

1 *Supplement of*

2 **Evaluation of linear regression techniques for**
3 **atmospheric applications: the importance of appropriate**
4 **weighting**

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

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

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

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

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

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

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

18 This document contains eight supporting tables, seven supporting figures, and discussion
19 on the impact of two primary sources.

20

21 A sampling site is often impated by multiple combustion sources in the real atmosphere.
22 In section 1 and 2 we evaluate the performance of OLS, DR, WODR and YR in scenarios
23 of two primary sources and arbitrarily dictate that the $(OC/EC)_{pri}$ of source 1 is lower than
24 source 2. By varying f_{EC1} (proportion of source 1 EC to total EC) from test to test, the effect
25 of different mixing ratios of the two sources can be examined. Two scenarios are
26 considered (Wu and Yu, 2016): two correlated primary sources and two independent
27 primary sources. Common configurations include: $EC_{total}=2 \mu\text{gC m}^{-3}$; f_{EC1} varies from 0 to
28 100%; ratio of the two OC/EC_{pri} values (γ_{pri}) vary in the range of 2~8. Studies by Chu
29 (2005) and Saylor et al. (2006) both suggest ROA being the best estimator of the expected
30 primary OC/EC ratio when SOC is zeroed. Since the overall OC/EC_{pri} from the two source
31 varies by γ_{pri} , ROA is considered as the reference OC/EC_{pri} to be compared with slope
32 regressed by of OLS, DR, WODR and YR. The abbreviations used for two primary sources
33 study are listted in Table S8.

34 **1 Impact of two correlated primary sources**

35 Simulations considering two correlated primary sources are performed, to examine the
36 effect on bias in the regression methods. The basic configuration is: $(OC/EC)_{pri1}=0.5$,
37 $(OC/EC)_{pri2}=5$, $\gamma_{unc}=30\%$, $N=8000$, $\text{intercept}=0$, and the following terms are compared:
38 ratio of average (ROA, which is considered as the true value of slope when $\text{intercept}=0$),
39 DR, WODR, WODR' (through origin) and OLS. As shown in Figure S5, when R^2 ($EC1$
40 vs. $EC2$) is very high, DR, WODR and WODR' can provide a result consistent with ROA.
41 If the R^2 decreases, the bias of the slope and intercept in DR and WODR is larger. OLS
42 constantly underestimate the slope.

43 **2 Impact of two independent primary sources**

44 Simulations of two independent primary sources are also conducted. If $RSD_{EC1}=RSD_{EC2}$,
45 slopes and intercepts may be either overestimated or underestimated (Figure S6), and the
46 degree of bias depends on the magnitude of RSD_{EC1} and RSD_{EC2} . Larger RSD results in

47 larger bias. Uneven RSD between two sources leads to even more bias (Figure S6 a&b).
48 The degree of bias also shows dependence on γ_{pri} . If γ_{pri} decreases, the bias becomes
49 smaller (FigureS6 c~f). These results indicate that the scenario with two independent
50 primary sources poses a challenge to $(OC/EC)_{pri}$ estimation by linear regression.

51 For the EC tracer method, if EC comes from two primary sources and contribution of the
52 two sources is comparable, the regression slope is no longer suitable for $(OC/EC)_{pri}$
53 estimation and the subsequent SOC calculation, and making EC a mixture that violates the
54 property of a tracer. For such a situation, pre-separation of EC into individual sources by
55 other tracers (if available) by the Minimum R Squared (MRS) method can provide unbiased
56 SOC estimation results (Wu and Yu, 2016).

57 **3 Igor programs for error in variables linear regression and simulated OC** 58 **EC data generation using MT**

59 An Igor Pro (WaveMetrics, Inc. Lake Oswego, OR, USA) based program (Scatter plot)
60 with graphical user interface (GUI) is developed to make the linear regression feasible and
61 user friendly (Figure 7). The program includes Deming and York algorithm for linear
62 regression, which consider uncertainties in both X and Y, that is more realistic for
63 atmospheric applications. It packed with lots of useful features for data analysis and graph
64 plotting, including batch plotting, data masking via GUI, color coding in Z axis, data
65 filtering and grouping by numerical values and strings.

66 Another program using MT can generate simulated OC and EC concentration through user
67 defined parameters via GUI as shown in Figure S7.

68 Both Igor programs and their operation manuals can be downloaded from the following
69 link: <https://sites.google.com/site/wuchengust>.

70

71

72 **References**

73 Chu, S. H.: Stable estimate of primary OC/EC ratios in the EC tracer method, *Atmos.*
74 *Environ.*, 39, 1383-1392, 10.1016/j.atmosenv.2004.11.038, 2005.

75 Saylor, R. D., Edgerton, E. S., and Hartsell, B. E.: Linear regression techniques for use in
76 the EC tracer method of secondary organic aerosol estimation, *Atmos. Environ.*, 40, 7546-
77 7556, 10.1016/j.atmosenv.2006.07.018, 2006.

78 Wu, C. and Yu, J. Z.: Determination of primary combustion source organic carbon-to-
79 elemental carbon (OC/EC) ratio using ambient OC and EC measurements: secondary OC-
80 EC correlation minimization method, *Atmos. Chem. Phys.*, 16, 5453-5465, 10.5194/acp-
81 16-5453-2016, 2016.

82

83

84

85 Table S1. Comparison of regression techniques using the data generation scheme of Chu
 86 (2005) with 5000 runs considering different LOD and α configurations. Non-linear relative
 87 measurement uncertainty follows Eqs. (18) & (19). Unbiased (difference <5% to true
 88 value) slopes and intercepts are highlighted in grey.

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 3	OLS	3.83	±0.08	0.95	±0.40
Slope=4, Intercept=0	DR $\lambda=1$	4.03	±0.09	-0.18	±0.44
$LOD_{POC}=0.5,$ $LOD_{EC}=0.5$ $\alpha_{POC}=0.5,$ $\alpha_{EC}=0.5$	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.00	±0.09	0.00	±0.44
	ODR	4.03	±0.09	-0.18	±0.44
$R^2(POC,EC) =$ 0.95 ± 0.01	WODR	4.00	±0.08	0.12	±0.37
	YR	4.00	±0.08	0.12	±0.37
Case 4	OLS	3.39	±0.15	3.34	±0.75
Slope=4, Intercept=0	DR1 $\lambda=1$	4.30	±0.21	-1.66	±1.06
$LOD_{POC}=1, LOD_{EC}=0.5$ $\alpha_{POC}=1, \alpha_{EC}=1$	DR1 $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.00	±0.19	-0.03	±0.99
	ODR	4.30	±0.21	-1.66	±1.06
$R^2(POC,EC) =$ 0.78 ± 0.02	WODR	4.00	±0.17	0.33	±0.81
	YR	4.00	±0.17	0.33	±0.81

89

90

91

92

93

94

95 Table S2. Comparison of regression techniques using the data generation scheme of Chu
 96 (2005) with 5000 runs with linear measurement uncertainty following Eqs. (23) & (24).
 97 Unbiased (difference <5% to true value) slopes and intercepts are highlighted in grey.

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 5 Slope=4, Intercept=0 $\gamma_{Unc}=30\%$ $R^2(POC,EC) = 0.69 \pm 0.04$	OLS	3.32	± 0.20	3.77	± 0.90
	DR $\lambda=1$	4.75	± 0.30	-4.14	± 1.36
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.01	± 0.25	-0.04	± 1.13
	ODR	4.75	± 0.30	-4.14	± 1.36
	WODR	4.00	± 0.18	-0.01	± 0.59
	YR	4.00	± 0.18	-0.01	± 0.59
Case 6 Slope=4, Intercept=3 $\gamma_{Unc}=30\%$ $R^2(POC,EC) = 0.66 \pm 0.04$	OLS	3.31	± 0.22	6.79	± 1.02
	DR1 $\lambda=1$	4.95	± 0.31	-2.26	± 1.48
	DR1 $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	3.99	± 0.26	3.05	± 1.22
	ODR	4.95	± 0.31	-2.26	± 1.48
	WODR	4.01	± 0.20	2.72	± 0.74
	YR	4.01	± 0.20	2.72	± 0.74

98

99

100 Table S3. Comparison of regression techniques using MT data generation scheme with
 101 5000 runs and a nonlinear measurement uncertainty following Eqs. (23) & (24).

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 9	OLS	0.43	±0.00	0.36	±0.02
Slope=0.5, Intercept=0	DR $\lambda=1$	0.46	±0.01	0.23	±0.03
$LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1.$	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	0.50	±0.01	-0.00	±0.03
$R^2(POC,EC) =$ 0.76±0.01	ODR	0.46	±0.01	0.23	±0.03
	WODR	0.50	±0.00	-0.00	±0.01
	YR	0.50	±0.00	-0.00	±0.01
Case 10	OLS	0.43	±0.01	3.36	±0.03
Slope=0.5, Intercept=3	DR $\lambda=1$	0.50	±0.01	3.02	±0.04
$LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1.$	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	0.49	±0.01	3.05	±0.04
$R^2(POC,EC) =$ 0.56±0.01	ODR	0.50	±0.01	3.02	±0.04
	WODR	0.51	±0.01	2.73	±0.03
	YR	0.51	±0.01	2.73	±0.03

102

103

104

105

106 Table S4. Comparison of regression techniques using MT data generation scheme with
 107 5000 runs, and a nonlinear measurement uncertainty following Eqs. (23) & (24).

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 11	OLS	0.87	±0.01	0.72	±0.05
Slope=1, Intercept=0	DR $\lambda=1$	1.00	±0.01	-0.00	±0.06
$LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1.$	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	1.00	±0.01	-0.00	±0.06
$R^2(POC,EC) =$ 0.76±0.01	ODR	1.00	±0.01	-0.00	±0.06
	WODR	1.00	±0.01	-0.00	±0.02
	YR	1.00	±0.01	-0.00	±0.02
Case 12	OLS	0.87	±0.01	3.72	±0.05
Slope=1, Intercept=3	DR $\lambda=1$	1.09	±0.01	2.52	±0.07
$LOD_{POC}=1, LOD_{EC}=1$ $a_{POC}=1, a_{EC}=1.$	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	0.99	±0.01	3.07	±0.06
$R^2(POC,EC) =$ 0.66±0.01	ODR	1.09	±0.01	2.52	±0.07
	WODR	1.01	±0.01	2.71	±0.04
	YR	1.01	±0.01	2.70	±0.04

108

109

110

111

112

113 Table S5. Comparison of regression techniques using MT data generation scheme with
 114 5000 runs and linear measurement uncertainty following Eqs. (23) & (24). Unbiased
 115 (difference <5% to true value) slopes and intercepts are highlighted with grey background.

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 13 Slope=4, Intercept=0 $\gamma_{Unc}=30\%$ $R^2(\text{POC,EC}) = 0.76 \pm 0.01$	OLS	3.48	± 0.04	2.87	± 0.18
	DR $\lambda=1$	4.53	± 0.05	-2.94	± 0.24
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	4.00	± 0.05	-0.00	± 0.22
	ODR	4.53	± 0.05	-2.94	± 0.24
	WODR	4.00	± 0.03	-0.00	± 0.09
	YR	4.00	± 0.03	-0.00	± 0.09
Case 14 Slope=4, Intercept=3 $\gamma_{Unc}=30\%$ $R^2(\text{POC,EC}) = 0.73 \pm 0.01$	OLS	3.48	± 0.04	5.87	± 0.19
	DR $\lambda=1$	4.67	± 0.05	-0.67	± 0.26
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	3.98	± 0.05	3.08	± 0.23
	ODR	4.67	± 0.05	-0.67	± 0.26
	WODR	4.02	± 0.03	2.68	± 0.11
	YR	4.02	± 0.03	2.68	± 0.11

116

117

118 Table S6. Comparison of regression techniques using MT data generation scheme with
 119 5000 runs. Linear measurement uncertainty follows Eqs. (23) & (24).

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 15 Slope=0.5 Intercept=0 $\gamma_{Unc}=30\%$ $R^2(POC,EC) = 0.54 \pm 0.01$	OLS	0.40	± 0.01	0.55	± 0.03
	DR $\lambda=1$	0.45	± 0.01	0.26	± 0.03
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	0.50	± 0.01	0.01	± 0.03
	ODR	0.45	± 0.01	0.26	± 0.03
	WODR	0.52	± 0.01	-0.23	± 0.02
	YR	0.52	± 0.01	-0.23	± 0.02
Case 16 Slope=0.5 Intercept=3 $\gamma_{Unc}=30\%$ $R^2(POC,EC) = 0.40 \pm 0.01$	OLS	0.40	± 0.01	3.54	± 0.04
	DR $\lambda=1$	0.50	± 0.01	2.98	± 0.04
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	0.50	± 0.01	3.00	± 0.04
	ODR	0.50	± 0.01	2.98	± 0.04
	WODR	0.52	± 0.01	2.65	± 0.04
	YR	0.52	± 0.01	2.65	± 0.04

120
 121
 122
 123

124 Table S7. Comparison of regression techniques using MT data generation scheme with
 125 5000 runs. Linear measurement uncertainty follows Eqs. (23)&(24).

Case	Regression scenario	Slope mean	Slope SD	Intercept mean	Intercept SD
Case 17 Slope=1 Intercept=0 $\gamma_{Unc}=30\%$ $R^2(POC,EC) = 0.65 \pm 0.01$	OLS	0.80	± 0.01	1.07	± 0.04
	DR $\lambda=1$	1.00	± 0.01	0.00	± 0.05
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	1.00	± 0.01	0.00	± 0.05
	ODR	1.00	± 0.01	0.00	± 0.05
	WODR	1.00	± 0.01	0.00	± 0.04
	YR	1.00	± 0.01	0.00	± 0.04
Case 18 Slope=1 Intercept=3 $\gamma_{Unc}=30\%$ $R^2(POC,EC) = 0.59 \pm 0.01$	OLS	0.80	± 0.01	4.07	± 0.05
	DR $\lambda=1$	1.07	± 0.01	2.62	± 0.07
	DR $\lambda = \frac{\omega(X_i)}{\omega(Y_i)}$	1.00	± 0.01	3.00	± 0.06
	ODR	1.07	± 0.01	2.62	± 0.07
	WODR	1.02	± 0.01	2.84	± 0.05
	YR	1.02	± 0.01	2.84	± 0.05

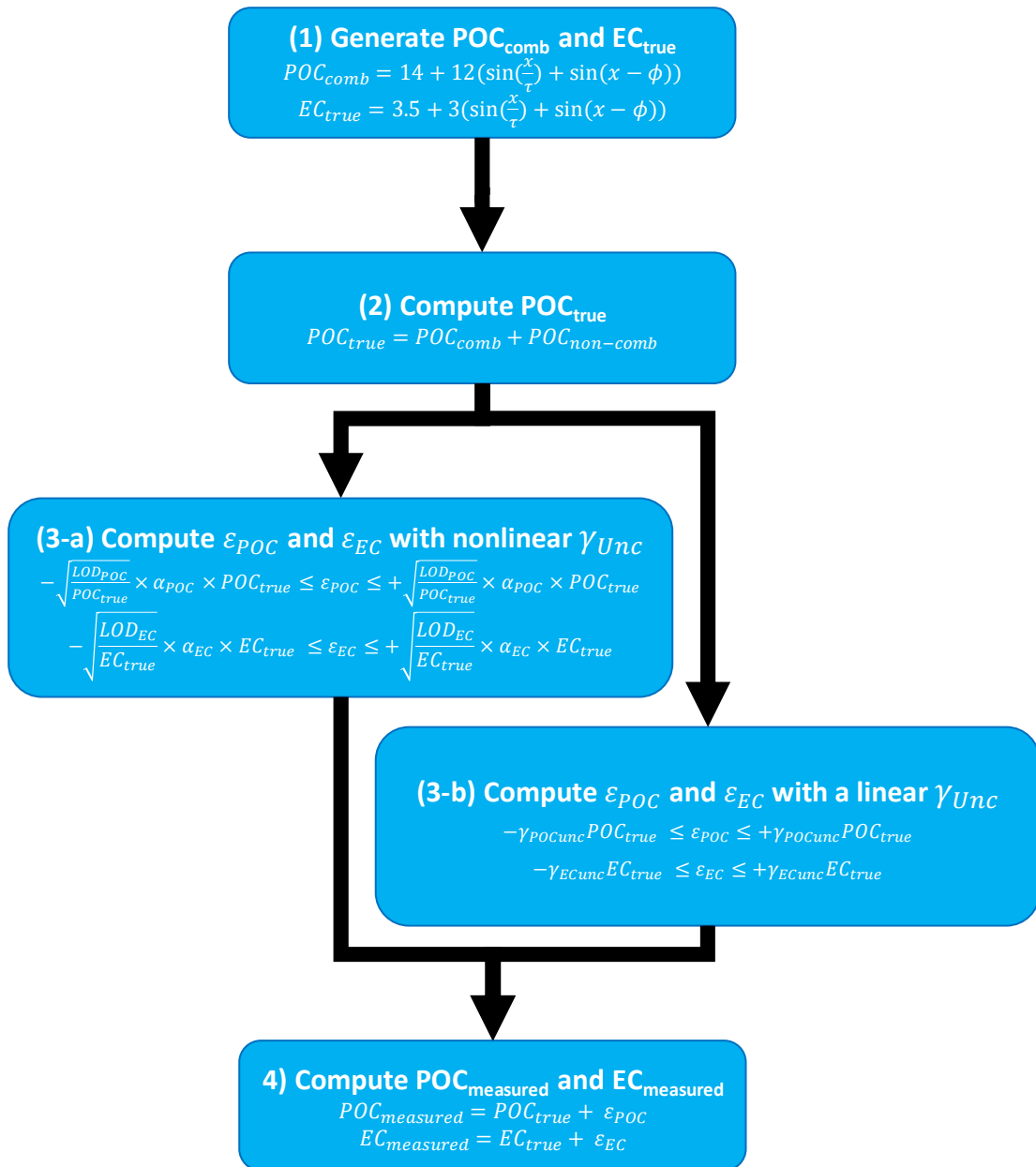
126

127 Table S8. Abbreviations used for two primary sources study.

Abbreviation	Definition
EC_1, EC_2	EC from source 1 and source 2 in the two sources scenario
f_{EC1}	fraction of EC from source 1 to the total EC
ROA	ratio of averages
γ_{pri}	ratio of the $(OC/EC)_{pri}$ of source 2 to source 1
RSD	relative standard deviation
RSD_{EC}	RSD of EC
$\epsilon_{EC}, \epsilon_{OC}$	measurement uncertainty of EC and OC
γ_{unc}	relative measurement uncertainty
γ_{RSD}	the ratio between the RSD values of $(OC/EC)_{pri}$ and EC

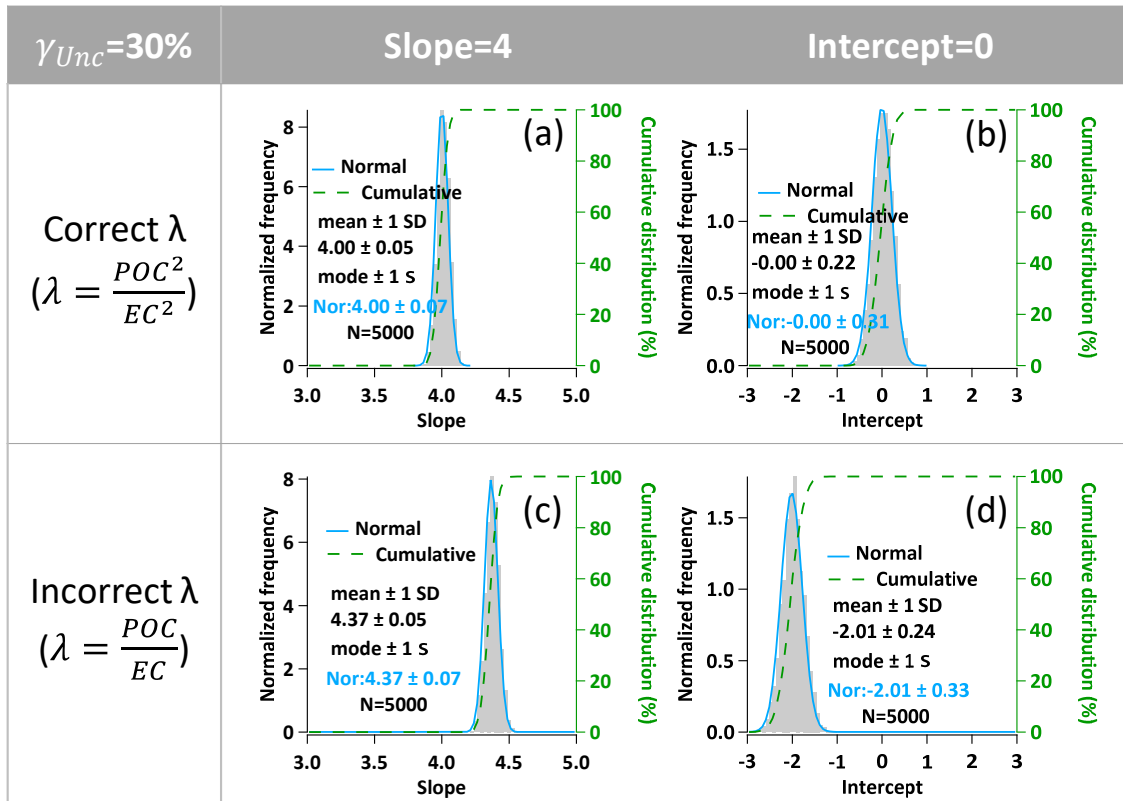
128

Data generations steps by the sine functions of Chu (2005)



129

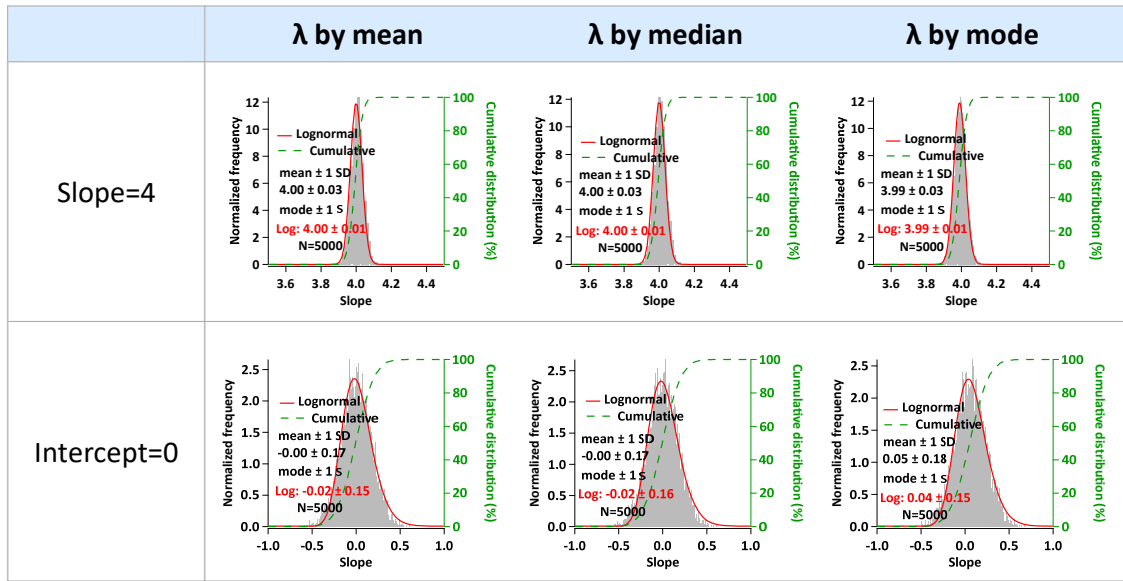
130 Figure S1. Flowchart of data generation steps using the sine functions of Chu (2005).



131

132 Figure S2. Example of bias in slope and intercept due to improper λ assignment. Data
 133 generation: Slope=4, Intercept=0; linear γ_{Unc} (30%). (a)&(b) Slopes and intercepts when
 134 proper λ is input following linear γ_{Unc} ($\lambda = \frac{POC^2}{EC^2}$); (c)&(d) Slopes and intercepts when
 135 improper λ is input following non-linear γ_{Unc} ($\lambda = \frac{POC}{EC}$).

136



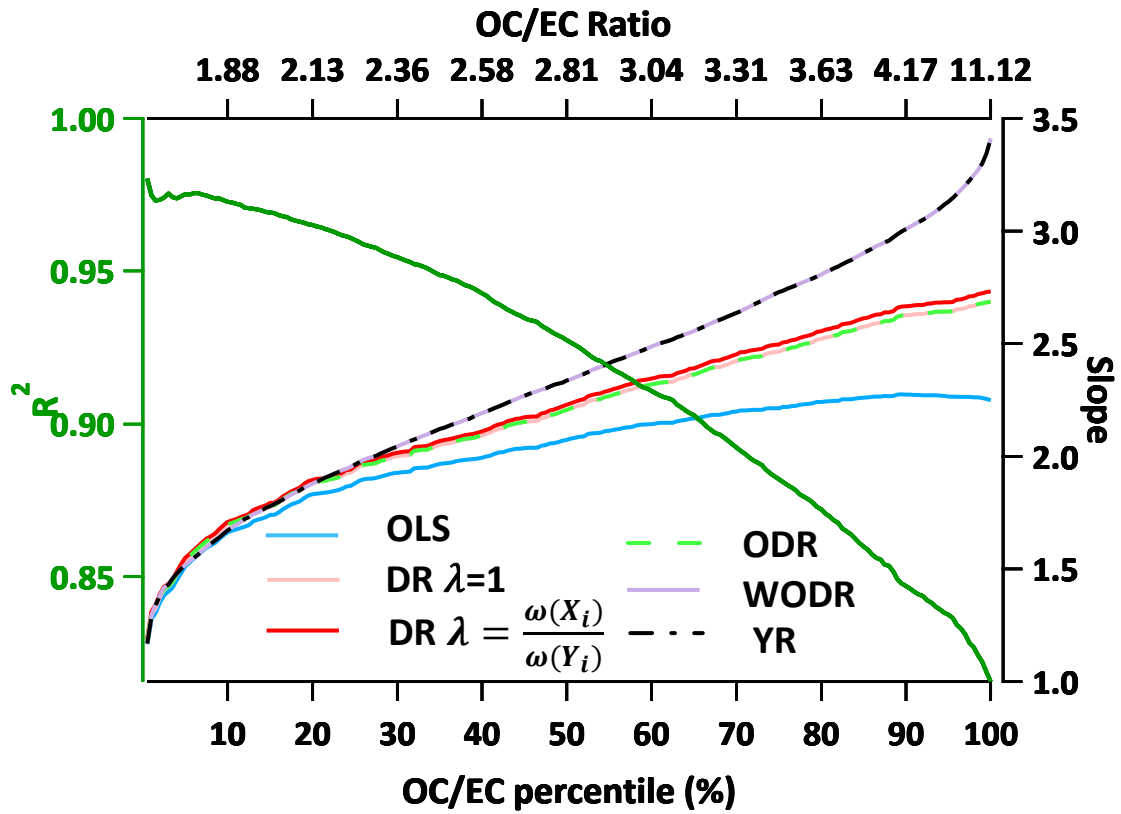
137

138 Figure S3. Sensitivity tests of λ calculated by mean, median and mode.

139

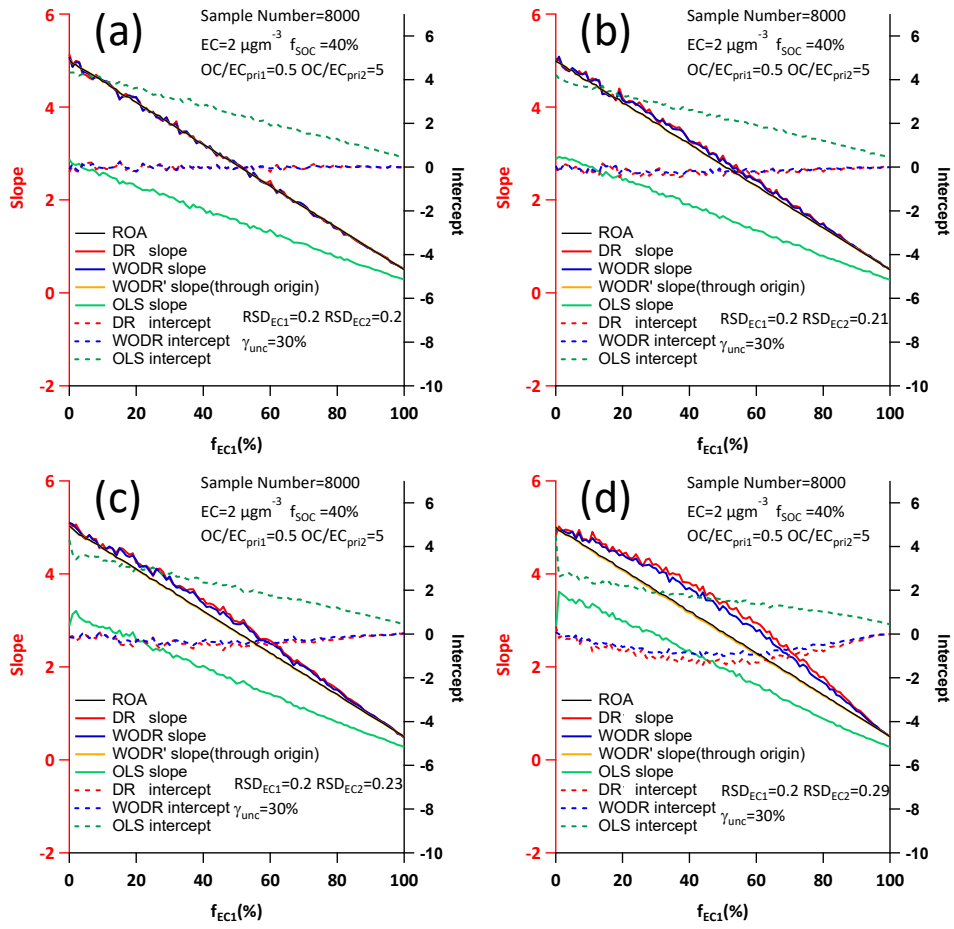
140

141



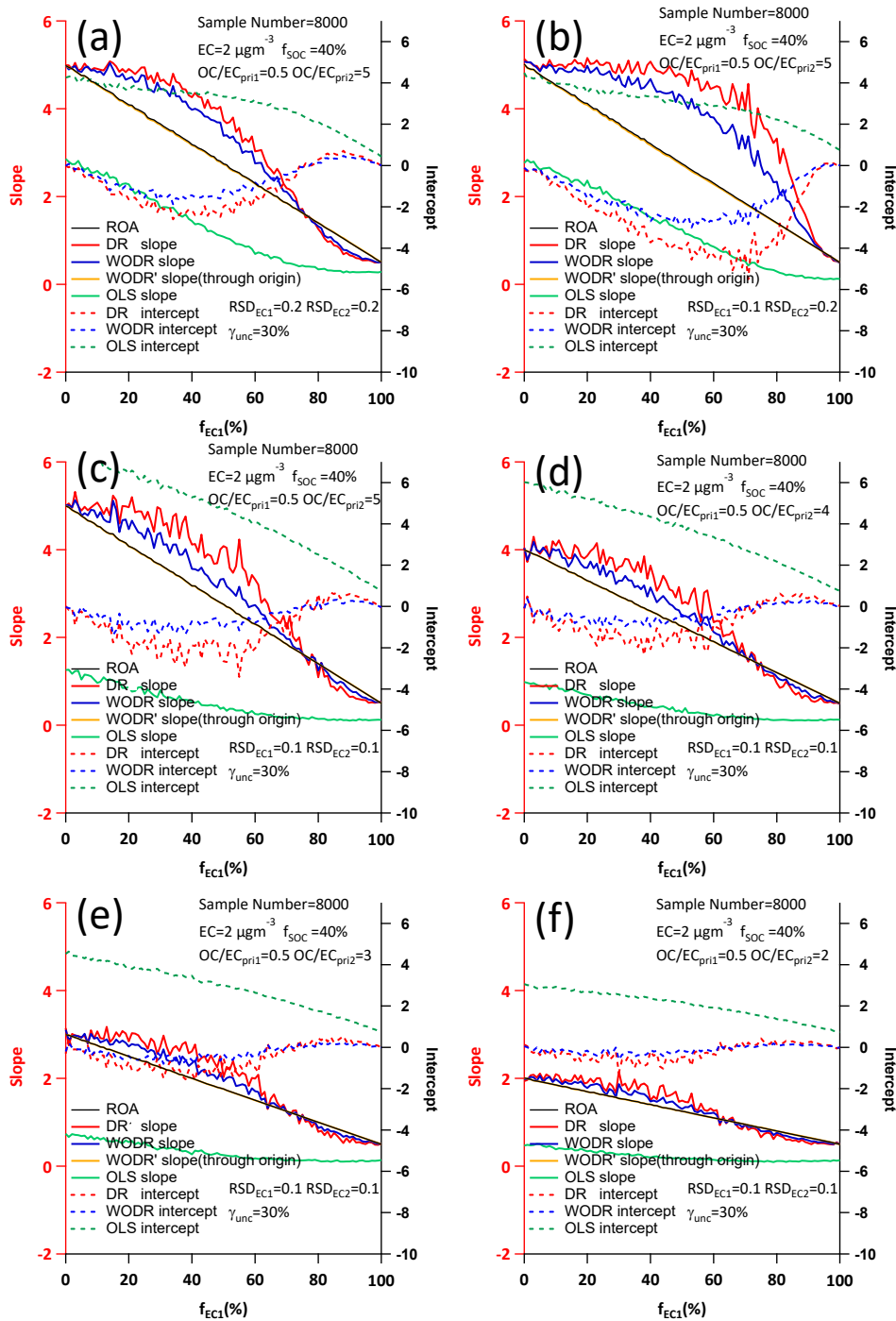
142

143 Figure S4. Regression slopes as a function of OC/EC percentile. OC/EC percentile range
 144 from 0.5% to 100%, with an interval of 0.5%.



145

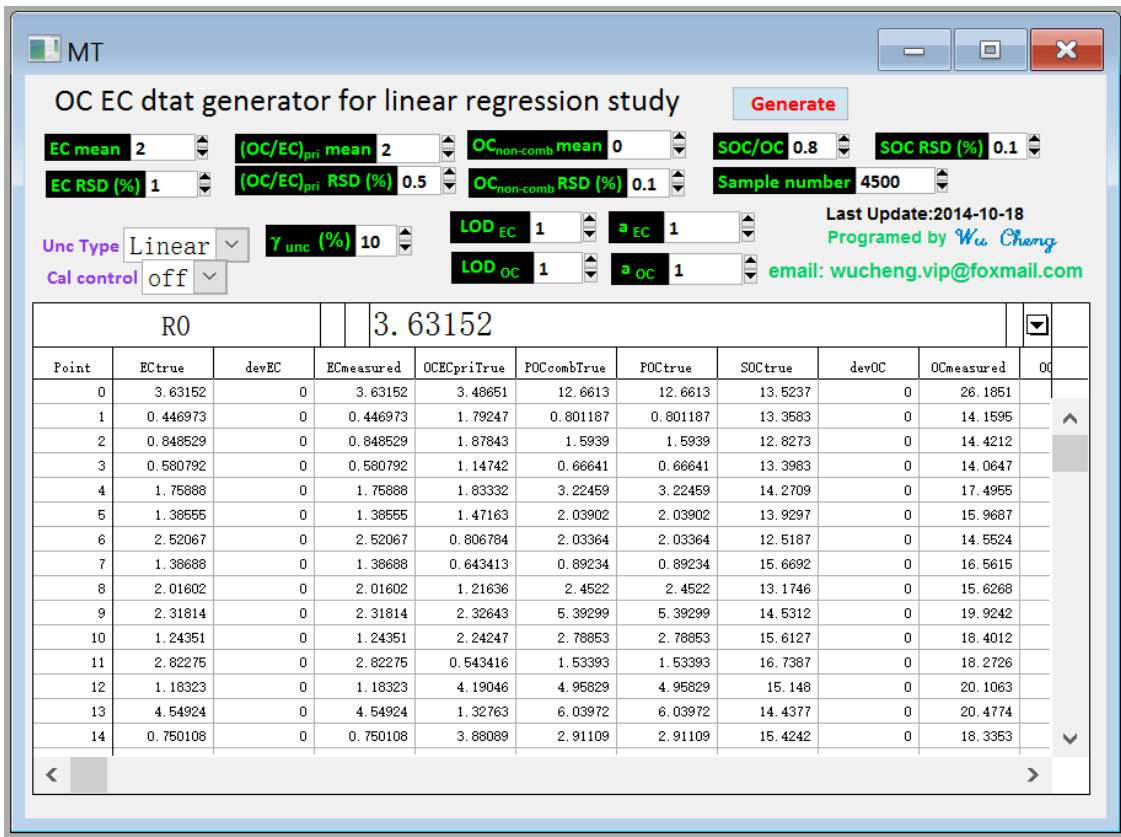
146 Figure S5. Study of two correlated sources scenario by different R^2 between the two
 147 sources. (a) $R^2 = 1$ (b) $R^2 = 0.86$ (c) $R^2 = 0.75$ (d) $R^2 = 0.49$



148

149 Figure S6. Study of two independent sources scenario by different parameters. (a) $\gamma_{pri}=10$,
 150 $RSD_{EC1}=0.2, RSD_{EC2}=0.2$ (b) $\gamma_{pri}=10$, $RSD_{EC1}=0.1, RSD_{EC2}=0.2$ (c) $\gamma_{pri}=10$,
 151 $RSD_{EC1}=0.1, RSD_{EC2}=0.1$ (d) $\gamma_{pri}=8$, $RSD_{EC1}=0.1, RSD_{EC2}=0.1$ (e) $\gamma_{pri}=6$, $RSD_{EC1}=0.1$,
 152 $RSD_{EC2}=0.1$ (f) $\gamma_{pri}=4$, $RSD_{EC1}=0.1, RSD_{EC2}=0.1$

153



154

155 Figure S7. MT Igor program. OC and EC data following log-normal distribution can be
 156 generated for statistical study purpose (no time series information). User can define mean
 157 and RSD of EC, $(OC/EC)_{pri}$, SOC/OC ratio, measurement uncertainty, sample size, etc.
 158 MT Igor program can be downloaded from the following link:
 159 <https://sites.google.com/site/wuchengust>.

160