

Atmos. Meas. Tech. Discuss., referee comment RC1
<https://doi.org/10.5194/amt-2022-99-RC1>, 2022
© Author(s) 2022. This work is distributed under
the Creative Commons Attribution 4.0 License.

Comment on amt-2022-99

Anonymous Referee #1

Referee comment on "Ch3MS-RF: a random forest model for chemical characterization and improved quantification of unidentified atmospheric organics detected by chromatography–mass spectrometry techniques" by Emily B. Franklin et al., Atmos. Meas. Tech. Discuss., <https://doi.org/10.5194/amt-2022-99-RC1>, 2022

This work presents a new method for predicting organic molecule properties (carbon number, mean oxidation state, oxygen-to-carbon ratio, vapor pressure) for compounds measured by gas chromatography and mass spectrometry but not listed in mass spectral databases. The novel idea is to use a statistical model trained on compounds listed in such databases together with parameters extracted from measurement, which is available regardless of the compound being listed in such a database. The authors include the caveat that this approach requires two-dimensional gas chromatography measurements that capture both volatility and polarity and on dimension is insufficient, which does require more complicated instrumentation than commonly deployed. Nonetheless, the general idea can be useful for the atmospheric science community, and is recommended for publication in Atmospheric Measurement Techniques. I note a few comments regarding the generalizability below.

The authors refer to their "extrapolation set" in they are not included in the training set, but in reality it appears that new samples span a subset of the feature domain spanned by the training set. The consideration of whether extrapolation in this sense is happening or likely to happen in new data sets is relevant because random forest is not capable of such extrapolations - and this would limit the model's utility substantially. Can the authors clarify this point?

I assume the results are solely applicable to samples run on the same instrument with the same protocol, as the retention time is dependent on the operating procedure. For any new protocol a new model would have to be trained. Can this model be used to generate predictions using measurements on similar instruments using the same protocol, or does a new model have to be trained on each instrument? For publication, the authors should include a statement regarding what is required for adaptation by other users.