Public Comment- Simon Ruske (simon.ruske@student.manchester.ac.uk)
Received and published: 3 June 2018

> Note regarding document formatting: black text shows original referee comment, blue text shows author response, and red text shows quoted manuscript text. Changes to manuscript text are shown as *italicized and underlined*. All line numbers refer to discussion/review manuscript.

[Public Comment] The study presented is an extremely well structured and written investigation into the use of Hierarchical Agglomerative Clustering for classification of biological aerosol using a UV-LIF sensor, and will make an excellent addition to the literature upon publication.

> [Author Response] Simon, thanks for taking the time to read and comment on the manuscript. We appreciate the useful comments, which will help improve the quality of the manuscript. We respond to each comment in detail below.

However, the authors may have made a small error [L161-L162] where they state that the conclusions for Ruske et al. (2017) were for ambient data, whereas in the abstract they correctly state that the study was on standardised laboratory particles [L19-L20]. Please could you correct this prior to final publication.

> I apologize for this mistake. I am not sure where this error came in our writing process, but I removed the incorrect statement, as requested: "Their conclusions, however, were based on ambient field data using unknown particle types and did not investigate laboratory generated particles of known origin."

In addition the authors may wish to consider the following comments prior to publication.
[L78-L79] Would it be possible to clarify the starting conditions for supervised learning you are referring to? Hyper-parameter selection is an extremely important consideration for neural networks, but other supervised techniques such as decision trees and ensemble methods do exist where low classification error can be attained without providing the algorithm with any initial conditions other than the training data.

> This may have been a bit of a miscommunication. We do not deal with any supervised learning methods in this manuscript. We trust your team as the experts in this area. Nicole simply wanted to provide a few sentences of general contrast between supervised and unsupervised methods. That is also why we pointed to your 2017 paper in this section. We have also included citation of your manuscript currently being reviewed in AMT.

[L84-L85] Is it necessary to apply unsupervised techniques to assess the advantages of supervised methods? Do you mean that supervised techniques require laboratory data of known types to assess their advantages? A very important disadvantage of supervised techniques is that they rely on adequate training data, and it is not clear at this point how much training data will be required to adequately represent an ambient environment, which is the point I think you are alluding to here.

> This is the way I understand some of the pros/cons of supervised and unsupervised. I agree that the community (probably you first) will continue to lean about how this all works together and how well lab-generated data can be useful to train supervised data algorithms. As you well know, the differences between nicely behaving lab particles and more complicated particles collected in the field confounds most areas of aerosol science to some degree. So these problems will not necessarily be trivial to solve, but I think collectively we are all learning little pieces that will help.

51
52 [L186 - 187] Does the z-score rely on the assumption of normality? The z-scores of a normal random
53 variable will be normally distributed whereas the z-scores of a non-normal random variable will be non-
54 normally distributed. **Applied to any data set, regardless of distribution, the resultant variables after**
55 **z-scoring will have mean of 0 and standard deviation of 1.** Is the purpose of standardising the data to
56 prevent one of the variables from dominating in the analysis or to produce normally distributed data?
57
58 Thanks to your prompting, we looked into these details and learned a bit more, which has been
59 helpful to us. You are right that the way we characterized the z-scoring process was not correct.
60 Talking back and forth with the university statistician, we now understand that values can indeed
61 be input scaled to a normal distribution or not. We chose to standardize our variables to a mean of
62 0 and a variance of 1 so that the output variables would be on comparable scales, but this is also
63 not the same as rigorously normalizing them in the rigorous sense. As a result, we have removed
64 the statement you correctly indicated was inaccurate and updated the sentence as follows:
65
66 Original text: "Standardization using the z-score method compares results to a normal (Gaussian)
67 population, ~~and therefore relies on the assumption that input data can be described by a normal~~
68 ~~distribution (Gordon, 2006).~~"
69
70 Updated text: "Standardization using the z-score method compares results to a normal (Gaussian)
71 population, and we have chosen to standardize our variables to a mean of 0 and a variance of 1 so
72 that the output variables would be on comparable scales."
73
74 [L203] It would be worth noting that in Crawford et al., 2015, there are particles for which negative
75 measurement of fluorescence was recorded. The option of logtransformations may have been overlooked,
76 as the logarithm is undefined for negative values. This was not intended to imply an assumption of
77 normality, although this assumption has been stated explicitly in Robinson et al., 2013. In these cases
78 would you recommend translating the fluorescence measurements to a range bounded below by 1, or
79 alternatively would it be more appropriate to reject measurements for which the fluorescence produced
80 was negative? It is also important to note that even if the data is log transformed, the data will still have a
81 finite range due to the saturation point on the detector, and hence the data will have a truncated normal
82 distribution rather than a normal distribution, and depending on how often saturation occurred there may
83 still be a peak to the right hand side of the distribution. It is however, perfectly acceptable to apply HAC
84 when the assumptions for best performance are not met as stated in Norusis, 2011.
85
86 My understanding is that negative fluorescence values can be observed after subtracting some
87 threshold value from the fluorescence intensity data. Instead of subtracting the data and looking
88 only at positive values, we did the same thing by filtering the data at several discreet thresholds.
89 This gets around the problem of negative values. In any case, we looked at three thresholding
90 scenarios (Table 3), i.e. no threshold, 3 sigma, and 9 sigma. The ultimate result is that we found
91 the most consistently positive results to be as a result of 3 sigma filtering, but this could be
92 different in other situations. You are correct about the fact that particles that exhibit saturation of
93 the detector in any channel will truncate a normal distribution.
94
95 [L222] How often did the CH index conclude that there were 2 clusters? When the CH index concluded a
96 number of clusters other than 2, how much of an impact did this have on the quality of the results? Were
97 the two cluster solutions always the best solution?
98
99 We did not explore solutions that had more than 2 solutions, simply as a matter of limited time.
100 There are certainly many scenarios in which individual bioparticle types (i.e. pollen, in many
101 instances) can split into two reasonable clusters by themselves, and so independently allowing 3

102 or more cluster solutions could significantly improve results in many cases. We just didn't have
103 the time to do this systematically, and so we chose to limit analysis to only 2 clusters in all cases.
104 To help clarify this point, we added text at:

105

106 L227: "In order to reduce the length and complexity of *discussion, analysis of results in Sections*
107 *4.1-4.3 was limited to using cluster products only from the 2-cluster solution. In some cases a 3-*
108 *cluster solution may have produced higher quality results, but these cases were not investigated.*"

109

110 [L267-270 & Figure 3] The HAC algorithm may not necessarily output clusters in the same order that
111 they were inputted as demonstrated in Figure 5. In Figure 3 for preparation strategy A for bacteria and
112 diesel for the 80:20 ratio, is it possible to attain 80% misclassification for a two cluster solution? Perhaps
113 I have misunderstood, but would this not mean that there were more diesel particles in the bacterial
114 cluster and more bacterial particles in the diesel cluster, and hence a better classification error could be
115 attained simply by swapping the labels on the clusters?

116

117 You are correct that the order of cluster numbering is unrelated to the order of particles input and
118 so the source of individual particles must be known already, but it is not possible to improve the
119 results by swapping labels in the way you suggest. We independently tracked the source of each
120 particle assigned to each cluster so we can rigorously calculate which particles were incorrectly
121 assigned. The numbering of the clusters is arbitrary and the naming was assigned simply as a
122 function of which particle was assigned in the largest concentration.

123

124 [Figure 3 & Table 2] Could you extend the results presented in Figure 3 to include at least one biological
125 versus biological matchup? I notice when considering matching ups which contained only biological
126 material the classification error is much higher. I believe that by not standardising the data this would
127 cause the fluorescence to dominate more in the analysis. In the case of attempting to discriminate between
128 fluorescent and non-fluorescent particles, this may be advantageous. However, in the case of attempting
129 to discriminate between two different types of biological particle, it may be advantageous to give the size
130 and shape measurements more weight, and hence it would be better in these cases to standardise the data.
131 In addition other instruments such as the WIBS-NEO will have fluorescence measurements over a much
132 larger range and fluorescent measurements are recorded often above 10000. What would the implication
133 then be when not standardising the data in this case?

134

135 This is another interesting idea, but it was beyond the scope of what we were able to accomplish
136 in the relatively short time we had available for this project. We chose to focus on the ability to
137 separate bio from non-bio particles. While we didn't explore all Scenarios (e.g. A-F) for
138 biological particles, we chose to look at bio-bio separations using Scenario B (i.e. Tables 2 and
139 3).