



1 **Evaluation of linear regression techniques for atmospheric**  
2 **applications: The importance of appropriate weighting**

3 **Cheng Wu<sup>1,2</sup> and Jian Zhen Yu<sup>3,4,5</sup>**

4 <sup>1</sup>Institute of Mass Spectrometer and Atmospheric Environment, Jinan University, Guangzhou 510632, China

5 <sup>2</sup>Guangdong Provincial Engineering Research Center for on-line source apportionment system of air  
6 pollution, Guangzhou 510632, China

7 <sup>3</sup>Division of Environment, Hong Kong University of Science and Technology, Clear Water Bay, Hong  
8 Kong, China

9 <sup>4</sup>Atmospheric Research Centre, Fok Ying Tung Graduate School, Hong Kong University of Science and  
10 Technology, Nansha, China

11 <sup>5</sup>Department of Chemistry, Hong Kong University of Science and Technology, Clear Water Bay, Hong  
12 Kong, China

13 *Corresponding to:* Cheng Wu ([wucheng.vip@foxmail.com](mailto:wucheng.vip@foxmail.com)) and Jian Zhen Yu ([jian.yu@ust.hk](mailto:jian.yu@ust.hk))

14



15 **Abstract**

16 Linear regression techniques are widely used in atmospheric science, but are often improperly applied due to  
17 lack of consideration or inappropriate handling of measurement uncertainty. In this work, numerical  
18 experiments are performed to evaluate the performance of five linear regression techniques, significantly  
19 extending previous works by Chu and Saylor. The regression techniques tested are Ordinary Least Square  
20 (OLS), Deming Regression (DR), Orthogonal Distance Regression (ODR), Weighted ODR (WODR), and  
21 York regression (YR). We first introduce a new data generation scheme that employs the Mersenne Twister  
22 (MT) pseudorandom number generator. The numerical simulations are also improved by: (a) refining the  
23 parameterization of non-linear measurement uncertainties, (b) inclusion of a linear measurement uncertainty,  
24 (c) inclusion of WODR for comparison. Results show that DR, WODR and YR produce an accurate slope,  
25 but the intercept by WODR and YR is overestimated and the degree of bias is more pronounced with a low  
26  $R^2$  XY dataset. The importance of a properly weighting parameter  $\lambda$  in DR is investigated by sensitivity tests,  
27 and it is found an improper  $\lambda$  in DR can leads to a bias in both the slope and intercept estimation. Because the  
28  $\lambda$  calculation depends on the actual form of the measurement error, it is essential to determine the exact form  
29 of measurement error in the XY data during the measurement stage. With the knowledge of an appropriate  
30 weighting, DR, WODR and YR are recommended for atmospheric studies when both x and y data have  
31 measurement errors.

32



## 33 1 Introduction

34 Linear regression is heavily used in atmospheric science to derive the slope and intercept of XY datasets.  
35 Examples of linear regression applications include primary OC (organic carbon) and EC (elemental carbon)  
36 ratio estimation (Turpin and Huntzicker, 1995), MAE (mass absorption efficiency) estimation from light  
37 absorption and EC mass (Moosmüller et al., 1998), source apportionment of polycyclic aromatic hydrocarbons  
38 using CO and NO<sub>x</sub> as combustion tracers (Lim et al., 1999), gas-phase reaction rate determination (Brauers  
39 and Finlayson-Pitts, 1997), inter-instrument comparison (Bauer et al., 2009; Cross et al., 2010; von Bobruzki  
40 et al., 2010; Zieger et al., 2011; Huang et al., 2014; Zhou et al., 2016), light extinction budget reconstruction  
41 (Malm et al., 1994; Watson, 2002), comparison between modeling and measurement (Petäjä et al., 2009),  
42 emission factor study (Janhäll et al., 2010), retrieval of shortwave cloud forcing (Cess et al., 1995), calculation  
43 of pollutant growth rate (Richter et al., 2005), estimation of ground level PM<sub>2.5</sub> from MODIS data (Wang and  
44 Christopher, 2003), distinguishing OC origin from biomass burning using K<sup>+</sup> as a tracer (Duan et al., 2004)  
45 and emission type identification by the EC/CO ratio (Chen et al., 2001).

46 Ordinary least squares (OLS) regression is the most widely used method due to its simplicity. In OLS, it is  
47 assumed that independent variables are error free. This is the case for certain applications, such as determining  
48 a calibration curve of an instrument in analytical chemistry. For example, a known amount of analyte (e.g.,  
49 through weighing) can be used to calibrate the instrument output response (e.g., voltage). Since the uncertainty  
50 of gravimetric analysis is much smaller compared to the uncertainty of the instrument response, the error free  
51 assumption in “x” is valid. However, in many other applications, such as inter-instrument comparison, x and  
52 y (from two instruments) may have a comparable degree of uncertainty. This deviation from the underlying  
53 assumption in OLS would produce biased slope and intercept when OLS is applied on the data set.

54 To overcome the drawback of OLS, a number of error-in-variable regression models (also known as bivariate  
55 fittings (Cantrell, 2008) or total least-squares methods (Markovsky and Van Huffel, 2007) arise. Deming  
56 (1943) proposed an approach by minimizing sum of squares of x and y residuals. A closed-form solution of  
57 Deming regression (DR) was provided by York (1966). Method comparison work of various regression  
58 techniques by Cornbleet and Gochman (1979) found significant error in OLS slope estimation when the  
59 relative standard deviation (RSD) of measurement error in “x” exceeded 20%, while DR was found to reach  
60 a more accurate slope estimation. In an early application of the EC tracer method, Turpin and Huntzicker  
61 (1995) realized the limitation of OLS since OC and EC have comparable measurement uncertainty, thus  
62 recommended the use of DR for (OC/EC)<sub>pri</sub> (primary OC to EC ratio) estimation. Ayers (2001) conducted a  
63 simple numerical experiment and concluded that reduced major axis regression (RMA) is more suitable for  
64 air quality data regression analysis. Linnet (1999) pointed out that when applying DR for inter-method (or



65 inter-instrument) comparison, special attention should be paid to the sample size. If the range ratio (max/min)  
66 is relatively small (e.g., less than 2), more samples are needed to obtain statistically significant results.

67 In principle, regression line should depend more on the precise data points rather than the less reliable data  
68 points. Chu (2005) performed a comparison study of OLS and DR specifically focusing on the EC tracer  
69 method application, and found the slope estimated by DR is closer to the ideal value than OLS but may still  
70 overestimate the ideal value. Saylor et al. (2006) extended the comparison work of Chu (2005) by including a  
71 regression technique developed by York et al. (2004). They found that the slope overestimation by DR in the  
72 study of Chu (2005) was due to improper configuration of the weighting parameter,  $\lambda$ . This  $\lambda$  is the key to  
73 handle the uneven errors between data points for the best fitting line calculation. This example demonstrates  
74 the importance of appropriate weighting in the calculation of fitting line for error-in-variable regression model,  
75 which is overlooked in many studies.

76 In this study, we extend the work by Saylor et al. (2006) to achieve four objectives. The first is to propose a  
77 new data generation scheme by applying the Mersenne Twister (MT) pseudorandom number generator for  
78 evaluation of linear regression techniques. In the study of Chu (2005), data generation is achieved by a variational  
79 sine function, which has limitations in sample size, sample distribution, and nonadjustable correlation ( $R^2$ )  
80 between X and Y. In comparison, the MT data generation provides more flexibility, permitting adjustable  
81 sample size, XY correlation and distribution. The second is to develop a non-linear measurement error  
82 parameterization scheme for use in the regression method. The third is to incorporate linear measurement  
83 errors in the regression methods. In the work by Chu (2005) and Saylor et al. (2006), the relative measurement  
84 uncertainty ( $\gamma_{Unc}$ ) is non-linear with concentration, but a constant  $\gamma_{Unc}$  is often applied on atmospheric  
85 instruments due to its simplicity. The fourth is to include weighted orthogonal distance regression (WODR)  
86 for comparison. Abbreviations used in this study are summarized in Table 1 for quick lookup.

## 87 **2 Description of regression techniques compared in this study**

88 **Ordinary least squares (OLS) method.** OLS only considers the errors in dependent variables (y). OLS  
89 regression is achieved by minimizing the sum of squares (S) in the y residuals:

$$90 \quad S = \sum_{i=1}^n (y_i - Y_i)^2 \quad (1)$$

91 where  $Y_i$  are observed y data points while  $y_i$  are regressed y data points of the regression line.

92 **Orthogonal distance regression (ODR).** ODR minimizes the sum of the squared orthogonal distances from  
93 all data points to the regressed line and considers equal error variances:

$$94 \quad S = \sum_{i=1}^n (x_i - X_i)^2 + (y_i - Y_i)^2 \quad (2)$$



95 **Weighted orthogonal distance regression (WODR).** Unlike ODR that considers even error in  $x$  and  $y$ ,  
96 weightings based on measurement errors in both  $x$  and  $y$  are considered in WODR when minimizing the sum  
97 of squared orthogonal distance from the data points to the regression line (Carroll and Ruppert, 1996):

$$98 \quad S = \sum_{i=1}^n (x_i - X_i)^2 + (y_i - Y_i)^2 / \eta \quad (3)$$

99 where  $\eta$  is error variance ratio. Implementation of ODR and WODR in Igor was done by the computer routine  
100 ODRPACK95 (Boggs et al., 1989; Zwolak et al., 2007).

101 **Deming regression (DR).** Deming (1943) proposed the following function to minimize both the  $x$  and  $y$   
102 residuals,

$$103 \quad S = \sum_{i=1}^n \omega(X_i)(x_i - X_i)^2 + \omega(Y_i)(y_i - Y_i)^2 \quad (4)$$

104 where  $X_i$  and  $Y_i$  are observed data points and  $x_i$  and  $y_i$  are regressed data points. Individual data points are  
105 weighted based on errors in  $X_i$  and  $Y_i$ ,

$$106 \quad \omega(X_i) = \frac{1}{\sigma_{X_i}^2}, \quad \omega(Y_i) = \frac{1}{\sigma_{Y_i}^2} \quad (5)$$

107 where  $\sigma_{X_i}$  and  $\sigma_{Y_i}$  are the standard deviation of the error in measurement of  $X_i$  and  $Y_i$  respectively. The closed  
108 form solutions for slope and intercept of DR are shown in Appendix A.

109 **York regression (YR).** The York method (York et al., 2004) introduces the correlation coefficient of errors  
110 in  $x$  and  $y$  into the minimization function.

$$111 \quad S = \sum_{i=1}^n \left[ \omega(X_i)(x_i - X_i)^2 - 2r_i \sqrt{\omega(X_i)\omega(Y_i)}(x_i - X_i)(y_i - Y_i) + \omega(Y_i)(y_i - Y_i)^2 \right] \frac{1}{1-r_i^2} \quad (6)$$

112 where  $r_i$  is the correlation coefficient between measurement errors in  $X_i$  and  $Y_i$ . The slope and intercept of YR  
113 are calculated iteratively through the formulas in Appendix A.

### 114 **3 Data description**

115 Two types of data are used for regression comparison. The first type is synthetic data generated by computer  
116 programs, which can be used in the EC tracer method (Turpin and Huntzicker, 1995) to demonstrate the  
117 regression application. The true “slope” and “intercept” are assigned during data generation, allowing  
118 quantitative comparison of the bias of each regression scheme. The second type of data comes from ambient  
119 measurement of light absorption, OC and EC in Guangzhou for demonstration of the real-world application.



### 120 3.1 Synthetic XY data generation

121 In this study, numerical simulations are conducted in Igor Pro (WaveMetrics, Inc. Lake Oswego, OR, USA)  
 122 through custom codes. Two types of generation schemes are employed, one is based on the MT pseudorandom  
 123 number generator (Matsumoto and Nishimura, 1998) and the other is based on the sine function described by  
 124 Chu (2005).

125 The general form of linear regression on XY data can be written as:

$$126 \quad Y = kX + b \quad (7)$$

127 Here  $k$  is the regressed slope and  $b$  is the intercept. The underlying meaning is that,  $Y$  can be decomposed  
 128 into two parts. One part is correlated with  $X$ , and the ratio is defined by  $k$ . The other part of  $Y$  is relatively  
 129 constant and independent of  $X$  and regarded as  $b$ .

130 To make the discussion easier to follow, we intentionally avoid discussion using the abstract general form and  
 131 instead opt to use a real-world application case in atmospheric science. Linear regression had been heavily  
 132 applied on OC and EC data, here we use OC and EC data as an example to demonstrate the regression  
 133 application in atmospheric science. In the EC tracer method, OC (mixture) is  $Y$  and EC (tracer) is  $X$ . OC can  
 134 be decomposed into three components based on their formation pathway:

$$135 \quad OC = POC_{comb} + POC_{non-comb} + SOC \quad (8)$$

136 Here  $POC_{comb}$  is primary OC from combustion.  $POC_{non-comb}$  is primary OC emitted from non-combustion  
 137 activities.  $SOC$  is secondary OC formed during atmospheric aging. Since  $POC_{comb}$  is co-emitted with EC and  
 138 well correlated with each other, their relationship can be parameterized as:

$$139 \quad POC_{comb} = (OC/EC)_{pri} \times EC \quad (9)$$

140 By carefully selecting an OC and EC subset when  $SOC$  is very low (considered as approximately zero), the  
 141 combination of Eqs. (8) & (9) become:

$$142 \quad POC = (OC/EC)_{pri} \times EC + POC_{non-comb} \quad (10)$$

143 The regressed slope of POC ( $Y$ ) against EC ( $X$ ) represents  $(OC/EC)_{pri}$  ( $k$  in Eq.(7)). The regressed intercept  
 144 become  $POC_{non-comb}$  ( $b$  in Eq. (7)). With known  $(OC/EC)_{pri}$  and  $POC_{non-comb}$ ,  $SOC$  can be estimated by:

$$145 \quad SOC = OC - ((OC/EC)_{pri} \times EC + POC_{non-comb}) \quad (11)$$

146 The data generation starts from EC ( $X$  values). Once EC is generated,  $POC_{comb}$  (correlated part of  $Y$ ) can be  
 147 obtained by multiplying EC with a preset constant,  $(OC/EC)_{pri}$  (slope  $k$ ). Then the other preset constant  
 148  $POC_{non-comb}$  is added on  $POC_{comb}$  and the sum becomes POC ( $Y$  values). To simulate the real-world situation,



149 measurement errors are added on X and Y values. Details of synthesized measurement error are discussed in  
 150 the next section. Implementation of data generation by two types of mathematical schemes are explained in  
 151 section 3.1.2 and 3.1.3 respectively.

### 152 3.1.1 Parameterization of synthesized measurement uncertainty

153 Weighting of variables is a crucial input for errors-in-variables linear regression methods such as DR, YR and  
 154 WODR. In practice, the weights are usually defined as the inverse of the measurement error variance (Eq.  
 155 (5)). When measurement errors are considered, measured concentrations ( $Conc_{measured}$ ) are simulated by  
 156 adding measurement uncertainties ( $\varepsilon_{Conc.}$ ) into the true concentrations ( $Conc_{true}$ ):

$$157 \quad Conc_{measured} = Conc_{true} + \varepsilon_{Conc.} \quad (12)$$

158 Here  $\varepsilon_{Conc.}$  is random error following an even distribution, the range of which is constrained by:

$$159 \quad -\gamma_{Unc} \times Conc_{true} \leq \varepsilon_{Conc.} \leq +\gamma_{Unc} \times Conc_{true} \quad (13)$$

160 The  $\gamma_{Unc}$  is a dimensionless factor that describes the fractional measurement uncertainties relative to the true  
 161 concentration ( $Conc_{true}$ ).  $\gamma_{Unc}$  could be a function of  $Conc_{true}$  or a constant. The term  $\gamma_{Unc} \times Conc_{true}$   
 162 defines the boundary of random measurement errors.

163 Two types of measurement error are considered in this study. The first type is non-linear  $\gamma_{Unc}$ . In the data  
 164 generation scheme of Chu (2005) for the measurement uncertainties ( $\varepsilon_{POC}$  and  $\varepsilon_{EC}$ ),  $\gamma_{Unc}$  is non-linearly  
 165 related to  $Conc_{true}$ :

$$166 \quad \gamma_{Unc} = \frac{1}{\sqrt{Conc_{true}}} \quad (14)$$

167 then Eq. (13) for POC and EC become:

$$168 \quad -\frac{1}{\sqrt{POC_{true}}} \times POC_{true} \leq \varepsilon_{POC} \leq +\frac{1}{\sqrt{POC_{true}}} \times POC_{true} \quad (15)$$

$$169 \quad -\frac{1}{\sqrt{EC_{true}}} \times EC_{true} \leq \varepsilon_{EC} \leq +\frac{1}{\sqrt{EC_{true}}} \times EC_{true} \quad (16)$$

170 The non-linear  $\gamma_{Unc}$  defined in Eq. (14) is more realistic than the constant approach, but there are two  
 171 limitations. First, the physical meaning of the uncertainty unit is lost. If the unit of OC is  $\mu\text{g m}^{-3}$ , then the unit  
 172 of  $\varepsilon_{OC}$  becomes  $\sqrt{\mu\text{g m}^{-3}}$ . Second, the concentration is not normalized by a consistent relative value, making  
 173 it sensitive to the x and y units used. For example, if  $POC_{true}=0.9 \mu\text{g m}^{-3}$ , then  $\varepsilon_{POC}=\pm 0.95 \mu\text{g m}^{-3}$  and  $\gamma_{Unc}$   
 174  $= 105\%$ , but by changing the concentration unit to  $POC_{true}=900 \text{ ng m}^{-3}$ , then  $\varepsilon_{OC}=\pm 30 \text{ ng m}^{-3}$  and  $\gamma_{Unc} = 3\%$ .  
 175 To overcome these deficiencies, we propose to modify Eq. (14) to:



$$176 \quad \gamma_{Unc} = \sqrt{\frac{LOD}{Conc.true}} \times \alpha \quad (17)$$

177 here LOD (limit of detection) is introduced to generate a dimensionless  $\gamma_{Unc}$ .  $\alpha$  is a dimensionless adjustable  
 178 factor to control the position of  $\gamma_{Unc}$  curve on the concentration axis, which is indicated by the value of  $\gamma_{Unc}$   
 179 at LOD level. As shown in Figure 1a, at different values of  $\alpha$  ( $\alpha=1, 0.5$  and  $0.3$ ), the corresponding  $\gamma_{Unc}$  at  
 180 the same LOD level would be 100%, 50% and 30% respectively. By changing  $\alpha$ , the location of the  $\gamma_{Unc}$   
 181 curve on x axis direction can be set, using the  $\gamma_{Unc}$  at LOD as the reference point. Then Eq. (17) for POC and  
 182 EC become:

$$183 \quad -\sqrt{\frac{LOD_{POC}}{POC_{true}}} \times \alpha_{POC} \times POC_{true} \leq \varepsilon_{POC} \leq +\sqrt{\frac{LOD_{POC}}{POC_{true}}} \times \alpha_{POC} \times POC_{true} \quad (18)$$

$$184 \quad -\sqrt{\frac{LOD_{EC}}{EC_{true}}} \times \alpha_{EC} \times EC_{true} \leq \varepsilon_{EC} \leq +\sqrt{\frac{LOD_{EC}}{EC_{true}}} \times \alpha_{EC} \times EC_{true} \quad (19)$$

185 With the improved non-linear  $\gamma_{Unc}$  parameterization, concentrations of POC and EC are normalized by a  
 186 corresponding LOD, which maintains unit consistency between  $POC_{true}$  and  $\varepsilon_{POC}$  and  $EC_{true}$  and  $\varepsilon_{EC}$ , and  
 187 eliminates dependency on the concentration unit.

188 Considering a uniform distribution of measurement error and defining the variance based on a single data  
 189 point as was done above for the standard deviation, the weights in DR and YR (inverse of variance) are  
 190 calculated from:

$$191 \quad \omega(X_i) = \frac{1}{\sigma_{X_i}^2} = \frac{3}{EC_{true} \times LOD_{EC} \times \alpha_{EC}^2} \quad (20)$$

$$192 \quad \omega(Y_i) = \frac{1}{\sigma_{Y_i}^2} = \frac{3}{POC_{true} \times LOD_{POC} \times \alpha_{POC}^2} \quad (21)$$

193 The parameter  $\lambda$  in Deming regression is then determined:

$$194 \quad \lambda = \frac{\omega(X_i)}{\omega(Y_i)} = \frac{POC_{true} \times LOD_{POC} \times \alpha_{POC}^2}{EC_{true} \times LOD_{EC} \times \alpha_{EC}^2} \quad (22)$$

195 Besides the non-linear  $\gamma_{Unc}$  discussed above, a second type measurement uncertainty with a constant  $\gamma_{Unc}$  is  
 196 very common in atmospheric applications:

$$197 \quad -\gamma_{POCunc} \times POC_{true} \leq \varepsilon_{POC} \leq +\gamma_{POCunc} \times POC_{true} \quad (23)$$

$$198 \quad -\gamma_{ECunc} \times EC_{true} \leq \varepsilon_{EC} \leq +\gamma_{ECunc} \times EC_{true} \quad (24)$$





199 where  $\gamma_{Unc}$  is the relative measurement uncertainty, e.g., for relative measurement uncertainty of 10%,  
 200  $\gamma_{Unc}=0.1$ . As a result, the measurement error is linearly proportional to the concentration. An example  
 201 comparison of non-linear  $\gamma_{Unc}$  and linear  $\gamma_{Unc}$  is shown in Figure 1b. For a linear  $\gamma_{Unc}$ , the weights become:

$$202 \quad \omega(X_i) = \frac{1}{\sigma_{X_i}^2} = \frac{3}{(\gamma_{ECunc} \times EC_{true})^2} \quad (25)$$

$$203 \quad \omega(Y_i) = \frac{1}{\sigma_{Y_i}^2} = \frac{3}{(\gamma_{POCunc} \times POC_{true})^2} \quad (26)$$

204 and  $\lambda$  for Deming regression can be determined:

$$205 \quad \lambda = \frac{\omega(X_i)}{\omega(Y_i)} = \frac{(\gamma_{POCunc} \times POC_{true})^2}{(\gamma_{ECunc} \times EC_{true})^2} \quad (27)$$

### 206 **3.1.2 XY data generation by Mersenne Twister (MT) generator following a specific** 207 **distribution**

208 The Mersenne twister (MT) is a pseudorandom number generator (PRNG) developed by Matsumoto and  
 209 Nishimura (1998). MT has been widely adopted by mainstream numerical analysis software (e.g., Matlab,  
 210 SPSS, SAS and Igor Pro) as well as popular programming languages (e.g., R, Python, IDL, C++ and PHP). Data  
 211 generation using MT provides a few advantages: (1) Frequency distribution can be easily assigned during the  
 212 data generation process, allowing straightforward simulation of the frequency distribution characteristics (e.g.,  
 213 Gaussian or Log-normal) observed in ambient measurements; (2) The inputs for data generation are simply  
 214 the mean and standard deviation of the data series and can be changed easily by the user; (3) The correlation  
 215 ( $R^2$ ) between X and Y can be manipulated easily during the data generation to satisfy various purposes; (4)  
 216 Unlike the sine function described by Chu (2005) that has a sample size limitation of 120, the sample size in  
 217 MT data generation is highly flexible.

218 In this section, we will use POC as Y and EC as X as an example to explain the data generation. Procedure of  
 219 applying MT to simulate ambient POC and EC data can be found in our previous study (Wu and Yu, 2016).  
 220 Details of the data generation steps are shown in Figure 2 and described below. The first step is generation of  
 221  $EC_{true}$  by MT. In our previous study, it was found that ambient POC and EC data follow a lognormal  
 222 distribution in various locations of the Pearl River Delta (PRD) region. Therefore, lognormal distributions are  
 223 adopted during  $EC_{true}$  generation. A range of average concentration and relative standard deviation (RSD)  
 224 from ambient samples are considered in formulating the lognormal distribution. The second step is to generate  
 225  $POC_{comb}$ . As shown in Figure 2,  $POC_{comb}$  is generated by multiplying  $EC_{true}$  with  $(OC/EC)_{pri}$ . Instead of having  
 226 a Gaussian distribution,  $(OC/EC)_{pri}$  in this study is a single value, which favors direct comparison between  
 227 the true value of  $(OC/EC)_{pri}$  and  $(OC/EC)_{pri}$  estimated from the regression slope. The third step is generation



228 of  $POC_{true}$  by adding  $POC_{non-comb}$  onto  $POC_{comb}$ . Instead of having a distribution,  $POC_{non-comb}$  in this study is a  
 229 single value, which favors direct comparison between the true value of  $POC_{non-comb}$  and  $POC_{non-comb}$  estimated  
 230 from the regression intercept. The fourth step is to compute  $\varepsilon_{POC}$  and  $\varepsilon_{EC}$ . As discussed in section 3.1, two  
 231 types of measurement errors are considered for  $\varepsilon_{POC}$  and  $\varepsilon_{EC}$  calculation: nonlinear  $\gamma_{Unc}$  and linear  $\gamma_{Unc}$ . In  
 232 the last step,  $POC_{measured}$  and  $EC_{measured}$  are calculated following Eq. (12), i.e., applying measurement errors  
 233 on  $POC_{true}$  and  $EC_{true}$ . Then  $POC_{measured}$  and  $EC_{measured}$  can be used as Y and X respectively to test the  
 234 performance of various regression techniques. An Igor Pro based program with graphical user interface (GUI)  
 235 is developed to facilitate the MT data generation for OC and EC. A brief introduction is given in SI.

### 236 3.1.3 XY data generation by the sine function of Chu (2005)

237 In this section, XY data generation by sine functions is demonstrated using POC as Y and EC as X. There are  
 238 four steps in POC and EC data generation as shown by the flowchart in Figure S1. Details are explained as  
 239 follows: (1) The first step is to generate POC and EC (Chu, 2005):

$$240 \quad POC_{comb} = 14 + 12\left(\sin\left(\frac{x}{\tau}\right) + \sin(x - \phi)\right) \quad (28)$$

$$241 \quad EC_{true} = 3.5 + 3\left(\sin\left(\frac{x}{\tau}\right) + \sin(x - \phi)\right) \quad (29)$$

242 Here x is the elapsed hour ( $x=1,2,3,\dots,n$ ;  $n \leq 120$ ),  $\tau$  is used to adjust the width of each peak, and  $\phi$  is used to  
 243 adjust the phase of the sine wave. The constants 14 and 3.5 are used to lift the sine wave to the positive range  
 244 of the y-axis. An example of data generation by the sine functions of Chu (2005) is shown in Figure 3. Dividing  
 245 Eq. (28) by Eq. (29) yields a value of 4. In this way the exact relation between POC and EC is defined clearly  
 246 as  $(OC/EC)_{pri} = 4$ . (2) With  $POC_{comb}$  and  $EC_{true}$  generated, the second step is to add  $POC_{non-comb}$  to  $POC_{comb}$  to  
 247 compute  $POC_{true}$ . As for  $POC_{non-comb}$ , a single value is assigned and added to all POC following Eq. (10). Then  
 248 the goodness of regression intercept can be evaluated by comparing the regressed intercept with preset  $POC_{non-}$   
 249  $comb$ . (3) The third step is to compute  $\varepsilon_{POC}$  and  $\varepsilon_{EC}$ , considering both nonlinear  $\gamma_{Unc}$  and linear  $\gamma_{Unc}$ . (4) The  
 250 last step is to apply measurement errors on  $POC_{true}$  and  $EC_{true}$  following Eq. (12). Then  $POC_{measured}$  and  
 251  $EC_{measured}$  can be used as Y and X respectively to evaluate the performance of various regression techniques.

### 252 3.2 Ambient measurement of $\sigma_{abs}$ and EC

253 Sampling was conducted from Feb 2012 to Jan 2013 at the suburban Nancun (NC) site ( $23^{\circ} 0'11.82''N$ ,  
 254  $113^{\circ}21'18.04''E$ ), which is situated on the top of the highest peak (141 m ASL) in the Panyu district of  
 255 Guangzhou. This site is located at the geographic center of Pearl River Delta region (PRD), making it a good  
 256 location for representing the average atmospheric mixing characteristics of city clusters in the PRD region.



257 Light absorption measurements were performed by a 7- $\lambda$  Aethalometer (AE-31, Magee Scientific Company,  
258 Berkeley, CA, USA). EC mass concentrations were measured by a real time ECOC analyzer (Model RT-4,  
259 Sunset Laboratory Inc., Tigard, Oregon, USA). Both instruments were equipped with a 2.5  $\mu\text{m}$  inlet. The  
260 algorithm of Weingartner et al. (2003) was adopted to correct the sampling artifacts (aerosol loading, filter  
261 matrix and scattering effect) (Coen et al., 2010) root in Aethalometer measurement. A customized computer  
262 program with graphical user interface, Aethalometer data processor (Wu et al., 2017), was developed to  
263 perform the data correction and detailed descriptions can be found in <https://sites.google.com/site/wuchengust>.  
264 More details of the measurements can be found in Wu et al. (2017).

#### 265 **4 Comparison study using synthetic data**

266 In the following comparisons, six regression scenarios (RS) are compared using two data generation schemes  
267 (Sine function and MT) separately, as illustrated in Figure 4. Each data generation scheme considers both  
268 nonlinear  $\gamma_{Unc}$  and linear  $\gamma_{Unc}$  in measurement error parameterization. The six RS includes OLS, DR ( $\lambda =$   
269 1), DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ), ODR, WODR and YR. In commercial software (e.g., Origin, Sigma Pro, Prism, etc),  $\lambda$  in  
270 DR is sets to 1 by default if not specified. As indicated by Saylor et al. (2006), the bias observed in the study  
271 of Chu (2005) is likely due to  $\lambda = 1$  in DR. The purpose of including DR ( $\lambda = 1$ ) in this study is to examine  
272 the potential bias using the default input in many software products. The six RS are considered to examine the  
273 sensitivity of regression results to different approaches for the various parameters. For each case, 5000 runs  
274 are performed to obtain statistically significant results, as recommended by Saylor et al. (2006). The mean  
275 slope and intercept from 5000 runs is compared with the true value assigned during data generation. If the  
276 difference is  $<5\%$ , the result is considered unbiased.

#### 277 **4.1 Comparison results using the data set of Chu (2005)**

278 In this section, the scheme of Chu (2005) is adopted for data generation to obtain a benchmark of 6 regression  
279 scenarios. With different setup of slope, intercept and  $\gamma_{Unc}$ , 6 cases (Case 1 ~ 6) are studied and the results  
280 are discussed below.

##### 281 **4.1.1 Results with nonlinear $\gamma_{Unc}$**

282 A comparison of the regression techniques results with a non-linear  $\gamma_{Unc}$  (following Eqs. (18) & (19)) are  
283 summarized in Table 2.  $LOD_{POC}$ ,  $LOD_{EC}$ ,  $\alpha_{POC}$  and  $\alpha_{EC}$  are all set to 1 to reproduce the data studied by Chu  
284 (2005) and Saylor et al. (2006). Two sets of true slope and intercept are considered (Case 1: Slope=4,  
285 Intercept=0; Case 2: Slope=4, Intercept=3) to examine if any results are sensitive to the non-zero intercept.  
286 The  $R^2$  (POC, EC) from 5000 runs for both case 1 and 2 are  $0.67 \pm 0.03$ .



287 As shown in Figure 5, for the zero-intercept case (Case 1), OLS significantly underestimates the slope  
288 ( $2.95 \pm 0.14$ ) while overestimates the intercept ( $5.84 \pm 0.78$ ). This result indicates that OLS is not suitable for  
289 errors-in-variables linear regression, consistent with similar analysis results from Chu (2005) and Saylor et al.  
290 (2006). With DR, if the  $\lambda$  is properly calculated by weights ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ), unbiased slopes ( $4.01 \pm 0.25$ ) and  
291 intercepts ( $-0.04 \pm 1.28$ ) are obtained, however, results from DR with  $\lambda=1$  shows obvious bias in the slopes  
292 ( $4.27 \pm 0.27$ ) and intercepts ( $-1.45 \pm 1.36$ ). ODR also produces biased slopes ( $4.27 \pm 0.27$ ) and intercepts ( $-$   
293  $1.45 \pm 1.36$ ), which are identical to results of DR when  $\lambda=1$ . With WODR, unbiased slopes ( $3.98 \pm 0.22$ ) are  
294 observed, but the intercepts are overestimated ( $1.12 \pm 1.02$ ). Results of YR are identical to WODR. For Case 2  
295 (slope=4, intercept=3), slopes from all six regression scenarios are consistent with Case 1 (Table 2). The Case  
296 2 intercepts are equal to the Case 1 intercepts plus 3, implying that all the regression methods are not sensitive  
297 to a non-zero intercept.

298 For case 3,  $LOD_{POC}=0.5$ ,  $LOD_{EC}=0.5$ ,  $\alpha_{POC}=0.5$ ,  $\alpha_{EC}=0.5$  are adopted (Table S1), leading to an offset to the  
299 left of  $\gamma_{Unc}$  (blue curve) compared to Case 1 and 2 (black curve) in Figure 1. As a result, for the same  
300 concentration of EC and OC in Case 3, the  $\gamma_{Unc}$  is smaller than in Case 1 and Case 2 as indicated by higher  
301 the  $R^2$  ( $0.95 \pm 0.01$  for Case 3, Table S1). With a smaller measurement uncertainty, the degree of bias in Case  
302 3 is smaller than Case 1. For example, OLS slopes are less biased in Case 3 ( $3.83 \pm 0.08$ ) compare to Case 1  
303 ( $2.94 \pm 0.14$ ). Similarly, the slopes ( $4.03 \pm 0.09$ ) and intercepts ( $-0.18 \pm 0.44$ ) of DR ( $\lambda=1$ ) exhibit a much smaller  
304 bias with a smaller measurement uncertainty, implying that the degree of bias by improperly weighting in DR,  
305 WODR and YR is associated with the degree of measurement uncertainty. A higher measurement uncertainty  
306 results in larger bias in slope and intercept.

307 An uneven  $LOD_{POC}$  and  $LOD_{EC}$  is tested in Case 4 with  $LOD_{POC}=1$ ,  $LOD_{EC}=0.5$ ,  $\alpha_{POC}=0.5$ ,  $\alpha_{EC}=0.5$ , which  
308 yield a  $R^2$ (POC, EC) of  $0.78 \pm 0.02$ . The results (Table S1) are similar to Case 1. For DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ) unbiased  
309 slopes and intercepts are obtained. For WODR and YR, unbiased slopes are reported with a small bias in the  
310 intercepts. Large bias values are observed in both the slopes and intercepts in Case 4 using OLS, DR ( $\lambda = 1$ )  
311 and ODR.

#### 312 **4.1.2 Results with linear $\gamma_{Unc}$**

313 Cases 5 and 6 represent the results from using a constant  $\gamma_{Unc}$  and are shown in Table S2.  $\gamma_{Unc}$  is set to be  
314 30% to achieve a  $R^2$  (POC, EC) of 0.7, a value close to the  $R^2$  in studies of Chu (2005) and Saylor et al.  
315 (2006). In Case 5 (slope=4, intercept=0), unbiased slopes and intercepts are determined by DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ),  
316 WODR and YR. OLS underestimates the slopes ( $3.32 \pm 0.20$ ) and overestimates intercepts ( $3.77 \pm 0.90$ ), while



317 DR ( $\lambda = 1$ ) and ODR overestimate the slopes ( $4.75 \pm 0.30$ ) and underestimates the intercepts ( $-4.14 \pm 1.36$ ). In  
318 Case 6 (slope=4, intercept=3), results similar to Case 5 are obtained. It is worth noting that although the mean  
319 intercept ( $3.05 \pm 1.22$ ) of DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ), is closest to the true value (intercept=3), the deviations are much  
320 larger than for WODR ( $2.72 \pm 0.74$ ).

## 321 4.2 Comparison results using data generated by MT

322 In this section, MT is adopted for data generation to obtain a benchmark of 6 regression scenarios. Both  
323 nonlinear  $\gamma_{Unc}$  and linear  $\gamma_{Unc}$  are considered. With different setup of slope, intercept and  $\gamma_{Unc}$ , 12 cases  
324 (Case 7 ~ Case 18) are studied and the results are discussed below.

### 325 4.2.1 Nonlinear $\gamma_{Unc}$ results

326 Cases 7 and 8 use data generated by MT and a non-linear  $\gamma_{Unc}$  with results shown in Table 3. In Case 7  
327 (slope=4, intercept=0,  $LOD_{POC}=1$ ,  $LOD_{EC}=1$ ,  $\alpha_{POC}=1$ ,  $\alpha_{EC}=1$ ), unbiased slopes ( $4.00 \pm 0.03$ ) and intercepts  
328 ( $0.00 \pm 0.17$ ) are estimated by DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ). WODR and YR yield unbiased slopes ( $3.96 \pm 0.03$ ) but  
329 overestimate the intercepts ( $1.21 \pm 0.13$ ). DR ( $\lambda = 1$ ) and ODR report slightly biased slopes ( $4.17 \pm 0.04$ ) with  
330 biased intercepts ( $-0.94 \pm 0.18$ ). OLS underestimates the slopes ( $3.22 \pm 0.03$ ) and overestimates the intercept  
331 ( $4.30 \pm 0.14$ ). In Case 8 (slope=4, intercept=3,  $LOD_{POC}=1$ ,  $LOD_{EC}=1$ ,  $\alpha_{POC}=1$ ,  $\alpha_{EC}=1$ ), DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ )  
332 provides unbiased slopes ( $4.00 \pm 0.03$ ) and intercepts ( $3.00 \pm 0.18$ ) estimations. WODR and YR report unbiased  
333 slopes ( $3.97 \pm 0.03$ ) and overestimate intercepts ( $4.11 \pm 0.13$ ). OLS, DR ( $\lambda = 1$ ) and ODR report biased slopes  
334 and intercepts.

335 To test the overestimation/underestimation dependency on the true slope, Case 9 (slope=0.5, intercept=0,  
336  $LOD_{POC}=1$ ,  $LOD_{EC}=1$ ,  $\alpha_{POC}=1$ ,  $\alpha_{EC}=1$ ) and case 10 (Case 10) (slope=0.5, intercept=3,  $LOD_{POC}=1$ ,  $LOD_{EC}=1$ ,  
337  $\alpha_{POC}=1$ ,  $\alpha_{EC}=1$ ) are conducted and the results are shown in Table S3. Unlike the overestimation observed in  
338 Case 1~Case 8, DR ( $\lambda = 1$ ) and ODR underestimate the slope ( $0.46 \pm 0.01$ ) in Case 9. In case 10, DR ( $\lambda = 1$ )  
339 DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ) and ODR report unbiased slopes and intercepts. Case 11 and case 12 test the bias when the  
340 true slope is 1 as shown in Table S4. In Case 11 (intercept=0), all RS except OLS can provide unbiased results.  
341 In Case 12, all RS report unbiased slope except OLS, but DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ) is the only RS that report unbiased  
342 intercept.

343 These results imply that if the true slope is less than 1, the improper weighting ( $\lambda = 1$ ) in Deming regression  
344 and ODR without weighting tends to underestimate slope. If the true slope is 1, these two estimators can



345 provide unbiased results. If the true slope is larger than 1, the improper weighting ( $\lambda = 1$ ) in Deming  
346 regression and ODR without weighting tends to overestimate slopes.

#### 347 **4.2.2 Fixed $\gamma_{Unc}$ results**

348 Cases 13 and 14 (results shown in Table S5) represent the results from using a linear  $\gamma_{Unc}$  (30%) and data  
349 generated from MT. For case 13 (slope=4, intercept=0), DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ), WODR and YR provide the best  
350 estimation of slopes and intercepts. DR ( $\lambda = 1$ ) and ODR overestimate slopes ( $4.53 \pm 0.05$ ) and underestimate  
351 intercepts ( $-2.94 \pm 0.24$ ). For case 14 (slope=4, intercept=3), DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ), WODR and YR provide an  
352 unbiased estimation of slopes. But DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ) is the only RS reports unbiased intercepts ( $3.08 \pm 0.23$ ).  
353 Cases 15 and 16 are tested to investigate whether the results are different if the true slope is smaller than 1.  
354 As shown in Table S6, the results are similar to Table S5 that DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ) can provide unbiased slopes and  
355 intercepts while WODR and YR can provide unbiased slopes but biased intercepts. Cases 17 and 18 are tested  
356 to see if the results are the same for a special case when the true slope is 1. As shown in Table S7, the results  
357 are similar to those shown in Table S5, implying that these results are not sensitive to the special case when  
358 the true slope is 1.

#### 359 **4.3 The importance of appropriate $\lambda$ input for Deming regression**

360 As discussed above, inappropriate  $\lambda$  assignment in the Deming regression (e.g.,  $\lambda=1$  by default for many  
361 commercial software) leads to biased slope and intercept. Beside  $\lambda=1$ , inappropriate  $\lambda$  input due to improper  
362 handling of measurement uncertainty can also result in bias for Deming regression. An example is shown in  
363 Figure S2. Data is generated by MT with following parameters: slope=4, intercept=0, and linear  $\gamma_{Unc}$  (30%).  
364 Figure S2 a&b demonstrates that when an appropriate  $\lambda$  is provided (following linear  $\gamma_{Unc}$ ,  $\lambda = \frac{POC^2}{EC^2}$ ),  
365 unbiased slopes and intercepts are obtained. If an improper  $\lambda$  is used due to a mismatched measurement  
366 uncertainty assumption (non-linear  $\gamma_{Unc}$ ,  $\lambda = \frac{POC}{EC}$ ), the slopes are overestimated (Figure S2c,  $4.37 \pm 0.05$ ) and  
367 intercepts are underestimated (Figure S2d,  $-2.01 \pm 0.24$ ). This result emphasizes the importance of  
368 determining the correct form of measurement uncertainty in ambient samples, since  $\lambda$  is a crucial parameter  
369 in Deming regression.

370 In the  $\lambda$  calculation, different representations for POC and EC, including mean, median and mode, are tested  
371 as shown in Figure S3. The results show that when x and y have a similar distribution (e.g., both are log-  
372 normal), any of mean, median or mode can be used for the  $\lambda$  calculation.



## 373 5 Regression applications to ambient data

374 This section demonstrates the application of the 6 regression approaches on a light absorption coefficient and  
375 EC dataset collected in a suburban site in Guangzhou. As mentioned in the last section, measurement  
376 uncertainties are crucial inputs for DR, YR and WODR. The measurement precision of Aethalometer is 5%  
377 (Hansen, 2005) while EC by RT-ECOC analyzer is 24% (Bauer et al., 2009). These measurement uncertainties  
378 are used in DR, YR and WODR calculation. The data-set contains 6926 data points with a  $R^2$  of 0.92.

379 As shown in Figure 6, y axis is light absorption at 520 nm ( $\sigma_{\text{abs}520}$ ) and the x axis is EC mass concentration.  
380 The regressed slopes represent the mass absorption efficiency (MAE) of EC at 520 nm, ranging from 13.66 to  
381  $15.94 \text{ m}^2\text{g}^{-1}$  by the six regression scenarios. OLS yields the lowest slope (13.66 as shown in Figure 6a) among  
382 all six regression scenarios, consistent with the results by synthetic data. This implies that OLS tends to  
383 underestimate regression slope when mean Y to X ratio is larger than 1. DR ( $\lambda = 1$ ) and ODR report the same  
384 slope (14.88) and intercept (5.54), this equivalency is also observed for the synthetic data. Similarly, WODR  
385 and YR yield identical slope (14.88) and intercept (5.54), in line with the synthetic data results. The regressed  
386 slope by DR ( $\lambda = 1$ ) is higher than DR ( $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$ ), and this relationship agrees well with the synthetic data  
387 results.

388 Regression comparison is also performed on hourly OC and EC data. Regression on OC/EC percentile subset  
389 is a widely used empirical approach for primary OC/EC ratio determination. Figure S4 shows the regression  
390 slopes as a function of OC/EC percentile. OC/EC percentile ranges from 0.5% to 100%, with an interval of  
391 0.5%. As the percentile increases, SOC contribution in OC increases as well, resulting decreased  $R^2$  between  
392 OC and EC. The deviations between six RS exhibit a dependency on  $R^2$ . When percentile is relatively small  
393 (e.g., <10%), the difference between the six RS is also small due to the high  $R^2$  (0.98). The deviations between  
394 the six RS become more pronounced as  $R^2$  decrease (e.g., <0.9). The deviations are expected to be even larger  
395 when  $R^2$  is less than 0.8. These results emphasize the importance of applying error-in-variables regression,  
396 since ambient XY data more likely has a  $R^2$  less than 0.9 in most cases.

397 As discussed in this section, the ambient data confirm the comparison results obtained with the synthetic data.  
398 The advantage of using the synthetic data for regression approaches evaluation is that the ideal slope and  
399 intercept are known values during the data generation, so the bias of each regression approach can be  
400 quantified.

## 401 6 Recommendations and conclusions

402 This study aims to provide a benchmark of commonly used linear regression algorithms using a new data  
403 generation scheme. The results show that OLS fails to estimate the correct slope and intercept when



404 measurement errors are expected in both x and y, consistent with previous studies. With appropriate weighting,  
405 DR, WODR and YR are recommended for atmospheric studies when both the x and y data have measurement  
406 errors. Sensitivity tests also reveal the importance of the weighting parameter  $\lambda$  in DR. An improper  $\lambda$  could  
407 lead to biased slope and intercept. Since the  $\lambda$  estimation depends on the form of the measurement errors, it is  
408 important to determine the measurement error during the experimentation stage rather than making  
409 assumptions. For ambient data with a XY  $R^2$  less than 0.9, error-in-variables regression is needed to minimize  
410 the bias in slope and intercept. Application of error-in-variables regression is often overlooked in atmospheric  
411 studies, partly due to the lack of a specified tool for the regression implementation. To facilitate the  
412 implementation of error-in-variables regression, a computer program (Scatter plot) with graphical user  
413 interface (GUI) in Igor Pro (WaveMetrics, Inc. Lake Oswego, OR, USA) is developed (Figure 7). It packed  
414 with lots of useful features for data analysis and graph plotting, including batch plotting, data masking via  
415 GUI, color coding in Z axis, data filtering and grouping by numerical values and strings. The Scatter plot  
416 program and user manual are available from <https://sites.google.com/site/wuchengust>.  
417





418 **Appendix A: Equations of regression techniques**

419 Ordinary Least Square (**OLS**) calculation steps.

420 First calculate average of observed  $X_i$  and  $Y_i$ .

421 
$$\bar{X} = \frac{\sum_{i=1}^N X_i}{N} \quad (\text{A1})$$

422 
$$\bar{Y} = \frac{\sum_{i=1}^N Y_i}{N} \quad (\text{A2})$$

423 Then calculate  $S_{xx}$  and  $S_{yy}$ .

424 
$$S_{xx} = \sum_{i=1}^N (X_i - \bar{X})^2 \quad (\text{A3})$$

425 
$$S_{yy} = \sum_{i=1}^N (Y_i - \bar{Y})^2 \quad (\text{A4})$$

426 OLS slope and intercept can be obtained from,

427 
$$k = \frac{S_{yy}}{S_{xx}} \quad (\text{A6})$$

428 
$$b = \bar{Y} - k\bar{X} \quad (\text{A7})$$

429

430 Deming regression (**DR**) calculation steps (York, 1966).

431 Besides  $S_{xx}$  and  $S_{yy}$  as shown above,  $S_{xy}$  can be calculated from,

432 
$$S_{xy} = \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y}) \quad (\text{A8})$$

433 DR slope and intercept can be obtained from,

434 
$$k = \frac{S_{yy} - \lambda S_{xx} + \sqrt{(S_{yy} - \lambda S_{xx})^2 + 4\lambda S_{xy}^2}}{2S_{xy}} \quad (\text{A9})$$

435 
$$b = \bar{Y} - k\bar{X} \quad (\text{A10})$$

436

437 York regression (**YR**) iteration steps (York et al., 2004).

438 Slope by OLS can be used as the initial  $k$  in  $W_i$  calculation.

439 
$$W_i = \frac{\omega(X_i)\omega(Y_i)}{\omega(X_i) + k^2\omega(Y_i) - 2kr_i\sqrt{\omega(X_i)\omega(Y_i)}} \quad (\text{A11})$$



440 
$$U_i = X_i - \bar{X} = X_i - \frac{\sum_{i=1}^N W_i X_i}{\sum_{i=1}^N W_i} \quad (\text{A12})$$

441 
$$V_i = Y_i - \bar{Y} = Y_i - \frac{\sum_{i=1}^N W_i Y_i}{\sum_{i=1}^N W_i} \quad (\text{A13})$$

442 Then calculate  $\beta_i$ .

443 
$$\beta_i = W_i \left[ \frac{U_i}{\omega(Y_i)} + \frac{kV_i}{\omega(X_i)} - [kU_i + V_i] \frac{r_i}{\sqrt{\omega(X_i)\omega(Y_i)}} \right] \quad (\text{A14})$$

444 Slope and intercept can be obtained from,

445 
$$k = \frac{\sum_{i=1}^n W_i \beta_i V_i}{\sum_{i=1}^n W_i \beta_i U_i} \quad (\text{A15})$$

446 
$$b = \bar{Y} - k\bar{X} \quad (\text{A16})$$

447 Since  $W_i$  and  $\beta_i$  are functions of  $k$ ,  $k$  must be solved iteratively by repeating A11 to A15. If the difference  
 448 between the  $k$  obtained from A15 and the  $k$  used in A11 satisfies the predefined tolerance, the calculation is  
 449 considered as converged. For the data set of Chu (2005), the iteration count is around 6.

450 **Acknowledgements**

451 This work is supported by the National Natural Science Foundation of China (41605002). The author would  
 452 like to thank Dr. Bin Yu Kuang at HKUST for discussion on mathematics and Dr. Stephen M Griffith at  
 453 HKUST for valuable comments.

454

455

456 **References**

- 457 Ayers, G. P.: Comment on regression analysis of air quality data, *Atmos. Environ.*, 35, 2423-2425, doi:  
458 10.1016/S1352-2310(00)00527-6, 2001.
- 459 Bauer, J. J., Yu, X.-Y., Cary, R., Laulainen, N., and Berkowitz, C.: Characterization of the sunset semi-  
460 continuous carbon aerosol analyzer, *J. Air Waste Manage. Assoc.*, 59, 826-833, doi: 10.3155/1047-  
461 3289.59.7.826, 2009.
- 462 Boggs, P. T., Donaldson, J. R., and Schnabel, R. B.: Algorithm 676: ODRPACK: software for weighted  
463 orthogonal distance regression, *ACM Trans. Math. Softw.*, 15, 348-364, doi: 10.1145/76909.76913, 1989.
- 464 Brauers, T. and Finlayson-Pitts, B. J.: Analysis of relative rate measurements, *Int. J. Chem. Kinet.*, 29, 665-  
465 672, doi: 10.1002/(SICI)1097-4601(1997)29:9<665::AID-KIN3>3.0.CO;2-S, 1997.
- 466 Cantrell, C. A.: Technical Note: Review of methods for linear least-squares fitting of data and application to  
467 atmospheric chemistry problems, *Atmos. Chem. Phys.*, 8, 5477-5487, doi: 10.5194/acp-8-5477-2008, 2008.
- 468 Carroll, R. J. and Ruppert, D.: The use and misuse of orthogonal regression in linear errors-in-variables  
469 models, *Am. Stat.*, 50, 1-6, doi: 10.1080/00031305.1996.10473533, 1996.
- 470 Cess, R. D., Zhang, M. H., Minnis, P., Corsetti, L., Dutton, E. G., Forgan, B. W., Garber, D. P., Gates, W. L.,  
471 Hack, J. J., Harrison, E. F., Jing, X., Kiehi, J. T., Long, C. N., Morcrette, J.-J., Potter, G. L., Ramanathan, V.,  
472 Subasilar, B., Whitlock, C. H., Young, D. F., and Zhou, Y.: Absorption of solar radiation by clouds:  
473 Observations versus models, *Science*, 267, 496-499, doi: 10.1126/science.267.5197.496, 1995.
- 474 Chen, L. W. A., Doddridge, B. G., Dickerson, R. R., Chow, J. C., Mueller, P. K., Quinn, J., and Butler, W.  
475 A.: Seasonal variations in elemental carbon aerosol, carbon monoxide and sulfur dioxide: Implications for  
476 sources, *Geophys. Res. Lett.*, 28, 1711-1714, doi: 10.1029/2000GL012354, 2001.
- 477 Chu, S. H.: Stable estimate of primary OC/EC ratios in the EC tracer method, *Atmos. Environ.*, 39, 1383-  
478 1392, doi: 10.1016/j.atmosenv.2004.11.038, 2005.
- 479 Coen, M. C., Weingartner, E., Apituley, A., Ceburnis, D., Fierz-Schmidhauser, R., Flentje, H., Henzing, J. S.,  
480 Jennings, S. G., Moerman, M., Petzold, A., Schmid, O., and Baltensperger, U.: Minimizing light absorption  
481 measurement artifacts of the Aethalometer: evaluation of five correction algorithms, *Atmos. Meas. Tech.*, 3,  
482 457-474, doi: 10.5194/amt-3-457-2010, 2010.
- 483 Cornbleet, P. J. and Gochman, N.: Incorrect least-squares regression coefficients in method-comparison  
484 analysis, *Clin. Chem.*, 25, 432-438, 1979.
- 485 Cross, E. S., Onasch, T. B., Ahern, A., Wrobel, W., Slowik, J. G., Olfert, J., Lack, D. A., Massoli, P., Cappa,  
486 C. D., Schwarz, J. P., Spackman, J. R., Fahey, D. W., Sedlacek, A., Trimborn, A., Jayne, J. T., Freedman, A.,  
487 Williams, L. R., Ng, N. L., Mazzoleni, C., Dubey, M., Brem, B., Kok, G., Subramanian, R., Freitag, S., Clarke,  
488 A., Thornhill, D., Marr, L. C., Kolb, C. E., Worsnop, D. R., and Davidovits, P.: Soot particle studies—  
489 instrument inter-comparison—project overview, *Aerosol. Sci. Technol.*, 44, 592-611, doi:  
490 10.1080/02786826.2010.482113, 2010.
- 491 Deming, W. E.: *Statistical Adjustment of Data*, Wiley, New York, 1943.
- 492 Duan, F., Liu, X., Yu, T., and Cachier, H.: Identification and estimate of biomass burning contribution to the  
493 urban aerosol organic carbon concentrations in Beijing, *Atmos. Environ.*, 38, 1275-1282, doi:  
494 10.1016/j.atmosenv.2003.11.037, 2004.
- 495 Hansen, A. D. A.: *The Aethalometer Manual*, Berkeley, California, USA, Magee Scientific, 2005.



- 496 Huang, X. H., Bian, Q., Ng, W. M., Louie, P. K., and Yu, J. Z.: Characterization of PM<sub>2.5</sub> major components  
497 and source investigation in suburban Hong Kong: A one year monitoring study, *Aerosol. Air. Qual. Res.*, 14,  
498 237-250, doi: 10.4209/aaqr.2013.01.0020, 2014.
- 499 Janhäll, S., Andreae, M. O., and Pöschl, U.: Biomass burning aerosol emissions from vegetation fires: particle  
500 number and mass emission factors and size distributions, *Atmos. Chem. Phys.*, 10, 1427-1439, doi:  
501 10.5194/acp-10-1427-2010, 2010.
- 502 Lim, L. H., Harrison, R. M., and Harrad, S.: The contribution of traffic to atmospheric concentrations of  
503 polycyclic aromatic hydrocarbons, *Environ. Sci. Technol.*, 33, 3538-3542, doi: 10.1021/es990392d, 1999.
- 504 Linnet, K.: Necessary sample size for method comparison studies based on regression analysis, *Clin. Chem.*,  
505 45, 882-894, 1999.
- 506 Malm, W. C., Sisler, J. F., Huffman, D., Eldred, R. A., and Cahill, T. A.: Spatial and seasonal trends in particle  
507 concentration and optical extinction in the United-States, *J. Geophys. Res.*, 99, 1347-1370, doi:  
508 10.1029/93JD02916, 1994.
- 509 Markovsky, I. and Van Huffel, S.: Overview of total least-squares methods, *Signal Process.*, 87, 2283-2302,  
510 doi: 10.1016/j.sigpro.2007.04.004, 2007.
- 511 Matsumoto, M. and Nishimura, T.: Mersenne twister: a 623-dimensionally equidistributed uniform pseudo-  
512 random number generator, *ACM Trans. Model. Comput. Simul.*, 8, 3-30, doi: 10.1145/272991.272995, 1998.
- 513 Moosmüller, H., Arnott, W. P., Rogers, C. F., Chow, J. C., Frazier, C. A., Sherman, L. E., and Dietrich, D. L.:  
514 Photoacoustic and filter measurements related to aerosol light absorption during the Northern Front Range Air  
515 Quality Study (Colorado 1996/1997), *J. Geophys. Res.*, 103, 28149-28157, doi: 10.1029/98jd02618, 1998.
- 516 Petäjä, T., Mauldin, I. R. L., Kosciuch, E., McGrath, J., Nieminen, T., Paasonen, P., Boy, M., Adamov, A.,  
517 Kotiaho, T., and Kulmala, M.: Sulfuric acid and OH concentrations in a boreal forest site, *Atmos. Chem.*  
518 *Phys.*, 9, 7435-7448, doi: 10.5194/acp-9-7435-2009, 2009.
- 519 Richter, A., Burrows, J. P., Nusz, H., Granier, C., and Niemeier, U.: Increase in tropospheric nitrogen dioxide  
520 over China observed from space, *Nature*, 437, 129-132, doi: 10.1038/nature04092, 2005.
- 521 Saylor, R. D., Edgerton, E. S., and Hartsell, B. E.: Linear regression techniques for use in the EC tracer method  
522 of secondary organic aerosol estimation, *Atmos. Environ.*, 40, 7546-7556, doi:  
523 10.1016/j.atmosenv.2006.07.018, 2006.
- 524 Turpin, B. J. and Huntzicker, J. J.: Identification of secondary organic aerosol episodes and quantitation of  
525 primary and secondary organic aerosol concentrations during SCAQS, *Atmos. Environ.*, 29, 3527-3544, doi:  
526 10.1016/1352-2310(94)00276-Q, 1995.
- 527 von Bobruzki, K., Braban, C. F., Famulari, D., Jones, S. K., Blackall, T., Smith, T. E. L., Blom, M., Coe, H.,  
528 Gallagher, M., Ghalaieny, M., McGillen, M. R., Percival, C. J., Whitehead, J. D., Ellis, R., Murphy, J.,  
529 Mohacsi, A., Pogany, A., Junninen, H., Rantanen, S., Sutton, M. A., and Nemitz, E.: Field inter-comparison  
530 of eleven atmospheric ammonia measurement techniques, *Atmos. Meas. Tech.*, 3, 91-112, doi: 10.5194/amt-  
531 3-91-2010, 2010.
- 532 Wang, J. and Christopher, S. A.: Intercomparison between satellite-derived aerosol optical thickness and  
533 PM<sub>2.5</sub> mass: Implications for air quality studies, *Geophys. Res. Lett.*, 30, 2095, doi: 10.1029/2003gl018174,  
534 2003.
- 535 Watson, J. G.: Visibility: Science and regulation, *J. Air Waste Manage. Assoc.*, 52, 628-713, doi:  
536 10.1080/10473289.2002.10470813, 2002.



- 537 Weingartner, E., Saathoff, H., Schnaiter, M., Streit, N., Bitnar, B., and Baltensperger, U.: Absorption of light  
538 by soot particles: determination of the absorption coefficient by means of aethalometers, *J. Aerosol. Sci.*, 34,  
539 1445-1463, doi: 10.1016/S0021-8502(03)00359-8, 2003.
- 540 Wu, C. and Yu, J. Z.: Determination of primary combustion source organic carbon-to-elemental carbon  
541 (OC/EC) ratio using ambient OC and EC measurements: secondary OC-EC correlation minimization method,  
542 *Atmos. Chem. Phys.*, 16, 5453-5465, doi: 10.5194/acp-16-5453-2016, 2016.
- 543 Wu, C., Wu, D., and Yu, J. Z.: Quantifying black carbon light absorption enhancement by a novel statistical  
544 approach, *Atmos. Chem. Phys. Discuss.*, 2017, 1-37, doi: 10.5194/acp-2017-582, 2017.
- 545 York, D.: Least-squares fitting of a straight line, *Can. J. Phys.*, 44, 1079-1086, doi: 10.1139/p66-090, 1966.
- 546 York, D., Evensen, N. M., Martinez, M. L., and Delgado, J. D. B.: Unified equations for the slope, intercept,  
547 and standard errors of the best straight line, *Am. J. Phys.*, 72, 367-375, doi: 10.1119/1.1632486, 2004.
- 548 Zhou, Y., Huang, X. H. H., Griffith, S. M., Li, M., Li, L., Zhou, Z., Wu, C., Meng, J., Chan, C. K., Louie, P.  
549 K. K., and Yu, J. Z.: A field measurement based scaling approach for quantification of major ions, organic  
550 carbon, and elemental carbon using a single particle aerosol mass spectrometer, *Atmos. Environ.*, 143, 300-  
551 312, doi: 10.1016/j.atmosenv.2016.08.054, 2016.
- 552 Zieger, P., Weingartner, E., Henzing, J., Moerman, M., de Leeuw, G., Mikkilä, J., Ehn, M., Petäjä, T., Clémer,  
553 K., van Roozendaal, M., Yilmaz, S., Frieß, U., Irie, H., Wagner, T., Shaiganfar, R., Beirle, S., Apituley, A.,  
554 Wilson, K., and Baltensperger, U.: Comparison of ambient aerosol extinction coefficients obtained from in-  
555 situ, MAX-DOAS and LIDAR measurements at Cabauw, *Atmos. Chem. Phys.*, 11, 2603-2624, doi:  
556 10.5194/acp-11-2603-2011, 2011.
- 557 Zwolak, J. W., Boggs, P. T., and Watson, L. T.: Algorithm 869: ODRPACK95: A weighted orthogonal  
558 distance regression code with bound constraints, *ACM Trans. Math. Softw.*, 33, 27, doi:  
559 10.1145/1268776.1268782, 2007.
- 560



561 **Table 1.** Abbreviations.

Abbreviation	Definition
$\alpha$	a dimensionless adjustable factor to control the position of $\gamma_{Unc}$ curve on the concentration axis
$\gamma_{Unc}$	fractional measurement uncertainties relative to the true concentration (%)
DR	Deming regression
$\epsilon_{EC}, \epsilon_{POC}$	absolute measurement uncertainties of EC and POC
EC	elemental carbon
$EC_{true}$	numerically synthesized true EC concentration without measurement uncertainty
$EC_{measured}$	EC with measurement error ( $EC_{true} + \epsilon_{EC}$ )
$\lambda$	$\omega(X_i)$ to $\omega(Y_i)$ ratio in Deming regression
LOD	limit of detection
MT	Mersenne twister pseudorandom number generator
OC	organic carbon
OC/EC	OC to EC ratio
$(OC/EC)_{pri}$	primary OC/EC ratio
$OC_{non-comb}$	OC from non-combustion sources
ODR	orthogonal distance regression
OLS	ordinary least squares regression
POC	primary organic carbon
$POC_{comb}$	numerically synthesized true POC from combustion sources (well correlated with $EC_{true}$ ), measurement uncertainty not considered
$POC_{non-comb}$	numerically synthesized true POC from non-combustion sources (independent of $EC_{true}$ ) without considering measurement uncertainty
$POC_{true}$	sum of $POC_{comb}$ and $POC_{non-comb}$ without considering measurement uncertainty
$POC_{measured}$	POC with measurement error ( $POC_{true} + \epsilon_{POC}$ )
$\sigma_{X_i}, \sigma_{Y_i}$	the standard deviation of the error in measurement of $X_i$ and $Y_i$
S	sum of squared residuals
SOC	secondary organic carbon
$\tau$	parameter in the sine function of Chu (2005) that adjust the width of each peak
$\phi$	parameter in the sine function of Chu (2005) that adjust the phase of the curve
WODR	weight orthogonal distance regression
YR	York regression
$\omega(X_i), \omega(Y_i)$	inverse of $\sigma_{X_i}$ and $\sigma_{Y_i}$ , used as weights in DR calculation.

562

563



564 **Table 2.** Comparison of regression techniques using the data generation scheme of Chu (2005) with 5000 runs  
 565 considering both zero intercept and non-zero intercept conditions. Non-linear relative measurement  
 566 uncertainty follows Eqs. (18) & (19). Unbiased (difference <5% to true value) slopes and intercepts are  
 567 highlighted in grey.

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 1 Slope=4, Intercept=0 $LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1$ . $R^2(POC,EC) = 0.67 \pm 0.03$	OLS	2.94	$\pm 0.14$	5.84	$\pm 0.78$
	DR $\lambda=1$	4.27	$\pm 0.27$	-1.45	$\pm 1.36$
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.01	$\pm 0.25$	-0.04	$\pm 1.28$
	ODR	4.27	$\pm 0.27$	-1.45	$\pm 1.36$
	WODR	3.98	$\pm 0.22$	1.12	$\pm 1.02$
	YR	3.98	$\pm 0.22$	1.12	$\pm 1.02$
Case 2 Slope=4, Intercept=3 $LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1$ . $R^2(POC,EC) = 0.67 \pm 0.04$	OLS	2.95	$\pm 0.15$	8.83	$\pm 0.80$
	DR1 $\lambda=1$	4.32	$\pm 0.28$	1.28	$\pm 1.43$
	DR1 $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.01	$\pm 0.26$	2.94	$\pm 1.34$
	ODR	4.32	$\pm 0.28$	1.28	$\pm 1.43$
	WODR	3.99	$\pm 0.23$	3.98	$\pm 1.05$
	YR	3.99	$\pm 0.23$	3.98	$\pm 1.05$

568



569 **Table 3.** Comparison of regression techniques using MT data generation scheme with 5000 runs with non-  
 570 linear relative measurement uncertainty following Eqs. (18) & (19). Unbiased (difference <5% to true value)  
 571 slopes and intercepts are highlighted with grey background.

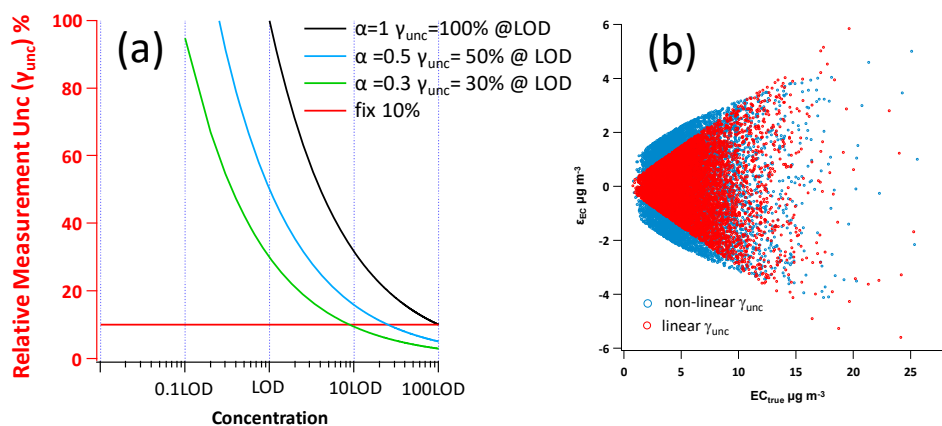
Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 7	OLS	3.22	±0.03	4.30	±0.14
Slope=4, Intercept=0	DR $\lambda=1$	4.17	±0.04	-0.94	±0.18
$LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1.$	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.00	±0.03	0.00	±0.17
$R^2(POC,EC) =$ 0.76±0.01	ODR	4.17	±0.04	-0.94	±0.18
	WODR	3.96	±0.03	1.21	±0.13
	YR	3.96	±0.03	1.21	±0.13
Case 8	OLS	3.22	±0.03	7.29	±0.14
Slope=4, Intercept=3	DR $\lambda=1$	4.20	±0.04	1.88	±0.18
$LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1.$	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.00	±0.03	3.00	±0.18
$R^2(POC,EC) =$ 0.75±0.01	ODR	4.20	±0.04	1.88	±0.18
	WODR	3.97	±0.03	4.11	±0.13
	YR	3.97	±0.03	4.11	±0.13

572





573

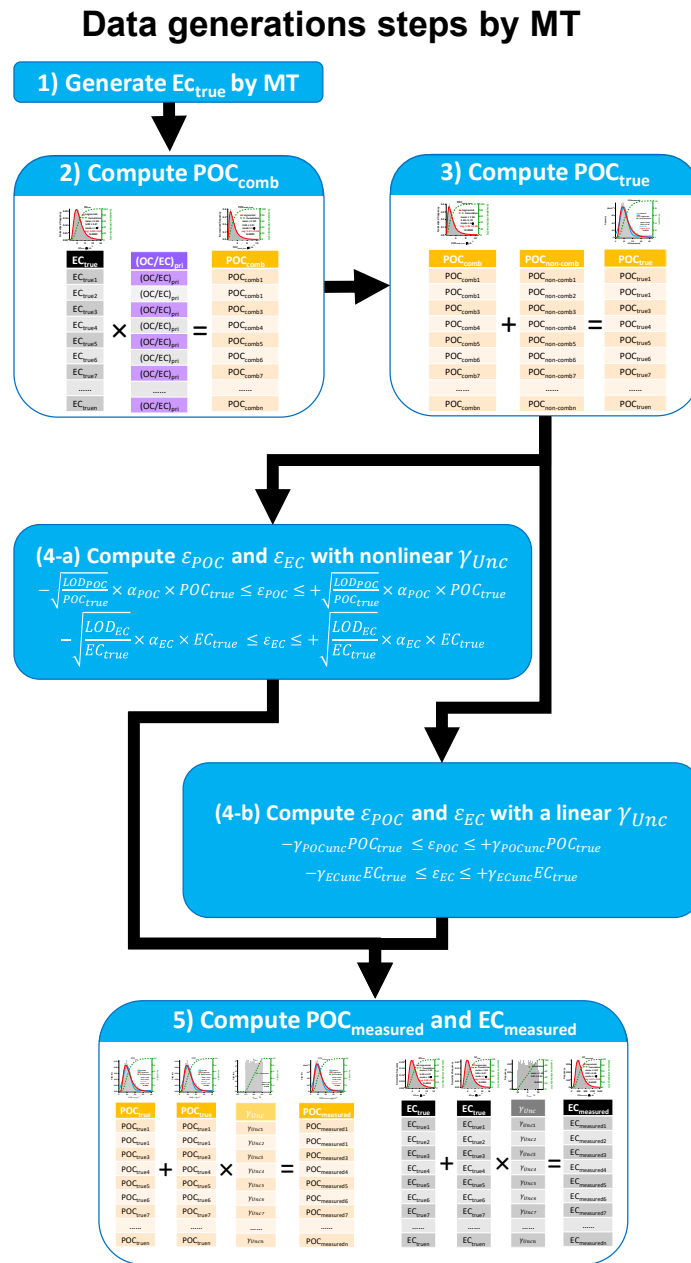


574

575 **Figure 1.** (a) Example  $\gamma_{unc}$  curves by different  $\alpha$  values (Eq. (17)). The x-axis is concentration (normalized  
 576 by LOD) and y axis is  $\gamma_{unc}$ . Black, blue and green line represent  $\alpha$  equal to 1, 0.5 and 0.3 respectively,  
 577 corresponding to the  $\gamma_{unc}$  at LOD level equals to 100%, 50% and 30% respectively. (b) Example of  
 578 measurement uncertainty generation with non-linear  $\gamma_{unc}$  and linear  $\gamma_{unc}$  constrain. The blue circles represent  
 579 non-linear  $\gamma_{unc}$  following Eq. (17) ( $LOD_{EC} = 1$ ,  $a_{EC} = 1$ ). The red circles represent linear  $\gamma_{unc}$  (30%).  
 580



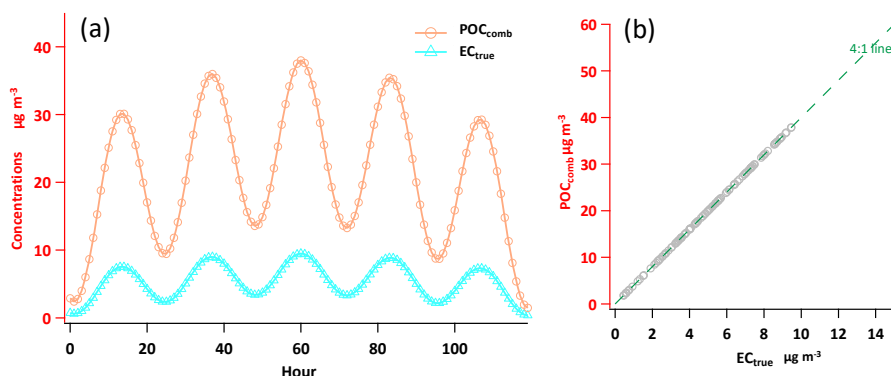
581



582

583 **Figure 2.** Flowchart of data generation steps using MT.

584



$$POC_{comb} = 14 + 12\left(\sin\left(\frac{x}{\tau}\right) + \sin(x - \phi)\right)$$

$$EC_{true} = 3.5 + 3\left(\sin\left(\frac{x}{\tau}\right) + \sin(x - \phi)\right)$$

585

586 **Figure 3.**  $POC_{comb}$  and  $EC_{true}$  data generated by the sine functions of Chu (2005). (a) Time series of the 120

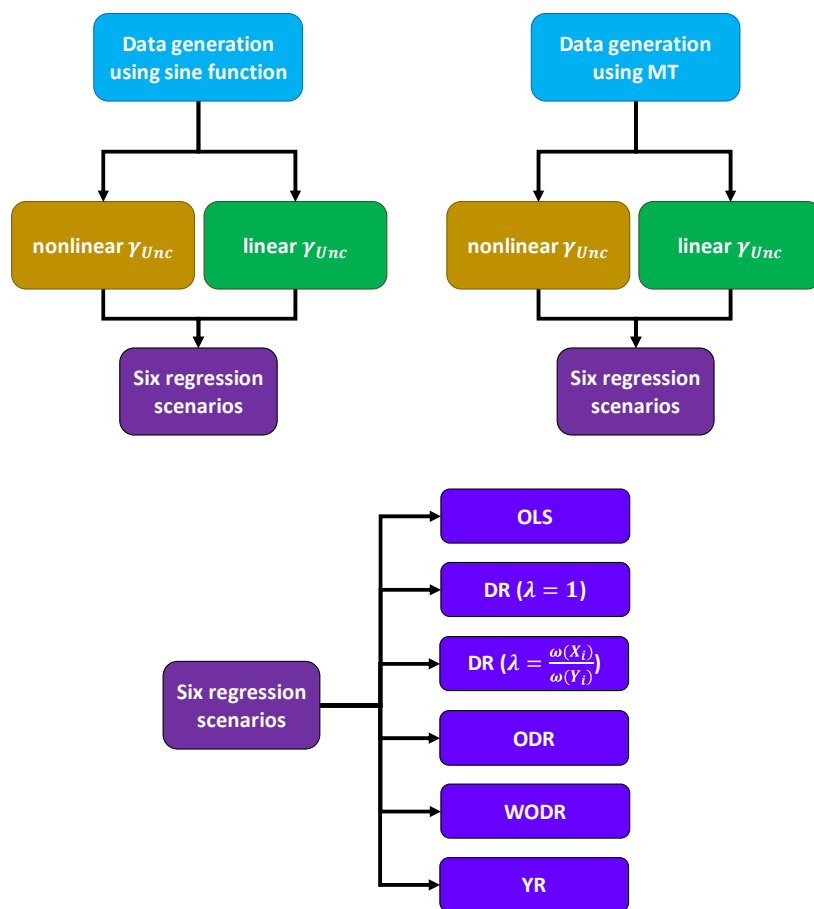
587 data points for  $POC_{comb}$  and  $EC_{true}$ . (b) Scatter plot of  $POC_{comb}$  vs.  $EC_{true}$

588

589



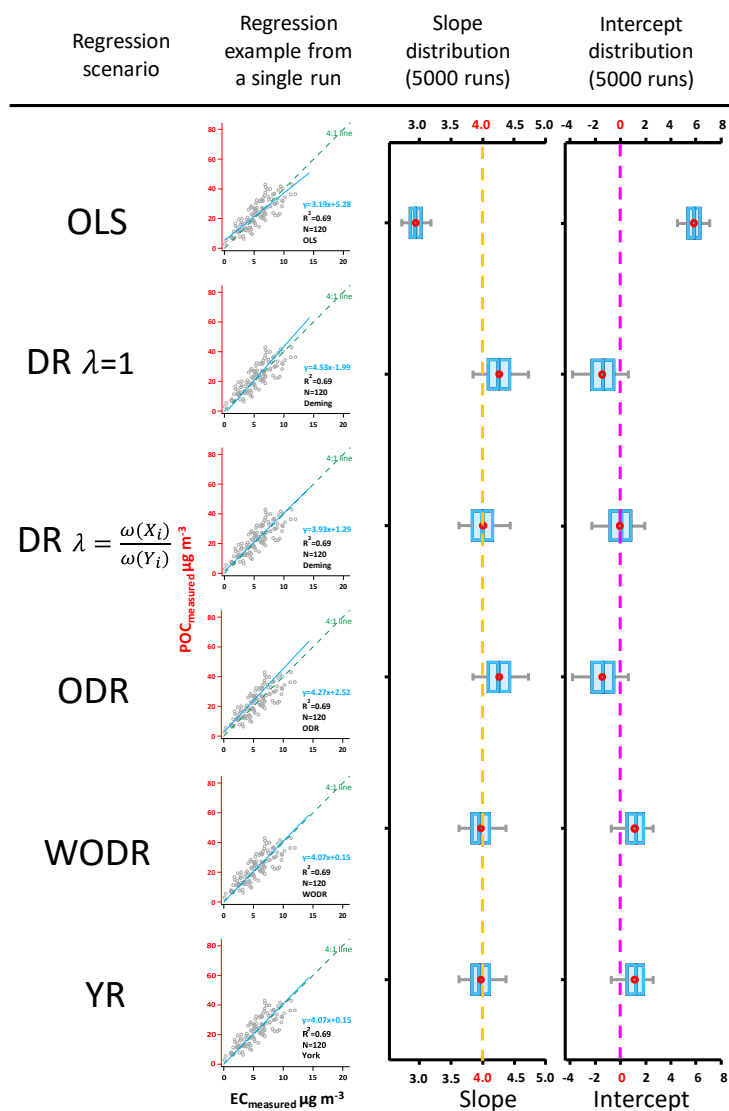
## Comparison study design



590

591 **Figure 4.** Overview of the comparison study design.

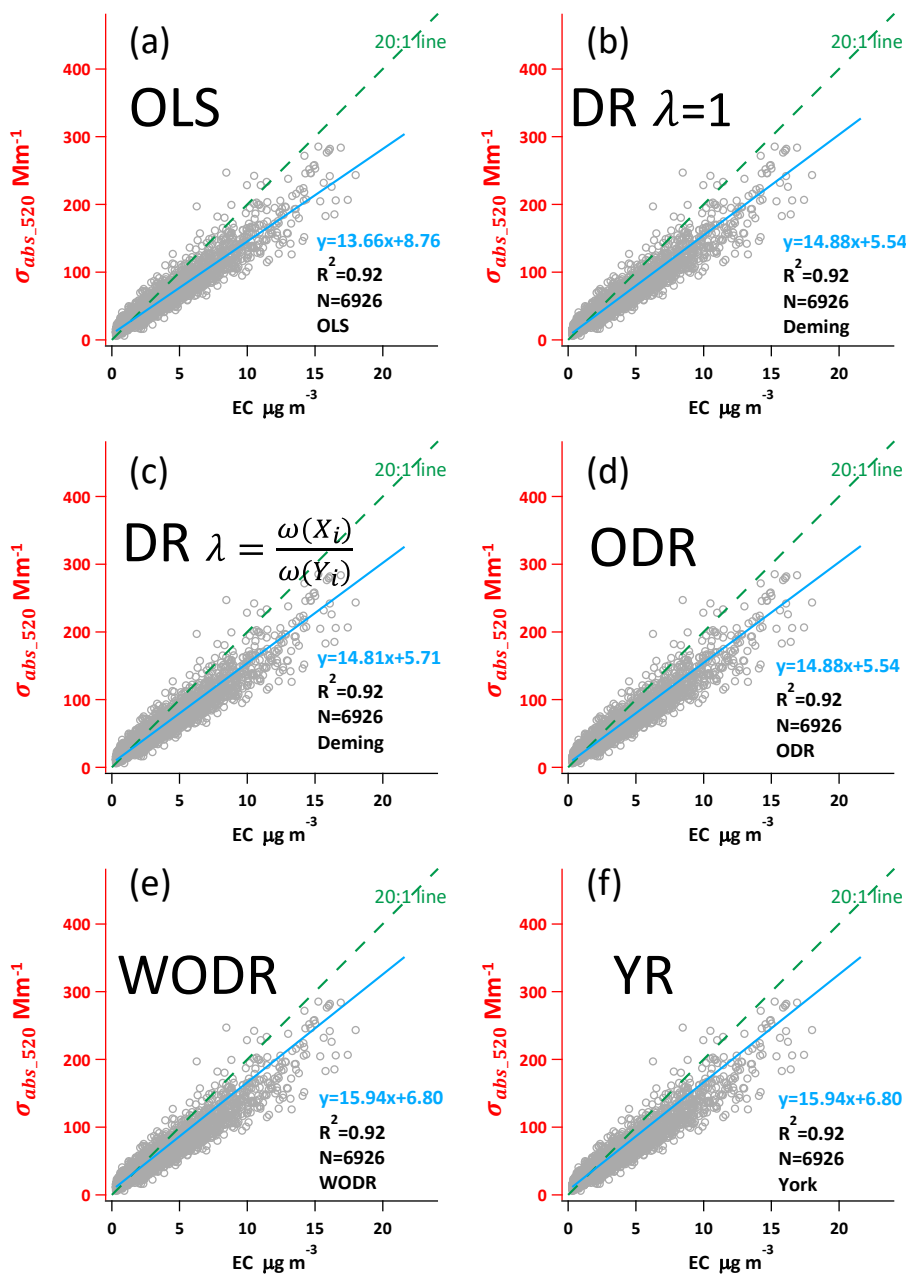
592



593

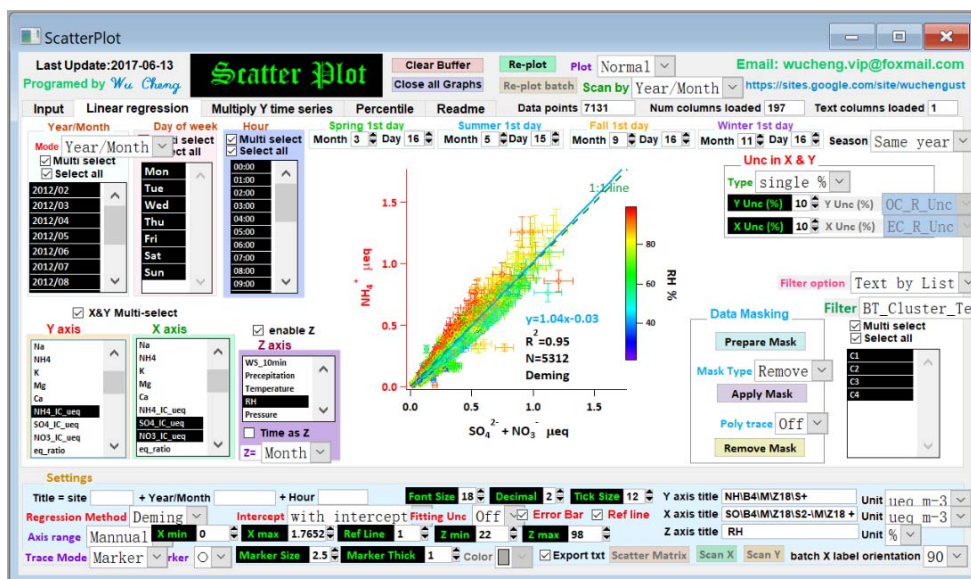
594 **Figure 5.** Regression results on synthetic data, case 1 (Slope=4, Intercept=0,  $LOD_{POC}=1$ ,  $LOD_{EC}=1$ ,  $a_{POC}=1$ ,  
 595  $a_{EC}=1$ ,  $R^2(POC, EC)=0.67\pm 0.03$ ). The scatter plots demonstrate regression examples from a single run. The  
 596 box plots show the distribution of regressed slopes and intercepts from 5000 runs of six regression scenarios.  
 597 The dashed line in orange and peachblow represent true slope and intercept respectively.

598



599

600 **Figure 6.** Regression results using ambient  $\sigma_{abs520}$  and EC data from a suburban site in Guangzhou, China.



601

602 **Figure 7.** The user interface of Scatter Plot Igor program. The program and its operation manual are available

603 from: <https://sites.google.com/site/wuchengust>.