the same behavior, except for datasets Forest Cover and Shuttle, where GIF outperforms EIF. Finally note that EIF performs slightly better than IF and GIF for the dataset Ionosphere. From these experiments, we conclude that the performances of EIF and GIF are globally similar. In order to appreciate the interest of GIF, we have compared the execution times of the different algorithms, i.e., the time required to produce the forest (for both EIF and GIF) and the average proportion of external nodes at the maximum depth among all the external nodes computed for all the trees of a forest. The results are shown in Figs. 8 and 9.

As one can see, the times to compute the forests are significantly smaller for GIF compared to EIF, with generally smaller confidence intervals. The mean proportion of limit nodes among all the external nodes shows the capability of the method to isolate data. Indeed, an external node is either due to a reach of the given maximal depth, or to an isolated data. Therefore, if this number is close to one, few data are isolated (and conversely, the lower the mean proportion of limit nodes, the more data are isolated by the method). As the purpose of IF methods is precisely to isolate data, this ratio should be as low as possible. As one can see in Fig. 9, GIF leads to smaller proportions of this ratio than EIF (except for the Aloi dataset), which was expected.

3.3. Anomaly scores

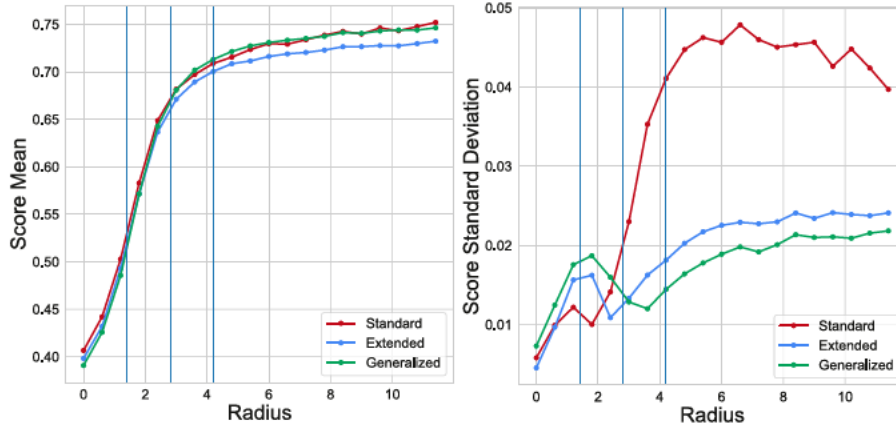
Additional experiments were conducted as presented in Hariri et al. [6] to analyze the anomaly scores of the different algorithms. More precisely, IF, EIF and GIF were trained on the 2D single blob synthetic dataset, and testing points were generated around constant radii, as shown in Fig. 10. This dataset allowed us to appreciate the behavior of the various algorithms to isotropic normal data. Indeed, for grey dots located around the same circle, the anomaly score should be approximately the same. The mean scores versus constant radius and the corresponding standard deviations are plotted for the various algorithms in Fig. 11a.

As one can see, the anomaly scores are equivalent for all the algorithms: there is a fast increase from zero to a value in the interval (2,3) when the radius increases, and slower variations after-

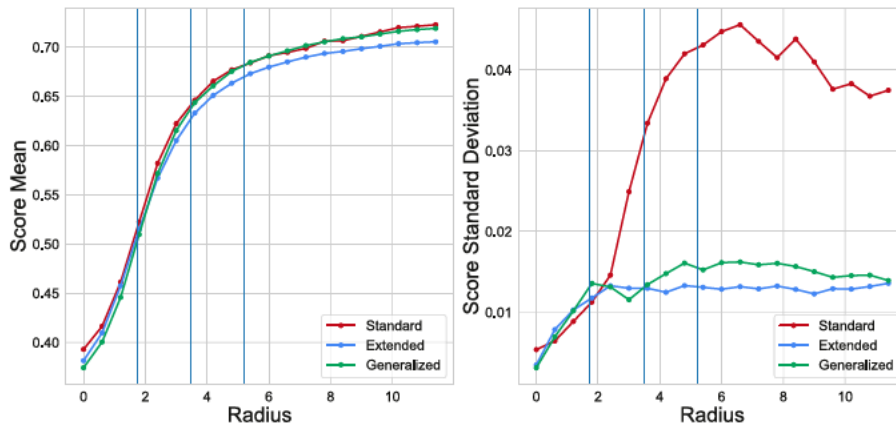
² <https://scikit-learn.org/stable/>

³ <https://github.com/sahandha/EIF>

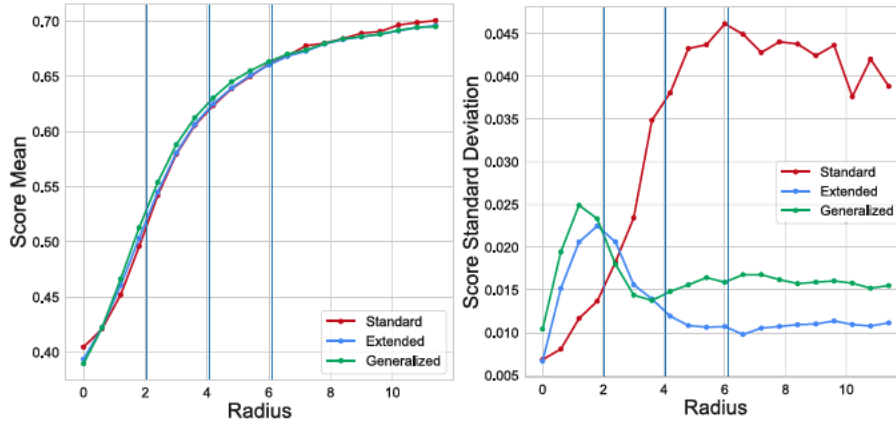
⁴ <http://perso.tesa.prd.fr/jlesouple/codes.html>



(a) 2D blob.



(b) 3D blob.



(c) 4D blob.

Fig. 11. Mean anomaly scores (left) and corresponding standard deviations (right) for the various algorithms versus the radius. The vertical blue lines represented the 1, 2 and 3 standard deviations. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

wards. One can observe that the standard deviations of the scores are significantly larger for IF than for EIF and GIF, which is explained by the absence of the “cross” effect for this dataset. These results were already shown in Hariri et al. [6] and are repeated here to show that the proposed GIF performs similarly to EIF, with the advantage of being faster, thanks to the absence of empty branches in the trees. The same experiments were run on a 3D blob and a 4D blob, as shown in Fig. 11b and c leading to the same conclusions.

The convergence of the mean anomaly scores was also studied, as in Hariri et al. [6]. The average anomaly scores for the inner an outer shell of each blob and the corresponding standard deviations were computed for each blob for various numbers of trees in the forest. The results are depicted in Fig. 12a–c for the 2D, 3D and 4D blobs respectively. As one can see, EIF and GIF provide similar results, with lower standard deviations when compared to the standard IF algorithm. Moreover, the anomaly scores for EIF and

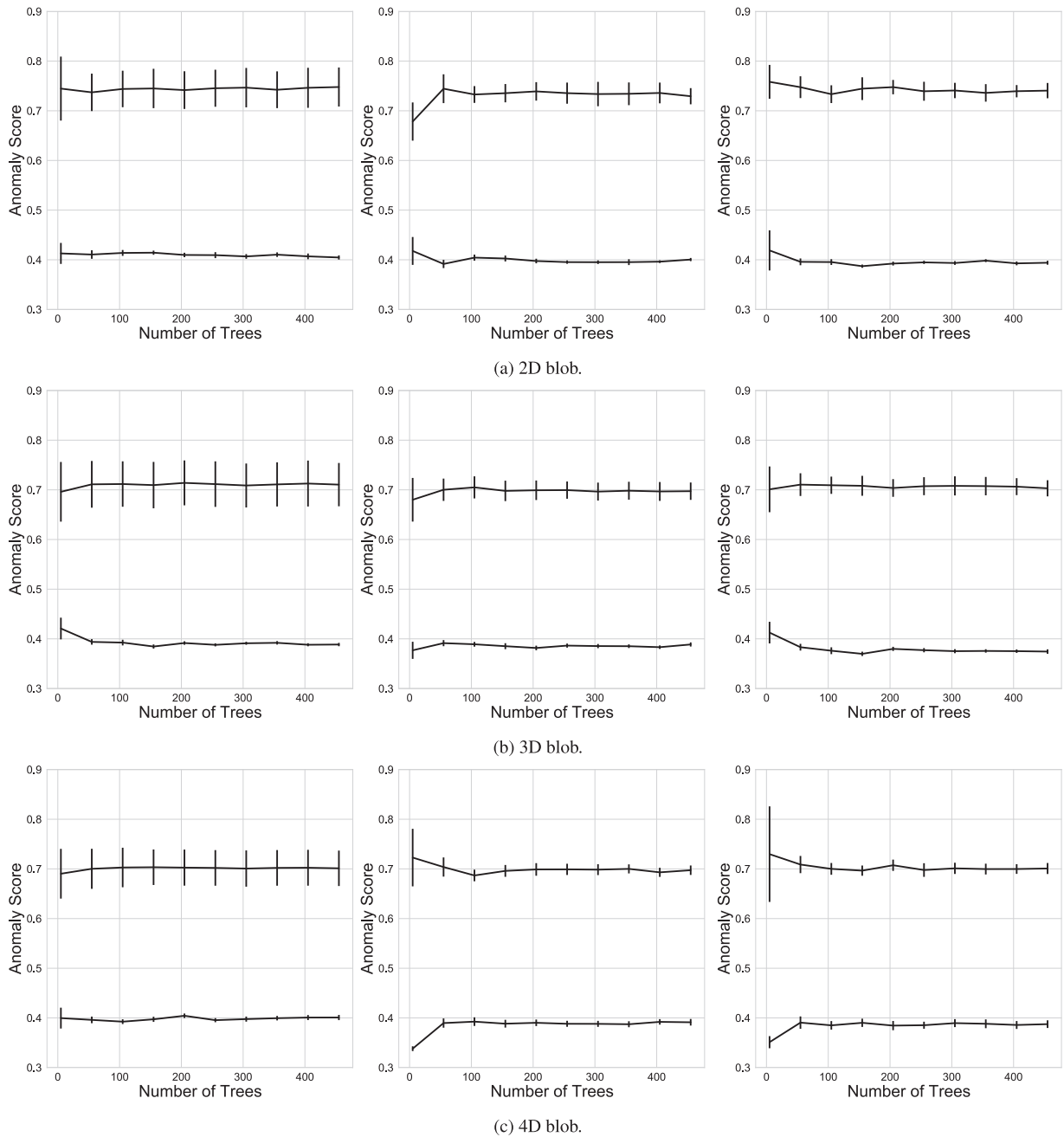


Fig. 12. Mean score for inner (bottom) and outer (top) shells versus the number of trees in the forest with corresponding standard deviations (vertical lines) for IF (left), EIF (center) and GIF (right).

GIF seem to converge to a constant value using a relatively small number of trees (around 100 trees in each forest).

4. Conclusion

This letter studied a new isolation forest algorithm referred to as generalized isolation forest for anomaly detection. This algorithm allows some artefacts of isolation forest to be bypassed and produces trees without empty branches, which is a drawback of the extended isolation forest (EIF) algorithm. Experimentations on both synthetic and benchmark datasets allowed us to evaluate the performance of the proposed method, which is similar to that obtained with EIF. However, the proposed algorithm has a significantly reduced execution time when compared to EIF, and requires few parameters to store (a threshold at each node for GIF versus an

intercept vector for each node for EIF). Future work will consider active learning and the injection of user feedback into the anomaly detectors to reduce the false alarm rate and improve anomaly detection.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- [1] R. Brause, T. Langsdorf, M. Hepp, Neural data mining for credit card fraud detection, in: *Proc. Int. Conf. on Tools with Artificial Intelligence*, 1999, pp. 103–106.

- [2] M.M. Breunig, H.-P. Kriegel, R.T. Ng, J. Sander, LOF: identifying density-based local outliers, in: Proc. Int. Conf. on Management of Data (SIGMOD), Dallas, TX, 2000, pp. 93–104.
- [3] V. Chandola, A. Banerjee, V. Kumar, Anomaly detection: a survey, *ACM Comput. Surv.* 41 (3) (2009) 15:1–15 58.
- [4] J.K. Dutta, B. Banerjee, Comparison of sparse coding-based versus traditional outlier detection methods, *Pattern Recognit. Lett.* 122 (2019) 99–105.
- [5] M. Goldstein, Unsupervised Anomaly Detection Benchmark, 2015, 10.7910/DVN/OPQMVF
- [6] S. Hariri, M.C. Kind, R.J. Brunner, Extended isolation forest, *IEEE Trans. Knowl. Data Eng.* 33 (4) (2021) 1.
- [7] T. İnkaya, S. Kayalgil, N.E. Özdemirel, An adaptive neighbourhood construction algorithm based on density and connectivity, *Pattern Recognit. Lett.* 52 (2015) 17–24.
- [8] H.-P. Kriegel, P. Kröger, E. Schubert, A. Zimek, LoOP: local outlier probabilities, in: Proc. Int. Conf. on Information and Knowledge Management (CIKM), Hong-Kong, China, 2009, pp. 1649–1652.
- [9] M.J.V. Leach, E.P. Sparks, N.M. Robertson, Contextual anomaly detection in crowded surveillance scenes, *Pattern Recognit. Lett.* 44 (2014) 71–79.
- [10] F.T. Liu, K.M. Ting, Z.-H. Zhou, Isolation forest, in: Proc. Int. Conf. on Data Mining (ICDM), Pisa, Italy, 2008, pp. 413–422.
- [11] M.E. Muller, A note on a method for generating points uniformly on N -dimensional spheres, *Commun. ACM* 2 (4) (1959) 19–20.
- [12] B. Pilastre, L. Boussof, S. d'Esquivan, J.-Y. Tournet, Anomaly detection in mixed telemetry data using a sparse representation and dictionary learning, *Signal Process.* 168 (2020) 107474.
- [13] B. Schölkopf, J.C. Platt, J.C. Shawe-Taylor, A.J. Smola, R.C. Williamson, Estimating the support of a high-dimensional distribution, *Neural Comput.* 13 (7) (2001) 1443–1471.
- [14] D.M.J. Tax, R.P.W. Duin, Support vector data description, *Mach. Learn.* 54 (2004) 45–66.
- [15] T. Yairi, N. Takeishi, T. Oda, Y. Nakajima, N. Nishimura, N. Takata, A data driven health monitoring method for satellite housekeeping data based on probabilistic clustering and dimensionality reduction, *IEEE Trans. Aerosp. Electron. Syst.* 53 (3) (2017) 1384–1401.