

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

4 **Cheng Wu^{1,2} and Jian Zhen Yu^{3,4,5}**

5 ¹Institute of Mass Spectrometer and Atmospheric Environment, Jinan University,
6 Guangzhou 510632, China

7 ²Guangdong Provincial Engineering Research Center for on-line source
8 apportionment system of air pollution, Guangzhou 510632, China

9 ³Division of Environment, Hong Kong University of Science and Technology, Clear
10 Water Bay, Hong Kong, China

11 ⁴Atmospheric Research Centre, Fok Ying Tung Graduate School, Hong Kong
12 University of Science and Technology, Nansha, China

13 ⁵Department of Chemistry, Hong Kong University of Science and Technology, Clear
14 Water Bay, Hong Kong, China

15 *Corresponding to:* Cheng Wu (wucheng.vip@foxmail.com) and Jian Zhen Yu
16 (jian.yu@ust.hk)

17

18 **Abstract**

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

41

42 **1 Introduction**

43 Linear regression is heavily used in atmospheric science to derive the slope and
44 intercept of XY datasets. Examples of linear regression applications include primary
45 OC (organic carbon) and EC (elemental carbon) ratio estimation (Turpin and
46 Huntzicker, 1995), MAE (mass absorption efficiency) estimation from light absorption
47 and EC mass (Moosmüller et al., 1998), source apportionment of polycyclic aromatic
48 hydrocarbons using CO and NO_x as combustion tracers (Lim et al., 1999), gas-phase
49 reaction rate determination (Brauers and Finlayson-Pitts, 1997), inter-instrument
50 comparison (Bauer et al., 2009; Cross et al., 2010; von Bobruzki et al., 2010; Zieger et
51 al., 2011; Huang et al., 2014; Zhou et al., 2016), light extinction budget reconstruction
52 (Malm et al., 1994; Watson, 2002), comparison between modeling and measurement
53 (Petäjä et al., 2009), emission factor study (Janhäll et al., 2010), retrieval of shortwave
54 cloud forcing (Cess et al., 1995), calculation of pollutant growth rate (Richter et al.,
55 2005), estimation of ground level PM_{2.5} from MODIS data (Wang and Christopher,
56 2003), distinguishing OC origin from biomass burning using K⁺ as a tracer (Duan et al.,
57 2004) and emission type identification by the EC/CO ratio (Chen et al., 2001).

58 Ordinary least squares (OLS) regression is the most widely used method due to its
59 simplicity. In OLS, it is assumed that independent variables are error free. This is the
60 case for certain applications, such as determining a calibration curve of an instrument
61 in analytical chemistry. For example, a known amount of analyte (e.g., through
62 weighing) can be used to calibrate the instrument output response (e.g., voltage).
63 However, in many other applications, such as inter-instrument comparison, X and Y
64 (from two instruments) may have comparable degrees of uncertainty. This deviation
65 from the underlying assumption in OLS would produce biased slope and intercept when
66 OLS is applied to the dataset.

67 To overcome the drawback of OLS, a number of error-in-variable regression models
68 (also known as bivariate fittings (Cantrell, 2008) or total least-squares methods
69 (Markovsky and Van Huffel, 2007) arise. Deming (1943) proposed an approach by
70 minimizing sum of squares of X and Y residuals. A closed-form solution of Deming
71 regression (DR) was provided by York (1966). Method comparison work of various
72 regression techniques by Cornbleet and Gochman (1979) found significant error in OLS

73 slope estimation when the relative standard deviation (RSD) of measurement error in
74 “X” exceeded 20%, while DR was found to reach a more accurate slope estimation. In
75 an early application of the EC tracer method, Turpin and Huntzicker (1995) realized
76 the limitation of OLS since OC and EC have comparable measurement uncertainty,
77 thus recommended the use of DR for $(OC/EC)_{pri}$ (primary OC to EC ratio) estimation.
78 Ayers (2001) conducted a simple numerical experiment and concluded that reduced
79 major axis regression (RMA) is more suitable for air quality data regression analysis.
80 Linnet (1999) pointed out that when applying DR for inter-method (or inter-instrument)
81 comparison, special attention should be paid to the sample size. If the range ratio
82 (max/min) is relatively small (e.g., less than 2), more samples are needed to obtain
83 statistically significant results.

84 In principle, a best-fit regression line should have greater dependence on the more
85 precise data points rather than the less reliable ones. Chu (2005) performed a
86 comparison study of OLS and DR specifically focusing on the EC tracer method
87 application, and found the slope estimated by DR is closer to the correct value than
88 OLS but may still overestimate the ideal value. Saylor et al. (2006) extended the
89 comparison work of Chu (2005) by including a regression technique developed by York
90 et al. (2004). They found that the slope overestimation by DR in the study of Chu (2005)
91 was due to improper configuration of the weighting parameter, λ . This λ value is the
92 key to handling the uneven errors between data points for the best-fit line calculation.
93 This example demonstrates the importance of appropriate weighting in the calculation
94 of best-fit line for error-in-variable regression model, which is overlooked in many
95 studies.

96 In this study, we extend the work by Saylor et al. (2006) to achieve four objectives.
97 The first is to propose a new data generation scheme by applying the Mersenne Twister
98 (MT) pseudorandom number generator for evaluation of linear regression techniques.
99 In the study of Chu (2005), data generation is achieved by a variational sine function,
100 which has limitations in sample size, sample distribution, and nonadjustable correlation
101 (R^2) between X and Y. In comparison, the MT data generation provides more
102 flexibility, permitting adjustable sample size, XY correlation and distribution. The
103 second is to develop a non-linear measurement error parameterization scheme for use
104 in the regression method. The third is to incorporate linear measurement errors in the

105 regression methods. In the work by Chu (2005) and Saylor et al. (2006), the relative
 106 measurement uncertainty (γ_{Unc}) is non-linear with concentration, but a constant γ_{Unc}
 107 is often applied on atmospheric instruments due to its simplicity. The fourth is to
 108 include weighted orthogonal distance regression (WODR) for comparison.
 109 Abbreviations and symbols used in this study are summarized in Table 1 for quick
 110 lookup.

111 **2 Description of regression techniques compared in this study**

112 **Ordinary least squares (OLS) method.** OLS only considers the errors in dependent
 113 variables (Y). OLS regression is achieved by minimizing the sum of squares (S) in the
 114 Y residuals:

$$115 \quad S = \sum_{i=1}^n (y_i - Y_i)^2 \quad (1)$$

116 where Y_i are observed Y data points while y_i are regressed Y data points of the
 117 regression line.

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

$$121 \quad S = \sum_{i=1}^n [(x_i - X_i)^2 + (y_i - Y_i)^2] \quad (2)$$

122 **Weighted orthogonal distance regression (WODR).** Unlike ODR that considers even
 123 error in X and Y, weightings based on measurement errors in both X and Y are
 124 considered in WODR when minimizing the sum of squared orthogonal distance from
 125 the data points to the regression line (Carroll and Ruppert, 1996):

$$126 \quad S = \sum_{i=1}^n [(x_i - X_i)^2 + (y_i - Y_i)^2 / \eta] \quad (3)$$

127 where η is error variance ratio. Implementation of ODR and WODR in Igor was done
 128 by the computer routine ODRPACK95 (Boggs et al., 1989; Zwolak et al., 2007).

129 **Deming regression (DR).** Deming (1943) proposed the following function to minimize
 130 both the X and Y residuals,

$$131 \quad S = \sum_{i=1}^n [\omega(X_i)(x_i - X_i)^2 + \omega(Y_i)(y_i - Y_i)^2] \quad (4)$$

132 where X_i and Y_i are observed data points and x_i and y_i are regressed data points.
 133 Individual data points are weighted based on errors in X_i and Y_i ,

134
$$\omega(X_i) = \frac{1}{\sigma_{X_i}^2}, \quad \omega(Y_i) = \frac{1}{\sigma_{Y_i}^2} \quad (5)$$

135 where σ_{X_i} and σ_{Y_i} are the standard deviation of the error in measurement of X_i and Y_i
 136 respectively. The closed form solutions for slope and intercept of DR are shown in
 137 Appendix A.

138 **York regression (YR).** The York method (York et al., 2004) introduces the correlation
 139 coefficient of errors in X and Y into the minimization function.

140
$$S = \sum_{i=1}^n \left[\omega(X_i)(x_i - X_i)^2 - 2r_i \sqrt{\omega(X_i)\omega(Y_i)}(x_i - X_i)(y_i - Y_i) + \omega(Y_i)(y_i - \right.$$

 141
$$\left. Y_i)^2 \right] \frac{1}{1-r_i^2} \quad (6)$$

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

145 **3 Data description**

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

153 **3.1 Synthetic XY data generation**

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

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

159
$$Y = kX + b \quad (7)$$

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

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

$$169 \quad \quad \quad OC = POC_{comb} + POC_{non-comb} + SOC \quad (8)$$

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

$$174 \quad \quad \quad POC_{comb} = (OC/EC)_{pri} \times EC \quad (9)$$

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

$$177 \quad \quad \quad POC = (OC/EC)_{pri} \times EC + POC_{non-comb} \quad (10)$$

178 The regressed slope of POC (Y) against EC (X) represents $(OC/EC)_{pri}$ (k in Eq.(7)). The
179 regressed intercept become $POC_{non-comb}$ (b in Eq. (7)). With known $(OC/EC)_{pri}$ and
180 $POC_{non-comb}$, SOC can be estimated by:

$$181 \quad \quad \quad SOC = OC - ((OC/EC)_{pri} \times EC + POC_{non-comb}) \quad (11)$$

182 The data generation starts from EC (X values). Once EC is generated, POC_{comb} (the part
183 of Y that is correlated with X) can be obtained by multiplying EC with a preset constant,
184 $(OC/EC)_{pri}$ (slope k). Then the other preset constant $POC_{non-comb}$ is added to POC_{comb}
185 and the sum becomes POC (Y values). To simulate the real-world situation,
186 measurement errors are added on X and Y values. Details of synthesized measurement
187 error are discussed in the next section. Implementation of data generation by two types
188 of mathematical schemes are explained in section 3.1.2 and 3.1.3 respectively.

189 **3.1.1 Parameterization of synthesized measurement uncertainty**

190 Weighting of variables is a crucial input for errors-in-variables linear regression
 191 methods such as DR, YR and WODR. In practice, the weights are usually defined as
 192 the inverse of the measurement error variance (Eq. (5)). When measurement errors are
 193 considered, measured concentrations ($Conc_{\text{measured}}$) are simulated by adding
 194 measurement uncertainties (ε_{Conc}) to the true concentrations ($Conc_{\text{true}}$):

195
$$Conc_{\text{measured}} = Conc_{\text{true}} + \varepsilon_{Conc}. \quad (12)$$

196 Here ε_{Conc} is the random error following an even distribution with an average of 0, the
 197 range of which is constrained by:

198
$$-\gamma_{Unc} \times Conc_{\text{true}} \leq \varepsilon_{Conc} \leq +\gamma_{Unc} \times Conc_{\text{true}} \quad (13)$$

199 The γ_{Unc} is a dimensionless factor that describes the fractional measurement
 200 uncertainties relative to the true concentration ($Conc_{\text{true}}$). γ_{Unc} could be a function of
 201 $Conc_{\text{true}}$ (Thompson, 1988) or a constant. The term $\gamma_{Unc} \times Conc_{\text{true}}$ defines the
 202 boundary of random measurement errors.

203 Two types of measurement error are considered in this study. The first type is
 204 $\gamma_{Unc\text{-nonlinear}}$. In the data generation scheme of Chu (2005) for the measurement
 205 uncertainties (ε_{POC} and ε_{EC}), $\gamma_{Unc\text{-nonlinear}}$ is non-linearly related to $Conc_{\text{true}}$:

206
$$\gamma_{Unc\text{-nonlinear}} = \frac{1}{\sqrt{Conc_{\text{true}}}} \quad (14)$$

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

208
$$-\frac{1}{\sqrt{POC_{\text{true}}}} \times POC_{\text{true}} \leq \varepsilon_{POC} \leq +\frac{1}{\sqrt{POC_{\text{true}}}} \times POC_{\text{true}} \quad (15)$$

209
$$-\frac{1}{\sqrt{EC_{\text{true}}}} \times EC_{\text{true}} \leq \varepsilon_{EC} \leq +\frac{1}{\sqrt{EC_{\text{true}}}} \times EC_{\text{true}} \quad (16)$$

210 In Eq. (14), the γ_{Unc} decreases as concentration increases, since low concentrations are
 211 usually more challenging to measure. As a result, the $\gamma_{Unc\text{-nonlinear}}$ defined in Eq.
 212 (14) is more realistic than the constant approach, but there are two limitations. First, the
 213 physical meaning of the uncertainty unit is lost. If the unit of OC is $\mu\text{g m}^{-3}$, then the
 214 unit of ε_{OC} becomes $\sqrt{\mu\text{g m}^{-3}}$. Second, the concentration is not normalized by a
 215 consistent relative value, making it sensitive to the X and Y units used. For example, if

216 $POC_{true}=0.9 \mu\text{g m}^{-3}$, then $\varepsilon_{POC}=\pm 0.95 \mu\text{g m}^{-3}$ and $\gamma_{Unc} = 105\%$, but by changing the
 217 concentration unit to $POC_{true}=900 \text{ ng m}^{-3}$, then $\varepsilon_{OC}=\pm 30 \text{ ng m}^{-3}$ and $\gamma_{Unc} = 3\%$. To
 218 overcome these deficiencies, we propose to modify Eq. (14) to:

$$219 \quad \gamma_{Unc} = \sqrt{\frac{LOD}{Conc.true}} \times \alpha \quad (17)$$

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

$$227 \quad -\sqrt{\frac{LOD_{POC}}{POC_{true}}} \times \alpha_{POC} \times POC_{true} \leq \varepsilon_{POC} \leq +\sqrt{\frac{LOD_{POC}}{POC_{true}}} \times \alpha_{POC} \times POC_{true}$$

$$228 \quad (18)$$

$$229 \quad -\sqrt{\frac{LOD_{EC}}{EC_{true}}} \times \alpha_{EC} \times EC_{true} \leq \varepsilon_{EC} \leq +\sqrt{\frac{LOD_{EC}}{EC_{true}}} \times \alpha_{EC} \times EC_{true} \quad (19)$$

230 With the modified $\gamma_{Unc-nonlinear}$ parameterization, concentrations of POC and EC are
 231 normalized by a corresponding LOD, which maintains unit consistency between
 232 POC_{true} and ε_{POC} and EC_{true} and ε_{EC} , and eliminates dependency on the concentration
 233 unit.

234 Uniform distribution had been used in previous studies (Cox et al., 2003; Chu, 2005;
 235 Saylor et al., 2006) and is adopted in this study to parameterize measurement error. For
 236 a uniform distribution in the interval [a,b], the variance is $\frac{1}{12}(a-b)^2$. Since ε_{POC} and
 237 ε_{EC} follows a uniform distribution in the interval as given by Eqs. (18) and (19), the
 238 weights in DR and YR (inverse of variance) become:

$$239 \quad \omega(X_i) = \frac{1}{\sigma_{X_i}^2} = \frac{3}{EC_{true} \times LOD_{EC} \times \alpha_{EC}^2} \quad (20)$$

$$240 \quad \omega(Y_i) = \frac{1}{\sigma_{Y_i}^2} = \frac{3}{POC_{true} \times LOD_{POC} \times \alpha_{POC}^2} \quad (21)$$

241 The parameter λ in Deming regression is then determined:

242
$$\lambda = \frac{\omega(X_i)}{\omega(Y_i)} = \frac{POC_{true} \times LOD_{POC} \times \alpha_{POC}^2}{EC_{true} \times LOD_{EC} \times \alpha_{EC}^2} \quad (22)$$

243 Besides the $\gamma_{Unc-nonlinear}$ discussed above, a second type measurement uncertainty
 244 parameterized by a constant proportional factor, $\gamma_{Unc-linear}$, is very common in
 245 atmospheric applications:

246
$$-\gamma_{POCunc} \times POC_{true} \leq \varepsilon_{POC} \leq +\gamma_{POCunc} \times POC_{true} \quad (23)$$

247
$$-\gamma_{ECunc} \times EC_{true} \leq \varepsilon_{EC} \leq +\gamma_{ECunc} \times EC_{true} \quad (24)$$

248 where γ_{POCunc} and γ_{ECunc} are the relative measurement uncertainties, e.g., for relative
 249 measurement uncertainty of 10%, $\gamma_{Unc}=0.1$. As a result, the measurement error is
 250 linearly proportional to the concentration. An example comparison of $\gamma_{Unc-nonlinear}$
 251 and $\gamma_{Unc-linear}$ is shown in Figure 1b. For $\gamma_{Unc-linear}$, the weights become:

252
$$\omega(X_i) = \frac{1}{\sigma_{X_i}^2} = \frac{3}{(\gamma_{ECunc} \times EC_{true})^2} \quad (25)$$

253
$$\omega(Y_i) = \frac{1}{\sigma_{Y_i}^2} = \frac{3}{(\gamma_{POCunc} \times POC_{true})^2} \quad (26)$$

254 and λ for Deming regression can be determined:

255
$$\lambda = \frac{\omega(X_i)}{\omega(Y_i)} = \frac{(\gamma_{POCunc} \times POC_{true})^2}{(\gamma_{ECunc} \times EC_{true})^2} \quad (27)$$

256 **3.1.2 XY data generation by Mersenne Twister (MT) generator following** 257 **a specific distribution**

258 The Mersenne twister (MT) is a pseudorandom number generator (PRNG) developed
 259 by Matsumoto and Nishimura (1998). MT has been widely adopted by mainstream
 260 numerical analysis software (e.g., Matlab, SPSS, SAS and Igor Pro) as well as popular
 261 programming languages (e.g., R, Python, IDL, C++ and PHP). Data generation using MT
 262 provides a few advantages: (1) Frequency distribution can be easily assigned during the
 263 data generation process, allowing straightforward simulation of the frequency
 264 distribution characteristics (e.g., Gaussian or Log-normal) observed in ambient
 265 measurements; (2) The inputs for data generation are simply the mean and standard
 266 deviation of the data series and can be changed easily by the user; (3) The correlation
 267 (R^2) between X and Y can be manipulated easily during the data generation to satisfy

268 various purposes; (4) Unlike the sine function described by Chu (2005) that has a
269 sample size limitation of 120, the sample size in MT data generation is highly flexible.

270 In this section, we will use POC as Y and EC as X as an example to explain the data
271 generation. Procedure of applying MT to simulate ambient POC and EC data can be
272 found in our previous study (Wu and Yu, 2016). Details of the data generation steps
273 are shown in Figure 2 and described below. The first step is generation of EC_{true} by MT.
274 In our previous study, it was found that ambient POC and EC data follow a lognormal
275 distribution in various locations of the Pearl River Delta (PRD) region. Therefore,
276 lognormal distributions are adopted during EC_{true} generation. A range of average
277 concentration and relative standard deviation (RSD) from ambient samples are
278 considered in formulating the lognormal distribution. The second step is to generate
279 POC_{comb} . As shown in Figure 2, POC_{comb} is generated by multiplying EC_{true} with
280 $(OC/EC)_{pri}$. Instead of having a Gaussian distribution, $(OC/EC)_{pri}$ in this study is a
281 single value, which favors direct comparison between the true value of $(OC/EC)_{pri}$ and
282 $(OC/EC)_{pri}$ estimated from the regression slope. The third step is generation of POC_{true}
283 by adding $POC_{non-comb}$ onto POC_{comb} . Instead of having a distribution, $POC_{non-comb}$ in
284 this study is a single value, which favors direct comparison between the true value of
285 $POC_{non-comb}$ and $POC_{non-comb}$ estimated from the regression intercept. The fourth step is
286 to compute ε_{POC} and ε_{EC} . As discussed in section 3.1, two types of measurement errors
287 are considered for ε_{POC} and ε_{EC} calculation: $\gamma_{Unc-nonlinear}$ and $\gamma_{Unc-linear}$. In the
288 last step, $POC_{measured}$ and $EC_{measured}$ are calculated following Eq. (12), i.e., applying
289 measurement errors on POC_{true} and EC_{true} . Then $POC_{measured}$ and $EC_{measured}$ can be used
290 as Y and X respectively to test the performance of various regression techniques. An
291 Igor Pro based program with graphical user interface (GUI) is developed to facilitate
292 the MT data generation for OC and EC. A brief introduction is given in the
293 Supplemental Information.

294 **3.1.3 XY data generation by the sine function of Chu (2005)**

295 Beside MT, the inclusion of the sine function data generation schemes in this study
296 mainly serves two purposes. First, the sine function scheme had been adopted by two
297 previous studies (Chu, 2005; Saylor et al., 2006), the inclusion of this scheme can help
298 to verify whether the codes in Igor for various regression approaches can yield the same

299 results from the two previous studies. Second, crosscheck between results from sine
300 function and MT can provides circumstantial evidence that the MT scheme works as
301 expected.

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

$$306 \quad POC_{comb} = 14 + 12(\sin(\frac{x}{\tau}) + \sin(x - \phi)) \quad (28)$$

$$307 \quad EC_{true} = 3.5 + 3(\sin(\frac{x}{\tau}) + \sin(x - \phi)) \quad (29)$$

308 Here x is the elapsed hour (x=1,2,3.....n; n≤120), τ is used to adjust the width of each
309 peak, and φ is used to adjust the phase of the sine wave. The constants 14 and 3.5 are
310 used to lift the sine wave to the positive range of the Y axis. An example of data
311 generation by the sine functions of Chu (2005) is shown in Figure 3. Dividing Eq. (28)
312 by Eq. (29) yields a value of 4. In this way the exact relation between POC and EC is
313 defined clearly as (OC/EC)_{pri} = 4. (2) With POC_{comb} and EC_{true} generated, the second
314 step is to add POC_{non-comb} to POC_{comb} to compute POC_{true}. As for POC_{non-comb}, a single
315 value is assigned and added to all POC following Eq. (10). Then the goodness of the
316 regression intercept can be evaluated by comparing the regressed intercept with preset
317 POC_{non-comb}. (3) The third step is to compute ε_{POC} and ε_{EC}, considering both
318 γ_{Unc-nonlinear} and γ_{Unc-linear}. (4) The last step is to apply measurement errors on
319 POC_{true} and EC_{true} following Eq. (12). Then POC_{measured} and EC_{measured} can be used as
320 Y and X respectively to evaluate the performance of various regression techniques.

321 **3.2 Ambient measurement of σ_{abs} and EC**

322 Sampling was conducted from Feb 2012 to Jan 2013 at the suburban Nancun (NC) site
323 (23° 0'11.82"N, 113°21'18.04"E), which is situated on the top of the highest peak (141
324 m ASL) in the Panyu district of Guangzhou. This site is located at the geographic center
325 of Pearl River Delta region (PRD), making it a good location for representing the
326 average atmospheric mixing characteristics of city clusters in the PRD region. Light
327 absorption measurements were performed by a 7-λ Aethalometer (AE-31, Magee

328 Scientific Company, Berkeley, CA, USA). EC mass concentrations were measured by
329 a real time ECOC analyzer (Model RT-4, Sunset Laboratory Inc., Tigard, Oregon,
330 USA). Both instruments utilized inlets with a 2.5 μm particle diameter cutoff. The algorithm
331 of Weingartner et al. (2003) was adopted to correct the sampling artifacts (aerosol
332 loading, filter matrix and scattering effect) (Coen et al., 2010) root in Aethalometer
333 measurement. A customized computer program with graphical user interface,
334 Aethalometer data processor (Wu et al., 2017), was developed to perform the data
335 correction and detailed descriptions can be found in
336 <https://sites.google.com/site/wuchengust>. More details of the measurements can be
337 found in Wu et al. (2017).

338 **4 Comparison study using synthetic data**

339 In the following comparisons, six regression approaches are compared using two data
340 generation schemes (Chu sine function and MT) separately, as illustrated in Figure 4.
341 Each data generation scheme considers both $\gamma_{Unc-nonlinear}$ and $\gamma_{Unc-linear}$ in
342 measurement error parameterization. In total, 18 cases are tested with different
343 combination of data generation schemes, measurement error parameterization schemes,
344 true slope and intercept settings. For each case, six regression approaches are tested,
345 including OLS, DR ($\lambda = 1$), DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$), ODR, WODR and YR. In commercial
346 software (e.g., Origin, SigmaPlot, GraphPad Prism, etc), λ in DR is set to 1 by default
347 if not specified. As indicated by Saylor et al. (2006), the bias observed in the study of
348 Chu (2005) is likely due to $\lambda = 1$ in DR. The purpose of including DR ($\lambda = 1$) in this
349 study is to examine the potential bias using the default input in many software products.
350 The six regression approaches are considered to examine the sensitivity of regression
351 results to various parameters used in data generation. For each case, 5000 runs are
352 performed to obtain statistically significant results, as recommended by Saylor et al.
353 (2006). The mean slope and intercept from 5000 runs is compared with the true value
354 assigned during data generation. If the difference is $<5\%$, the result is considered
355 unbiased.

356 4.1 Comparison results using the data set of Chu (2005)

357 In this section, the scheme of Chu (2005) is adopted for data generation to obtain a
358 benchmark of six regression approaches. With different setup of slope, intercept and
359 γ_{Unc} , 6 cases (Case 1 ~ 6) are studied and the results are discussed below.

360 4.1.1 Results with $\gamma_{Unc-nonlinear}$

361 A comparison of the regression techniques results with $\gamma_{Unc-nonlinear}$ (following Eqs.
362 (18) & (19)) are summarized in Table 2. LOD_{POC} , LOD_{EC} , α_{POC} and α_{EC} are all set to
363 1 to reproduce the data studied by Chu (2005) and Saylor et al. (2006). Two sets of true
364 slope and intercept are considered (Case 1: Slope=4, Intercept=0; Case 2: Slope=4,
365 Intercept=3) to examine if any results are sensitive to the non-zero intercept. The R^2
366 (POC, EC) from 5000 runs for both case 1 and 2 are 0.67 ± 0.03 .

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

380 For case 3, $LOD_{POC}=0.5$, $LOD_{EC}=0.5$, $\alpha_{POC}=0.5$, $\alpha_{EC}=0.5$ are adopted (Table 2),
381 leading to an offset to the left of $\gamma_{Unc-nonlinear}$ (blue curve) compared to Case 1 and 2
382 (black curve) in Figure 1. As a result, for the same concentration of EC and OC in Case
383 3, the $\gamma_{Unc-nonlinear}$ is smaller than in Case 1 and Case 2 as indicated by higher the R^2
384 (0.95 ± 0.01 for Case 3, Table 2). With a smaller measurement uncertainty, the degree
385 of bias in Case 3 is smaller than Case 1. For example, OLS slope is less biased in Case

386 3 (3.83 ± 0.08) compare to Case 1 (2.94 ± 0.14). Similarly, the slope (4.03 ± 0.09) and
387 intercept (-0.18 ± 0.44) of DR ($\lambda=1$) exhibit a much smaller bias with a smaller
388 measurement uncertainty, implying that the degree of bias by improperly weighting in
389 DR, WODR and YR is associated with the degree of measurement uncertainty. A higher
390 measurement uncertainty results in larger bias in slope and intercept.

391 An uneven LOD_{POC} and LOD_{EC} is tested in Case 4 with $LOD_{POC}=1$, $LOD_{EC}=0.5$,
392 $\alpha_{POC}=0.5$, $\alpha_{EC}=0.5$, which yield a $R^2(\text{POC, EC})$ of 0.78 ± 0.02 . The results are similar
393 to Case 1. For DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) unbiased slope and intercept are obtained. For WODR
394 and YR, unbiased slopes are reported with a small bias in the intercepts. Large bias
395 values are observed in both the slopes and intercepts in Case 4 using OLS, DR ($\lambda = 1$)
396 and ODR.

397 **4.1.2 Results with $\gamma_{Unc-linear}$**

398 Cases 5 and 6 represent the results from using $\gamma_{Unc-linear}$ and are shown in Table 2.
399 γ_{Unc} is set to be 30% to achieve a $R^2(\text{POC, EC})$ of 0.7, a value close to the R^2 in studies
400 of Chu (2005) and Saylor et al. (2006). In Case 5 (slope=4, intercept=0), unbiased
401 slopes and intercepts are determined by DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$), WODR and YR. OLS
402 underestimates the slope (3.32 ± 0.20) and overestimates intercept (3.77 ± 0.90), while
403 DR ($\lambda = 1$) and ODR overestimate the slopes (4.75 ± 0.30) and underestimates the
404 intercepts (-4.14 ± 1.36). In Case 6 (slope=4, intercept=3), results similar to Case 5 are
405 obtained. It is worth noting that although the mean intercept (3.05 ± 1.22) of DR ($\lambda =$
406 $\frac{\omega(X_i)}{\omega(Y_i)}$), is closest to the true value (intercept=3), the deviations are much larger than for
407 WODR (2.72 ± 0.74).

408 **4.2 Comparison results using data generated by MT**

409 In this section, MT is adopted for data generation to obtain a benchmark of six
410 regression approaches. Both $\gamma_{Unc-nonlinear}$ and $\gamma_{Unc-linear}$ are considered. With
411 different configuration of slope, intercept and γ_{Unc} , 12 cases (Case 7 ~ Case 18) are
412 studied and the results are discussed below.

413 **4.2.1 $\gamma_{Unc-nonlinear}$ results**

414 Cases 7 and 8 use data generated by MT and $\gamma_{Unc-nonlinear}$ with results shown in Table
415 2. In Case 7 (slope=4, intercept=0, $LOD_{POC}=1$, $LOD_{EC}=1$, $\alpha_{POC}=1$, $\alpha_{EC}=1$), unbiased
416 slope (4.00 ± 0.03) and intercept (0.00 ± 0.17) is estimated by DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$). WODR
417 and YR yield unbiased slopes (3.96 ± 0.03) but overestimate the intercepts (1.21 ± 0.13).
418 DR ($\lambda = 1$) and ODR report slightly biased slopes (4.17 ± 0.04) with biased intercepts
419 (-0.94 ± 0.18). OLS underestimates the slope (3.22 ± 0.03) and overestimates the
420 intercept (4.30 ± 0.14). In Case 8 (slope=4, intercept=3, $LOD_{POC}=1$, $LOD_{EC}=1$, $\alpha_{POC}=1$,
421 $\alpha_{EC}=1$), DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) provides unbiased slope (4.00 ± 0.03) and intercept (3.00 ± 0.18)
422 estimations. WODR and YR report unbiased slopes (3.97 ± 0.03) and overestimate
423 intercepts (4.11 ± 0.13). OLS, DR ($\lambda = 1$) and ODR report biased slopes and intercepts.
424 To test the overestimation/underestimation dependency on the true slope, Case 9
425 (slope=0.5, intercept=0, $LOD_{POC}=1$, $LOD_{EC}=1$, $\alpha_{POC}=1$, $\alpha_{EC}=1$) and case 10
426 (slope=0.5, intercept=3, $LOD_{POC}=1$, $LOD_{EC}=1$, $\alpha_{POC}=1$, $\alpha_{EC}=1$) are conducted and the
427 results are shown in Table 2. Unlike the overestimation observed in Case 1~Case 8, DR
428 ($\lambda = 1$) and ODR underestimate the slopes (0.46 ± 0.01) in Case 9. In case 10, DR ($\lambda =$
429 1), DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) and ODR report unbiased slopes and intercepts. Case 11 and case
430 12 test the bias when the true slope is 1 as shown in Table 2. In Case 11 (intercept=0),
431 all regression approaches except OLS can provide unbiased results. In Case 12, all
432 regression approaches report unbiased slopes except OLS, but DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) is the
433 only regression approach that report unbiased intercept.
434 These results imply that if the true slope is less than 1, the improper weighting ($\lambda = 1$)
435 in Deming regression and ODR without weighting tends to underestimate slope. If the
436 true slope is 1, these two estimators can provide unbiased results. If the true slope is
437 larger than 1, the improper weighting ($\lambda = 1$) in Deming regression and ODR without
438 weighting tends to overestimate slope.

439 **4.2.2 $\gamma_{Unc-linear}$ results**

440 Cases 13 and 14 (Table 2) represent the results from using $\gamma_{Unc-linear}$ (30%) and data
441 generated from MT. For case 13 (slope=4, intercept=0), DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$), WODR and
442 YR provide the best estimation of slopes and intercepts. DR ($\lambda = 1$) and ODR
443 overestimate slopes (4.53 ± 0.05) and underestimate intercepts (-2.94 ± 0.24). For case
444 14 (slope=4, intercept=3), DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$), WODR and YR provide an unbiased
445 estimation of slopes. But DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) is the only regression approaches reports
446 unbiased intercept (3.08 ± 0.23). Cases 15 and 16 are tested to investigate whether the
447 results are different if the true slope is smaller than 1. As shown in Table 2, the results
448 are similar to case 13&14 that DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) can provide unbiased slope and intercept
449 while WODR and YR can provide unbiased slopes but biased intercepts. Cases 17 and
450 18 are tested to see if the results are the same for a special case when the true slope is
451 1. As shown in Table 2, the results are similar to case 13&14, implying that these results
452 are not sensitive to the special case when the true slope is 1.

453 **4.3 The importance of appropriate λ input for Deming regression**

454 As discussed above, inappropriate λ assignment in the Deming regression (e.g., $\lambda=1$ by
455 default for many commercial software) leads to biased slope and intercept. Beside $\lambda=1$,
456 inappropriate λ input due to improper handling of measurement uncertainty can also
457 result in bias for Deming regression. An example is shown in Figure S2. Data is
458 generated by MT with following parameters: slope=4, intercept=0, and $\gamma_{Unc-linear}$
459 (30%). Figure S2 a&b demonstrates that when an appropriate λ is provided (following
460 $\gamma_{Unc-linear}$, $\lambda = \frac{POC^2}{EC^2}$), unbiased slopes and intercepts are obtained. If an improper λ is
461 used due to a mismatched measurement uncertainty assumption ($\gamma_{Unc-nonlinear}$, $\lambda =$
462 $\frac{POC}{EC}$), the slopes are overestimated (Figure S2c, 4.37 ± 0.05) and intercepts are
463 underestimated (Figure S2d, -2.01 ± 0.24). This result emphasizes the importance of
464 determining the correct form of measurement uncertainty in ambient samples, since λ
465 is a crucial parameter in Deming regression.

466 In the λ calculation, different representations for POC and EC, including mean, median
467 and mode, are tested as shown in Figure S3. The results show that when X and Y have
468 a similar distribution (e.g., both are log-normal), any of mean, median or mode can be
469 used for the λ calculation.

470 **4.4 Caveats of regressions with unknown X and Y uncertainties**

471 When applying linear regression on real world data, it happens that a priori error in
472 one of the variables is unknown, or the measurement error described cannot be trusted.
473 In other words, that would be certain degree of discrepancy between the measurement
474 error used for linear regression and measurement error embed in the data. It is common
475 that measurement error cannot be determined due to the lack of duplicated or
476 collocated measurements and an arbitrarily assumed uncertainty is used. For example,
477 Flanagan et al. (2006) found that the whole-system uncertainty retrieved by data from
478 collocated sampler is different from the arbitrarily assumed 5% uncertainty, which is
479 previously used by the Speciation Trends Network (STN). In addition, the degree of
480 discrepancy between the actual uncertainty by collocated samples and arbitrarily
481 assumed uncertainty also varied by different chemical species. To investigate the
482 impact of such cases on different regression approaches, two tests are conducted. In
483 Test A, the actual measurement error for X is fixed at 30% while γ_{Unc} for Y varied
484 from 1% to 50%. The assumed measurement error for regression is 10% for both X
485 and Y. Results of Test A are shown in Figure 6 a&b. For OLS, the slopes are
486 underestimated (-14 ~ -12%) and intercepts are overestimated (90 ~ 103%). The biases
487 in OLS slope and intercept are independent of variations in γ_{Unc_Y} . ODR and DR ($\lambda =$
488 1) yield similar results with overestimated slopes (0 ~ 44%) and underestimated
489 intercepts (-330 ~ 0%). The degree of bias in slopes and intercepts depends on γ_{Unc_Y} .
490 WODR, DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) and YR performed much better than other regression
491 approaches in Test A, with a smaller bias in both slopes (-8 ~ 12%) and intercepts -98
492 ~ 55%).

493 The results of Test B are shown in Figure 6 c&d. which has a fixed γ_{Unc_Y} of 30% and
494 γ_{Unc_X} varied between 1 ~ 50%. The assumed measurement error for regression is 10%
495 for both X and Y. OLS underestimates slopes (-29 ~ -0.2%) and overestimates

496 intercepts (2 ~ 209%) in Test B. In contrast to Test A which slope and intercept biases
497 are independent of variations in γ_{Unc_Y} , the OLS slope and intercept biases in Test B
498 exhibit dependency on γ_{Unc_X} . The reason behind is because OLS only considers
499 errors in Y, while X is assumed to be error free. ODR and DR ($\lambda = 1$) yield similar
500 results with overestimated slopes (11 ~ 18%) and underestimated intercepts (-144 ~ -
501 87%). The degree of biases in slopes and intercepts is relatively independent to the
502 γ_{Unc_X} . WODR, DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) and YR performed much better than other regression
503 approaches in Test B, with a smaller bias in both slopes (-14 ~ 8%) and intercepts (-
504 59 ~ 106%).

505 The results from these two tests suggest that, in case of one of the measurement error
506 described cannot be trusted or a priori error in one of the variables is unknown, WODR,
507 DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$) and YR should be used instead of ODR, DR ($\lambda = 1$) and OLS. This
508 conclusion also agrees with section 4.1 and 4.2. The results also suggest that, in general,
509 the magnitude of bias in slope estimation by these regression approaches are smaller
510 than those for intercept. In other words, slope is a more reliable quantity compare to
511 intercept when extracting quantitative information from linear regressions.

512 **5 Regression applications to ambient data**

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

520 As shown in Figure 7, Y axis is light absorption at 520 nm (σ_{abs520}) and the X axis is
521 EC mass concentration. The regressed slopes represent the mass absorption efficiency
522 (MAE) of EC at 520 nm, ranging from 13.66 to 15.94 m^2g^{-1} by the six regression
523 approaches. OLS yields the lowest slope (13.66 as shown in Figure 7a) among all six
524 regression approaches, consistent with the results using synthetic data. This implies that
525 OLS tends to underestimate regression slope when mean Y to X ratio is larger than 1.

526 DR ($\lambda = 1$) and ODR report the same slope (14.88) and intercept (5.54), this
527 equivalency is also observed for the synthetic data. Similarly, WODR and YR yield
528 identical slope (14.88) and intercept (5.54), in line with the synthetic data results. The
529 regressed slope by DR ($\lambda = 1$) is higher than DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$), and this relationship
530 agrees well with the synthetic data results.

531 Regression comparison is also performed on hourly OC and EC data. Regression on
532 OC/EC percentile subset is a widely used empirical approach for primary OC/EC ratio
533 determination. Figure S4 shows the regression slopes as a function of OC/EC percentile.
534 OC/EC percentile ranges from 0.5% to 100%, with an interval of 0.5%. As the
535 percentile increases, SOC contribution in OC increases as well, resulting decreased R^2
536 between OC and EC. The deviations between six regression approaches exhibit a
537 dependency on R^2 . When percentile is relatively small (e.g., <10%), the differences
538 between the six regression approaches are also small due to the high R^2 (0.98). The
539 deviations between the six regression approaches become more pronounced as R^2
540 decreases (e.g., <0.9). The deviations are expected to be even larger when R^2 is less
541 than 0.8. These results emphasize the importance of applying error-in-variables
542 regression, since ambient XY data more likely has a R^2 less than 0.9 in most cases.

543 As discussed in this section, the ambient data confirm the results obtained in comparing
544 methods with the synthetic data. The advantage of using the synthetic data for
545 regression approaches evaluation is that the ideal slope and intercept are known values
546 during the data generation, so the bias of each regression approach can be quantified.

547 **6 Recommendations and conclusions**

548 This study aims to provide a benchmark of commonly used linear regression algorithms
549 using a new data generation scheme (MT). Six regression approaches are tested,
550 including OLS, DR ($\lambda = 1$), DR ($\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$), ODR, WODR and YR. The results show
551 that OLS fails to estimate the correct slope and intercept when both X and Y have
552 measurement errors. This result is consistent with previous studies. For ambient data
553 with R^2 less than 0.9, error-in-variables regression is needed to minimize the biases in
554 slope and intercept. If measurement uncertainties in X and Y are determined during the
555 measurement, measurement uncertainties should be used for regression. With

556 appropriate weighting, DR, WODR and YR can provide the best results among all
557 tested regression techniques. Sensitivity tests also reveal the importance of the
558 weighting parameter λ in DR. An improper λ could lead to biased slope and intercept.
559 Since the λ estimation depends on the form of the measurement errors, it is important
560 to determine the measurement errors during the experimentation stage rather than
561 making assumptions. If measurement errors are not available from the measurement
562 and assumptions are made on measurement errors, DR, WODR and YR are still the
563 best option that can provide the least bias in slope and intercept among all tested
564 regression techniques. For these reasons, DR, WODR and YR are recommended for
565 atmospheric studies when both X and Y data have measurement errors.

566 Application of error-in-variables regression is often overlooked in atmospheric studies,
567 partly due to the lack of a specified tool for the regression implementation. To facilitate
568 the implementation of error-in-variables regression (including DR, WODR and YR), a
569 computer program (Scatter plot) with graphical user interface (GUI) in Igor Pro
570 (WaveMetrics, Inc. Lake Oswego, OR, USA) is developed (Figure 8). It packed with
571 many useful features for data analysis and plotting, including batch plotting, data
572 masking via GUI, color coding in Z axis, data filtering and grouping by numerical
573 values and strings. The Scatter plot program and user manual are available from
574 <https://sites.google.com/site/wuchengust> and <https://doi.org/10.5281/zenodo.832417>.

575

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

577 Ordinary Least Square (OLS) calculation steps.

578 First calculate average of observed X_i and Y_i .

579
$$\bar{X} = \frac{\sum_{i=1}^N X_i}{N} \quad (A1)$$

580
$$\bar{Y} = \frac{\sum_{i=1}^N Y_i}{N} \quad (A2)$$

581 Then calculate S_{xx} and S_{yy} .

582
$$S_{xx} = \sum_{i=1}^N (X_i - \bar{X})^2 \quad (A3)$$

583
$$S_{yy} = \sum_{i=1}^N (Y_i - \bar{Y})^2 \quad (A4)$$

584 OLS slope and intercept can be obtained from,

585
$$k = \frac{S_{yy}}{S_{xx}} \quad (A6)$$

586
$$b = \bar{Y} - k\bar{X} \quad (A7)$$

587

588 Deming regression (DR) calculation steps (York, 1966).

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

590
$$S_{xy} = \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y}) \quad (A8)$$

591 DR slope and intercept can be obtained from,

592
$$k = \frac{S_{yy} - \lambda S_{xx} + \sqrt{(S_{yy} - \lambda S_{xx})^2 + 4\lambda S_{xy}^2}}{2S_{xy}} \quad (A9)$$

593
$$b = \bar{Y} - k\bar{X} \quad (A10)$$

594

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

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

597
$$W_i = \frac{\omega(X_i)\omega(Y_i)}{\omega(X_i) + k^2\omega(Y_i) - 2kr_i\sqrt{\omega(X_i)\omega(Y_i)}} \quad (A11)$$

598
$$U_i = X_i - \bar{X} = X_i - \frac{\sum_{i=1}^N W_i X_i}{\sum_{i=1}^N W_i} \quad (\text{A12})$$

599
$$V_i = Y_i - \bar{Y} = Y_i - \frac{\sum_{i=1}^N W_i Y_i}{\sum_{i=1}^N W_i} \quad (\text{A13})$$

600 Then calculate β_i .

601
$$\beta_i = W_i \left[\frac{U_i}{\omega(Y_i)} + \frac{kV_i}{\omega(X_i)} - [kU_i + V_i] \frac{r_i}{\sqrt{\omega(X_i)\omega(Y_i)}} \right] \quad (\text{A14})$$

602 Slope and intercept can be obtained from,

603
$$k = \frac{\sum_{i=1}^n W_i \beta_i V_i}{\sum_{i=1}^n W_i \beta_i U_i} \quad (\text{A15})$$

604
$$b = \bar{Y} - k\bar{X} \quad (\text{A16})$$

605 Since W_i and β_i are functions of k , k must be solved iteratively by repeating A11 to
 606 A15. If the difference between the k obtained from A15 and the k used in A11 satisfies
 607 the predefined tolerance ($\frac{k_{i+1}-k_i}{k_i} < e^{-15}$), the calculation is considered as converged. The
 608 calculation is straightforward and usually converged in 10 iterations. For example, the
 609 iteration count on the data set of Chu (2005) is around 6.

610 **Acknowledgements**

611 This work is supported by the National Natural Science Foundation of China (Grant
 612 No. 41605002 and 21607056), NSFC of Guangdong Province (Grant No.
 613 2015A030313339), Guangdong Province Public Interest Research and Capacity
 614 Building Special Fund (Grant No. 2014B020216005). The author would like to thank
 615 Dr. Bin Yu Kuang at HKUST for discussion on mathematics and Dr. Stephen M
 616 Griffith at HKUST for valuable comments.

617

618

619 **References**

- 620 Ayers, G. P.: Comment on regression analysis of air quality data, *Atmos. Environ.*, 35,
621 2423-2425, doi: 10.1016/S1352-2310(00)00527-6, 2001.
- 622 Bauer, J. J., Yu, X.-Y., Cary, R., Laulainen, N., and Berkowitz, C.: Characterization of
623 the sunset semi-continuous carbon aerosol analyzer, *J. Air Waste Manage. Assoc.*, 59,
624 826-833, doi: 10.3155/1047-3289.59.7.826, 2009.
- 625 Boggs, P. T., Donaldson, J. R., and Schnabel, R. B.: Algorithm 676: ODRPACK:
626 software for weighted orthogonal distance regression, *ACM Trans. Math. Softw.*, 15,
627 348-364, doi: 10.1145/76909.76913, 1989.
- 628 Brauers, T. and Finlayson-Pitts, B. J.: Analysis of relative rate measurements, *Int. J.*
629 *Chem. Kinet.*, 29, 665-672, doi: 10.1002/(SICI)1097-4601(1997)29:9<665::AID-
630 KIN3>3.0.CO;2-S, 1997.
- 631 Cantrell, C. A.: Technical Note: Review of methods for linear least-squares fitting of
632 data and application to atmospheric chemistry problems, *Atmos. Chem. Phys.*, 8, 5477-
633 5487, doi: 10.5194/acp-8-5477-2008, 2008.
- 634 Carroll, R. J. and Ruppert, D.: The use and misuse of orthogonal regression in linear
635 errors-in-variables models, *Am. Stat.*, 50, 1-6, doi: 10.1080/00031305.1996.10473533,
636 1996.
- 637 Cess, R. D., Zhang, M. H., Minnis, P., Corsetti, L., Dutton, E. G., Forgan, B. W.,
638 Garber, D. P., Gates, W. L., Hack, J. J., Harrison, E. F., Jing, X., Kiehi, J. T., Long, C.
639 N., Morcrette, J.-J., Potter, G. L., Ramanathan, V., Subasilar, B., Whitlock, C. H.,
640 Young, D. F., and Zhou, Y.: Absorption of solar radiation by clouds: Observations
641 versus models, *Science*, 267, 496-499, doi: 10.1126/science.267.5197.496, 1995.
- 642 Chen, L. W. A., Doddridge, B. G., Dickerson, R. R., Chow, J. C., Mueller, P. K., Quinn,
643 J., and Butler, W. A.: Seasonal variations in elemental carbon aerosol, carbon monoxide
644 and sulfur dioxide: Implications for sources, *Geophys. Res. Lett.*, 28, 1711-1714, doi:
645 10.1029/2000GL012354, 2001.
- 646 Chu, S. H.: Stable estimate of primary OC/EC ratios in the EC tracer method, *Atmos.*
647 *Environ.*, 39, 1383-1392, doi: 10.1016/j.atmosenv.2004.11.038, 2005.
- 648 Coen, M. C., Weingartner, E., Apituley, A., Ceburnis, D., Fierz-Schmidhauser, R.,
649 Flentje, H., Henzing, J. S., Jennings, S. G., Moerman, M., Petzold, A., Schmid, O., and
650 Baltensperger, U.: Minimizing light absorption measurement artifacts of the
651 Aethalometer: evaluation of five correction algorithms, *Atmos. Meas. Tech.*, 3, 457-
652 474, doi: 10.5194/amt-3-457-2010, 2010.
- 653 Cornbleet, P. J. and Gochman, N.: Incorrect least-squares regression coefficients in
654 method-comparison analysis, *Clin. Chem.*, 25, 432-438, 1979.
- 655 Cox, M., Harris, P., and Siebert, B. R.-L.: Evaluation of Measurement Uncertainty
656 Based on the Propagation of Distributions Using Monte Carlo Simulation,
657 *Measurement Techniques*, 46, 824-833, doi: 10.1023/B:METE.0000008439.82231.ad,
658 2003.
- 659 Cross, E. S., Onasch, T. B., Ahern, A., Wrobel, W., Slowik, J. G., Olfert, J., Lack, D.
660 A., Massoli, P., Cappa, C. D., Schwarz, J. P., Spackman, J. R., Fahey, D. W., Sedlacek,

661 A., Trimborn, A., Jayne, J. T., Freedman, A., Williams, L. R., Ng, N. L., Mazzoleni,
662 C., Dubey, M., Brem, B., Kok, G., Subramanian, R., Freitag, S., Clarke, A., Thornhill,
663 D., Marr, L. C., Kolb, C. E., Worsnop, D. R., and Davidovits, P.: Soot particle studies—
664 instrument inter-comparison—project overview, *Aerosol. Sci. Technol.*, 44, 592-611,
665 doi: 10.1080/02786826.2010.482113, 2010.

666 Deming, W. E.: *Statistical Adjustment of Data*, Wiley, New York, 1943.

667 Duan, F., Liu, X., Yu, T., and Cachier, H.: Identification and estimate of biomass
668 burning contribution to the urban aerosol organic carbon concentrations in Beijing,
669 *Atmos. Environ.*, 38, 1275-1282, doi: 10.1016/j.atmosenv.2003.11.037, 2004.

670 Flanagan, J. B., Jayanty, R. K. M., Rickman, J. E. E., and Peterson, M. R.: PM2.5
671 Speciation Trends Network: Evaluation of Whole-System Uncertainties Using Data
672 from Sites with Collocated Samplers, *J. Air Waste Manage. Assoc.*, 56, 492-499, doi:
673 10.1080/10473289.2006.10464516, 2006.

674 Hansen, A. D. A.: *The Aethalometer Manual*, Berkeley, California, USA, Magee
675 Scientific, 2005.

676 Huang, X. H., Bian, Q., Ng, W. M., Louie, P. K., and Yu, J. Z.: Characterization of
677 PM2.5 major components and source investigation in suburban Hong Kong: A one
678 year monitoring study, *Aerosol. Air. Qual. Res.*, 14, 237-250, doi:
679 10.4209/aaqr.2013.01.0020, 2014.

680 Janhäll, S., Andreae, M. O., and Pöschl, U.: Biomass burning aerosol emissions from
681 vegetation fires: particle number and mass emission factors and size distributions,
682 *Atmos. Chem. Phys.*, 10, 1427-1439, doi: 10.5194/acp-10-1427-2010, 2010.

683 Lim, L. H., Harrison, R. M., and Harrad, S.: The contribution of traffic to atmospheric
684 concentrations of polycyclic aromatic hydrocarbons, *Environ. Sci. Technol.*, 33, 3538-
685 3542, doi: 10.1021/es990392d, 1999.

686 Linnet, K.: Necessary sample size for method comparison studies based on regression
687 analysis, *Clin. Chem.*, 45, 882-894, 1999.

688 Malm, W. C., Sisler, J. F., Huffman, D., Eldred, R. A., and Cahill, T. A.: Spatial and
689 seasonal trends in particle concentration and optical extinction in the United-States, *J.*
690 *Geophys. Res.*, 99, 1347-1370, doi: 10.1029/93JD02916, 1994.

691 Markovsky, I. and Van Huffel, S.: Overview of total least-squares methods, *Signal*
692 *Process.*, 87, 2283-2302, doi: 10.1016/j.sigpro.2007.04.004, 2007.

693 Matsumoto, M. and Nishimura, T.: Mersenne twister: a 623-dimensionally
694 equidistributed uniform pseudo-random number generator, *ACM Trans. Model.*
695 *Comput. Simul.*, 8, 3-30, doi: 10.1145/272991.272995, 1998.

696 Moosmüller, H., Arnott, W. P., Rogers, C. F., Chow, J. C., Frazier, C. A., Sherman, L.
697 E., and Dietrich, D. L.: Photoacoustic and filter measurements related to aerosol light
698 absorption during the Northern Front Range Air Quality Study (Colorado 1996/1997),
699 *J. Geophys. Res.*, 103, 28149-28157, doi: 10.1029/98jd02618, 1998.

700 Petäjä, T., Mauldin, I. R. L., Kosciuch, E., McGrath, J., Nieminen, T., Paasonen, P.,
701 Boy, M., Adamov, A., Kotiaho, T., and Kulmala, M.: Sulfuric acid and OH
702 concentrations in a boreal forest site, *Atmos. Chem. Phys.*, 9, 7435-7448, doi:
703 10.5194/acp-9-7435-2009, 2009.

704 Richter, A., Burrows, J. P., Nusz, H., Granier, C., and Niemeier, U.: Increase in
705 tropospheric nitrogen dioxide over China observed from space, *Nature*, 437, 129-132,
706 doi: 10.1038/nature04092, 2005.

707 Saylor, R. D., Edgerton, E. S., and Hartsell, B. E.: Linear regression techniques for use
708 in the EC tracer method of secondary organic aerosol estimation, *Atmos. Environ.*, 40,
709 7546-7556, doi: 10.1016/j.atmosenv.2006.07.018, 2006.

710 Thompson, M.: Variation of precision with concentration in an analytical system,
711 *Analyst*, 113, 1579-1587, doi: 10.1039/AN9881301579, 1988.

712 Turpin, B. J. and Huntzicker, J. J.: Identification of secondary organic aerosol episodes
713 and quantitation of primary and secondary organic aerosol concentrations during
714 SCAQS, *Atmos. Environ.*, 29, 3527-3544, doi: 10.1016/1352-2310(94)00276-Q, 1995.

715 von Bobruzki, K., Braban, C. F., Famulari, D., Jones, S. K., Blackall, T., Smith, T. E.
716 L., Blom, M., Coe, H., Gallagher, M., Ghalaieny, M., McGillen, M. R., Percival, C. J.,
717 Whitehead, J. D., Ellis, R., Murphy, J., Mohacsi, A., Pogany, A., Junninen, H.,
718 Rantanen, S., Sutton, M. A., and Nemitz, E.: Field inter-comparison of eleven
719 atmospheric ammonia measurement techniques, *Atmos. Meas. Tech.*, 3, 91-112, doi:
720 10.5194/amt-3-91-2010, 2010.

721 Wang, J. and Christopher, S. A.: Intercomparison between satellite-derived aerosol
722 optical thickness and PM_{2.5} mass: Implications for air quality studies, *Geophys. Res.
723 Lett.*, 30, 2095, doi: 10.1029/2003gl018174, 2003.

724 Watson, J. G.: Visibility: Science and regulation, *J. Air Waste Manage. Assoc.*, 52, 628-
725 713, doi: 10.1080/10473289.2002.10470813, 2002.

726 Weingartner, E., Saathoff, H., Schnaiter, M., Streit, N., Bitnar, B., and Baltensperger,
727 U.: Absorption of light by soot particles: determination of the absorption coefficient by
728 means of aethalometers, *J. Aerosol. Sci.*, 34, 1445-1463, doi: 10.1016/S0021-
729 8502(03)00359-8, 2003.

730 Wu, C. and Yu, J. Z.: Determination of primary combustion source organic carbon-to-
731 elemental carbon (OC/EC) ratio using ambient OC and EC measurements: secondary
732 OC-EC correlation minimization method, *Atmos. Chem. Phys.*, 16, 5453-5465, doi:
733 10.5194/acp-16-5453-2016, 2016.

734 Wu, C., Wu, D., and Yu, J. Z.: Quantifying black carbon light absorption enhancement
735 by a novel statistical approach, *Atmos. Chem. Phys. Discuss.*, 2017, 1-37, doi:
736 10.5194/acp-2017-582, 2017.

737 York, D.: Least-squares fitting of a straight line, *Can. J. Phys.*, 44, 1079-1086, doi:
738 10.1139/p66-090, 1966.

739 York, D., Evensen, N. M., Martinez, M. L., and Delgado, J. D. B.: Unified equations
740 for the slope, intercept, and standard errors of the best straight line, *Am. J. Phys.*, 72,
741 367-375, doi: 10.1119/1.1632486, 2004.

742 Zhou, Y., Huang, X. H. H., Griffith, S. M., Li, M., Li, L., Zhou, Z., Wu, C., Meng, J.,
743 Chan, C. K., Louie, P. K. K., and Yu, J. Z.: A field measurement based scaling approach
744 for quantification of major ions, organic carbon, and elemental carbon using a single
745 particle aerosol mass spectrometer, *Atmos. Environ.*, 143, 300-312, doi:
746 10.1016/j.atmosenv.2016.08.054, 2016.

747 Zieger, P., Weingartner, E., Henzing, J., Moerman, M., de Leeuw, G., Mikkilä, J., Ehn,
748 M., Petäjä, T., Clémer, K., van Roozendaal, M., Yilmaz, S., Frieß, U., Irie, H., Wagner,
749 T., Shaiganfar, R., Beirle, S., Apituley, A., Wilson, K., and Baltensperger, U.:
750 Comparison of ambient aerosol extinction coefficients obtained from in-situ, MAX-
751 DOAS and LIDAR measurements at Cabauw, *Atmos. Chem. Phys.*, 11, 2603-2624,
752 doi: 10.5194/acp-11-2603-2011, 2011.

753 Zwolak, J. W., Boggs, P. T., and Watson, L. T.: Algorithm 869: ODRPACK95: A
754 weighted orthogonal distance regression code with bound constraints, *ACM Trans.*
755 *Math. Softw.*, 33, 27, doi: 10.1145/1268776.1268782, 2007.

756

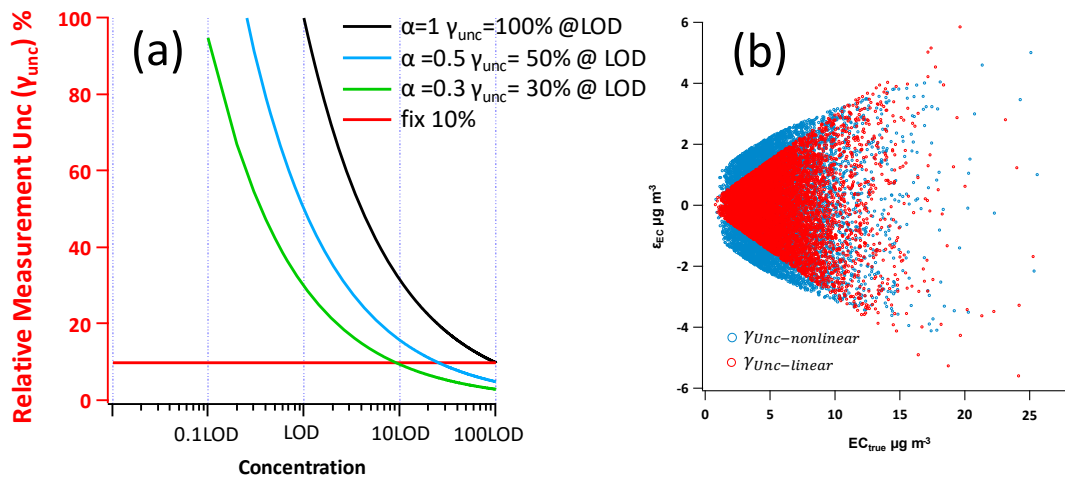
757 **Table 1.** Summary of abbreviations and symbols.

Abbreviation/symbol	Definition
α	a dimensionless adjustable factor to control the position of γ_{Unc} curve on the concentration axis
b	intercept in linear regression
β_i, U_i, V_i, W_i	intermediates in York regression calculations
γ_{Unc}	fractional measurement uncertainties relative to the true concentration (%)
DR	Deming regression
$\varepsilon_{EC}, \varepsilon_{POC}$	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 ($EC_{true} + \varepsilon_{EC}$)
λ	$\omega(X_i)$ to $\omega(Y_i)$ ratio in Deming regression
k	slope in linear 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 ($POC_{true} + \varepsilon_{POC}$)
$\sigma_{X_i}, \sigma_{Y_i}$	the standard deviation of the error in measurement of X_i and Y_i
r_i	correlation coefficient between errors in X_i and Y_i in YR
S	sum of squared residuals
SOC	secondary organic carbon
τ	parameter in the sine function of Chu (2005) that adjust the width of each peak
ϕ	parameter in the sine function of Chu (2005) that adjust the phase of the curve
WODR	weight orthogonal distance regression
\bar{X}, \bar{Y}	average of X_i and Y_i
YR	York regression
$\omega(X_i), \omega(Y_i)$	inverse of σ_{X_i} and σ_{Y_i} , used as weights in DR calculation.

758

Table 2. Summary of six regression approaches comparison with 5000 runs for 18 cases.

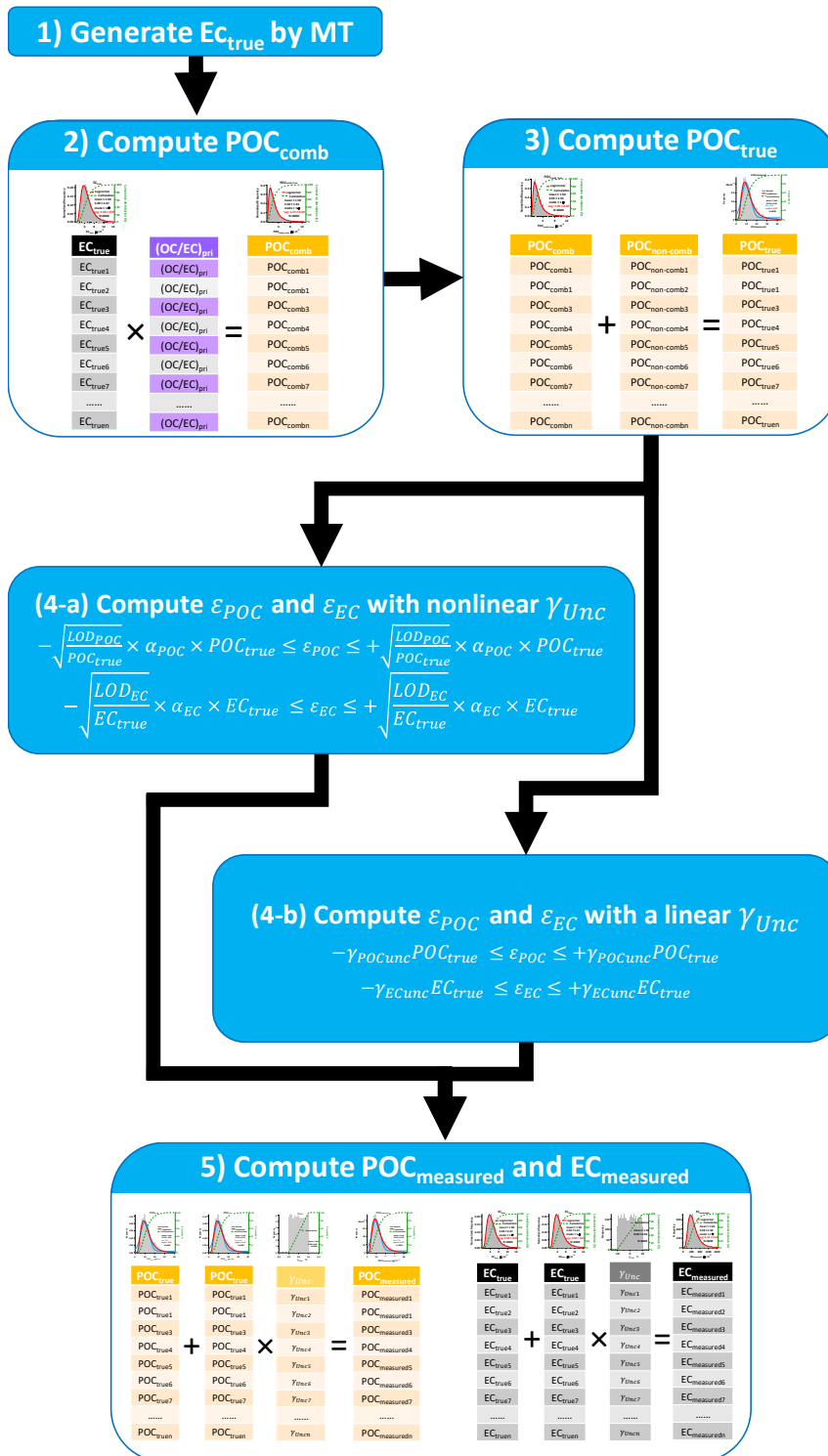
Case	Data generation					Results by different regression approaches											
	Data scheme	True Slope	True Intercept	R ² (X, Y)	Measurement error	OLS		DR $\lambda=1$		DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$		ODR		WODR		YR	
						Slope	Intercept	Slope	Intercept	Slope	Intercept	Slope	Intercept	Slope	Intercept	Slope	Intercept
1	Chu	4	0	0.67±0.03	$LOD_{POC}=1, LOD_{EC}=1$	2.94±0.14	5.84±0.78	4.27±0.27	-1.45±1.36	4.01±0.25	-0.04±1.28	4.27±0.27	-1.45±1.36	3.98±0.22	1.12±1.02	3.98±0.22	1.12±1.02
2		4	3	0.67±0.04	$a_{POC}=1, a_{EC}=1$	2.95±0.15	8.83±0.80	4.32±0.28	1.28±1.43	4.01±0.26	2.94±1.34	4.32±0.28	1.28±1.43	3.99±0.23	3.98±1.05	3.99±0.23	3.98±1.05
3		4	0	0.95±0.01	$LOD_{POC}=0.5, LOD_{EC}=0.5, \alpha_{POC}=0.5, \alpha_{EC}=0.5$	3.83±0.08	0.95±0.40	4.03±0.09	-0.18±0.44	4±0.09	0±0.44	4.03±0.09	-0.18±0.44	4±0.08	0.12±0.37	4±0.08	0.12±0.37
4		4	0	0.78±0.02	$LOD_{POC}=1, LOD_{EC}=0.5, \alpha_{POC}=1, \alpha_{EC}=1$	3.39±0.15	3.34±0.75	4.3±0.21	-1.66±1.06	4±0.19	-0.03±0.99	4.3±0.21	-1.66±1.06	4±0.17	0.33±0.81	4±0.17	0.33±0.81
5		4	0	0.69±0.04	$\gamma_{Unc}=30\%$	3.32±0.20	3.77±0.90	4.75±0.30	-4.14±1.36	4.01±0.25	-0.04±1.13	4.75±0.30	-4.14±1.36	4±0.18	-0.01±0.59	4±0.18	-0.01±0.59
6		4	3	0.66±0.04		3.31±0.22	6.79±1.02	4.95±0.31	-2.26±1.48	3.99±0.26	3.05±1.22	4.95±0.31	-2.26±1.48	4.01±0.20	2.72±0.74	4.01±0.20	2.72±0.74
7	MT	4	0	0.76±0.01	$LOD_{POC}=1, LOD_{EC}=1, a_{POC}=1, a_{EC}=1, \gamma_{Unc}=30\%$	3.22±0.03	4.3±0.14	4.17±0.04	-0.94±0.18	4±0.03	0±0.17	4.17±0.04	-0.94±0.18	3.96±0.03	1.21±0.13	3.96±0.03	1.21±0.13
8		4	3	0.75±0.01		3.22±0.03	7.29±0.14	4.2±0.04	1.88±0.18	4±0.03	3±0.18	4.2±0.04	1.88±0.18	3.97±0.03	4.11±0.13	3.97±0.03	4.11±0.13
9		0.5	0	0.76±0.01		0.43±0.00	0.36±0.02	0.46±0.01	0.23±0.03	0.5±0.01	0±0.03	0.46±0.01	0.23±0.03	0.5±0.00	0±0.01	0.5±0.00	0±0.01
10		0.5	3	0.56±0.01		0.43±0.01	3.36±0.03	0.5±0.01	3.02±0.04	0.49±0.01	3.05±0.04	0.5±0.01	3.02±0.04	0.51±0.01	2.73±0.03	0.51±0.01	2.73±0.03
11		1	0	0.76±0.01		0.87±0.01	0.72±0.05	1±0.01	0±0.06	1±0.01	0±0.06	1±0.01	0±0.06	1±0.01	0±0.02	1±0.01	0±0.02
12		1	3	0.66±0.01		0.87±0.01	3.72±0.05	1.09±0.01	2.52±0.07	0.99±0.01	3.07±0.06	1.09±0.01	2.52±0.07	1.01±0.01	2.71±0.04	1.01±0.01	2.7±0.04
13		4	0	0.76±0.01		3.48±0.04	2.87±0.18	4.53±0.05	-2.94±0.24	4±0.05	0±0.22	4.53±0.05	-2.94±0.24	4±0.03	0±0.09	4±0.03	0±0.09
14		4	3	0.73±0.01		3.48±0.04	5.87±0.19	4.67±0.05	-0.67±0.26	3.98±0.05	3.08±0.23	4.67±0.05	-0.67±0.26	4.02±0.03	2.68±0.11	4.02±0.03	2.68±0.11
15		0.5	0	0.54±0.01		0.4±0.01	0.55±0.03	0.45±0.01	0.26±0.03	0.5±0.01	0.01±0.03	0.45±0.01	0.26±0.03	0.52±0.01	-0.23±0.02	0.52±0.01	-0.23±0.02
16		0.5	3	0.40±0.01		0.4±0.01	3.54±0.04	0.5±0.01	2.98±0.04	0.5±0.01	3±0.04	0.5±0.01	2.98±0.04	0.52±0.01	2.65±0.04	0.52±0.01	2.65±0.04
17		1	0	0.65±0.01		0.8±0.01	1.07±0.04	1±0.01	0±0.05	1±0.01	0±0.05	1±0.01	0±0.05	1±0.01	0±0.04	1±0.01	0±0.04
18		1	3	0.59±0.01		0.8±0.01	4.07±0.05	1.07±0.01	2.62±0.07	1±0.01	3±0.06	1.07±0.01	2.62±0.07	1.02±0.01	2.84±0.05	1.02±0.01	2.84±0.05



763

764 **Figure 1.** (a) Example $\gamma_{Unc-nonlinear}$ curves by different α values (Eq. (17)). The X
 765 axis is concentration (normalized by LOD) in log scale and Y axis is γ_{Unc} . Black, blue
 766 and green line represent α equal to 1, 0.5 and 0.3 respectively, corresponding to the
 767 $\gamma_{Unc-nonlinear}$ at LOD level equals to 100%, 50% and 30% respectively. The red line
 768 represents $\gamma_{Unc-linear}$ of 10%. (b) Example of measurement uncertainty generation of
 769 $\gamma_{Unc-nonlinear}$ and $\gamma_{Unc-linear}$. The blue circles represent $\gamma_{Unc-nonlinear}$ following
 770 Eq. (17) ($LOD_{EC} = 1$, $a_{EC} = 1$). The red circles represent $\gamma_{Unc-linear}$ (30%).
 771

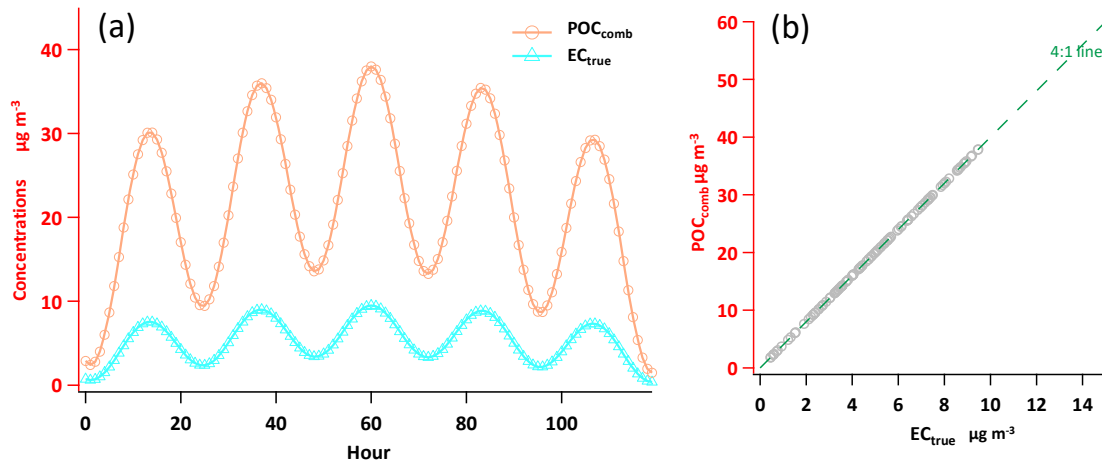
Data generations steps by MT



773

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

775



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

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

776

777 **Figure 3.** POC_{comb} and EC_{true} data generated by the sine functions of (Chu (2005)). (a)

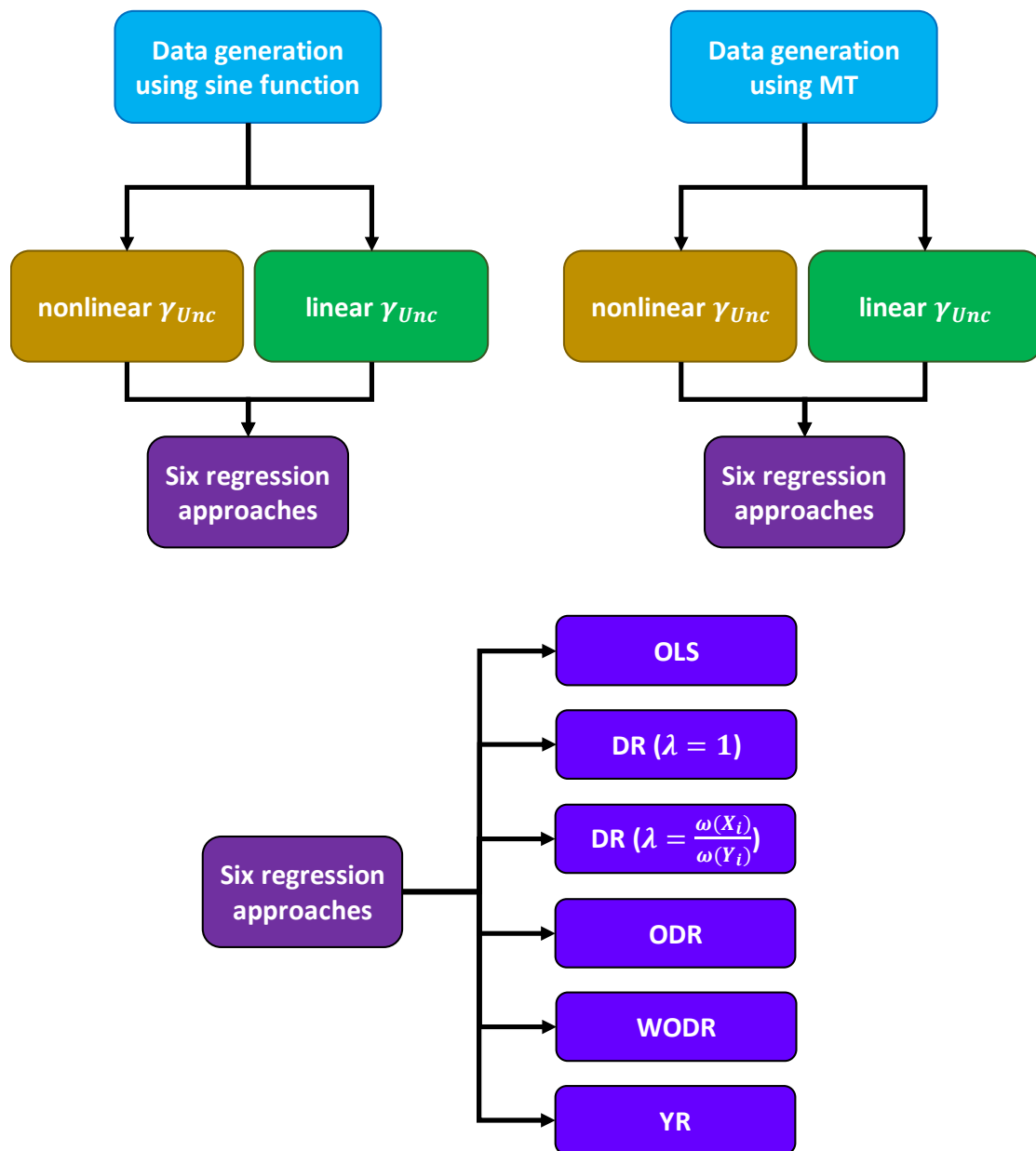
778 Time series of the 120 data points for POC_{comb} and EC_{true}. (b) Scatter plot of POC_{comb}

779 vs. EC_{true}

780

781

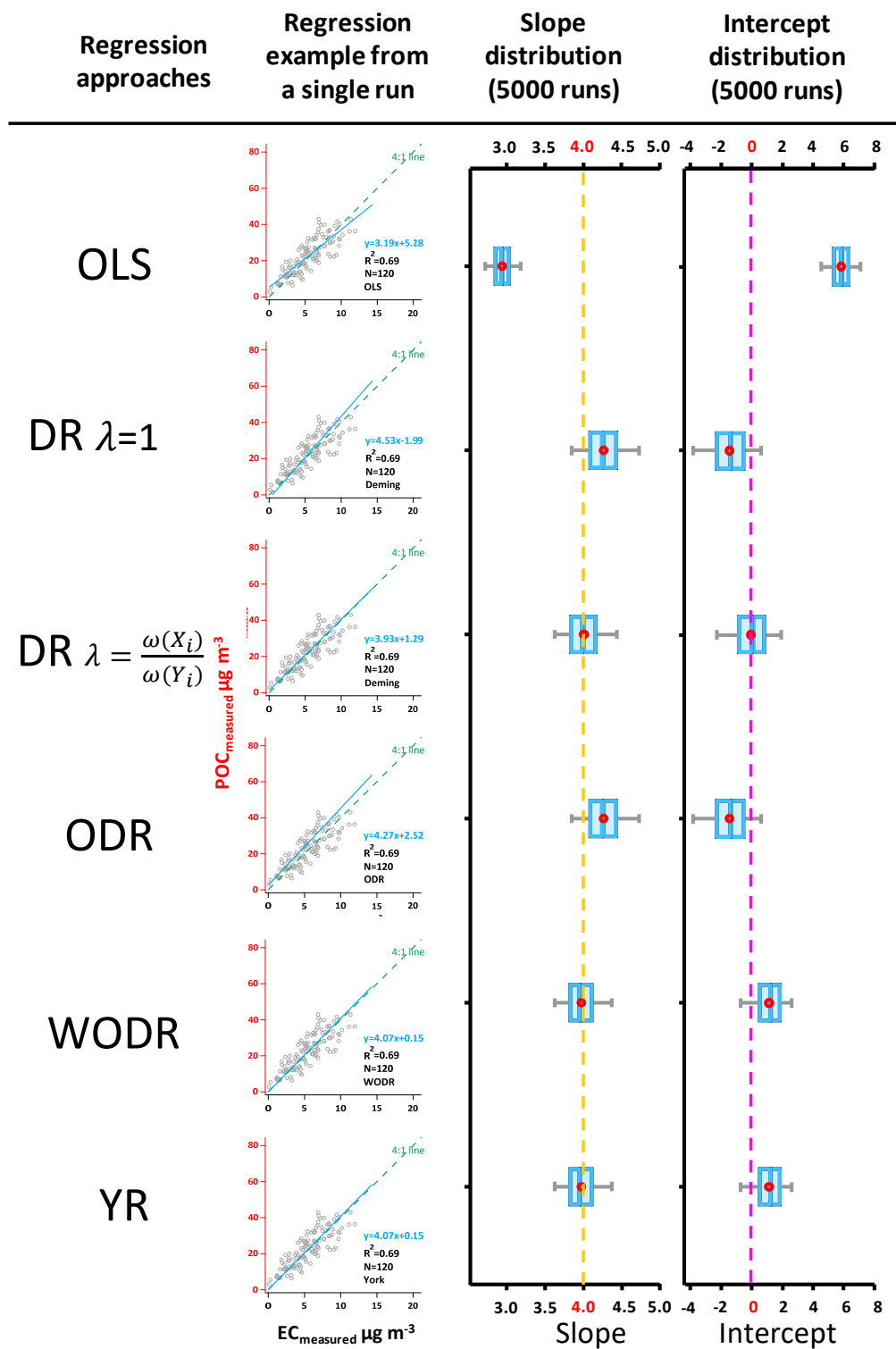
Comparison study design



782

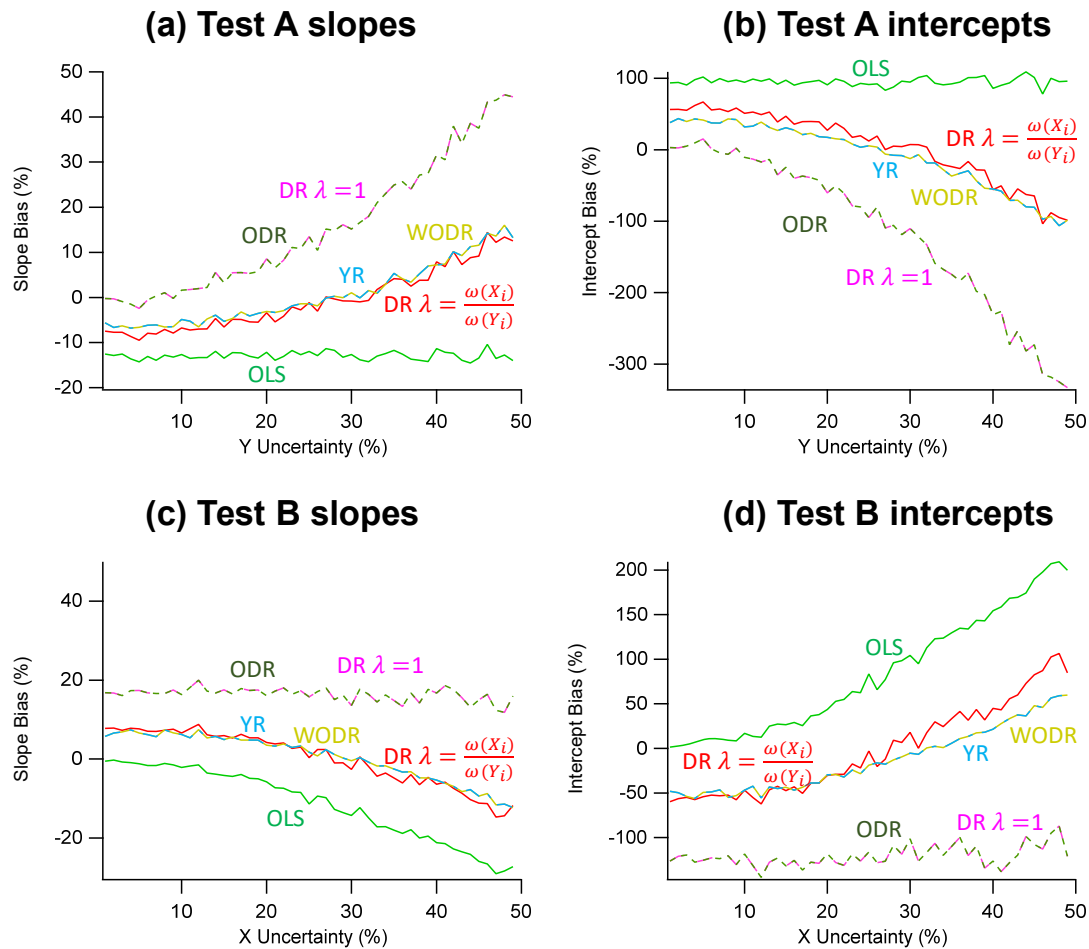
783 **Figure 4.** Overview of the comparison study design.

784



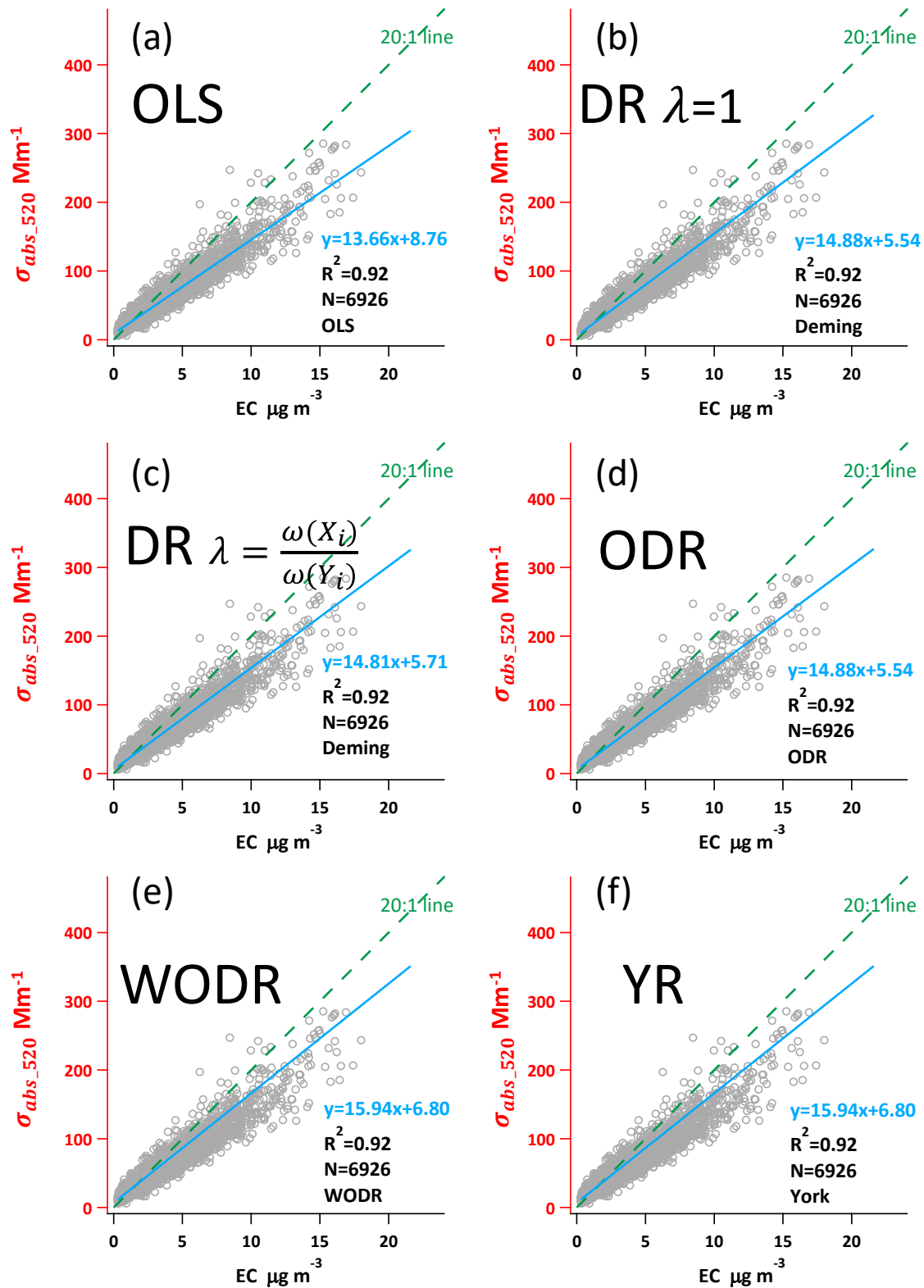
785

786 **Figure 5.** Regression results on synthetic data, case 1 (Slope=4, Intercept=0,
 787 $LOD_{POC}=1, LOD_{EC}=1, a_{POC}=1, a_{EC}=1, R^2(POC, EC) = 0.67 \pm 0.03$). The scatter plots
 788 demonstrate regression examples from a single run. The box plots show the distribution
 789 of regressed slopes and intercepts from 5000 runs of six regression approaches. The
 790 dashed line in orange and peachblow represent true slope and intercept respectively.



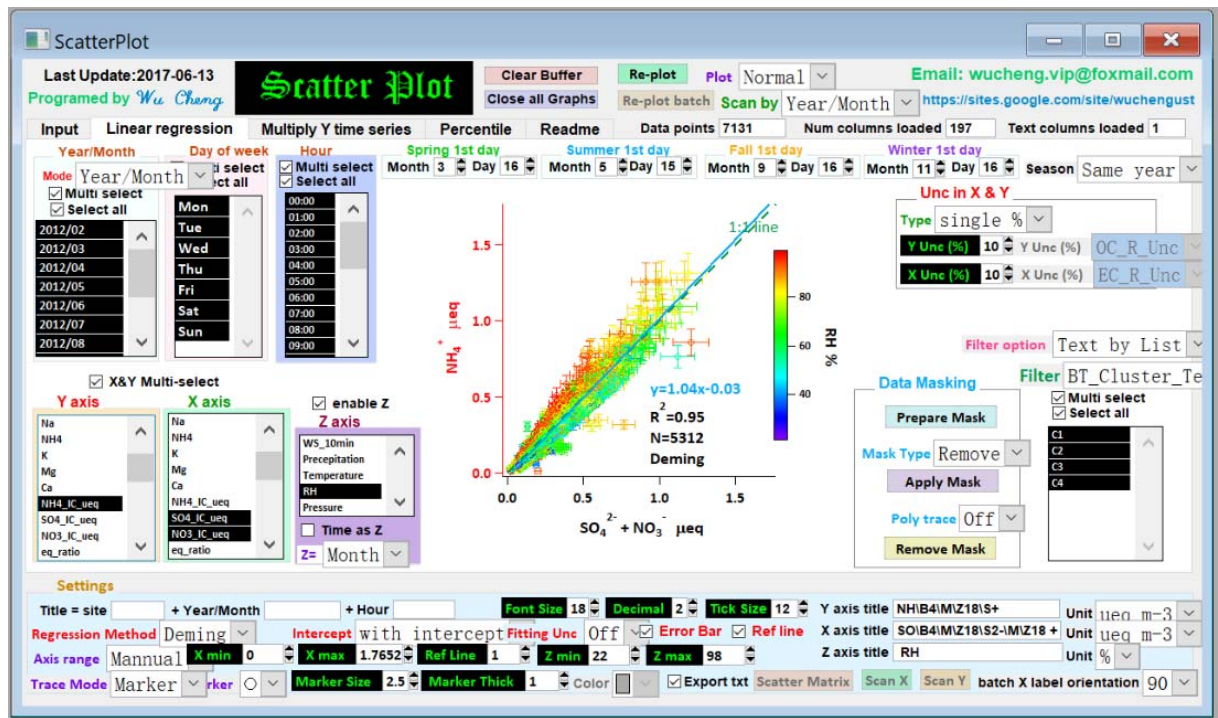
791

792 **Figure 6.** Slope and intercept biases due to the inconsistency between measurement error of
 793 data and measurement error used in regression. In Test A data generation, γ_{Unc_X} is fixed at
 794 30% and γ_{Unc_Y} varied between 1 ~ 50%. In Test B, γ_{Unc_X} varied between 1 ~ 50% and γ_{Unc_Y}
 795 is fixed at 30%. The assumed measurement error for regression is 10% for both X and Y. (a)
 796 Slopes biases as a function of γ_{Unc_Y} in Test A. (b) Intercepts biases as a function of γ_{Unc_Y} in
 797 Test A. (c) Slopes biases as a function of γ_{Unc_X} in Test B. (d) Intercepts biases as a function
 798 of γ_{Unc_X} in Test B.



799

800 **Figure 7.** Regression results using ambient σ_{abs520} and EC data from a suburban site in
 801 Guangzhou, China.



802

803 **Figure 8.** The user interface of Scatter Plot Igor program. The program and its operation

804 manual are available from: <https://doi.org/10.5281/zenodo.832417>.