

Detection/removal of outliers is critical because – as the author correctly states – the raw data are altered thus impacting any subsequent statistical analysis. A widely accepted description of outliers as “an observation (or subset of observations) which APPEARS to be inconsistent with the remainder of that set of data” (Barnett&Lewis 1994) stresses the notorious difficulty to unambiguously identify a data point as an outlier and there is always the choice between correcting the data OR correcting the model. The problem is most eminent if there are few data – then there is a great chance that data are removed as outliers just because the assumed model is misspecified (and we must fear that this practice is wide-spread). The author considers the problem of outlier detection in the context of time series where the amount of data is typically relatively large. This opens the chance that we have a lot of information about the underlying distribution (the model) and we can exploit that knowledge to identify inconsistent observations. The approach chosen by the author addresses exactly this situation and the proposed modification of the original boxplot-idea seems to be very promising and robust. However, in the present form the implementation of the idea is not consistent. Roughly speaking, the author compares his approach with several older suggestions which are more or less inspired by the Gaussian case and he replaces these by a suggestion ($k=0.6$) which is based on an average(median) of reasonable models. This will almost certainly lead to an improvement since the Gaussian case is an extreme case in the Skewness-Kurtosis plane. But in any particular application we don't have an average of distributions – we have one particular distribution and we should use an outlier detection method which is tailored to this particular distribution as good as possible. The author notes in his discussion of the precipitation use case: “... $k=0.6$ is insufficient here. A value of $k \sim 5$ is optimum in this case ...”. As the authors goal is to provide a “generic pre-processing procedure implemented in R” I would strongly advise to choose the value of k in a data-driven way instead of relying on a default value of $k=0.6$ thus introducing a quasi-standard even if there is the opportunity for the user to change this value. As far as I see, all the required knowledge is there. My idea would be to robustly estimate skewness and kurtosis from the data and then choose a value for k derived from the corresponding Pearson distribution. This would probably lead to a substantial improvement in robustness of the ctbi method and would also justify the publication of the LogBox approach in a more statistically oriented journal.

I am not a specialist in robust estimation but Kim&White (2004) could probably provide some ideas on how to estimate the kurtosis.

More detailed comments:

- The fact, that the Pearson family provides a distribution for any theoretically possible combination of skewness and kurtosis could be more explicitly stated. Also the fact that all the mentioned families (Beta, Gamma, ...) are Pearson families and only later got their current names could be worth mentioning.
- The argument with the 3σ , 4σ and 5σ convention seems relatively arbitrary. Only together with the requirement that the absolute number of erroneously flagged outliers should be restricted this makes sense. This important information is hidden only in the caption to Table 1. Perhaps you could argue the other way around starting from the expected number of erroneously flagged outliers. This would of course change the regression equations (If you go for exactly one erroneously flagged outlier in the Gaussian case the equation should roughly be $-0.5 + 0.322 \log(n)$)
- I don't see the need to present the results for the generalized Gamma distribution as the bounded support precludes its use a priori.
- You are looking for methods which are adapted to particular distributions. A general procedure to deal with non-normal data, which is also popular in the time-series context, is

the transformation (log, square root, Box-Cox, ...) of the data before statistical analysis. It would be worth comparing the transformation approach with the LogBox approach when detecting outliers.

- I personally find ctbi a bit cryptic. How about CTBin which at least makes the main ingredient – the bins – visible.

Barnett, V. and T. Lewis (1994) Outliers in Statistical Data. 3rd ed. Wiley&Sons

Kim, T.-H. and H. White (2004) On more robust estimation of skewness and kurtosis. Finance Research Letters 1 (1)