

Review of “Evaluation of linear regression techniques for atmospheric applications: The importance of appropriate weighting” by Wu and Yu

General comments:

This manuscript evaluates five linear regression techniques, ranging from standard (ordinary) least squares to those that account for errors in both variables. Described is a technique to generate data with desired properties for analysis by the regression techniques. The proper accounting for uncertainties, and thus the appropriate weighting, is emphasized. Approaches are recommended that retrieve slopes and intercepts of datasets with uncertainties in x and y variables that have minimal bias in slope and intercept.

The analysis is systematic and apparently carefully done. It does surprise this reviewer, however, that none of the regression techniques precisely recover the input slope and intercept (for example, results from Figure 5), particularly for the more sophisticated methods. Other papers have shown that the York method retrieves correct slopes and intercepts for a wide variety of conditions. It seems that with 5000 (or more) runs, regression with proper weighting should yield average slopes and intercepts very close to the input values. Suggest making use of Pearson’s data with York’s weights (for which the slope and intercept are known with high accuracy) to verify the coding used to perform the regression, as there may be some coding errors that remain and are affecting the results. The coding for data generation should also be checked carefully to ensure that this is not the problem. It is just stated that the  $r^2$  value is 0.67. The situation at the top of Figure S2 is what I would expect for properly generated data with proper accounting for uncertainties in x and y, namely that the average slope and intercept are precisely the input values.

The data generation schemes presented need more explanation. To test data regression schemes, it is not necessary that the data behave precisely like ambient atmospheric data. While not stated, it appears that the Chu 2005) method is attempted to reproduce the diurnal behavior of species concentration. This reviewer does not see that the use of this method adds to the comparison of the various regression methods, and probably only adds confusion. Suggest either providing a better explanation and justification of using this approach, or remove it from the paper.

Specific comments:

Several cases are considered in the paper and the supplemental material. For clarity, suggest presenting all the cases in a single table, showing the input slopes and intercepts as well as the linear and non-linear uncertainties of the x and y variables. Yes, values are shown for some of the cases, but they are split between the main paper and the supplement, and are hard to directly compare. There is also inconsistency between Figure 4, which indicates that there are 12 scenarios, and the various tables that go up to Case 18 (Table S7). Suggest describing the various scenarios in the text earlier in the paper than page 14 where they are discussed.

Page 3, line 57. Suggest “...is much smaller than the uncertainty...”. Suggest making this discussion more quantitative. In other words, give a precise value and to the how large the relative uncertainty must be to require use of techniques beyond OLS.

Page 3, line 59. Suggest “...may have comparable degrees of uncertainty.”

Page 3, line 61. Suggest “...applied to the dataset.”

Page 4, line 78. Suggest “In principle, a best-fit regression line should have greater dependence on the more precise data points rather than the less reliable ones.”

Page 4, line 81. Suggest “...is closer to the correct value than OLS, but may...”

Page 4, line 84. Suggest “This  $\lambda$  value is the key to handling the...”

Page 4, line 85. Suggest “...for the best-fit line calculations.”

Page 4, line 86. Suggest “...in the calculation of the best-fit line for an error-in-variable...”

Page 5, equations 2, 3 and 4. It appears that brackets or parentheses are needed to include both the  $x_i$  and  $y_i$  containing terms in the summation (such as done in equation 6).

Page 6, line 136. Suggest “...for demonstration in a real-world application.”

Page 6, line 146. Not sure why the word “relatively” was added. Suggest removing it.

Page 7, line 166. Suggest “... $POC_{comb}$  (the part of Y that is correlated with X)...”

Page 7, line 168. Suggest “...is added to  $POC_{comb}$ ...”

Page 7, line 178. Suggest “...uncertainties ( $\epsilon_{comb}$ ) to the true...”. Also, suggest indicating (somewhere) that the uncertainties are both positive and negative with a defined distribution, and an average of 0.

Page 8, line 199. The modification of the definition of  $Y_{unc}$  is stated, but no references are given, and the justification is not clear. Does this formula represent the uncertainties in an appropriate way?

Page 9, line 209-211. Related to the previous comment, this is asserted, but not really proven.

Page 9, line 212. Does “uniform distribution” mean “flat distribution” (also used on page 8, line 180)? In other words, is the distribution variance (and thus the weight) constant with deviation from the mean (rather than Gaussian or some other distribution). If so, why was this chosen?

Page 9, equations 20 and 21. The origin of these equations is not clear. Why is  $EC_{true}$  multiplied by  $LOD_{EC}$ ? Why is the factor of 3 included?

Page 9, line 223. Suggest “...where  $Y_{POC_{unc}}$  and  $Y_{EC_{nc}}$  are the relative measurement uncertainties...”

Page 10, line 239. Have you done analyses of the fitting accuracy with various frequency distributions? Since ambient data is typically log-normal distributed, its use might make sense, if it does make a difference.

Page 10, line 253. Suggest “...in this study is a single value...”

Page 11, line 260. It might be useful to have separate symbols for the non-linear and linear parts of the uncertainties (e.g.  $Y_{unc-linear}$  and  $Y_{unc-nonlinear}$ ).

Page 11, line 264. Suggest "...is given in the Supplemental Information."

Page 11, Section 3.1.3. Suggest a statement indicating why Chu (2005) used this method to generate data (if this remains in the paper per earlier comment).

Page 11, line 278. Suggest "...goodness of the regression intercept..".

Page 12, line 291. Suggest "...instruments utilized inlets with a 2.5  $\mu\text{m}$  particle diameter cutoff."

Page 12, line 300. Do you mean "SigmaPlot" rather than "Sigma Pro"? Also suggest "...DR is set to 1..."

Page 15, line 369. Suggest "...unbiased slope...and intercept..."

Page 17, line 426. Suggest "...the results using synthetic data."

Page 17, line 434, 435, and 437. Suggest changing "percentile" to "percentage".

Page 17, line 438. Suggest "...the differences between the six RS are also small..."

Page 17, line 439. Suggest "...as  $r^2$  decreases..."

Page 17, line 443. Suggest "...confirm the results obtained in comparing methods with the..."

Page 18, line 455. Suggest "...the measurement errors during..."

Page 18, line 456. Suggest "...data with  $r^2$  less than..."

Page 18, line 457. Suggest "...to minimize biases in the slope..."

Page 18, line 461-2. Suggest "...packed with many use features for data analysis and plotting,..."

Page 18, line 464. Is the program planned to be archived at the given site for a long time? Check the journal's policies regarding links to download sites.

Page 20, line 498. Suggest that for the York iterative method, that a relative tolerance between successive iterations be calculated, and that convergence be considered when this tolerance is reached. While 6 iterations could be sufficient for some datasets, it may not be enough for others.

Page 24. There are a few abbreviations missing from the table (e.g.  $W_i$ ,  $\beta_i$ ,  $r_i$ ).

Supplemental Material

Page 2, line 20. Suggest "...is often impacted by..."