

# Check your outliers! An introduction to identifying statistical outliers in R with *easystats*

Rémi Thériault<sup>1</sup>, Mattan S. Ben-Shachar<sup>2</sup>, Indrajeet Patil<sup>3</sup>, Daniel Lüdtke<sup>4</sup>, Brenton M. Wiernik<sup>5</sup>, and Dominique Makowski<sup>6</sup>

<sup>1</sup> Department of Psychology, Université du Québec à Montréal, Montréal, Québec, Canada <sup>2</sup> Independent Researcher <sup>3</sup> Center for Humans and Machines, Max Planck Institute for Human Development, Berlin, Germany <sup>4</sup> Institute of Medical Sociology, University Medical Center Hamburg-Eppendorf, Germany <sup>5</sup> Independent Researcher, Tampa, FL, USA <sup>6</sup> School of Psychology, University of Sussex, Brighton, UK

DOI:

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Submitted:

Published:

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

## Summary

Beyond the challenge of keeping up-to-date with current best practices regarding the diagnosis and treatment of outliers, an additional difficulty arises concerning the mathematical implementation of the recommended methods. In this paper, we provide an overview of current recommendations and best practices and demonstrate how they can easily and conveniently be implemented in the R statistical computing software, using the *{performance}* package of the *easystats* ecosystem. We cover univariate, multivariate, and model-based statistical outlier detection methods, their recommended threshold, standard output, and plotting methods. We conclude with recommendations on the handling of outliers: the different theoretical types of outliers, whether to exclude or winsorize them, and the importance of transparency.

## Statement of Need

Real-life data often contain observations that can be considered *abnormal* when compared to the main population. The cause of it—be it because they belong to a different distribution (originating from a different generative process) or simply being extreme cases, statistically rare but not impossible—can be hard to assess, and the boundaries of “abnormal” difficult to define.

Nonetheless, the improper handling of these outliers can substantially affect statistical model estimations, biasing effect estimations and weakening the models’ predictive performance. It is thus essential to address this problem in a thoughtful manner. Yet, despite the existence of established recommendations and guidelines, many researchers still do not treat outliers in a consistent manner, or do so using inappropriate strategies (Leys et al., 2013; Simmons et al., 2011).

One possible reason is that researchers are not aware of the existing recommendations, or do not know how to implement them using their analysis software. In this paper, we show how to follow current best practices for automatic and reproducible statistical outlier detection (SOD) using R and the *{performance}* package (Lüdtke et al., 2021), which is part of the *easystats* ecosystem of packages that build an R framework for easy statistical modeling, visualization, and reporting (Lüdtke et al., 2019/2023). Installation instructions can be found on [GitHub](#) or its [website](#), and its list of dependencies on [CRAN](#).

The instructional materials that follow is aimed at an audience of researchers who want to follow good practices, and is appropriate for advanced undergraduate students, graduate students, professors, or professionals having to deal with the nuances of outlier treatment.

## Identifying Outliers

Although many researchers attempt to identify outliers with measures based on the mean (e.g.,  $z$  scores), those methods are problematic because the mean and standard deviation themselves are not robust to the influence of outliers and they assume normally distributed data (i.e., a Gaussian distribution). Therefore, current guidelines recommend using robust methods to identify outliers, such as those relying on the median as opposed to the mean (Leys et al., 2013, 2018, 2019).

Nonetheless, which exact outlier method to use depends on many factors. In some cases, eye-gauging odd observations can be an appropriate solution, though many researchers will favour algorithmic solutions to detect potential outliers, for example, based on a continuous value expressing the observation stands out from the others.

One of the factors to consider when selecting an algorithmic outlier detection method is the statistical test of interest. When using a regression model, relevant information can be found by identifying observations that do not fit well with the model. This approach, known as model-based outliers detection (as outliers are extracted after the statistical model has been fit), can be contrasted with distribution-based outliers detection, which is based on the distance between an observation and the “center” of its population. Various quantification strategies of this distance exist for the latter, both univariate (involving only one variable at a time) or multivariate (involving multiple variables).

When no method is readily available to detect model-based outliers, such as for structural equation modelling (SEM), looking for multivariate outliers may be of relevance. For simple tests ( $t$  tests or correlations) that compare values of the same variable, it can be appropriate to check for univariate outliers. However, univariate methods can give false positives since  $t$  tests and correlations, ultimately, are also models/multivariable statistics. They are in this sense more limited, but we show them nonetheless for educational purposes.

Importantly, whatever approach researchers choose remains a subjective decision, which usage (and rationale) must be transparently documented and reproducible (Leys et al., 2019). Researchers should commit (ideally in a preregistration) to an outlier treatment method before collecting the data. They should report in the paper their decisions and details of their methods, as well as any deviation from their original plan. These transparency practices can help reduce false positives due to excessive researchers’ degrees of freedom (i.e., choice flexibility throughout the analysis). In the following section, we will go through each of the mentioned methods and provide examples on how to implement them with R.

### Univariate Outliers

Researchers frequently attempt to identify outliers using measures of deviation from the center of a variable’s distribution. One of the most popular such procedure is the  $z$  score transformation, which computes the distance in standard deviation (SD) from the mean. However, as mentioned earlier, this popular method is not robust. Therefore, for univariate outliers, it is recommended to use the median along with the Median Absolute Deviation (MAD), which are more robust than the interquartile range or the mean and its standard deviation (Leys et al., 2013, 2019).

Researchers can identify outliers based on robust (i.e., MAD-based)  $z$  scores using the `check_outliers()` function of the `{performance}` package, by specifying `method = "zscore_robust"`.<sup>1</sup> Although Leys et al. (2013) suggest a default threshold of 2.5 and Leys et al. (2019) a threshold of 3, `{performance}` uses by default a less conservative

<sup>1</sup>Note that `check_outliers()` only checks numeric variables.

threshold of  $\sim 3.29$ .<sup>2</sup> That is, data points will be flagged as outliers if they go beyond  $\pm \sim 3.29$  MAD. Users can adjust this threshold using the `threshold` argument, as demonstrated below.

```
library(performance)

# Create some artificial outliers and an ID column
data <- rbind(mtcars[1:4], 42, 55)
data <- cbind(car = row.names(data), data)

outliers <- check_outliers(data, method = "zscore_robust", ID = "car")
outliers

#> 2 outliers detected: cases 33, 34.
#> - Based on the following method and threshold: zscore_robust (3.291).
#> - For variables: mpg, cyl, disp, hp.
#>
#> -----
#>
#> The following observations were considered outliers for two or more
#>   variables by at least one of the selected methods:
#>
#>   Row car n_Zscore_robust
#> 1  33  33                2
#> 2  34  34                2
#>
#> -----
#> Outliers per variable (zscore_robust):
#>
#> $mpg
#>   Row car Distance_Zscore_robust
#> 33  33  33                3.709699
#> 34  34  34                5.848328
#>
#> $cyl
#>   Row car Distance_Zscore_robust
#> 33  33  33                12.14083
#> 34  34  34                16.52502
```

The row numbers of the detected outliers can be obtained by using `which()` on the output object, which can be used for exclusions for example:

```
which(outliers)

#> [1] 33 34

data_clean <- data[-which(outliers), ]
```

All `check_outliers()` output objects possess a `plot()` method, meaning it is also possible to visualize the outliers using the generic `plot()` function on the resulting outlier object after loading the `{see}` package.

Other univariate methods are available, such as using the interquartile range (IQR), or based on different intervals, such as the Highest Density Interval (HDI) or the Bias Corrected and Accelerated Interval (BCI). These methods are documented and described in the function's [help page](#).

<sup>2</sup>3.29 is an approximation of the two-tailed critical value for  $p < .001$ , obtained through `qnorm(p = 1 - 0.001 / 2)`. We chose this threshold for consistency with the thresholds of all our other methods.

## Multivariate Outliers

Univariate outliers can be useful when the focus is on a particular variable, for instance the reaction time, as extreme values might be indicative of inattention or non-task-related behavior<sup>3</sup>.

However, in many scenarios, variables of a data set are not independent, and an abnormal observation will impact multiple dimensions. For instance, a participant giving random answers to a questionnaire. In this case, computing the  $z$  score for each of the questions might not lead to satisfactory results. Instead, one might want to look at these variables together.

One common approach for this is to compute multivariate distance metrics such as the Mahalanobis distance. Although the Mahalanobis distance is very popular, just like the regular  $z$  scores method, it is not robust and is heavily influenced by the outliers themselves. Therefore, for multivariate outliers, it is recommended to use the Minimum Covariance Determinant, a robust version of the Mahalanobis distance (MCD, Leys et al., 2018, 2019).

In `{performance}`'s `check_outliers()`, one can use this approach with `method = "mcd"`.<sup>4</sup>

```
outliers <- check_outliers(data, method = "mcd")
outliers
```

```
#> 9 outliers detected: cases 7, 15, 16, 17, 24, 29, 31, 33, 34.
#> - Based on the following method and threshold: mcd (20).
#> - For variables: mpg, cyl, disp, hp.
```

Other multivariate methods are available, such as another type of robust Mahalanobis distance that in this case relies on an orthogonalized Gnanadesikan-Kettenring pairwise estimator (Gnanadesikan & Kettenring, 1972). These methods are documented and described in the function's [help page](#).

## Model-Based Outliers

Working with regression models creates the possibility of using model-based SOD methods. These methods rely on the concept of *leverage*, that is, how much influence a given observation can have on the model estimates. If few observations have a relatively strong leverage/influence on the model, one can suspect that the model's estimates are biased by these observations, in which case flagging them as outliers could prove helpful (see next section, "Handling Outliers").

In `{performance}`, two such model-based SOD methods are currently available: Cook's distance, for regular regression models, and Pareto, for Bayesian models. As such, `check_outliers()` can be applied directly on regression model objects, by simply specifying `method = "cook"` (or `method = "pareto"` for Bayesian models).<sup>5</sup>

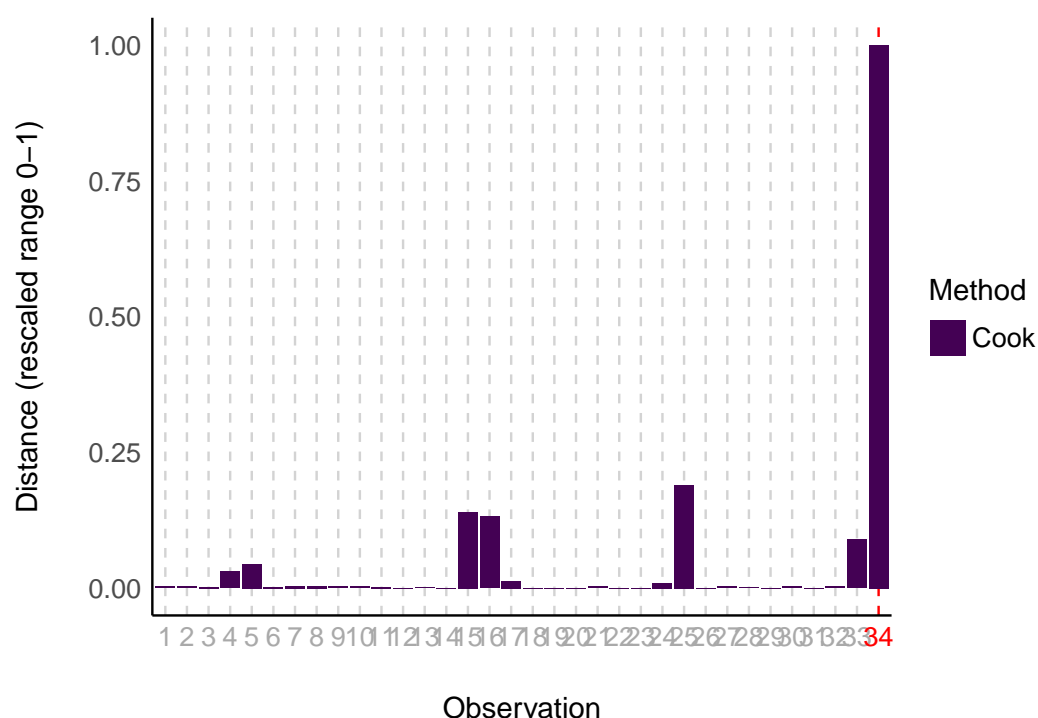
```
model <- lm(disp ~ mpg * disp, data = data)
outliers <- check_outliers(model, method = "cook")
outliers
```

```
#> 1 outlier detected: case 34.
#> - Based on the following method and threshold: cook (0.708).
```

<sup>3</sup>Note that they might not be the optimal way of treating reaction time outliers (Ratcliff, 1993; Van Zandt & Ratcliff, 1995)

<sup>4</sup>Our default threshold for the MCD method is defined by `stats::qchisq(p = 1 - 0.001, df = ncol(x))`, which again is an approximation of the critical value for  $p < .001$  consistent with the thresholds of our other methods.

<sup>5</sup>Our default threshold for the Cook method is defined by `stats::qf(0.5, ncol(x), nrow(x) - ncol(x))`, which again is an approximation of the critical value for  $p < .001$  consistent with the thresholds of our other methods.



**Figure 1:** Visual depiction of outliers based on Cook's distance (leverage and standardized residuals), based on the fitted model.

```
#> - For variable: (Whole model).
```

```
plot(outliers)
```

Table 1 below summarizes which methods to use in which cases, and with what threshold.

**Table 1:** Summary of Statistical Outlier Detection Methods Recommendations.

Statistical Test	Diagnosis Method	Recommended Threshold
Supported regression model	<b>Model-based:</b> Cook (or Pareto for Bayesian models)	$qf(0.5, ncol(x), nrow(x) - ncol(x))$ (or 0.7 for Pareto)
Structural Equation Modeling (or other unsupported model)	<b>Multivariate:</b> Minimum Covariance Determinant (MCD)	$qchisq(p = 1 - 0.001, df = ncol(x))$
Simple test with few variables ( <i>t</i> test, correlation, etc.)	<b>Univariate:</b> robust <i>z</i> scores (MAD)	$qnorm(p = 1 - 0.001 / 2), \sim 3.29$

### Cook's Distance vs. MCD

Leys et al. (2018) report a preference for the MCD method over Cook's distance. This is because Cook's distance removes one observation at a time and checks its corresponding influence on the model each time (Cook, 1977), and flags any observation that has a large influence. In the view of these authors, when there are several outliers, the process of removing a single outlier at a time is problematic as the model remains "contaminated" or influenced by other possible outliers in the model, rendering this method suboptimal in the presence of multiple outliers.

However, distribution-based approaches are not a silver bullet either, and there are cases where the usage of methods agnostic to theoretical and statistical models of interest might be problematic. For example, a very tall person would be expected to also be much heavier than average, but that would still fit with the expected association between height and weight (i.e., it would be in line with a model such as `weight ~ height`). In contrast, using multivariate outlier detection methods there may flag this person as being an outlier—being unusual on two variables, height and weight—even though the pattern fits perfectly with our predictions.

Finally, unusual observations happen naturally: extreme observations are expected even when taken from a normal distribution. While statistical models can integrate this “expectation”, multivariate outlier methods might be too conservative, flagging too many observations despite belonging to the right generative process. For these reasons, we believe that model-based methods are still preferable to the MCD when using supported regression models. Additionally, if the presence of multiple outliers is a significant concern, regression methods that are more robust to outliers should be considered—like *t* regression or quantile regression—as they render their precise identification less critical (McElreath, 2020).

## Composite Outlier Score

The `{performance}` package also offers an alternative, consensus-based approach that combines several methods, based on the assumption that different methods provide different angles of looking at a given problem. By applying a variety of methods, one can hope to “triangulate” the true outliers (those consistently flagged by multiple methods) and thus attempt to minimize false positives.

In practice, this approach computes a composite outlier score, formed of the average of the binary (0 or 1) classification results of each method. It represents the probability that each observation is classified as an outlier by at least one method. The default decision rule classifies rows with composite outlier scores superior or equal to 0.5 as outlier observations (i.e., that were classified as outliers by at least half of the methods). In `{performance}`’s `check_outliers()`, one can use this approach by including all desired methods in the corresponding argument.

```
outliers <- check_outliers(model, method = c("zscore_robust", "mcd", "cook"))
which(outliers)
```

```
#> [1] 33 34
```

Outliers (counts or per variables) for individual methods can then be obtained through attributes. For example:

```
attributes(outliers)$outlier_var$zscore_robust
```

```
#> $mpg
#>   Row Distance_Zscore_robust
#> 33   33                3.709699
#> 34   34                5.848328
```

An example sentence for reporting the usage of the composite method could be:

Based on a composite outlier score (see the ‘`check_outliers()`’ function in the ‘performance’ R package, Lüdtke et al., 2021) obtained via the joint application of multiple outliers detection algorithms ((a) median absolute deviation (MAD)-based robust *z* scores, Leys et al., 2013; (b) Mahalanobis minimum covariance determinant (MCD), Leys et al., 2019; and (c) Cook’s distance, Cook, 1977), we excluded two participants that were classified as outliers by at least half of the methods used.

## Handling Outliers

The above section demonstrated how to identify outliers using the `check_outliers()` function in the *{performance}* package. But what should we do with these outliers once identified? Although it is common to automatically discard any observation that has been marked as “an outlier” as if it might infect the rest of the data with its statistical ailment, we believe that the use of SOD methods is but one step in the get-to-know-your-data pipeline; a researcher or analyst’s *domain knowledge* must be involved in the decision of how to deal with observations marked as outliers by means of SOD. Indeed, automatic tools can help detect outliers, but they are nowhere near perfect. Although they can be useful to flag suspect data, they can have misses and false alarms, and they cannot replace human eyes and proper vigilance from the researcher. If you do end up manually inspecting your data for outliers, it can be helpful to think of outliers as belonging to different types of outliers, or categories, which can help decide what to do with a given outlier.

### Error, Interesting, and Random Outliers

Leys et al. (2019) distinguish between error outliers, interesting outliers, and random outliers. *Error outliers* are likely due to human error and should be corrected before data analysis or outright removed since they are invalid observations. *Interesting outliers* are not due to technical error and may be of theoretical interest; it might thus be relevant to investigate them further even though they should be removed from the current analysis of interest. *Random outliers* are assumed to be due to chance alone and to belong to the correct distribution and, therefore, should be retained.

It is recommended to *keep* observations which are expected to be part of the distribution of interest, even if they are outliers (Leys et al., 2019). However, if it is suspected that the outliers belong to an alternative distribution, then those observations could have a large impact on the results and call into question their robustness, especially if significance is conditional on their inclusion, so should be removed.

On the other hand, there are also outliers that cannot be detected by statistical tools, but should be found and removed. For example, if we are studying the effects of X on Y among teenagers and we have one observation from a 20-year-old, this observation might not be a *statistical outlier*, but it is an outlier in the *context* of our research, and should be discarded to allow for valid inferences.

### Winsorization

*Removing* outliers can in this case be a valid strategy, and ideally one would report results with and without outliers to see the extent of their impact on results. This approach however can reduce statistical power. Therefore, some propose a *recoding* approach, namely, winsorization: bringing outliers back within acceptable limits (e.g., 3 MADs, Tukey & McLaughlin, 1963). However, if possible, it is recommended to collect enough data so that even after removing outliers, there is still sufficient statistical power without having to resort to winsorization (Leys et al., 2019).

The *easystats* ecosystem makes it easy to incorporate this step into your workflow through the `winsorize()` function of *{datawizard}*, a lightweight R package to facilitate data wrangling and statistical transformations (Patil et al., 2022). This procedure will bring back univariate outliers within the limits of ‘acceptable’ values, based either on the percentile, the *z* score, or its robust alternative based on the MAD.



## The Importance of Transparency

Finally, it is a critical part of a sound outlier treatment that regardless of which SOD method used, it should be reported in a reproducible manner. Ideally, the handling of outliers should be specified *a priori* with as much detail as possible, and preregistered, to limit researchers' degrees of freedom and therefore risks of false positives (Leys et al., 2019). This is especially true given that interesting outliers and random outliers are often times hard to distinguish in practice. Thus, researchers should always prioritize transparency and report all of the following information: (a) how many outliers were identified (including percentage); (b) according to which method and criteria, (c) using which function of which R package (if applicable), and (d) how they were handled (excluded or winsorized, if the latter, using what threshold). If at all possible, (e) the corresponding code script along with the data should be shared on a public repository like the Open Science Framework (OSF), so that the exclusion criteria can be reproduced precisely.

## References

- Cook, R. D. (1977). Detection of influential observation in linear regression. *Technometrics*, 19(1), 15–18. <https://doi.org/10.1080/00401706.1977.10489493>
- Gnanadesikan, R., & Kettenring, J. R. (1972). Robust estimates, residuals, and outlier detection with multiresponse data. *Biometrics*, 81–124.
- Leys, C., Delacre, M., Mora, Y. L., Lakens, D., & Ley, C. (2019). How to classify, detect, and manage univariate and multivariate outliers, with emphasis on pre-registration. *International Review of Social Psychology*. <https://doi.org/10.5334/irsp.289>
- Leys, C., Klein, O., Dominicy, Y., & Ley, C. (2018). Detecting multivariate outliers: Use a robust variant of the mahalanobis distance. *Journal of Experimental Social Psychology*, 74, 150–156. <https://doi.org/10.1016/j.jesp.2017.09.011>
- Leys, C., Ley, C., Klein, O., Bernard, P., & Licata, L. (2013). Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. *Journal of Experimental Social Psychology*, 49(4), 764–766. <https://doi.org/10.1016/j.jesp.2013.03.013>
- Lüdtke, D., Ben-Shachar, M. S., Patil, I., Waggoner, P., & Makowski, D. (2021). performance: An R package for assessment, comparison and testing of statistical models. *Journal of Open Source Software*, 6(60), 3139. <https://doi.org/10.21105/joss.03139>
- Lüdtke, D., Makowski, D., Ben-Shachar, M. S., Patil, I., Wiernik, B. M., Bacher, E., & Thériault, R. (2023). *easystats: Streamline model interpretation, visualization, and reporting*. <https://easystats.github.io/easystats/> (Original work published 2019)
- McElreath, R. (2020). *Statistical rethinking: A bayesian course with examples in R and stan*. CRC press.
- Patil, I., Makowski, D., Ben-Shachar, M. S., Wiernik, B. M., Bacher, E., & Lüdtke, D. (2022). datawizard: An R package for easy data preparation and statistical transformations. *Journal of Open Source Software*, 7(78), 4684. <https://doi.org/10.21105/joss.04684>
- Ratcliff, R. (1993). Methods for dealing with reaction time outliers. *Psychological Bulletin*, 114(3), 510. <https://doi.org/10.1037/0033-2909.114.3.510>
- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, 22(11), 1359–1366. <https://doi.org/10.1177/0956797611417632>
- Tukey, J. W., & McLaughlin, D. H. (1963). Less vulnerable confidence and significance procedures for location based on a single sample: Trimming/winsorization 1. *Sankhyā: The Indian Journal of Statistics, Series A*, 331–352.
- Van Zandt, T., & Ratcliff, R. (1995). Statistical mimicking of reaction time data: Single-process models, parameter variability, and mixtures. *Psychonomic Bulletin & Review*, 2(1), 20–54. <https://doi.org/10.3758/BF03214411>