

# High Dimensional Outliers

Leland Wilkinson \*

June 22, 2016

## Abstract

Outliers have more than two centuries' history in the field of statistics. Recently, they have become a focal topic because of their relevance to terrorism, network intrusions, financial fraud, and other areas where rare events are critical to understanding a process. This paper presents a new algorithm, called `hdoutliers`, for detecting multidimensional outliers. It is unique for a) dealing with a mixture of categorical and continuous variables, b) dealing with the curse of dimensionality (many columns of data), c) dealing with many rows of data, d) dealing with outliers that mask other outliers, and e) dealing consistently with unidimensional and multidimensional datasets. Unlike ad hoc methods found in many machine learning papers, `hdoutliers` is based on a distributional model that allows outliers to be tagged with a probability.

---

\*H2O.ai and Department of Computer Science, University of Illinois at Chicago.

E-mail: leland@h2o.ai.

This research was supported by NSF/DHS grant DMS-FODAVA-0808860

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Related Work</b>	<b>5</b>
2.1	Univariate Outliers . . . . .	5
2.1.1	The Distance from the Center Rule . . . . .	5
2.1.2	The Box Plot Rule . . . . .	6
2.1.3	The Gaps Rule . . . . .	8
2.2	Multivariate Outliers . . . . .	9
2.2.1	Mahalanobis Distance . . . . .	9
2.2.2	Multivariate Gap Tests . . . . .	10
2.2.3	Clustering . . . . .	10
<b>3</b>	<b>A New Multivariate Outlier Algorithm</b>	<b>11</b>
3.1	The Algorithm . . . . .	11
3.2	Comments on the Algorithm . . . . .	12
3.3	Validation . . . . .	14
3.3.1	False Positives . . . . .	14
3.3.2	False Negatives . . . . .	14
<b>4</b>	<b>Visualization</b>	<b>16</b>
4.1	Visualizing Unidimensional Outliers . . . . .	16
4.2	Low-dimensional visualizations cannot be used to discover multivariate outliers . . . . .	17
4.3	Using statistical algorithms to highlight outliers in visualizations . . . . .	18
4.3.1	Parallel Coordinates . . . . .	18
4.3.2	Regression Residuals . . . . .	19
4.3.3	Time Series Outliers . . . . .	20
4.3.4	Ipsative Outliers . . . . .	21
4.3.5	Text Outliers . . . . .	22
4.4	Graph Outliers . . . . .	23
4.4.1	Scagnostics Outliers . . . . .	25
4.4.2	Geographic Outliers . . . . .	26
<b>5</b>	<b>Conclusions</b>	<b>27</b>

## List of Figures

1	Why ordinary box plots cannot be used to discover probable outliers. . . .	7
2	Dot plots of small batches of data with comparable means and standard deviations. . . . .	9
3	Mahalanobis Distance. . . . .	10
4	hdoutliers applied to data shown in Figure 2. . . . .	15
5	Inlier datasets. . . . .	16
6	German Research Center for Artificial Intelligence (DFKI) dataset. . . .	16
7	Dot plot and normal probability plot of residuals from a two-factor experiment. . . . .	17
8	Why lower-dimensional projections cannot be used to discern outliers. . .	18
9	Parallel coordinates plot of five variables from the Adult dataset in the UCI data repository. . . . .	19
10	Ordinary Least Squares regression and Least Median of Squares regression on attributes of stars. . . . .	20
11	Snow cover at a Greenland weather station. . . . .	21
12	US unemployment series outliers. . . . .	21
13	Baseball player profile. . . . .	22
14	Document outliers. . . . .	23
15	Les Misérables characters network graph. . . . .	24
16	Karate Club graph and comparably scaled random minimum spanning tree graphs. . . . .	25
17	Scatterplot outliers based on Scagnostics. . . . .	26
18	Marriage and Divorce rates in the US. . . . .	26

# 1 Introduction

According to Hawkins [28], “An outlier is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism”. The modern history of outlier detection began with concerns over combining astronomical observations [56, 4]. The prevailing concern among many scientists from the earliest times was over how much outliers could bias estimates of location and spread. Statistically-based rules for outlier rejection emerged as early as the late 18th century.

There are two predominant reasons for analysts’ longstanding interest in outliers. The first is to identify cases that can bias estimates of statistical models. The second is to locate extreme cases in a distribution because such values may be especially (or solely) interesting. The first reason is not a good approach to the estimation problem. Outliers should not be eliminated from model fits unless a clear reason for their occurrence is available. Furthermore, there are numerous robust versions of classical models that automatically downweight outliers without introducing substantial bias [25]. The second reason overlooks the possibility of inliers [29]. These are unusual cases found in the middle of mixtures of distributions.

This paper is concerned with the interplay of visual methods and outlier detection methods. It is not an attempt to survey the vast field of outlier detection or to cover the full range of currently available methods. For general introductions, see the references at the beginning of the Related Work section below.

The contributions in this paper are:

- We demonstrate why the classical definition of an outlier (a large distance of a point from a central location estimate (mean, median, etc.) is unnecessarily restrictive and often involves a circularity.
- We introduce a new algorithm, called `hdoutliers`, for multidimensional outliers on  $n$  rows by  $p$  columns of data that addresses the curse of dimensionality (large  $p$ ), scalability (large  $n$ ), categorical variables, and non-normal distributions. This algorithm is designed to be paired with visualization methods that can help an analyst explore unusual features in data.
- We demonstrate why visual analytic tools cannot be used to detect multidimensional outliers.
- We introduce some novel applications of outlier detection and accompanying visualizations based on `hdoutliers`.

## 2 Related Work

There are several excellent books on outliers written by statisticians [4, 28, 52, 57]. Statisticians have also written survey papers [24, 34, 2]. Computer scientists have written books and papers on this topic as well [1, 11, 31]. The latter include surveys of the statistical sources.

### 2.1 Univariate Outliers

The detection of outliers in the observed distribution of a single variable spans the entire history of outlier detection. It spans this history not only because it is apparently the oldest formulation of the problem, but also is the focus of relatively recent research on outliers.

#### 2.1.1 The Distance from the Center Rule

The word *outlier* implies lying at an extreme end of a set of ordered values – far away from the center of those values. The modern history of outlier detection emerged with methods that depend on a measure of centrality and a measure of distance from that measure of centrality. As early as the 1860’s, Chauvenet (cited in [4]) judged an observation to be an outlier if it lies outside the lower or upper  $1/(4n)$  points of the Normal distribution. Barnett and Lewis [4] document many other early rules that depend on the Normal distribution but fail to distinguish between population and sample variance.

Grubbs [23], in contrast, based his rule on the sample moments of the Normal:

$$G = \frac{\max_{1 \leq i \leq n} |x_i - \bar{x}|}{s}$$

where  $\bar{x}$  and  $s$  are the sample mean and standard deviation, respectively.

Grubbs referenced  $G$  against the  $t$  distribution in order to spot an upper or lower outlier:

$$G > \frac{n-1}{\sqrt{n}} \sqrt{\frac{t_{\alpha/(2n), n-2}^2}{n-2 + t_{\alpha/(2n), n-2}^2}}$$

If one knows that the values on a variable are sampled randomly from a Normal distribution and if the estimates of location and scale are unbiased and if one wishes to detect only the largest absolute outlier, it is a valid test.

Unfortunately, the usual sample estimates of the mean and standard deviation are not robust against outliers. So we have a circularity problem. We assume a null distribution (say, the Normal), estimate its parameters, and then use those estimates to test whether a point could have plausibly come from that distribution. But if our alternative hypothesis is that it doesn't (the usual case), then the outlier should not be included in the estimation. Barnett and Lewis [4] discuss this problem in more detail, where they distinguish *inclusive* and *exclusive* methods. They, as well as [52], also discuss robust estimation methods for overcoming this circularity problem.

Barnett and Lewis discuss other detection methods for non-Normal distributions. The same principals apply in these cases, namely, that the sample is random, the population distributions are known and that the parameter estimates are unbiased.

### 2.1.2 The Box Plot Rule

A box-plot graphically depicts a batch of data using a few summary statistics called *letter values* [58, 19]. The letter values in Tukey's original definition are the median and the *hinges* (medians of the upper and lower halves of the data). The hinge values correspond closely, but not necessarily, to the lower quartile (Q1) and the upper quartile (Q3). Tukey called the difference between the hinges the *Hspread*, which corresponds closely to the quantity  $Q3 - Q1$ , or the Inter Quartile Range (IQR). In Tukey's version of the box-plot (see the upper panel of Figure 1), a box is drawn to span the *Hspread*. The median is marked inside the box. Whiskers extend from the edges of the box to the farthest upper and lower data points (*Adjacent values*) inside the so-called *inner fences*. The upper inner fence is the

$$upperhinge + 1.5 \times Hspread$$

and the lower inner fence is the

$$lowerhinge - 1.5 \times Hspread$$

Any data point beyond the *Adjacent values* is plotted as an outlying point.<sup>1</sup>

---

<sup>1</sup>Few statistics packages produce box plots according to Tukey's definition [19]. Surprisingly, the `boxplot` function in the core *R* package does not, despite its ancestry inside Tukey's group at Bell Laboratories.

Tukey designed the box plot (he called it a *schematic plot*) to be drawn by hand on a small batch of numbers. The whiskers were designed not to enable outlier detection, but to locate the display on the interval that supports the bulk of the values. Consequently, he chose the *Hspread* to correspond roughly to three standard deviations on normally distributed data. This choice led to two consequences: 1) it doesn't apply to skewed distributions, which constitute the instance many advocates think is the best reason for using a box plot in the first place, and 2) it doesn't include sample size in its derivation, which means that the box plot will falsely flag outliers on larger samples. As Dawson [14] shows, "regardless of size, at least 30% of samples drawn from a normally-distributed population will have one or more data flagged as outliers." The top panel of Figure 1 illustrates this problem for a sample of 100,000 Normally distributed numbers. Thousands of points are denoted as outliers in the display.

To deal with the skewness problem, Hubert and Vandervieren [33] and others have suggested modifying the fences rule by using a robust estimate of skewness. By contrast, Tukey's approach for this problem involved transforming the data through his *ladder of powers* [58] before drawing the box plot.

The Letter-Value-Box-Plot [32] was designed to deal with the second problem. The authors compute additional letter values (splitting the splits) until a statistical measure of fit is satisfied. Each letter-value region is represented by a rectangle. The lower panel of Figure 1 shows the result. On the same 100,000 Normal variables, only two points are identified as outliers.

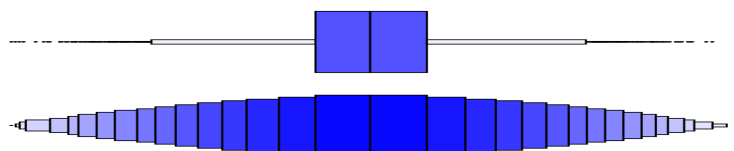


Figure 1: Outliers revealed in a box plot [58] and letter values box plot [32]. These plots are based on 100,000 values sampled from a Gaussian (Standard Normal) distribution. By definition, the data contain no probable outliers, yet the ordinary box plot shows thousands of outliers. This example illustrates why ordinary box plots cannot be used to discover probable outliers.

### 2.1.3 The Gaps Rule

Suppose that we do not know the population distribution and suppose, further, that our idea of outliers is that they do not belong to the generating distribution we suspect underlies our data. Figure 2 shows two dotplots of batches of data that have the same mean and standard deviation. Absent knowledge of the parametric distribution, we cannot apply the usual outlier detection algorithms. Furthermore, we are more inclined to say the the largest point in the right dot plot is an outlier, whereas the largest point in the left plot, having the same score, is not.

A simple example emphasizes this point. Suppose we give a test to 100 students and find the mean score is 50 and the standard deviation is 5. Among these students, we find one perfect score of 100. The next lower score is 65. We might be inclined to suspect the student with a score of 100 is a genius or a cheat. And if there were three students with perfect scores in this overall distribution, we might suspect cheating even more. On the other hand, if the perfect score is at the top of a chain of scores spaced not more than 5 points apart, we might be less suspicious. Classical outlier tests would not discriminate among these possibilities.

These considerations and others led to a different criterion for discovering outliers. Namely, we should look for *gaps* (*spacings*) between the ordered values rather than extremities. A consequence of this point of view is that we can identify unusual scores in the *middle* of distributions as well as in the extremes, as long as they are separated from other scores by a large gap.

Dixon [15] headed in this direction by developing an outlier test based on the gap between the largest point and the second largest point, standardized by the range of scores. His test was originally based on a normal distribution, but in subsequent publications, he developed nonparametric variations. Dixon tabulated percentage points for a range of  $Q$  statistics.

$$Q = \frac{x_n - x_{n-1}}{x_n - x_1}$$

Tukey [58] considered the more general question of identifying gaps anywhere in a batch of scores. Wainer and Schacht [60] adapted Tukey's gapping idea for a version of the test that weighted extreme values more than middle ones. They derived an approximate  $z$  score that could be used to test the significance of gaps.

Burridge and Taylor [9] developed an outlier test based on the extreme-value distribution of gaps between points sampled from the Exponential family of distributions:



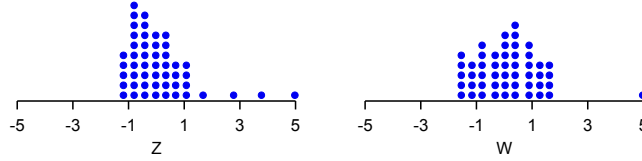


Figure 2: Dot plots of small batches of data with comparable means and standard deviations.

$$f(x_i; \theta_i, \phi) = \exp \left[ \frac{x\theta - a(\theta)}{b(\phi)} + c(x, \phi) \right]$$

where  $x$  is a scalar or vector,  $\theta$  is a scalar or vector of parameters,  $\phi$  is a scale parameter, and  $a(\cdot), b(\cdot), c(\cdot)$  are functions. This family of mathematical distributions is quite large (including the Normal, Exponential, Gamma, Beta, Bernoulli, Poisson, and many others).

## 2.2 Multivariate Outliers

### 2.2.1 Mahalanobis Distance

The squared Mahalanobis distance ( $D^2$ ) of a multidimensional point  $\mathbf{x}$  from the centroid of a multivariate Normal distribution described by covariance matrix  $\Sigma$  and centroid  $\mu$  is

$$D^2 = (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu)$$

Figure 3 shows how this works in two dimensions. The left panel shows a bivariate normal distribution with level curves inscribed at different densities. The right panel shows the same level curves as horizontal slices through this mountain. Each is an ellipse. Distance to the centroid of the ellipses is measured differently for different directions. The weights are determined by the covariance matrix  $\Sigma$ . If  $\Sigma$  is an identity matrix, then  $D^2$  is equivalent to squared Euclidean distance.

The squared distance in the above formula is a chi-square variate (a member of the Gamma distribution family). This means that, if the assumption of Normality is met,  $D^2$  can be tested against a chi-square distribution with  $p$  degrees of freedom. As with univariate outlier tests based on a Normality assumption, this test is valid if the assumption of multivariate Normality is met. Unfortunately,

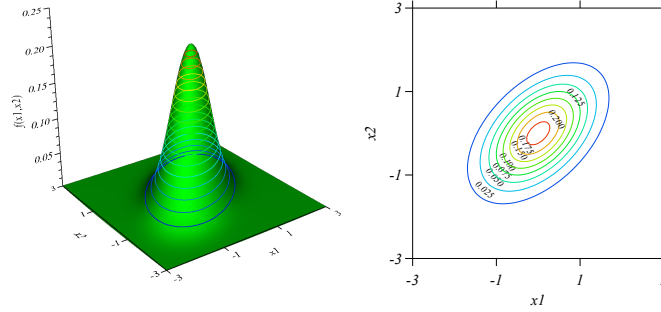


Figure 3: Mahalanobis Distance. The left panel shows a bivariate Normal distribution. The right panel shows level curves for that distribution. Each curve corresponds to a value of  $D^2$ .

this is seldom true for real data and, furthermore, estimates of the covariance matrix and centroid are far from robust. Consequently, this outlier test has limited applicability.

Rousseeuw and Van Zomeren [51] introduce a robust Mahalanobis Distance estimator that can be used to overcome some of these problems. Ram Gnanadesikan [21] discusses applications of Gamma probability plots to these multivariate problems. They can be interpreted similarly to the way univariate probability plots are interpreted.

### 2.2.2 Multivariate Gap Tests

Multivariate data do not have a simple ordering for computing gaps between adjacent points. There have been several attempts at getting around this problem. Rohlf [49] proposed using the edge lengths of the geometric minimum spanning tree (MST) as a single distribution measure. Assuming these edges follow a gamma distribution, one could construct a gamma probability plot on them or examine the upper tail for judgments on outliers. There are problems with this method, however, when variates are correlated [10]. Similar methods based on the MST have been proposed [44, 48], but they suffer from the same problem.

### 2.2.3 Clustering

A popular multivariate outlier detection method has been to cluster the data and then look for any points that are far from their nearest cluster centroids [66, 35,

47, 36]. This method works reasonably well for moderate-size datasets with a few singleton outliers. Most clustering algorithms do not scale well to larger datasets, however.

A related approach, called Local Outlier Factor (LOF) [8], is similar to density-based clustering. Like DBSCAN clustering [17], it is highly sensitive to the choice of input parameter values.

Most clustering methods are not based on a probability model (see [18] for an exception) so they are susceptible to false negatives and false positives. We will show one remedy in Section 3.3.2.

### 3 A New Multivariate Outlier Algorithm

The new algorithm `hdoutliers` is designed to meet several criteria at once:

- It allows us to identify outliers in a mixture of categorical and continuous variables.
- It deals with the curse of dimensionality by exploiting random projections for large  $p$  (number of dimensions).
- It deals with large  $n$  (number of points) by exploiting a one-pass algorithm to compress the data.
- It deals with the problem of *masking* [4], in which clusters of outlying points can elude detection by traditional methods.
- It works for both single-dimensional and multi-dimensional data.

#### 3.1 The Algorithm

1. If there are any categorical variables in the dataset, convert each categorical variable to a continuous variable by using Correspondence Analysis [22].
2. If there are more than 10,000 columns, use random projections to reduce the number of columns to  $p = 4 \log n / (\epsilon^2/2 - \epsilon^3/3)$ , where  $\epsilon$  is the error bound on squared distances.
3. Normalize the columns of the resulting  $n$  by  $p$  matrix  $X$ .
4. Let  $\text{row}(i)$  be the  $i$ th row of  $X$ .

5. Let  $radius = .1/(\log n)^{1/p}$ .
6. Initialize *exemplars*, a list of exemplars with initial entry  $[row(1)]$ .
7. Initialize *members*, a list of lists with initial entry  $[1]$ ; each exemplar will eventually have its own list of affiliated member indices.
8. Now do one pass through  $X$ :

```

forall the  $row(i), i = 1, \dots, n$  do
     $d$  = distance to closest exemplar, found in exemplars( $j$ )
    if  $d < radius$  then
        | add  $i$  to members( $j$ )list
    else
        | add  $row(i)$  to exemplars
        | add new list to members, initialized with  $[i]$ 
    end
end

```

9. Now compute nearest-neighbor distances between all pairs of exemplars in the *exemplars* list.
10. Fit an Exponential distribution to the upper tail of the nearest-neighbor distances and compute the upper  $1 - \alpha$  point of the fitted cumulative distribution function (CDF).
11. For any exemplar that is significantly far from all the other exemplars based on this cutpoint, flag all entries of *members* corresponding to *exemplar* as outliers.

### 3.2 Comments on the Algorithm

1. Correspondence Analysis (CA) begins by representing a categorical variable with a set of dummy codes, one code (1 or 0) for each category. These codes comprise a matrix of 1's and 0's with as many columns as there are categories on that variable. We then compute a principal components decomposition of the covariance matrix of the dummy codes. This analysis is done separately for each of  $k$  categorical variables in a dataset. CA scores

on the rows are computed on each categorical variable by multiplying the dummy codes on that row's variable times the eigenvectors of the decomposition for that variable. <sup>2</sup>

2. The Johnson-Lindenstrauss lemma [37] states that if a metric on  $X$  results from an embedding of  $X$  into a Euclidean space, then  $X$  can be embedded in  $R^p$  with distortion less than  $1 + \epsilon$ , where  $p \sim O(\epsilon^2 \log n)$ . Remarkably, this embedding is achieved by projecting onto a  $p$ -dimensional subspace using random Gaussian coefficients. Because our algorithm depends only on a similarity transformation of Euclidean distances, we can logarithmically reduce the complexity of the problem through random projections and avoid the curse of dimensionality. The number of projected columns based on the formula in this step was based on  $\epsilon = .2$  for the analyses in this paper. The value 10,000 is the lower limit for the formula's effectiveness in reducing the number of dimensions when  $\epsilon = .2$ .
3.  $X$  is now bounded by the unit (hyper) cube.
4. A *row* represents a  $p$ -dimensional vector in a finite vector space.
5. The value of *radius* is designed to be well below the expected value of the distances between  $n(n - 1)/2$  pairs of points distributed randomly in a  $p$  dimensional space.
6. The *exemplars* list contains a list of row values representing clusters of points.
7. The *members* list of lists contains one list of indices for each exemplar that point to rows represented by that exemplar.
8. The Leader algorithm [27] in this step creates clusters in one pass through the data. It is equivalent to centering balls in  $p$  dimensional space on points considered to be exemplars. Unlike  $k$ -means clustering, the Leader algorithm centers clusters on actual data points rather than on centroids and it involves only one pass through the data. In rare instances, the resulting clusters can be dependent on the order of the data, but not enough to affect the identification of outliers because of the large number of clusters produced. We are characterizing a high-dimensional density by covering it with many small balls.

9. The number of clusters resulting from *radius* applied even to large numbers of data points is small enough to allow the simple brute-force algorithm for finding nearest neighbors.
10. We use a modification of the Burridge and Taylor [9] algorithm due to Schwarz [55]. For all examples in this paper,  $\alpha$  (the critical value) was set to .05.
11. Flagging all members of an outlying cluster means that this algorithm can identify outlying sets of points as well as outlying singletons.

### 3.3 Validation

We validate `hdoutliers` by examining its performance with regard to 1) false positives and 2) false negatives. If the claims for the algorithm are true, then we should expect it 1) to find outliers in random data not more than  $100\alpha$  percent of the time and 2) not to miss outliers when they are truly due to mixtures of distributions or anomalous instances.

#### 3.3.1 False Positives

- Table 1 contains results of a simulation using random distributions. The entries are based on 1,000 runs of `hdoutliers` on normally distributed variables with  $\alpha$  (the critical value) set to .05. The entries show that `hdoutliers` is generally conservative.
- The results were similar for random Exponential and Uniform variables.

#### 3.3.2 False Negatives

- Figure 4 is based on the dataset in Figure 2. The `hdoutliers` identifies the outlier in the right dot plot but finds none in the left.
- Figure 5 shows that `hdoutliers` correctly identifies the inlier in the center of both one-dimensional and two-dimensional configurations.

---

<sup>2</sup>Computing the decomposition separately for each categorical variable is equivalent to doing an MCA separately for each variable instead of pooling all the categorical variable dummy codes into one matrix.

Table 1: Empirical level of `hdoutliers` test based on null model with Gaussian variables and critical value  $\alpha = .05$ .

	p=1	p=5	p=10	p=100
n=100	.011	.040	.018	.012
n=500	.015	.035	.027	.020
n=1000	.017	.045	.027	.024

- Figure 6 is based on the `dfki` dataset in [20]. The left panel shows what the authors consider to be outliers. The right panel shows the result of implementing `hdoutliers` inside a k-means clustering. On each iteration of the k-means algorithm, we apply `hdoutliers` to the points assigned to each cluster in order to determine if any points belonging to their nearest cluster should be treated as outliers. The outliers are then left out of the computation of the cluster centroids.
- Table 2 shows that `hdoutliers` correctly identifies the outlier in a table defined by two categorical variables. The data consist of two columns of strings, one for  $\{A,B,C,W\}$  and one for  $\{A,B,C,X\}$ . There is only one row with the tuple  $\langle W,X \rangle$ . The `hdoutliers` also handles mixtures of categorical and continuous variables.

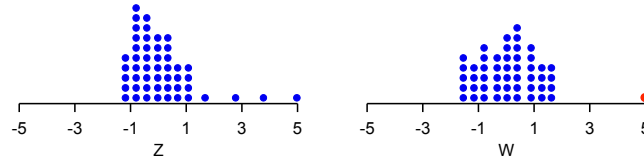


Figure 4: The `hdoutliers` algorithm applied to data shown in Figure2.

Table 2: Crosstab with an outlier (red entry)

	A	B	C	X
A	100	0	0	0
B	0	100	0	0
C	0	0	100	0
W	0	0	0	1

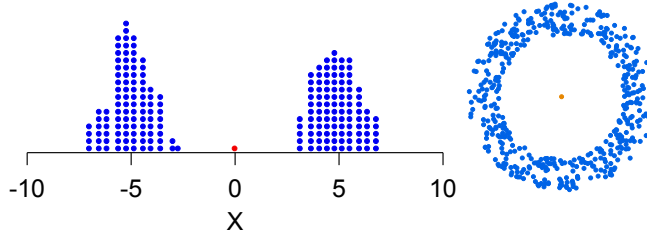


Figure 5: Inlier datasets; `hdoutliers` correctly identifies the inliers.

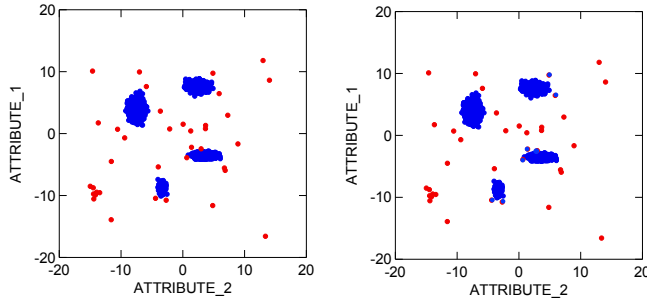


Figure 6: Test dataset from [20]. The left plot shows what the authors consider to be outliers and the right plot is the result produced by `hdoutliers` inside a k-means clustering. The outliers are colored red in both plots.

## 4 Visualization

This section covers various probability-based methods for visualizing outliers. The main point in all these examples is that a statistical algorithm based on probability theory is necessary for reliably discovering outliers but visualizations are necessary for interpreting the results of these discoveries.

### 4.1 Visualizing Unidimensional Outliers

For univariate outlier detection, histograms, probability plots [12], and dot plots [61] are most useful. Figure 7 shows a dot plot and normal probability plot of residuals from a two-factor experiment. In these probability plots, we look for major distortions from a straight line. A probability plot can be constructed from any parametric distribution for which a cumulative distribution function can be



computed. They are widely used in experimental design and analysis of residuals.

Even though these univariate displays can be helpful in exploratory analysis to detect outliers, they do not yield the kind of risk estimate that `hdoutliers` or the parametric methods described in the Related Work sections provide. Without a risk estimate, the chance of false discoveries is uncontrolled. In practical terms, we might see terrorists in groups where none exist. Thus, as in Figure 7, it is helpful to highlight outliers using a statistical method like `hdoutliers`. This approach will also help with false negatives, where significant outliers may not be visually salient.

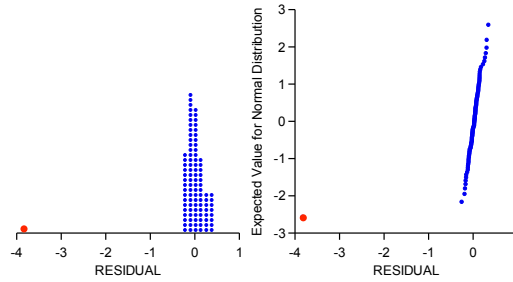


Figure 7: Dot plot and normal probability plot of residuals from a two-factor experiment. One lower outlier is evident.

## 4.2 Low-dimensional visualizations cannot be used to discover multivariate outliers

There have been many outlier identification proposals based on looking at axis-parallel views or low-dimensional projections (usually 2D) that are presumed to reveal high-dimensional outliers (e.g., [39, 30, 38, 42]). This approach is infeasible. Figure 8 shows why. The data are samples from a multivariate Normal distribution. The left panel plot illustrates the problem for two dimensions. The figure incorporates a 95 percent joint confidence ellipse based on the sample distribution of points. Two points are outside this ellipse. The red point on the left is at the extreme of both marginal histograms. But the one on the right is well inside both histograms. Examining the separate histograms would fail to identify that point.

The right panel plot shows the situation for three dimensions. The three marginal 2D plots are shown as projections onto the facets of the 3D cube. Each

confidence ellipse is based on the pairwise plots. The red outlying point in the joint distribution is inside all three marginal ellipses. The 2D scatterplots fail to reveal the 3D outlier. The situation gets even worse in higher dimensions.

Some authors have proposed methods for finding low-dimensional views based on projection pursuit or ad hoc projections (e.g., [7, 53]). This approach is relatively ineffective for visualizing outliers in higher-dimensional datasets because many projections are required to discriminate outliers. Furthermore, most outlier-seeking projection methods are impractical on large datasets.

Parallel coordinates, another 2D display, have been advocated for the purpose of outlier detection [46]. Figure 9 shows why this is infeasible.

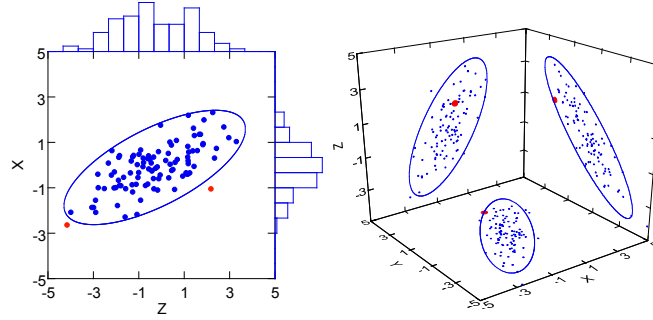


Figure 8: 2D (left) and 3D (right) joint outliers. The figures show why lower-dimensional projections cannot be used to discern outliers.

### 4.3 Using statistical algorithms to highlight outliers in visualizations

While visualizations cannot be used to detect multidimensional outliers, they are invaluable for inspecting and understanding outliers detected by statistical methods. This section covers a variety of visualizations that lend themselves to outlier description.

#### 4.3.1 Parallel Coordinates

As mentioned in the last section, parallel coordinates cannot be used to discover outliers. Figure 9 shows parallel coordinates on four variables from the Adult dataset in the UCI dataset repository [41]. The algorithm discovered two outliers

out of 32,561 cases. The profiles appear to run through the middle of the densities even though they are multivariate outliers. Although parallel coordinates are generally useless for discovering outliers, they can be useful for inspecting outlier profiles detected by a statistical algorithm.

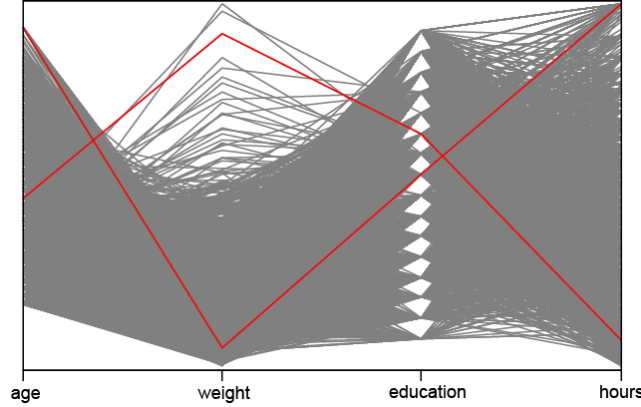


Figure 9: Parallel coordinates plot of five variables from the Adult dataset in the UCI data repository. The red profiles are multivariate outliers.

#### 4.3.2 Regression Residuals

The conventional statistical wisdom for dealing with outliers in a regression context is to examine residuals using a variety of diagnostic graphics and statistics [3, 5, 13]. Following this advice is critical before promoting any particular regression model on a dataset. It is a necessary but not sufficient strategy, however. The reason is that some outliers have a high influence on the regression and can pull the estimates so close to them that they are masked.

Figure 10, derived from an example in [52], shows how this can happen in even the simplest bivariate regression. The data are measurements of light intensity and temperature of a sample of stars. In the left panel, the ordinary least squares (OLS) regression line is pulled down by the four outliers in the lower right corner, leaving a bad fit to the bulk of the points. We would detect most, but not all, of the outliers in a residual plot. The right pane, based on a least median of squares regression (LMS) [50], shows six red points as regression outliers. They are, in fact, dwarf stars.

There are numerous robust regression models, but LMS has the lowest breakdown point against outliers [16]. Therefore, the most prudent approach to regression modeling is to compute the fit both ways and see if the regression coefficients and residual plots differ substantially. If they do, then LMS should be the choice. Otherwise, the simpler OLS model is preferable.

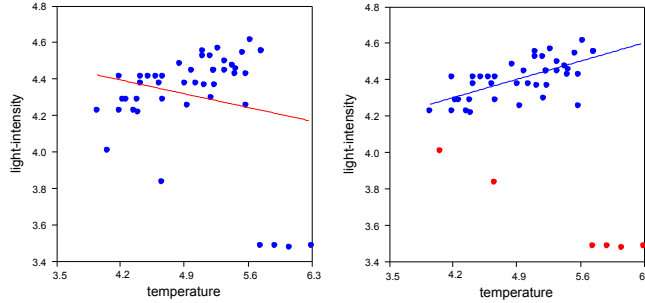


Figure 10: Ordinary Least Squares (OLS) regression (left panel) and Least Median of Squares (LMS) regression (right panel) on attributes of stars. Data are from [52]

### 4.3.3 Time Series Outliers

Detecting time series outliers requires some pre-processing. In particular, we need to fit a time series model and then examine residuals. Fitting parametric models like ARIMA [6] can be useful for this purpose, but appropriate model identification can be complicated. A simpler approach is to fit a nonparametric smoother. The example in Figure 11 was fit by a kernel smoother with a biweight function on the running mean. The data are measurements of snowfall at a Greenland weather station, used in [62]. The outliers (red dots) are presumably due to malfunctions in the recording equipment.

Computing outlying series for multiple time series is straightforward with the `hdoutliers` algorithm. We simply treat each series as a row in the data matrix. For  $n$  series on  $p$  time points, we have a  $p$ -dimensional outlier problem. Figure 12 shows series for 20 years of the Bureau of Labor Statistics Unemployment data. The red series clearly indicate the consequences of the Great Recession. This example illustrates why a probability-based outlier method is so important. We could rank the series by their average levels of unemployment or use one of the other ad-hoc multidimensional outlier detectors, but we would have no way of knowing how many at the top are significant outliers.

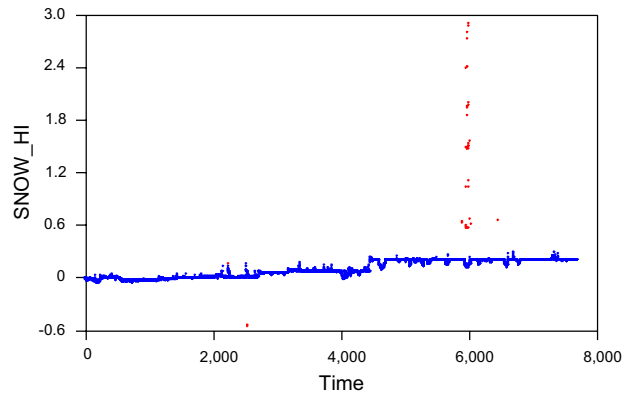


Figure 11: Outlying measurements of snow cover at a Greenland weather station.

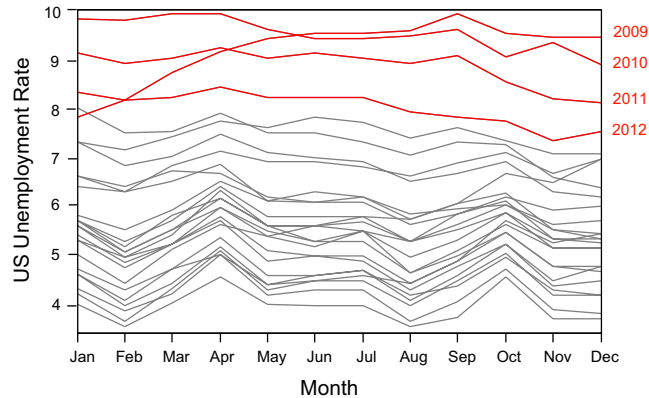


Figure 12: US unemployment series outliers. The shock and ensuing recovery from the Great Recession is clearly indicated in the outliers.

#### 4.3.4 Ipsative Outliers

An *ipsative* outlier is a case that is an outlier with respect to itself. That is, we standardize values within each case (row) and then look for outliers in each standardized profile. Any profile with an outlier identified by `hdoutliers` is considered noteworthy; in other words, we can characterize a person simply by referring to his outliers. It is easiest to understand this concept by examining a graphic. Figure 13 shows an outlying profile for a baseball player who is hit by pitches more frequently than we would expect from looking at his other characteristics. This player may not be hit by pitches significantly more than other players, however. We are instead interested in a player with a highly unusual profile that can

be described simply by his outlier(s). In every other respect, the player is not necessarily noteworthy. This method should not be used, of course, unless there are enough features to merit computing the statistical outlier model on a case.

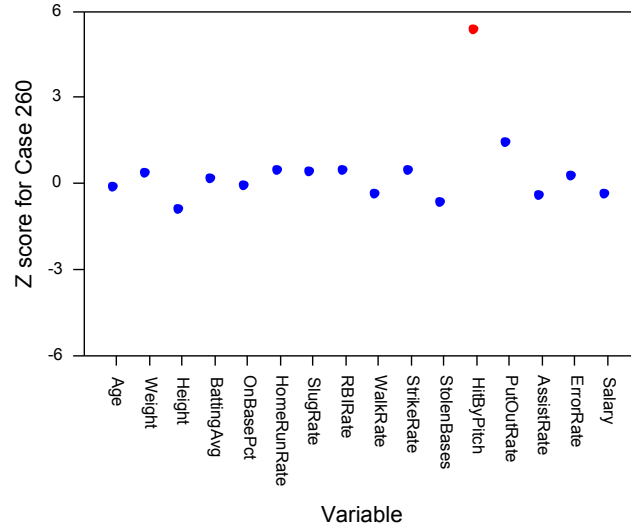


Figure 13: One baseball player’s profile showing an outlier (hit by pitch) that deviates significantly from his other features.

#### 4.3.5 Text Outliers

An important application for multivariate outlier detection involves document analysis. Given a collection of documents (Twitter messages, Wikipedia pages, emails, news pages, etc.), one might want to discover any document that is an outlier with respect to the others. The simplest approach to this problem is to use a bag-of-words model. We collect all the words in the documents, stem them to resolve variants, remove stopwords and punctuation, and then apply the tf-idf measure [54] on the words within each document. The resulting vectors for each document are then submitted to `hdoutliers`.

Figure 14 shows the results for an analysis of 21 novels from the Gutenberg Web site [26]. This problem requires the use of random projections. Before projection, there are 21,021 columns (tf-idf measures) in the dataset. After projection there are 653. Not surprisingly, *Ulysses* stands out as an outlier. Distinctively, it contains numerous neologisms.

*Tristram Shandy* was identified by `hdoutliers` as the second largest, but not significant, outlier. It too contains numerous neologisms. These two novels lie outside most of the points in Figure 14. Not all multivariate outliers will fall on the periphery of 2D projections, however, as we showed in Section 4.2.

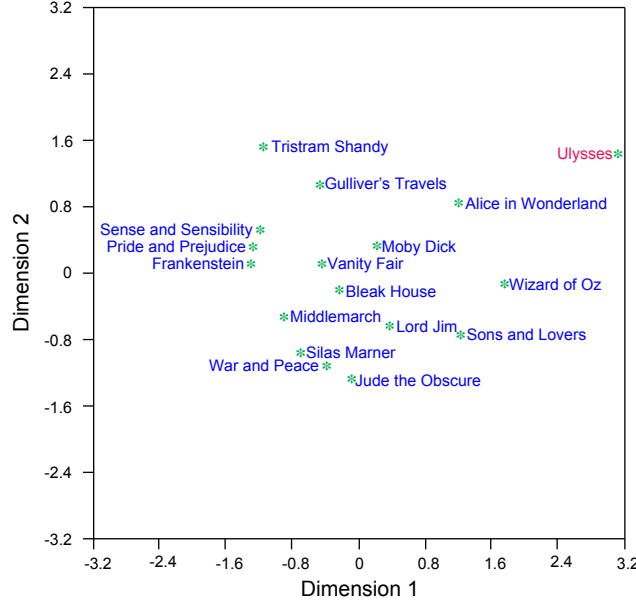


Figure 14: Document outliers. Nonmetric multidimensional scaling on matrix of Spearman correlations computed on tfidf scores. The stress for this solution is .163 and one document (*Ulysses*) is flagged as an outlier by `hdoutliers`.

## 4.4 Graph Outliers

There are several possibilities related to finding outliers in graphs. One popular application is the discovery of outliers among nodes of a network graph. The best way to exploit `hdoutliers` in this context is to featurize the nodes. Common candidates are Prominence, Transitivity (Watts-Strogatz Clustering Coefficient), Closeness Centrality, Betweenness Centrality, Node Degree, Average Degree of Neighbors, and Page Rank [45]. Figure 15 shows an example for the *Les Misérables* dataset [40]. The nodes were featurized for Betweenness Centrality in order to discover any extraordinarily influential characters. Not surprisingly, Valjean is connected to significantly more characters than anyone else in the book.

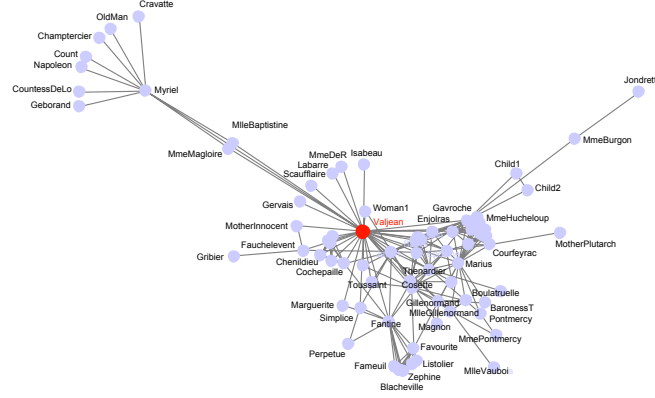


Figure 15: Les Misérables characters network graph. Valjean is identified as outlying on Betweenness Centrality.

An alternative application involves discovering outlying graphs in a collection of graphs. For this problem, we need to find a way to characterize a graph and to derive a distance measure that can be fed to `hdoutliers`. This application depends on assuming the collection of graphs is derived from a common population model and that any outliers involve a contamination from some alternative model. We need a measure of the distance between two graphs to do this. Unfortunately, graph matching and related graph edit distance calculations have impractical complexities. Approximate distances are easier to calculate, however [59]. The approach we take is as follows:

First, we compute the adjacency matrix for each graph. We then convert the adjacencies above the diagonal to a single binary string. When doing that, however, we have to reorder the adjacency matrix to a canonical form; otherwise, arbitrary input orderings could affect distance calculations on the string. A simple way to do this is to compute the eigendecomposition of the related Laplacian matrix and permute the adjacencies according to the ordering of the values of the eigenvector corresponding to the smallest nonzero eigenvalue. After permuting and encoding the adjacency matrices into strings, we compute the Levenshtein distances [43] between pairs of strings. Finally, we assemble the nearest-neighbor distances from the resulting distance matrix and subject them to the `hdoutliers` algorithm.

Figure 16 shows an example of this approach using the Karate Club graph [65]. We generated 15 random minimum spanning tree graphs having the same



number of nodes as the Karate Club graph. Then we applied the above procedure to identify outliers. The Karate Club graph was strongly flagged as an outlier by the algorithm.

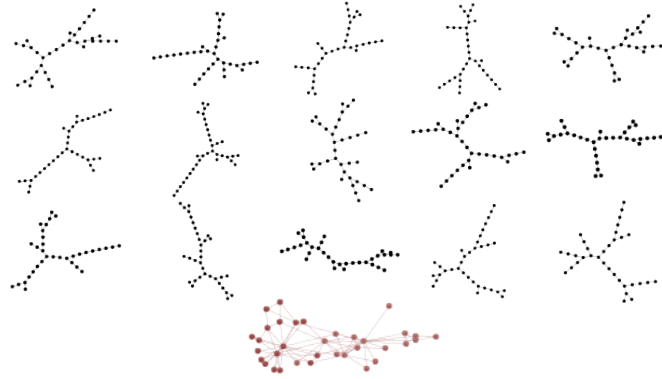


Figure 16: Karate Club graph (red) is an outlier with respect to comparably scaled random minimum spanning tree graphs.

#### 4.4.1 Scagnostics Outliers

Scagnostics [63] can be used to identify outlying scatterplots. Because the calculations are relatively efficient, these measures can be computed on many thousands of plots in practical time. This outlier application is multivariate, because there are nine scagnostics for each scatterplot, so a multivariate detection algorithm like `hdoutliers` is required.

Figure 17 shows two outlying scatterplots identified by `hdoutliers` when applied to a dataset of baseball player characteristics featured in [64]. While the left plot in the figure is clearly unusual, the surprising result is to see an evidently bivariate Normal scatterplot of Weight against Height in the right plot. Although the dataset includes many physical and performance features of real baseball players, the type of Normal bivariate distribution found in many introductory statistics books is an outlier among the 120 scatterplots considered in this example. This result should motivate authors writing tutorials on data analysis to include examples beyond Normal distributions.

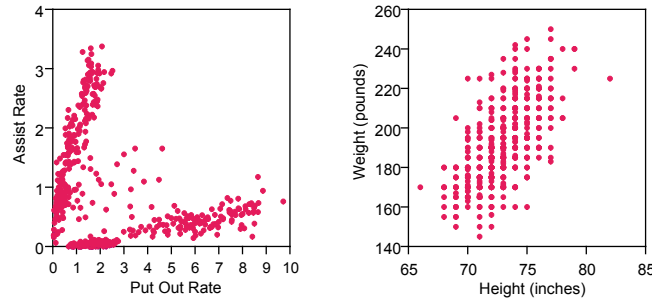


Figure 17: Scatterplot outliers based on Scagnostics computed on 120 scatterplots of baseball player features.

#### 4.4.2 Geographic Outliers

We can compute spatial outliers using the `hdoutliers` algorithm. More frequently, however, maps are a convenient way to display the results of outlier detection on other variables. Figure 18 shows an example of outlier detection on marriage and divorce rates by US state. Nevada is clearly an outlier. Despite the simplicity of this example, analyses at the State level are usually too coarse to be useful. Outliers displayed at a higher resolution (e.g., counties) are often preferable.

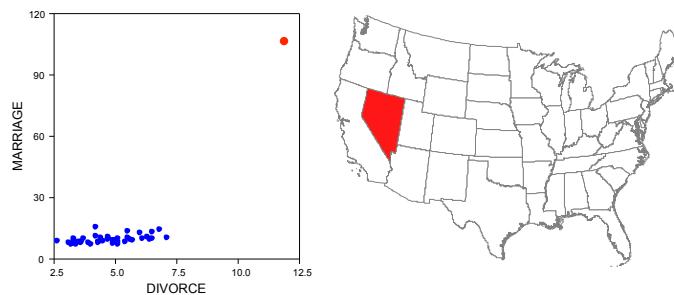


Figure 18: Marriage and Divorce rates in the US. There is one state that is an outlier.

## 5 Conclusions

There is a huge assortment of papers on outlier detection in the machine learning community; only a fraction is cited here. While many of these approaches are ingenious, few rest on a statistical foundation that takes risk into account. If we label something as an outlier, we had better be able to quantify or control our risk.

Outliers are anomalies. An anomaly is not a thing; literally, *anomaly* means lack of a law. It is a judgment based on evidence. Sometimes evidence is a collection of facts. Sometimes it is a collection of indications that cause us to modify our prior belief that what we are observing is not unusual.

The statistical detection of outliers is concerned with the latter case. Lacking certainty of the process that generated what we think might be an outlier, we must derive a judgment that an observation is inconsistent with our belief in that process.

Many statistical outlier detection algorithms assume a generating process derives from a parametric distribution. The more interesting cases are when we cannot presume such a distribution. The most useful cases, ones that are more relevant to real applications, involve the broadest class of prior beliefs in possible generating processes.

In order to be consistent in our behavior, we need to assign a probability to the strength of our belief that we are looking at an outlier. Methods that do not do this, that simply rank discrepancies or flag observations above an arbitrary threshold (like most of the algorithms in the Related Work section), can lead to inconsistent results.

The `hdoutliers` algorithm reduces the risk of making a false outlier discovery for a broad class of prior beliefs. Even for unusual applications such as the graph outlier problem, this algorithm provides a foundation for framing the judgment concerning an outlier. And importantly for the applications in this paper, `hdoutliers` is designed specifically to guide, protect, and deepen our visual analysis of data.

## References

- [1] C. Aggarwal. *Outlier Analysis*. Springer Verlag, 2013.
- [2] F. Anscombe. Rejection of outliers. *Technometrics*, 2:123–147, 1960.

- [3] A. Atkinson. *Plots, Transformations and Regression: An Introduction to Graphical Methods of Diagnostic Regression Analysis*. Oxford University Press, 1985.
- [4] V. Barnett and T. Lewis. *Outliers in Statistical Data*. John Wiley & Sons, 1994.
- [5] D. A. Belsley, E. Kuh, and R. E. Welsch. *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. John Wiley & Sons, 1980.
- [6] G. E. P. Box and G. M. Jenkins. *Time Series Analysis: Forecasting and Control (rev. ed.)*. Holden-Day, Oakland, CA, 1976.
- [7] M. Breaban and H. Luchian. Outlier detection with nonlinear projection pursuit. *International Journal of Computers Communications & Control*, 8(1):30–36, 2013.
- [8] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. LOF: Identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, SIGMOD '00, pages 93–104, New York, NY, USA, 2000. ACM.
- [9] P. Burrridge and A. Taylor. Additive outlier detection via extreme-value theory. *Journal of Time Series Analysis*, 27:685–701, 2006.
- [10] C. Caroni and P. Prescott. On Rohlf's method for the detection of outliers in multivariate data. *Journal of Multivariate Analysis*, 52:295–307, 1995.
- [11] V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM Comput. Surveys*, 41:15:1–15:58, July 2009.
- [12] W. S. Cleveland. *The Elements of Graphing Data*. Hobart Press, Summit, NJ, 1985.
- [13] R. D. Cook and S. Weisberg. *Residuals and Influence in Regression*. Chapman & Hall, London, 1982.
- [14] R. Dawson. How significant is a boxplot outlier? *Journal of Statistics Education*, 19, 2011.

- [15] W. Dixon. Ratios involving extreme values. *Annals of Mathematical Statistics*, 22:68–78, 1951.
- [16] D. Donoho and P. Huber. The notion of breakdown point. In P. Bickel, K. Doksum, and J. Hodges, editors, *A Festschrift for Erich L. Lehman*, pages 157–184. Wadsworth, Belmont, CA, 1983.
- [17] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In E. Simoudis, J. Han, and U. M. Fayyad, editors, *Second International Conference on Knowledge Discovery and Data Mining*, pages 226–231. AAAI Press, 1996.
- [18] C. Fraley and A. Raftery. Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association*, 97:611–631, 2002.
- [19] M. Frigge, D. Hoaglin, and B. Iglewicz. Some implementations of the box-plot. *The American Statistician*, 43:50–54, 1989.
- [20] German Research Center for Artificial Intelligence (DFKI). dataset: dfki-artificial-3000-unsupervised-ad.csv. <http://madm.dfki.de/downloads>. Accessed: 2016-02-08.
- [21] R. Gnanadesikan. *Methods for statistical data analysis of multivariate observations*. John Wiley & Sons, New York, 1977.
- [22] M. Greenacre. *Theory and Applications of Correspondence Analysis*. Academic Press, 1984.
- [23] F. Grubbs. Sample criteria for testing outlying observations. *The Annals of Mathematical Statistics*, 21:27–58, 1950.
- [24] A. Hadi and J. Simonoff. Procedures for the identification of multiple outliers in linear models. *Journal of the American Statistical Association*, 88:1264–1272, 1993.
- [25] F. R. Hampel, E. M. Ronchetti, P. J. Rousseeuw, and W. A. Stahel. *Robust Statistics: The Approach Based on Influence Functions*. John Wiley & Sons, 2005.
- [26] M. Hart. Project Gutenberg. <https://www.gutenberg.org>.

- [27] J. Hartigan. *Clustering Algorithms*. John Wiley & Sons, New York, 1975.
- [28] D. Hawkins. *Identification of Outliers*. Chapman & Hall/CRC, 1980.
- [29] R. Hayden. A dataset that is 44% outliers. *Journal of Statistics Education*, 13, 2005.
- [30] A. Hinneburg, D. Keim, and M. Wawryniuk. HD-Eye: Visual mining of high-dimensional data. *IEEE Computer Graphics and Applications*, 19(5):22–31, Sept. 1999.
- [31] V. Hodge. *Outlier and Anomaly Detection: A Survey of Outlier and Anomaly Detection Methods*. LAP LAMBERT Academic Publishing, 2011.
- [32] H. Hofmann, K. Kafadar, and H. Wickham. Letter-value plots: Boxplots for large data. Technical report, had.co.nz, 2011.
- [33] M. Hubert and E. Vandervieren. An adjusted boxplot for skewed distributions. *Computational Statistics & Data Analysis*, 52:5186–5201, 2008.
- [34] I. Iglewicz and D. Hoaglin. How to detect and handle outliers. In E. Mykytka, editor, *The ASQC Basic References in Quality Control: Statistical Techniques*. ASQC, 1993.
- [35] M. Jiang, S. Tseng, and C. Su. Two-phase clustering process for outliers detection. *Pattern Recognition Letters*, 22:691–700, 2001.
- [36] J. Jobe and M. Pokojovy. A cluster-based outlier detection scheme for multivariate data. *Journal of the American Statistical Association*, 110:1543–1551, 2015.
- [37] W. B. Johnson and J. Lindenstrauss. Lipschitz mapping into Hilbert space. *Contemporary Mathematics*, 26:189–206, 1984.
- [38] E. Kandogan. Visualizing multi-dimensional clusters, trends, and outliers using star coordinates. In *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '01, pages 107–116, New York, NY, USA, 2001. ACM.
- [39] E. Kandogan. Just-in-time annotation of clusters, outliers, and trends in point-based data visualizations. In *Proceedings of the 2012 IEEE Conference on Visual Analytics Science and Technology (VAST)*, VAST '12, pages 73–82, Washington, DC, USA, 2012. IEEE Computer Society.

- [40] D. Knuth. *The Stanford GraphBase: A Platform for combinatorial computing*. Addison-Wesley, Reading, MA, 1993.
- [41] R. Kohavi and B. Becker. Adult data set. <http://archive.ics.uci.edu/ml/datasets/Adult>, 1996.
- [42] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek. Outlier detection in axis-parallel subspaces of high dimensional data. In *13th Pacific-Asia Conf. on Knowledge Discovery and Data Mining (PAKDD 2009)*, Tangkok, Thailand, 2009.
- [43] V. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10:707–710, 1966.
- [44] J. Lin, D. Ye, C. Chen, and M. Gao. Minimum spanning tree based spatial outlier mining and its applications. In *Proceedings of the 3rd International Conference on Rough Sets and Knowledge Technology*, pages 508–515, Berlin, Heidelberg, 2008. Springer-Verlag.
- [45] M. E. J. Newman, A.-L. Barabasi, and D. J. Watts. *The Structure and Dynamics of Networks*. Princeton University Press, 2006.
- [46] M. Novotny and H. Hauser. Outlier-preserving focus+context visualization in parallel coordinates. *IEEE Transactions on Visualization and Computer Graphics*, 12(5):893–900, Sept. 2006.
- [47] R. Pamula, J. Deka, and S. Nandi. An outlier detection method based on clustering. In *International Conference on Emerging Applications of Information Technology, EAIT*, pages 253–256, 2011.
- [48] S. Peter and S. Victor. Hybrid – approach for outlier detection using minimum spanning tree. *International Journal of Computational & Applied Mathematics*, 5:621, 2010.
- [49] F. Rohlf. Generalization of the gap test for the detection of multivariate outliers. *Biometrics*, 31:93–101, 1975.
- [50] P. Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 79:871–880, 1984.

- [51] P. Rousseeuw and B. V. Zomeren. Unmasking multivariate outliers and leverage points. *Journal of the American Statistical Association*, 85:633–651, 1990.
- [52] P. J. Rousseeuw and A. Leroy. *Robust Regression & Outlier Detection*. John Wiley & Sons, 1987.
- [53] A. Ruiz-Gazen, S. L. Marie-Sainte, and A. Berro. Detecting multivariate outliers using projection pursuit with particle swarm optimization. In *COMPSTAT2010*, pages 89–98, 2010.
- [54] G. Salton and C. Buckley. Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5):513–523, Aug. 1988.
- [55] K. Schwarz. *Wind Dispersion of Carbon Dioxide Leaking from Underground Sequestration, and Outlier Detection in Eddy Covariance Data using Extreme Value Theory*. PhD thesis, Physics, University of California, Berkeley, 2008.
- [56] S. Stigler. *The History of Statistics: The Measurement of Uncertainty Before 1900*. Harvard University Press, 1986.
- [57] H. Thode. *Testing For Normality*. Taylor & Francis, 2002.
- [58] J. W. Tukey. *Exploratory Data Analysis*. Addison-Wesley Publishing Company, Reading, MA, 1977.
- [59] S. Umeyama. An eigendecomposition approach to weighted graph matching problems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10:695–703, 1988.
- [60] H. Wainer and S. Schacht. Gapping. *Psychometrika*, 43:203–212, 1978.
- [61] L. Wilkinson. Dot plots. *The American Statistician*, 53:276–281, 1999.
- [62] L. Wilkinson. *The Grammar of Graphics*. Springer-Verlag, New York, 1999.
- [63] L. Wilkinson, A. Anand, and R. Grossman. Graph-theoretic scagnostics. In *Proceedings of the IEEE Information Visualization 2005*, pages 157–164. IEEE Computer Society Press, 2005.



- [64] L. Wilkinson, A. Anand, and R. Grossman. High-dimensional visual analytics: Interactive exploration guided by pairwise views of point distributions. *IEEE Transactions on Visualization and Computer Graphics*, 12(6):1363–1372, 2006.
- [65] W. Zachary. An information flow model for conflict and fission in small groups. *Journal of Anthropological Research*, 33:452–473, 1977.
- [66] C. T. Zahn. Graph-theoretical methods for detecting and describing Gestalt clusters. *IEEE Transactions on Computers*, C-20:68–86, 1971.