

Atmos. Meas. Tech. Discuss., author comment AC1
<https://doi.org/10.5194/amt-2021-154-AC1>, 2021
© Author(s) 2021. This work is distributed under
the Creative Commons Attribution 4.0 License.

Reply on RC1

Jeffrey Bean

Author comment on "Evaluation methods for low-cost particulate matter sensors" by
Jeffrey K. Bean, Atmos. Meas. Tech. Discuss.,
<https://doi.org/10.5194/amt-2021-154-AC1>, 2021

The author thanks the reviewer for the comments. The manuscript has been updated and more information on the updates is provided below.

This paper presents and discusses an improved air sensor evaluation method the prediction interval. Overall, this appears to be a valuable addition to statistics more commonly used to evaluate air sensor performance. However, the methods are lacking to understand how to implement this method in future work. Some of the discussion and conclusions are also lacking.

Major:

It seems like this paper would be more helpful to the air sensor community if the sensor type was provided.

Response: The sensor models are withheld in part to allow focus on evaluation methods, rather than specific sensors, but models are also withheld to avoid giving the impression that the author endorses or disparages one sensor company over another. Additional details on sensor measurements have been added to the manuscript, as requested, but the specific makes and models are still withheld for the reasons stated.

Line 98: I think this method is super helpful to think about determining how to pick a threshold to remove outliers based on precision moving forward! Did you consider whether percentage and absolute concentration criteria would be more useful than just 50 ug/m3 by itself (as used in <https://doi.org/10.5194/amt-14-4617-2021>)? What was the range of concentrations experienced by the sensors during this period (as reported by the sensors since that is the criteria you are using to exclude)?

Response: That is a good suggestion about using a percentage data agreement criteria. This was considered, but in the end we chose to only use an absolute criteria for simplicity and because it worked well. Some discussion on this has been added to the manuscript (line 108).

Figure 1: How much does removing disagreeing datapoints clean up the bias? Or does it only clean up the R2 as shown in this figure? It would be helpful to include a 3rd panel on Fig 1 that would be average concentration reported by the sensor versus the allowed difference and also the average of the FEM as well.

Response. Good question. Improvements in RMSE are also observed as disagreeing datapoints are removed. These results have been added to Figure 1.

Line 120: How was no correlation defined? How did you decide which parameters to include in the figures and which not to? What are the correlations of the parameters in the figures? It would be helpful to also show the correlation between PM10 or PM10-2.5 and the error since as you say in the paper particle size can influence the accuracy.

Response: "no correlation" was a strong choice of words and that has been changed. The key message was that extremely close correlation was observed between sensors, even at times when bias was very high, but despite this there was no indication that this was caused by factors such as humidity, temperature or any other measured parameter. Unfortunately, PM₁₀ concentrations were not available during this measurement campaign.

Figure 2 seems overly complicated and challenging to interpret. It seems like 2A would be easier to interpret if it just showed sensor 1 versus sensor 2 (not the bias) and then if figure 2d showed sensor 1 versus the BAM concentration. I think this would be valuable even if you wanted to add a figure showing the basic plots and then a second figure showing all the bias plots if you feel you need both.

Response: The reason we focus on bias in the panels of the figure is to draw attention to the fact that these sensors often read well above or below the actual concentration, but in a very consistent way from one sensor to another. We believe that 2A and 2E draw attention to the fact that sensors will read higher or lower than they should in a way that suggests that it should be correlated with an external factor, though that factor was not observed in the measurements available during this campaign.

Prediction Interval: I don't understand how to calculate this based on the information you have included in the text. Please cite additional references and include the equations in the SI if needed. Please include a citation for the R package if one was used. What are the 3 diagonal lines shown on Figure 4? Can you label them?

Response: The three diagonal lines are the best fit and upper/lower prediction intervals. A prediction interval can be provided by the standard linear model (lm) package in R. For example:

```
sensor_model <- lm(data, formula = reference ~ sensor)
```

```
prediction <- predict(sensor_model, data, interval="predict")
```

Line 186: "Through examination it was found that residual trends were best eliminated by raising both the sensor and reference data to the 0.4 power." More scientific explanation is needed. How would we repeat this method in the future would others always just use 0.4 since that is what you "examined" and found or would they need to examine their datasets and come to different conclusions? Your figure doesn't show transformed data so it is unclear how you would determine this. Please provide additional explanation.

Response: Some additional clarification has been added to the manuscript. There are many transformation methods that could be applied but the key criteria before accepting a linear regression is that the residuals should be evenly spaced across the sampling domain. In our case raising both reference and sensor data to the 0.4 power eliminated any trends in residuals. Raising data to a power is just one option for any data that has uneven residuals. Future applications of this type of method might find other ways to ensure consistent residuals.

Past work has shown that sensors may respond nonlinearly at high concentrations

(<https://doi.org/10.1111/ina.12621>, <https://doi.org/10.1016/j.envpol.2018.11.065>).
Would a prediction interval still be appropriate in a case like this? Would you be able to remove the residual trends in a dataset like this?

Response: A prediction interval is especially relevant if there is nonlinearity under different concentration regimes. The method in the manuscript would probably work, but if it didn't then an equally thorough alternative would be to split the data into different ranges and fit separate prediction intervals for each.

Line 189: "This did not change the outcome significantly". Please define what you mean by this. How did you decide to use 70 ug/m3 as your split to have equal number of data points above and below? Does Figure 4 show the full dataset or the subset? How important is having a well-balanced dataset to getting accurate results?

Response: 70 $\mu\text{g}/\text{m}^3$ was chosen because it approximately split the 5% of data during high concentration events from the rest of the data. The random sampling below this line meant that the model and prediction intervals were a result of 50% of data below 70 $\mu\text{g}/\text{m}^3$ and 50% of data above that point. If the model instead used 95% of data below 70 $\mu\text{g}/\text{m}^3$ and only 5% above then the model is fit in a way that is more weighted towards lower concentration. If the transformation is done well and the residuals are the same across the domain then the weight of lower/higher concentrations does not matter. However since residual distribution is not perfect it is best to fit the model an equal amount of data across the range of observed concentrations.

The need to balance a dataset seems like a limitation of this method as compared to R2 or RMSE that has not been discussed.

Response: Balancing a dataset can improve the robustness of any model, as it ensures the model is built equally using the entire domain over which it will predict. This can improve R2 and RMSE in a standard linear model as well. However as discussed here, the result of doing this is small if residuals are correctly accounted for, so doing that should be the priority in building any model that predicts PM predictions from sensors.

Figure 4: I would recommend coloring the arrows uniquely and including them next to the text in the bottom right corner as a legend so the figure can be more quickly interpreted.

Response: Great suggestion. Additional labels have been added to make this quicker to interpret.

Figure 5: Is this only for T640 data? If so, why not also show BAM data as shown in the other figure?

Response: Figure 5 was built only for T640 data as this allowed 5-minute resolution in comparison with 1-hour or 24-hour. The BAM is limited to 1-hour resolution.

Figure 5: Did you consider whether uncertainty as a % would be more stable?

Response: Uncertainty as a percentage is an interesting idea. When applied to 5-minute data in Figure 5 it results in something like an exponential decay with uncertainty as high as 400% for low concentrations (1 $\mu\text{g}/\text{m}^3$) that decays towards $\sim 45\%$ uncertainty at 100 $\mu\text{g}/\text{m}^3$. We believe that absolute concentrations are a little easier/quicker to interpret, but percent uncertainty is an idea that could be interesting to explore in the future as well.

Line 219: "It allows for better comparison between sensors, as the evaluation results are not biased by the range of concentrations observed during evaluation." I think more explanation is needed here. It doesn't seem to me this is one of the findings of your

analysis since you only showed the results from one dataset covering the full range. It might be interesting to show another data subset with a different concentration range to understand how the range doesn't impact the results but maybe there is another way to explain since as I said above I don't really understand how you are calculating the prediction interval.

Response: More explanation has been added to that section to clarify: "Uncertainty at any given concentration can be compared from one brand of sensor to another and is not impacted by the range of concentrations observed, in contrast to RMSE or R2. In other words, the uncertainty of a sensor at 35 µg/m3 does not change depending on whether concentrations of 100 µg/m3 were also measured during evaluation, though the overall R2 or RMSE of that evaluation can be influenced by the 100 µg/m3 measurements, as shown in Fig. 3."

Lines 306-308: "Two of the most popular evaluation metrics, R2 and RMSE, can be influenced by averaging time, choice of reference instrument, and the range of concentrations observed (see Fig. 3). This study shows how a prediction interval can be used as a more statistically thorough evaluation tool." Figure 5 shows that prediction interval is also influenced by averaging time. Are you saying it isn't influenced by reference instrument? If so, I think you need more results to show that. Overall, this statement seems misleading.

Response: The reviewer is correct that the statement here was misleading. A few additional sentences have been added to clarify that the choice of reference instrument still needs to be standardized and that evaluations would be more simple if averaging time was also standardized (starting line 357).

Figure 6: With much of the data below 5 ug/m3 did you consider how LOD of the sensor and reference influence your results?

Response: A LOD was not provided for the sensors being evaluated and was not explored in this work. Future work could consider the connection between uncertainty, as measured here, and LOD, which is a similar concept.

Figure 2 seems to show that the bias is much more variable at high RH. How can you take that into account using prediction interval?

Response: Yes, a prediction interval can also be found for a multiple linear regression. The approach to do so in R is similar as with a single linear regression, as shown below. Exploring a prediction interval based on many predicting variables would be interesting, but would greatly complicate the figures in this work. To keep the concept simple we focus just on one prediction variable.

```
sensor_model <- lm(data, formula = reference ~ sensor + RH)
```

```
prediction <- predict(sensor_model, data, interval="predict")
```

Did you consider how the precision of the sensor influences prediction interval? I'm assuming that Figure 4 is for all the sensors but if it is for a group of sensors that would be helpful to clarify in the caption/text.

Response: Figure 4 is from a single sensor. The Figure 4 analysis could either be repeated for replicas of a sensor or data from multiple sensors could be combined to create an analysis such as Figure 4. Additional numbers are now included for a repeat of this analysis of a second sensor and more discussion has been added about how one might approach this for multiple sensors and how sensor precision will impact this (starting line

251).

It would be helpful to add the prediction interval for all of the sensors you tested not just the best sensor so that readers could compare the R2/RMSE/PI more closely across devices and understand how they could use this in the future.

Response: Good suggestion. The analysis has been added for a 2nd sensor to show how they compare and more discussion has been added about how one might approach doing this analysis for multiple sensors. Comparisons between sensor uncertainty for different sensors have been added as well (line 264).

Have you thought about how you could report this PI as something more easily to compare across papers than a plot (which may have different axis labels etc.)? For example, fitting a function or reporting the 95% uncertainty at various AQI breakpoints, etc?

Response: This is a nice suggestion and the following text has been added to the manuscript: "Picking a single comparison point allows users to quickly compare measurement uncertainty between different sensor types, as they might currently using R2 or RMSE. The breakpoints in the United States Air Quality Index (AQI) could be considered as standard comparison points. For example, the United States AQI transitions to "Unhealthy for Sensitive Groups" at 35 µg/m3."

This work is missing relevant citations. Examples: Giordano 2021 calibration review paper <https://doi.org/10.1016/j.jaerosci.2021.105833>, Zheng 2018 similar discussion of averaging interval and the precision of the reference <https://doi.org/10.5194/amt-11-4823-2018>, some others included in my other responses.

Response: These works have now been appropriately cited and the author appreciates the reviewer pointing them out.

Minor:

Figure 1: Could you include the averaging interval you are using for exclusion in the figure caption?

Response: They were averaged to 1-hour intervals for this figure and this has now been included in the caption.

Line 45: "The root of a calibration for low-cost particulate matter sensors is simple: sensors and reference instruments measure the same mass of air for a period and then adjustments are made to better align sensor measurements.". I'm not sure "root" is the clearest way to express this.

Response: This sentence has been adjusted to read: "During a calibration, low-cost sensors and reference instruments measure the same mass of air for a period and then adjustments are made to better align sensor measurements."

Line 255: US EPA recommends at least 30 days for their PM2.5 sensor evaluations.

https://cfpub.epa.gov/si/si_public_file_download.cfm?p_download_id=542106&Lab=CEM
M

Response: This has been added to the manuscript.