# Hunting inside n-quantiles of outliers (Hino)

Jessy Colonval $^{1[0000-0001-6792-0264]}$  and Fabrice Bouquet $^{1[0000-0001-9181-1172]}$ 

Université de Franche-Comté, Institut Femto-ST, CNRS, 16 route de Gray, Besançon, France {jessy.colonval,fabrice.bouquet}@femto-st.fr

**Abstract.** This paper presents an approach (called **Hino**) to detect outliers present in a data set, also called aberrations or anomalies. These data may reduce the quality of data analysis and lead to erroneous results. In the case of learning algorithms, they can deviate their behavior, i.e. reduce their efficiency. Thus, outlier detection is crucial where it improves performance by providing better data quality and reduces the influence of outliers. Inter Quartiles Range (IQR) is a popular statistical detection method, which has the advantage of being simple and fast in calculation time. It is based on the distribution quartiles of a data set and considers the most extreme values as outliers. This means that this method only searches for **point outliers**, which is a restrictive and naive approach. Indeed, nothing prevents an element from having an extreme value while remaining consistent with the rest of the elements. The proposed method is also a statistical detection method based on quantiles, but it looks for contextual outliers instead of point outliers and consider the context of a point to determine whether it is an outlier or not. The effectiveness of **Hino** is compared with the original **IQR** method and other approaches, including Isolation Forest, SVM, and LOF, using 16 real and 278 synthetic data sets.

**Keywords:** Outliers detection  $\cdot$  Data filtering  $\cdot$  Interquartile range  $\cdot$  Machine learning.

### 1 Introduction

An important principle in the use of artificial intelligence algorithms (AIs) is the quality of the learning data. So, the implementation of tools to help filter the data is a critical point [10]. However, data sets may contain outliers, which can be defined as "an observation (or subset of observations) which appears to be inconsistent with the remainder of that set of data" [3]. Outliers can reduce the quality of data analysis and lead to erroneous results. And in the case of AIs, they can deviate behavior of AIs and reduce their efficiency with two effects:

- 1. Weaken the predictive power of the model obtained after the learning phase.
- 2. Weaken the score obtained from the model during its validation phase.

Therefore, it is essential to identify and eliminate them directly. Thus, the quality of the data is improving and outlier's influence is reducing [9].

Three categories of outliers are identified [8]:

- 1. **Points outliers** are isolated individuals, they are far from other data.
- 2. Contextual outliers are isolated individuals in a specific context, i.e. according to their relations between their attributes and their behavioral values.
- 3. Collective outliers are a set of individuals consistent with each other but isolated from the rest of the data.

The IQR approach considers the distribution of a point to determine if it is isolated which makes it a statistical detection method. And since the context is not taken into consideration, then this method detects **point outliers**. Moreover, this method considers that a point with at least one of these values isolated is sufficient to consider it as an outlier. However, we will see that this method is ineffective for detecting aberrations. The objective is to extend the IQR approach so that it is able to detect **contextual outliers** and tolerate that the values of a point can be isolated a minimum number of times before considering the point fully isolated. This new approach is called **Hunting inside** n-quantiles of outliers (Hino).

Since **Hino** is looking for contextual outliers and the data sets used in the experimentation are spatial (see Section 5), then the terminology defines two types of attributes [1, Chapter 11]:

- 1. **Behavioral attributes** are attributes of interest measured for each point, commonly called *class*. For example: the type of glass, the presence of a heart abnormality, or the description of an image. A behavioral value is a value of this attribute. For example, a behavioral attribute that describes 7 different glass categories, here **1** would be one of the existing behavioral values.
- 2. **Contextual attributes** is an attribute expressing the characteristics of a point, commonly called *feature*. For example: the composition of a glass pane, the number of heartbeats per minute, or the color of a pixel.

The Section 2 present the **IQR** method of detection by quantiles, its functioning and the problems it raises. Then, the Section 3 present the **Hino** approach, its functioning and its meta-parameters. Follow-up of the Section 4 which proposes a methodology to generate the synthetic data sets used (see Section 4.1) and equations to calculate the **Hino** meta-parameters (see Section 4.2). Finally, the performance detection of **Hino** is compared with three methods (**IQR**, **Isolation Forest**, **SVM** and **LOF**) on a panel of 278 synthetic data sets (see Section 5.1) and on a panel of 16 real data sets (see Section 5.2) to conclude with a discussion of some biases present in the evaluation of real data sets and on the values used for the meta-parameters (see Section 5.3).

# 2 Detection method based on the interquartile range

In general, the median divides a data set into two more or less equal parts (depending on if the number of elements is even or odd), so approximately half of the elements are below the median and the rest are above it. On the same

principle, it is possible to make this division into two unequal parts such that a percentage p of the data is less than a certain number and a percentage 1-p is greater than that number. This number is called the 100p empirical percentile or the pth empirical quantile and is denoted by  $q_n(p)$  [5].

To calculate the interquartile range, the data set is divide into four parts. Using the previous definition we define: the 25th empirical percentile  $q_n(0.25)$  is called the lower quartile (Q1) and the 75th empirical percentile  $q_n(0.75)$  is called the upper quartile (Q3). Together with the median, they allow the division of a data set into four parts more or less equal in number of elements. The distance between these two quartiles gives an indication of the skewness of the data set. This distance is called **IQR** (Inter Quartile Range), and it equals to:

$$IQR = Q3 - Q1 = q_n(0.75) - q_n(0.25)$$
(1)

The interquartile range is used to describe the spread of a distribution. Thus, it is possible to detect outliers by considering the most scattered elements as outliers. For this, it is necessary to define thresholds that will cover the majority of the points so that only the most scattered ones are considered as outliers. The lower fence is equal to Q1-1.5\*IQR and the upper fence is equal to Q3+1.5\*IQR. The constant 1.5 allows controlling the sensitivity of the interval and thus the rule of detection. A larger one would cause outliers to be considered healthy, while a smaller one would cause some points to be detected as outliers. Thus, with a Gaussian distribution, then a constant equal to 1.5 gives fences covering  $-2.698\sigma$  (see equation 2) and  $2.698\sigma$  or 99.3% of the points.

$$Q1 - 1.5 * IQR = Q1 - 1.5 * (Q3 - Q1)$$

$$= -0.6745\sigma - 1.5 * (0.6745\sigma - (-0.6745\sigma))$$

$$= -0.6745\sigma - 1.5 * 1.349\sigma = -2.698\sigma$$
(2)

Figure 1, gives a visualization of quartiles and fences as a box plot and relates them to a Gaussian distribution. The majority of the points are included between these fences and only a small part of the points, the most extreme, will be excluded. Vinutha et al. [11] use the interquartile range as an outlier detection method. For each contextual attribute, this range is calculated and all points within are considered as outliers.

This method have three problems that make hazardous the use of the interquartile range as an outlier detection method. First, it removes the points that have extreme values who

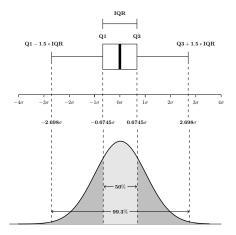


Fig. 1: Box plot and the Gaussian probability density of interquartile range

#### J. Colonval et al.

4

aren't always outliers. Indeed, nothing prevents an inlier to be extreme. Second, it doesn't consider the individual distribution of a behavioral value. One of them can be essentially present on the extreme values which would be suppressed with this approach. Third, there is no **tolerance limit** to determine if a point is an outlier. With only one extreme values among the contextual attribute, the large data sets are disadvantaged and risk to have a high false positivity rate because a point will have more chance to have an extreme value.

# 3 Extension of the IQR method

As a reminder, an outlier can be defined as an abrupt change. **IQR** detects only **point outliers** and considers as outliers the points with at least one extreme values. As explained previously, this approach is problematic. Thus, our objective is to extend **IQR** to be able to detect **contextual outliers** and tolerate the isolation of a point for few attributes before being an outlier.

For this purpose, we propose to separate the data set into several quantiles for each contextual attribute, instead of only 4. For each quantile, we check that all the behavioral values that constitute it are also present in adjacent quantiles. Otherwise, all representatives of absent values are considered isolated for this attribute. Except in the case where all the representatives of a value are present in a single quantile. A point isolated for few attributes is not necessarily an outlier. Indeed, depending on the total number of attributes, an isolation for just one of them does not have the same impact. For example, an isolation for 1 attribute out of 10 does not have the same significance for 100. Our approach tolerates a certain level of isolation before identifying a point as an outlier. To summarize, an outlier is a point "too often" isolated.

Our extension is summarized in the Algorithm 1. It takes as parameters the data set (P), the number of quantiles used  $(n\_qtils)$ , the tolerance limit (limit) and the maximum percentage of points that can be considered as outliers  $(p_{max})$  to preserve the consistency of data set. And it returns the mapping between the points in P and a Boolean to indicate if they are outliers. Except P, the other arguments strongly influence the detection result. However, it is difficult to determine the best values for each of them to obtain an optimal outlier detection.

- $-p_{max}$  is set at 20%, greater, we assume the data set will be distorted.
- n\_qtils is calculated according to the number of points they must contain.
   They must be small enough so that the points they contain share similar characteristics and make it easier to detect abrupt changes.
- limit is calculated according to the number of contextual attributes and behavioral values. When the number of contextual attributes and/or behavioral values are high, there's a greater chance of obtaining many isolated values.

This algorithm has been implemented in Python and its time complexity is equal in the worst case to O(p(ac + a + 1)), while in the best case to O(p(ac + a + 1)), where p is the number of points, a is the number of attributes and

c is the number of behavioral values. Plus, its memory complexity is equal to  $O(p(a+\frac{a}{a+1})).$ 

```
Algorithm 1: Hino.
   input: Let P be a set of n points
             Let n qtils \ge 2 be the number of quantiles
             Let limit > 0 be the maximum of breaking rule
             Let p_{max} be the maximum percentage of outliers
   output: Let result be an mapping indicating if a point is an outlier
 1 A_c be the set of contextual attributes labels in P
   /* Counters that track the frequency of isolated points.
 \mathbf{2} \ n \ cdt \leftarrow \{(p,0)|p \in P\}
 з for a \in A_c do
      qtils \leftarrow A list of list that groups the points divided by n qtils
       quantiles on the attribute a
       for i \in [0, card(qtils)] do
 5
           /* Gets the behavioral values of the previous, next
              and current quantile, respectively.
          prev \ cls \leftarrow \emptyset
 6
          if i > 0 then
             prev\_cls \leftarrow Behavioral values in <math>qtils[i-1]
 8
          next \ cls \leftarrow \emptyset
          if i < card(qtils) - 2 then
10
            | next \ cls \leftarrow Behavioral \ values \ in \ qtils[i+1]
11
          cur cls \leftarrow Behavioral values in <math>qtils[i]
12
          /* Missing behavioral values are those in the current
              quantile, but absent from the two adjacent ones. */
          missing\_cls \leftarrow cur\_cls \setminus (prev\_cls \cup next\_cls)
13
          /* If one is missing, then the counters of their
              points in the current quantile are increased.
          if card(missing\_cls) > 0 then
14
              for pq \in qtils[i] do
15
                  pq\_cls \leftarrow \text{Behavioral value of } pq
16
                  all\_here \leftarrow True if all of pq\_cls are in qtils[i]
17
                  /* In that case their counters are not
                                                                             */
                      incremented.
                  if pq cls \in missing cls and not all here then
18
                  n \ cdt[pq] + +
   /* Repeat with a limit incremented if a behavioral value is
       fully identified as an outlier or if more than p_{\max}
      points are outliers.
                                                                             */
20 repeat
21
       result \leftarrow \{(p, False) | p \in P\}
       for (p, count) \in n\_cdt do
22
         result[p] = count > limit
23
      limit + +
\mathbf{24}
25 until (countTrue(result) \le p_{max} and not viableCls(result))
```

## 4 Determination of meta-parameters

The first meta-parameters is the number of quantiles and it equals to  $\frac{p}{c+1}$  to ensure that each quantile is large enough so that each behavioral value is represented at least once. The second meta-parameter is the maximum percentage of outliers set at 20%, we assume that over this value a data set may become inconsistency. The third meta-parameter is the **tolerance limit**, which is determined by the number of attributes and behavioral values. However, building an equation from these values is difficult, so it's estimated empirically using a regression study based on synthetic data sets (see Section 4.2).

### 4.1 Generation of synthetics data sets and these outliers

Synthetic data sets are created with the same algorithm designed to generate Madelon [6]. They all have only one behavioral attribute whose values are homogeneously distributed. And all contextual attributes are useful for establishing the behavioral value. The original generation does not contain any outliers, but they are added after by using different approaches of  $machine\ learning\ algorithms\ (ML)$  from the Python library  $scikit\-learn^1$ :  $SVC^2$ , K  $Neighbors\ Classifier^2$ ,  $Random\ Forest\ Classifier^3$ ,  $Extra\ Trees\ Classifier^3$ ,  $Gradient\ Boosting\ Classifier^2$  and  $Logistic\ Regression^4$ .

The goal is to detect the most useful points for predictions in order to change their behavioral value. Thus, the chance of having strong outliers which are harmful to the predictions is increased. While a purely random method wouldn't prevent the selection of points that are not very useful for the prediction. For this purpose, the data sets are randomly divided 10 times into a pair of training  $\mathbf{T}$  (70%) and test  $\mathbf{V}$  (30%) sets. For each pairs, the  $\mathbf{ML}$  are trained with  $\mathbf{T}$  and predict the points in  $\mathbf{V}$ . The number of times each point is wrongly predicted are counted. Thus, each point will have an associated counter with a value between 0 and 60 = 10 \* 6 ( $\mathbf{ML}$  algorithms). Those with the lowest values will be the points that have been most frequently correctly predicted. Finally, the n points that will become outliers are the first n with the lowest scores. Where there are more than 2 behavioral values, the new one is chosen randomly.

#### 4.2 Tolerance limit estimation

The tolerance limit is determined by a regression study based on 346 data sets which contain 5% of outliers. Plus, they have distinct characteristics in order to cover a large set of attribute and behavioral value number (see Table 1). The number of points is small in order to perform a large amount of computation in a reasonable time. Since **Hino** works with quantiles of fixed sizes, so the number of points has no influence on the outlier detection.

<sup>&</sup>lt;sup>1</sup> https://scikit-learn.org/stable/index.html

<sup>&</sup>lt;sup>2</sup> All meta-parameters have the default values.

<sup>&</sup>lt;sup>3</sup> All meta-parameters have the default values, except n\_estimator set to 200.

<sup>&</sup>lt;sup>4</sup> All meta-parameters have the default values, except max\_iter set to 1000.

This study consists of determining what is the optimal **tolerance limit** for each data set. It is based on the assumption that there is a relationship between this limit, the numbers of contextual attributes and behavioral values. More concretely, all possible **tolerance limits** are used in order to determine which one is the best. A detection is performed for each limit and a **sensitivity** and **specificity** score is

Points $(p)$	Attributes (a)	Behavioral values (c)		
250	1050, 25	25		
500	1070, 25, 75	25		
1 000	$10 \dots 100, 25, 75$	27		
2500	$10 \dots 100, 25, 75$	28		
5 000	$ 10\dots 100, 25, 75, 125 $	211		

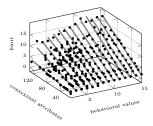
Table 1: data sets used in regression study with a step of 10 for a and 1 for c

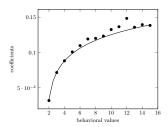
calculated. These scores are determined with the numbers of true positive (**TP**), true negative (**TN**), false positive (**FP**) and false negative (**FN**).

$$sensitivity = \frac{TP}{TP + FN}$$
 (3) 
$$specificity = \frac{TN}{TN + FP}$$
 (4) The **sensitivity** measures the ability of a test to give a positive result when

The **sensitivity** measures the ability of a test to give a positive result when a hypothesis is verified, and it is calculated using equation 3. It corresponds to the percentage of correctly detected outliers. While the **specificity** measures the ability of a test to give a negative result when the hypothesis is not verified, and it is calculated using equation 4. It corresponds to the percentage of undetected points that are actually healthy. These metrics are interpreted together, a good **tolerance limit** must give these scores as close to 100% as possible.

Plus, a good **tolerance limit** must not remove all points of a behavioral value, and it must not remove too many points, e.g. 20%. The optimal limit is selected from the remaining ones, and is the one with the highest sum of **sensitivity** and **specificity**. Thus, those found for each data set are visualized on a 3D graph according to the number of contextual attributes and behavioral values in Figure 2a. This limit is shown to increase with the number of attributes and behavioral values, which confirms the relationship between them.





(a) Graph of linear regression computed on optimal limits obtained

(b) Logarithmic regression of linear's coefficients

(c) Linear function for each behavioral value

Fig. 2: Regression study

From this relation, a second regression is performed in order to determine the slope of each linear line in Figure 2a according to the number behavioral values. Note that the equations of these lines are shown in Figure 2c. As shown in Figure 2b, these slopes do increase with the characteristics of the data set, but logarithmically. Thus, this regression gives the following equation:

$$0.0205 * log_2(-1.623730 + c) + 0.062579$$
(5)

Finally, it is possible to establish an equation that will approach the optimal **tolerance limit** with the number of attributes and behavioral values:

$$[(0.0205 * log_2(-1.623730 + c) + 0.062579) * a]$$
(6)

The result is only an estimation and is not guaranteed to be the optimal **tolerance limit**. This is due to the choices made on the characteristics of the synthetic data sets and the approximations made during the regressions. However, we can be confident that this estimation gives a viable solution.

## 5 Experimentation

The efficiency of the **IQR** (see Section 2), **Isolation Forest** [7], **SVM** [2], **LOF** [4] and **Hino** (see Section 3) methods are compared to their ability to correctly detect outliers on two kinds of data sets: synthetics (different from those in Table 1) and real. Computations have been performed on the supercomputer facilities of the *Mésocentre de calcul de Franche-Comté*. The results obtains are present in a GitHub directory<sup>5</sup>.

### 5.1 On synthetic data sets

An additional 278 data sets with their outliers are created with the same way as those used in the Section 4.1, which has the characteristics present in Table 2. **LOF**, **Isolation Forest** and **SVM** come from the Python library *scikit-learn* and performed with the default settings. Meta-parameters used for **Hino** are already established in Section 4.2.

Outliers detected are compared with the real ones in order to determine a sensitivity and specificity.

Points $(p)$	Attributes (a)	Behavioral values $(c)$			
375	1050	25			
1250	10100	27			
3 000	10100	28			
6000	10100	29			
10 000	10, 50, 75, 100	[2, 3, 5, 7, 9, 11]			
25000	10, 50, 75, 100	[2, 3, 5, 7, 9, 11]			

Table 2: Data sets used in evaluation with a step of 10 for a and 1 for c

Results are given in Figure 3 and the best method is the one that has these two scores closest to 100%. The results conclude:

<sup>&</sup>lt;sup>5</sup> https://github.com/JessyColonval/Hino

- IQR have a linear relationship between these scores, it means that the number of false positives is proportional to the number of outliers detected.
- The whole **sensitivity** of **Isolation forest** remains < 20% (the majority is < 10% or = 0%) while its **specificity** is between 80% and 100%. It means few or no outliers are removed, but it's always removing inliers.
- SVM and LOF preserve correctly the inliers with a score of specificity
   96%. However, they remove a fewer percent of outlier, i.e. < 5%.</li>
- The sensitivity of Hino are < 60% (the majority are between 10 and 40%), while the specificity are > 70% (the majority are > 80%). It means that generally, it remove more outlier than inliers unlike IQR but also a larger number than the Isolation Forest.

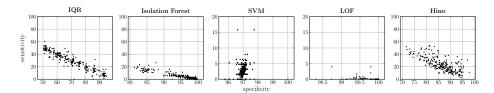


Fig. 3: Specificity and sensitivity on the 5 methods

#### 5.2 On real data sets

In order to confirm the effectiveness of **Hino** compared to state-of-the-art methods, the outlier detection is computed on 16 real data sets from *UCI Machine Learning Repository*<sup>6</sup>. Few are used, due to the difficulty of finding ones with expected characteristics. Their contextual attributes are numerical values and they rep-

		Attributes	Behavioral values	Data sets	Points	Attributes	Behavioral values
annthyroid breastW	7 200 683	21 9	3 2	multiple features	2 000	619	10
cardio	2 126	21	10	musk	6 598	166	2
glass	214	9	6	parkinson	756	752	2
ionosphere	351	33	2	pendigits	10992	16	10
isolet	7 797	617	26	satimage2	6435	36	6
letter	20 000	16	26	shuttle	58 000	9	7
recognition	20 000	10	20	wine	178	13	3
mammo- graphy	11 183	6	2	wine- quality	1 599	11	6

Table 3: Characteristics of real data sets.

resents a panel of different characteristics (see Table 3) in order to observe the behavior of these methods according to these different characteristics.

Two detections are analyzed with **Hino**: one with the meta-parameters calculation described in Section 4.2, and other one with the best solution among all possible **tolerance limits**. The goal is to validate the **tolerance limit** provide

<sup>&</sup>lt;sup>6</sup> https://archive.ics.uci.edu/ml/index.php

by equation 6 and observe the variation with the optimal one. Once detected, these outliers are removed from the data set, and the relevance of this removal is measured using two metrics by using 6 ML algorithms described in Section 4.1:

- Cross-validation measures the prediction performance of these algorithms.
   We assume that the presence of outliers reduces the performance of these algorithms, then it determines which removal was the most beneficial.
- False positive measure the performance of these ML algorithms to correctly predict the outliers detected. We assume that the outliers removed from the training set cannot be correctly predicted by these algorithms..

Table 4 summarizes and presents these results in 2 separate columns:

- Cross-validation cells contain the percentage of cross-validation (1<sup>st</sup> line); the standard deviation (2<sup>nd</sup> line) and tolerance limit (only for Hino) | percentage of outliers | ✓ if all behavioral values are kept, ✗ otherwise (3<sup>rd</sup> line).
   Tolerance limits annotated with a '\*' are those that have been modified to obtain a percentage of outliers < 20%.</li>
- **False positive** cells contain the percentage of points falsely detected as outliers, and the standard deviation of this percentage.

Cases where no outliers have been detected are indicated by the value **None** with the **cross validation** of the full data set. Cells with **Error** imply that the computations of the metrics could not be done because of an overly distorted data set. For example, when a behavioral value has only one representative.

A method is considered better when: it doesn't suppress any behavioral value; the number of points removed isn't too high; the cross-validation is higher and the false positive is lower. Sometimes the choice is not obvious, and the method that is considered better depends on the weight given to these different metrics. These choices are symbolized by two color: **gray** when one solution is better; **light gray** when multiple solutions are closed and better than the others.

According to Table 4, **Hino** gives equivalent or better results for almost all data sets with the estimation of the meta-parameters proposed in Section 4 (column **Hino**). Except for *wine*, where the best solutions are given by **SVM** and **LOF**. But by manually changing the **tolerance limits**, it is possible to get better results (column **Best of Hino**). Thus, out of 16 data sets, the results are clearly improved for 8, remain equivalent for 5 and do not change for the last 3.

#### 5.3 Discussion

The use of **ML** algorithms to compute **false positives** introduces a bias. Indeed, a point detected as an outlier, but correctly predicted, does not necessarily mean that the detection has failed. In a case where too many inliers have been removed from the training set, then it is possible that the **ML** algorithms no longer have enough information to correctly predict the outliers. Thus, these outliers would be wrongly predicted because the training set is too distorted. This situation is less likely when few points are removed, but it should be read carefully and in addition to the other metrics.

	IQR		Isolate	plate Forest SVM		LOF		Hino		Best of Hino		
Data sets	Cross	False	Cross	False	Cross	False	Cross	False	Cross	False	Cross	False
Data sets	validation		validation	positive	validation	positive	validation	positive	validation	positive	validation	positive
	Erre	or	97.68%	94.22%	97.12%	97.98%	98.12%	94.72%	97.65%	94.43%		
annthyroid			±0.19	±0.07	±0.22	$\pm 0.06$	±0.22	±0.12	±0.19	±0.06	No change	
	53.6%		5.54%   🗸		3.01%   🗸		11.22%   🗸		2   10.4%   🗸			
	97.13%	91.66%	98.61%	94.22%	97.09%	99.67%	96.73%	95.08%	98.24%	89.66%	97.61%	81.48%
breastW	±1.12	$\pm 0.36$	±0.29	±0.90	±0.91	±0.0	±1.19	±0.25	±0.86	±0.17	±1.18	±0.00
	28.1%   🗸		37.63%   ✓		7.32%   🗸		24.60%   🗸		2*   12.7%   ✓		3   3.95%   🗸	
	81.4%	60.35%	81.23%	72.71%	82.06%	71.67%	82.45%	71.67%	83.16%	39.72%	84.76%	43.77%
cardio	±1.94	±0.19	±1.26	±0.35	±1.5	±0.3	±1.69	±3.15	±1.30	±0.59	±1.3	±0.57
	56.3%   X	00.4007	16.04%   🗸	40 807	3.15%   🗸	40.4807	1.69%   🗸	10.0007	3   5.9%   🗸	20.2807	2   12.61%   🗸	0.0007
,	73.46% ±4.87	28.12% ±1.09	74.75% ±5.05	48.7% ±1.71	72.28% ±4.66	49.17% ±1.76	73.49% ±5.15	40.98% ±1.25	73.95% ±4.21	29.25% ±1.29	72.16% ±5.55	0.00%
glass		T1.09	16.82%   ✓	T1./1	1.87%   ✓	±1.70	15.89%   🗸	T1.20	2*   9.35%   <b>✓</b>	T1.29	4   0.93%   🗸	±0.00
	36.4%   X 95.12%	59.47%	90.81%	86.54%	91.21%	91.67%	94.23%	55.61%	92.46%	78.02%	92.57%	63.54%
ionosphere	±2.58	±0.43	±1.97	±0.52	±1.88	±0.00	±2.00	±0.72	92.40% ±2.19	±0.47	±2.26	±0.00
ionosphere	48.1%   ✓	TU.43	32.48%	±0.32	2.85%   🗸	±0.00	34.76% ↓ ✓	±0.72	6*   16.8%   ✓	±0.47	9   4.56%   🗸	±0.00
		O.P.	93.1%	93.4%	93.53%	90.55%	93.60%	75.00%	95.08%	76.61%	9   4.30%   🗸	
isolet	Error		±0.49	±0.14	±0.44	±0.25	±0.50	±0.00	95.08% 76.61% ±0.40 ±0.23		No chan	umo.
isolet	96.6	1%	14.21%	10.14	2.98%   ✓	10.20	0.05%   🗸	±0.00	97   13.7%   🗸		No chan	ige
	90,04%	68.37%	91.48%	88.15%	91.64%	92.75%	91.48%	94.52%	None - 91.	72%	91.93%	64.47%
letter	±0.38	±0.23	±0.27	±0.09	±0.32	±0.12	±0.25	±0.23			±0.24	±0.15
recognition	47.4%   🗸		13.57%   🗸		3.01%   🗸		0.19%   🗸		3   0.0%		0   8.16%   🗸	
	99.41%	94.33%	99.26%	91.06%	98.69%	91.23%	98.66%	90.98%	99.85%	83.75%	99.24%	27.28%
mammo-	±0.04	$\pm 0.03$	±0.03	±0.2	±0.12	$\pm 0.16$	±0.00	±0.19	$\pm 0.03$	±0.05	$\pm 0.15$	±0.47
graphy	37.6%   ✓		19.88%   🗸		3.08%   🗸	İ	1.78%   🗸		1*   10.7%   <		3   1.11%   🗸	
multiple	Error	35.3%	99.0%	91.15%	98.29%	90.14%	98.33%	3.33%	99.43%	85.2%	98.7%	25.6%
features		$\pm 0.74$	±0.45	±0.2	±0.45	$\pm 0.54$	±0.47	$\pm 10.54$	$\pm 0.25$	±0.37	$\pm 0.34$	±2.18
icacures	97.65%   X		35.45%   ✓		3.00%   🗸		0.05%   🗸		78   17.9%   🗸		119   0.7%   🗸	
	91.38%	59.22%	96.39%	98.78%	96.45%	99.31%	96.54%	93.17%	98.31%	43.46%	97.47%	22.91%
musk	±4.61	$\pm 0.37$	±0.34	±0.06	±0.27	$\pm 0.25$	±0.35	±0.41	±0.23	±0.31	±0.36	$\pm 0.72$
	97.4%   🗸		6.68%   🗸		2.99%   🗸		0.32%   🗸		6   4.5%   🗸		11   1.55%   🗸	
parkinson	Error		87.04%	86.19%	86.45%	87.03%	86.67%	83.26%	96.93%	9.48%	96.94%	9.47%
			±1.78	±1.14	±1.9	$\pm 1.08$	±1.67	±0.79	±0.80	±0.29	±1.06	±0.29
	100% 98.58%   85.18%		2.78%   🗸	*0 1*07	3.04%	010007	3.04%   🗸	WO 0 107	95*   19.8%   ✓	OH OH07	96   19.71%   🗸	0.0007
10.00	98.58%		98.87%	50.45%	98.52%	94.99%	98.67%	76.34%	98.63%	97.27%	98.58%	0.00%
pendigits	±0.20 4.7%   ✓	$\pm 0.17$	±0.23 51.91%   X	±0.17	±0.16 2.98%   ✓	±0.16	±0.20 1.56%   ✓	±0.45	±0.20 2   16.7%   ✓	±0.04	±0.12 9   0.01%   ✓	±0.00
	89.19%	96.65%	88.48%	58.24%	89.28%	97.56%	89.89%	76.35%	90.04%	88.20%	90.10%	93.89%
satimage2	±0.46	±0.16	±0.75	±4.04	±0.52	±0.06	±0.52	±0.49	±0.48	±0.10	±0.48	±1.76
Satimage2	8.4%   🗸	10.10	19.40%	14.04	2.97%	10.00	1.54%   🗸	10.40	9*   18.9%   🗸	10.10	18   0.05%   🗸	11.70
	99.83%	81.91%	99.78%	60.8%	Erre	or	Erro	or	99.36%	56.97%	99.36%	51.67%
shuttle	±0.05	±0.0	±0.03	±0.04					±0.03	±1.28	±0.03	±4.56
	79.6%   X		15.62%   X	10.01	2.79	%	19.00	0%	6*   0.02%   🗸	1	7   0.01%   🗸	1.00
	98.03%	90.2%	95.96%	95.83%	98.3%	88.12%	98.11%	70.67%	96.94%	95.56%	98.14%	91.67%
wine	±1.66	±0.0	±2.81	±0.00	±1.45	±1.01	±1.54	±1.41	±2.37	±0.0	±1.35	±0.00
	9.6%		13.48%		4.49%   🗸		2.81%		2*   8.4%   🗸		3   1.12%   🗸	
	64.42%	51.85%	63.99%	46.05%	63.48%	35.72%	63.66%	41.31%	63.08%	0.00%		
winequality	±2.17	±0.61	±1.95	±0.57	$\pm 1.85$	$\pm 1.16$	±1.69	±1.02	±2.15	±0.00	No chan	ige
	25.2%   ✓		16.45%   🗸		2.88%   ✓		2.94%   ✓		7*   1.06%   $\checkmark$			

Table 4: Detection methods comparison on real data sets

The way the **tolerance limit** equation is established in Section 4.2 is strongly related to the maximum percentage of inlier that is agreed to be removed (i.e. 20%). A different hypothesis would change this equation and the resulting detection. Plus, for the preservation of data sets, the algorithm shifts these limit when the maximum number of outliers is reached or when a behavioral value is completely removed. In practical, in Table 4, the limits of 9 detections has been shifted: 7 cause the maximum number of outliers and 2 cause the integrity of behavioral values. Despite this precaution, the **tolerance limit** doesn't always provide the optimal solution. As Table 4 shows, there is a much better solution in 7 over 16 cases. Thus, the **tolerance limit** suit better as a meta-parameter.

# 6 Conclusion

This paper presented a parameterized method of outlier detection based on the quantile principle and a calculation to determine the parameters to be used. This method was compared with a state-of-the-art method (**IQR**, **Isolation** 

Forest, SVM and LOF) on 278 synthetic data sets and 16 real data sets. It was clearly shown that the **Hino** approach was the more efficient. The computational efficiency of meta-parameters for **Hino** has been discussed and can be studied to adapt it to take into account the context of data sets. Finally, this approach can be easily integrated into standard machine learning library (as scikit-Learn or Tensor Flow) to help preparation of data sets.

### References

- Aggarwal, C.C.: Outlier Analysis. Springer International Publishing, Cham, 2 edn. (2017). https://doi.org/10.1007/978-3-319-47578-3
- Amer, M., Goldstein, M., Abdennadher, S.: Enhancing one-class support vector machines for unsupervised anomaly detection. In: Proceedings of the ACM SIGKDD Workshop on Outlier Detection and Description. pp. 8–15. ODD '13, Association for Computing Machinery, New York, NY, USA (Aug 2013). https://doi.org/10.1145/2500853.2500857
- 3. Barnett, V., Lewis, T., others: Outliers in statistical data, vol. 3. Wiley New York (1994). https://doi.org/10.1057/jors.1995.142
- Breunig, M., Kriegel, H.P., Ng, R., Sander, J.: LOF: Identifying Density-Based Local Outliers. vol. 29, pp. 93–104 (Jun 2000). https://doi.org/10.1145/342009.335388
- Dekking, F.M., Kraaikamp, C., Lopuhaä, H.P., Meester, L.E.: Exploratory data analysis: numerical summaries. In: Dekking, F.M., Kraaikamp, C., Lopuhaä, H.P., Meester, L.E. (eds.) A Modern Introduction to Probability and Statistics: Understanding Why and How, pp. 231–243. Springer, London (2005). https://doi.org/10.1007/1-84628-168-7 16
- 6. Guyon, I.: Design of experiments for the NIPS 2003 variable selection benchmark. In: undefined (2003). https://doi.org/10.1007/978-3-540-35488-8\_10, https://www.semanticscholar.org/paper/Design-of-experiments-forthe-NIPS-2003-variable-Guyon/b979fa88ca448fb08633f961131f45214b1cf109
- Liu, F.T., Ting, K.M., Zhou, Z.H.: Isolation-Based Anomaly Detection. ACM Transactions on Knowledge Discovery from Data 6(1), 3:1–3:39 (Mar 2012). https://doi.org/10.1145/2133360.2133363
- 8. Sadik, S., Gruenwald, L.: Research issues in outlier detection for data streams. ACM SIGKDD Explorations Newsletter **15**(1), 33–40 (Mar 2014). https://doi.org/10.1145/2594473.2594479
- Souiden, I., Omri, M.N., Brahmi, Z.: A survey of outlier detection in high dimensional data streams. Computer Science Review 44, 100463 (May 2022). https://doi.org/10.1016/j.cosrev.2022.100463, https://www.sciencedirect.com/science/article/pii/S1574013722000107
- Thung, F., Wang, S., Lo, D., Jiang, L.: An Empirical Study of Bugs in Machine Learning Systems. In: 2012 IEEE 23rd International Symposium on Software Reliability Engineering. pp. 271–280 (Nov 2012). https://doi.org/10.1109/ISSRE.2012.22, iSSN: 2332-6549
- 11. Vinutha, H.P., Poornima, B., Sagar, B.M.: Detection of Outliers Using Interquartile Range Technique from Intrusion Dataset. In: Satapathy, S.C., Tavares, J.M.R., Bhateja, V., Mohanty, J.R. (eds.) Information and Decision Sciences. pp. 511–518. Springer, Singapore (2018). https://doi.org/10.1007/978-981-10-7563-6 53