

Review of *Effects of undetected data quality issues on climatological analyses* by Hunziker et al.

### **General comments**

This paper uses data from the central Andes to illustrate the impact of systematic data issues on climate analyses. By applying standard and enhanced quality control procedures before data homogenisation, the study shows that systematic biases – which might not be detected using standard QC methods – can have a big impact on the analysis of temperature and rainfall trends.

This paper is a nice contribution to the field of climate data quality, and raises some important considerations that are often overlooked when large-scale climate trend analysis is conducted. The study is succinct and well written, and is a good complement to this group's previous paper (Hunziker et al. 2017).

Some suggestions for possible improvements to the paper are given below for the authors' consideration.

### **Specific comments**

1. The main concern I have about the results from the study is how much data are removed in the enhanced QC procedure, and what impact that would have on the final trends. Around 40% is a significant amount of data, and you'd expect that removing that much information would have an effect on the homogenisation and analyses whether or not the data have issues.  
It'd be great if the authors could repeat their analysis using a synthetic dataset, or a dataset that does not have systematic data issues and see the impact of removing 40% of the data. That way they could ascribe some statistical significance to their findings. If this is not plausible, then at least a discussion on the role of missing data in the results, OR more detail on where and when the removed data are.
2. In section 3.1.2, I think a little more information is needed about the QC issues mentioned. I know it's in Hunziker et al. 2017, but the two papers are independent and should be understandable on their own.

### **Technical corrections**

- Introduction line 2: replace "most" with "many national"
- Intro lines 20–30: This part confused me a little. Are you saying that trends actually do vary between neighbouring stations (due to climate factors), or that issues with data quality cause the differences?
- Page 3, line 1: as "the" enhanced approach
- Page 3, lines 5–9: The Central Andean region is also a good case study because of its complex terrain, as quality issues/homogenisation is notoriously difficult in such conditions.
- Page 3, lines 9–12: Change chapter to sections
- Page 3, line 20: add "network" after weather observation
- Page 4, line 5: data "were" quality-controlled
- Page 4, lines 10–15: I'd add that the Durre et al. tests include spatial inconsistency tests.

- Page 4, lines 27–29: Why doesn't the GHCN-Daily QC approach flag these clearly erroneous values?
- Page 4, line 30: add % after 0.35
- Page 5, line 6: on "an" annual time scale
- Page 5, line 26: add % after 0.26
- Page 5, lines 28 and 30 [and throughout the manuscript]: I think you need "a" before monthly and daily time scale.
- Page 6, line 5–6: I suggest you reword to: Note that Hunziker et al. (2017) further suggest *the inclusion of* additional information derived from 5 metadata into the QC process. This allows *the removal of* station records that were generated under inappropriate conditions such as poor station siting or severe...
- Page 6, line 15: one day "a week"
- Page 6, line 30: I'd add reference to Peterson et al. 1998 when discussing using first-differences: Peterson, T. C., T. R. Karl, P. F. Jamason, R. Knight, and D. R. Easterling (1998), First difference method: Maximizing station density for the calculation of long-term global temperature change, *J. Geophys. Res.*, 103(D20), 25967–25974, doi:10.1029/98JD01168.
- Page 6, line 31: Spearman (with a capital S)
- Page 7, line 1: add "the" before monthly time scale
- Section 3.4.1. Is there a way to graphically show the different clusters? Perhaps some additional panels in Figure 1 using polygons?
- Section 3.4.2. It sounds like you've used ACMANT3 because it is fully automated. Perhaps mention this earlier in the section to make that point stronger.
- Page 8, line 14: remove "of"
- Page 8, line 15: remove the hyphen between DECADE and dataset.
- Page 8, line 16: add "and" after requirements
- Page 8, lines 20–24: You need an extra line or two here to define/explain the Theil-Sen estimator, and how you have pre-whitened the data.
- Page 9, lines 5-15: I understand what you're trying to say here, but it's a bit confusing. Consider revising the text.
- Page 13, line 13: "and" Switzerland
- Page 13, line 14: Why are you hypothesising this? Why would the tropics have more UDQI? Do you also mean that there are more likely to be UDQI in developing countries?
- Page 15, line 5: add "a" before few climate change indices
- Table 1 caption: were, not where
- Table 2: can you add any information about the spread of the data in this table?
- Figure 2 caption: I'd add the word "show" after the words green triangles and red triangles.